











# Effect of sequence depth and length in long-read assembly of the maize inbred NC358

Shujun Ou<sup>1</sup>, Jianing Liu<sup>2</sup>, Kapeel M. Chougule<sup>3</sup>, Arkarachai Fungtammasan<sup>4</sup>, Arun S. Seetharam <sup>1,5</sup>, Joshua C. Stein<sup>3</sup>, Victor Llaca <sup>6</sup>, Nancy Manchanda<sup>1</sup>, Amanda M. Gilbert<sup>7</sup>, Sharon Wei<sup>3</sup>, Chen-Shan Chin<sup>4</sup>, David E. Hufnagel<sup>1</sup>, Sarah Pedersen<sup>1</sup>, Samantha J. Snodgrass<sup>1</sup>, Kevin Fengler<sup>6</sup>, Margaret Woodhouse <sup>8</sup>, Brian P. Walenz <sup>9</sup>, Sergey Koren <sup>9</sup>, Adam M. Phillippy <sup>9</sup>, Brett T. Hannigan<sup>4</sup>, R. Kelly Dawe <sup>2</sup>✉, Candice N. Hirsch <sup>7</sup>✉, Matthew B. Hufford <sup>1</sup>✉ & Doreen Ware <sup>3,10</sup>✉

Improvements in long-read data and scaffolding technologies have enabled rapid generation of reference-quality assemblies for complex genomes. Still, an assessment of critical sequence depth and read length is important for allocating limited resources. To this end, we have generated eight assemblies for the complex genome of the maize inbred line NC358 using PacBio datasets ranging from 20 to 75 × genomic depth and with N50 subread lengths of 11–21 kb. Assemblies with ≤30 × depth and N50 subread length of 11 kb are highly fragmented, with even low-copy genic regions showing degradation at 20 × depth. Distinct sequence-quality thresholds are observed for complete assembly of genes, transposable elements, and highly repetitive genomic features such as telomeres, heterochromatic knobs, and centromeres. In addition, we show high-quality optical maps can dramatically improve contiguity in even our most fragmented base assembly. This study provides a useful resource allocation reference to the community as long-read technologies continue to mature.

<sup>1</sup>Department of Ecology, Evolution, and Organismal Biology, Iowa State University, Ames, Iowa 50011, USA. <sup>2</sup>Department of Genetics, University of Georgia, Athens, Georgia 30602, USA. <sup>3</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, New York 11724, USA. <sup>4</sup>DNAnexus, Inc., Mountain View, San Francisco, California 94040, USA. <sup>5</sup>Genome Informatics Facility, Iowa State University, Ames, Iowa 50011, USA. <sup>6</sup>Genomics Technologies, Applied Science and Technology, Corteva Agriscience TM, Johnston, Iowa 50131, USA. <sup>7</sup>Department of Agronomy and Plant Genetics, University of Minnesota, St. Paul, Minnesota 55108, USA. <sup>8</sup>USDA ARS Corn Insects and Crop Genetics Research Unit, Ames, Iowa 50011, USA. <sup>9</sup>Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland 20892, USA. <sup>10</sup>USDA ARS Robert W. Holley Center for Agriculture and Health, Agricultural Research Service, Ithaca, New York 14853, USA.

✉email: [kdawe@uga.edu](mailto:kdawe@uga.edu); [cnhirsch@umn.edu](mailto:cnhirsch@umn.edu); [mhufford@iastate.edu](mailto:mhufford@iastate.edu); [ware@cshl.edu](mailto:ware@cshl.edu)

During the two decades following publication of the first larger eukaryotic genomes (i.e., *Drosophila melanogaster*<sup>1</sup> and *Homo sapiens*<sup>2</sup>), considerable progress has been made in sequencing technology and assembly methods, improving our basic knowledge of genome complexity across the tree of life. For example, we now understand that genome composition (e.g., gene complement, the extent of intergenic space, and the landscape of transposable elements (TEs)) varies substantially at both the inter- and intra-specific levels.

Throughout the genomic era, continual improvements have been made to both the data underlying assemblies and assembly algorithms. The first reference genomes employed bac-by-bac approaches with long-read Sanger sequencing generated at great expense over several years<sup>2,3</sup>. Next-generation assemblies initially relied on short-read data due to cost and technological limitations. These data have often been used for assembly in a combination of paired-end and mate pair formats with local assembly of contigs relying on paired-end sequences, which are then scaffolded using mate pair information<sup>4–6</sup>. While these assemblies represented genes reasonably well, repetitive regions containing TEs and tandem repeats were either omitted or highly fragmented<sup>7</sup>. Newly developed long-read sequencing technologies now enable contiguous assembly of even the repetitive fraction of eukaryotic genomes<sup>8</sup> with, e.g., a complete telomere-to-telomere human X chromosome recently being assembled<sup>9</sup>. Highly contiguous long-read assemblies have now been published for a wide range of species<sup>10–13</sup>. Recent long-read assemblies in maize also show considerable improvement in both completeness and contiguity relative to previous efforts<sup>13–16</sup>, suggesting these data are particularly useful for plant species like maize with genomes that are large (2.3 Gb), complex (paleopolyploid and comprised primarily of TEs), and highly repetitive (extensive tandem sequence arrays in heterochromatic knobs and centromeres).

Assembly algorithms are also continually being benchmarked and improved to better utilize long-read data with substantial variability in performance of particular methods observed across species<sup>10,17–19</sup>.

The cost of long-read sequence data can still be prohibitive for species with larger genomes, and the critical target for average read length and read depth remains unclear. A full assessment of the impacts of varying sequence read length and depth on the contiguity and completeness of assemblies in genic and repetitive regions is therefore essential for informed allocation of finite resources. Here we conduct a comprehensive assembly experiment using subsets of a high-depth, long-read (PacBio) dataset for the maize inbred line NC358 to evaluate critical inflection points of quality during the assembly of a complex, repeat-rich genome.

## Results

**Genome sequence subsampling and assembly statistics.** We sequenced the NC358 genome to 75× depth (based on a ~2.27 Gb genome size<sup>20</sup>) using the PacBio Sequel platform, which generated a subread N50 of 21.2 kb (Table 1, Supplementary Table 1, and Supplementary Fig. 1). To identify an optimal assembly approach for this study, the complete data from NC358 were each assembled using Falcon<sup>21</sup>, Canu<sup>22</sup>, WTDG2<sup>23</sup>, and hybrid approaches in which Falcon was used for error correction and Canu, Flye<sup>24</sup>, and Peregrine<sup>25</sup> were used for assembly (Supplementary Table 2). A subset of methods (Supplementary Table 2) were confirmed to be robust using data from the B73 v4 genome assembly (68× depth)<sup>13</sup>. All assembled contigs were superscaffolded with a de novo Bionano optical map (Supplementary Fig. 2) and pseudomolecules were constructed based on maize GoldenGate genetic markers<sup>26</sup> and high-density maize pan-genome markers<sup>27</sup> (Methods).

**Table 1 Summary statistics for NC358 assemblies.**

Experiment <sup>a</sup>	21k_20×	21k_30×	21k_40×	21k_50×	21k_60×	21k_75×	11k_50×	16k_50×
Subreads size (Gb)	45.62	68.16	91.01	113.89	136.80	171.08	113.63	113.60
Subread coverage	20×	30×	40×	50×	60×	75×	50×	50×
Max read length (kb)	89.6	103.3	103.3	103.3	103.3	103.3	88.3	69.8
Subread N25 (kb)	30.1	30.1	30.1	30.1	30.1	30.1	14.5	21.6
Subread N50 (kb)	21.2	21.2	21.2	21.2	21.2	21.2	11.1	16.8
Corrected reads (Gb)	25.11	48.13	66.05	82.96	88.93	100.90	79.26	80.22
Corrected coverage	11×	21×	29×	37×	39×	44×	35×	35×
Corrected read N50 (kb)	18.42	17.13	17.10	17.25	18.80	20.05	10.37	14.48
Contig number	10,563	2015	641	407	360	327	5683	1036
Contig total (Gb)	1.60	2.11	2.12	2.12	2.13	2.13	2.10	2.12
Longest contig (Mb)	1.06	11.50	47.89	76.00	79.68	78.40	4.37	21.45
Contig N50 (Mb)	0.18	1.82	7.48	16.27	22.12	24.54	0.56	4.24
Longest scaffold (Mb)	198.5	198.7	237.1	237.2	237.1	237.3	205.4	237.6
Superscaffold N50 (Mb)	95.3	96.9	99.2	98.5	99.4	99.2	98.5	99.4
Assembled (%) <sup>b</sup>	70.4%	92.8%	93.3%	93.3%	93.7%	93.7%	92.4%	93.2%
Assembly gaps (%)	24.50%	0.90%	0.43%	0.34%	0.31%	0.31%	2.01%	0.48%
Effective assembly size (Gb) <sup>c</sup>	1.33	1.67	1.70	1.72	1.74	1.75	1.68	1.70
Optical map conflict <sup>d</sup>	594	125	56	31	22	21	386	107
Complete BUSCOs <sup>e</sup>	68.0%	95.5%	96.5%	96.4%	96.2%	96.3%	95.7%	96.7%
LTR Assembly Index (LAI)	12.2	19.8	20.4	20.2	20.4	20.6	19.1	21.0
Falcon CPU hour	1563	4162	6363	10,693	12,386	32,950	9721	9224
Falcon RAM (Gb)	75	75	75	75	75	75	75	75
Canu CPU hour	1860	4036	5959	7914	8849	11,520	6400	7174
Canu RAM (Gb)	61	112	149	177	201	120	183	174

<sup>a</sup>Each dataset was assembled only once with the Falcon–Canu hybrid approach (see Methods).

<sup>b</sup>Calculated based on total contig size and the estimated genome size of 2.2724 Gb.

<sup>c</sup>Sum of unique 150-mers.

<sup>d</sup>The optical map was generated using the Direct Label and Stain (DLS) approach with enzyme DLE-1.

<sup>e</sup>Pilon-polished assemblies were used to calculate Benchmarking Universal Single-Copy Orthologs (BUSCO) scores. CPU central processing.

The Falcon–Canu hybrid assemblies of both the NC358 and B73 genomes showed consistently higher quality in terms of contig length, Bionano conflicts, Benchmarking Universal Single-Copy Orthologs (BUSCOs)<sup>28</sup>, and LTR Assembly Index (LAI)<sup>8</sup> (Supplementary Table 2); thus, this method was used for all subsequent assemblies performed on subsets of the data.

NC358 subreads were downsampled from 75× to 60×, 50×, 40×, 30×, and 20×, while maintaining a 21 kb subread N50 and to 50× depth with a subread N50 of 11 and 16 kb. For these latter assemblies, we based read-length distributions on empirical datasets generated for the human HG002<sup>29</sup> and maize B73 v4<sup>13</sup> genome assemblies (Supplementary Fig. 3). NC358 read subsets were error-corrected and assembled independently using the hybrid assembly approach described above (Methods and Supplementary Note 1). These processes were resource-intensive and were accelerated through cloud computing. The central processing unit (CPU) time required for both Falcon error correction and Canu assembly increased substantially as read depth increased, whereas the required maximum memory was fairly similar (Fig. 1h and Table 1).

Most assemblies had a total contig size covering >92% of the flow-cytometry estimated genome size of NC358 (2.27 Gb<sup>20</sup>), with the notable exception of the 21k\_20x assembly (70.4% covered; Table 1). Contig length metrics were positively correlated with both read length and sequence coverage (Fig. 1b), with the highest contig N50 (24.54 Mb) and the longest contig (79.68 Mb) observed in the 21k\_75x and 21k\_60x assembly, respectively (Table 1). A dramatic drop in contiguity was observed for both the lowest depth (21k\_20x) and shortest sequence length (11k\_50x) assemblies, where the number of contigs was 17×–32× more than the complete 21k\_75x dataset (Table 1 and Fig. 1e).

For each assembly, superscaffolds were generated from the contigs using a common Bionano optical map. Even the most fragmented Falcon–Canu assembly could be scaffolded to high contiguity using this optical map due to the high density of labels in the map (Fig. 1a–c). The resulting assemblies all had scaffold N50s at ~98 Mb (Table 1), highlighting the utility of a high-quality optical map when sequencing quality is suboptimal. In fact, chromosome 3 (~237 Mb) consisted of a single scaffold in five out of eight assemblies (Table 1). However, conflicts vs. the Bionano map were much higher in the assemblies with 20× coverage and a subread N50 of 11 kb (Table 1 and Fig. 1e), suggesting assembly error increased with lower coverage and shorter read length. Assemblies with shorter read length contained many more deletions relative to the optical map (Fig. 1i; Supplementary Table 3), which may be due to the collapse of repetitive sequences. We did not observe a clear pattern between read length and deletion size (Fig. 1i). Assembly misjoins were reduced with both longer reads and higher coverage, as shown by the relative number of insertions (Fig. 1i and Supplementary Table 3).

For each of the assemblies, pseudomolecules were constructed using the GoldenGate and pan-genome genetic markers, which placed >99% of the total assembled bases into pseudomolecules (Supplementary Table 4 and Supplementary Fig. 4). The resulting NC358 pseudomolecules were highly syntenic across assemblies and to the B73 v4 genome (Supplementary Fig. 5).

**Evaluation of the gene assembly space.** We evaluated the completeness of gene-rich regions in each of the assemblies using BUSCO<sup>28</sup>. The percentage of complete BUSCO genes increased from 68.0% to 96.3% from the 21k\_20x to the 21k\_75x assembly (Table 1, Fig. 1d, and Supplementary Table 5). Minimal improvement in BUSCO scores was achieved at depths higher

than 30× (95.5% complete BUSCO genes), indicating this depth provides satisfactory gene-space assembly.

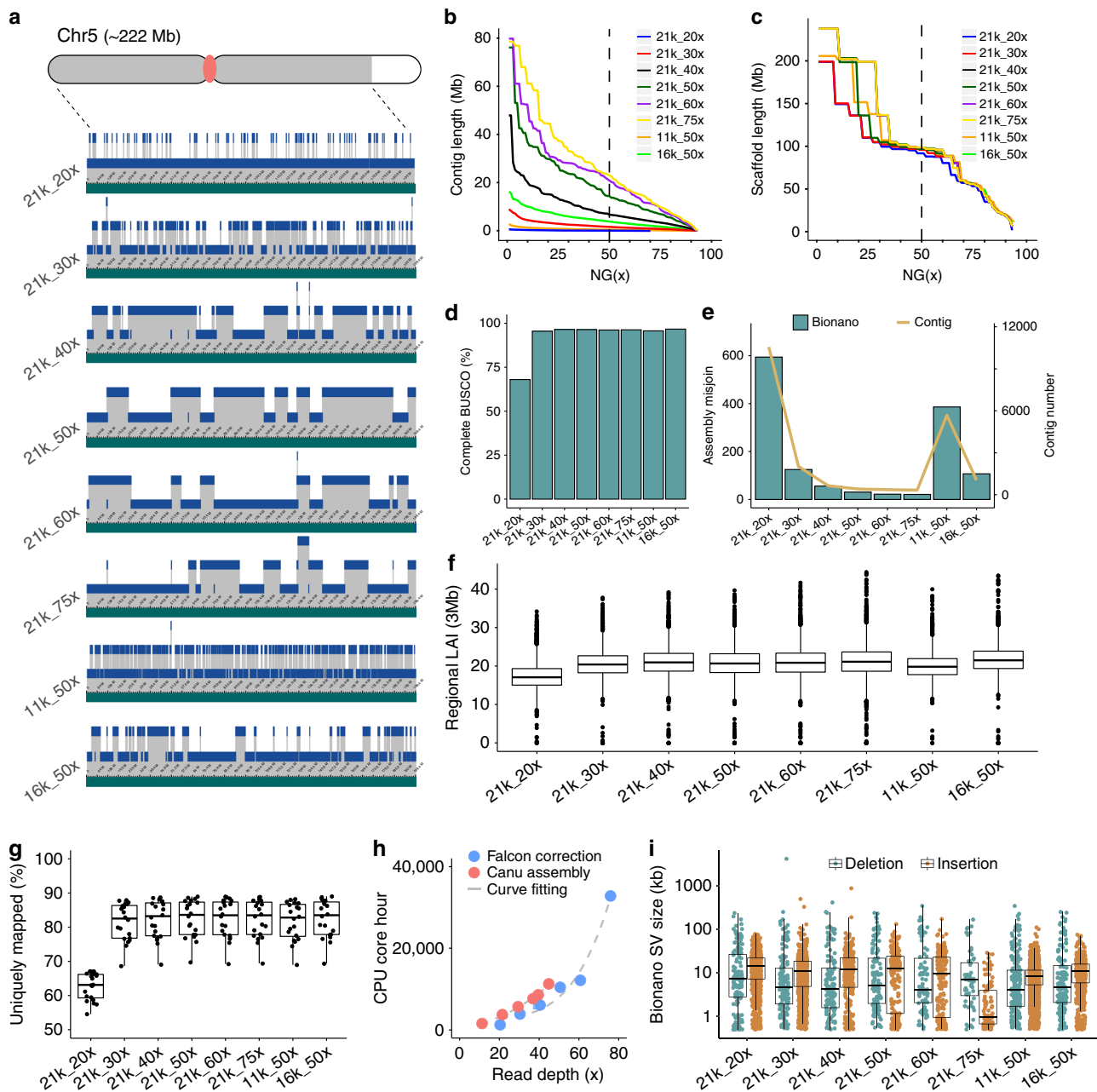
To further evaluate the assembly of genic regions, we annotated gene models in the 21k\_20x and the 21k\_75x assemblies (Methods) and obtained a total of 28,275 and 39,578 genes, respectively (Supplementary Table 6). More than 99% of the 21k\_75x genes could be mapped to all but the 21k\_20x assembly, to which only 90% were mapped (Supplementary Table 7). Exon and intron lengths of the annotated genes were similar across the assemblies (Supplementary Table 6). However, genes in the 21k\_20x assembly contained 50 times more gaps than those in the 21k\_75x assembly (Supplementary Table 7), suggesting higher coverage contributed to the contiguity of genic regions. Longer reads similarly improved the gene-space contiguity (Supplementary Table 7).

In addition, we sequenced RNA libraries from ten tissues with two biological replicates (Methods). On average, 80% of reads in these libraries could be uniquely mapped to the various NC358 assemblies (Fig. 1g). The 21k\_20x assembly was a notable exception with only 63% of reads uniquely mapped (Fig. 1g and Supplementary Fig. 6). We extracted the reads that did not map to the 21k\_20x assembly and remapped them to the 21k\_75x assembly, obtaining a unique mapping rate of 36% (Supplementary Table 8). These reads mapped to 3184 genes in the 21k\_75x assembly (Supplementary Table 9). Of these genes, 20% are present in the 21k\_20x assembly but had assembly errors that prevented the RNA sequencing (RNA-seq) reads from mapping, while the remaining 80% were within sequence gaps (Supplementary Table 9).

In addition to metrics of gene completeness, we also examined each assembly for its ability to capture two notable maize tandem gene arrays, *Rp1-D*<sup>30</sup> and *zein*<sup>31</sup>. The total length of these gene arrays was estimated at 536 kb and 62 kb in NC358, respectively, based on the optical map. Both the *Rp1-D* and *zein* loci were completely assembled in all, except for the 21k\_20x assembly, where only 70% and 91% of the loci were assembled, respectively (Fig. 2g and Supplementary Table 10).

**Evaluation of the TE assembly space.** The completeness of transposon-rich regions of the genome was assessed through the assembly index of LTR retrotransposons, called LAI<sup>8</sup>. A higher LAI score is indicative of a more complete assembly in TE-rich regions. The 21k\_20x assembly had a substantially lower LAI score compared to other assemblies (LAI = 12.2; Table 1). As sequence depth increased a substantial improvement in LAI was observed, while the effect of sequence length on LAI was minimal (Fig. 1f). This is likely due to the fact that the length of LTR retrotransposons is ~10 kb on average (Supplementary Fig. 7), which could be spanned by even the 11 kb reads. The assemblies that were generated from ≥40× genomic depth achieved gold quality (LAI ≥ 20 (ref. <sup>8</sup>)) (Table 1 and Fig. 1f), which was comparable to the B73 v4 genome and much higher than many previously published maize genome assemblies generated with short-read data (Supplementary Fig. 8).

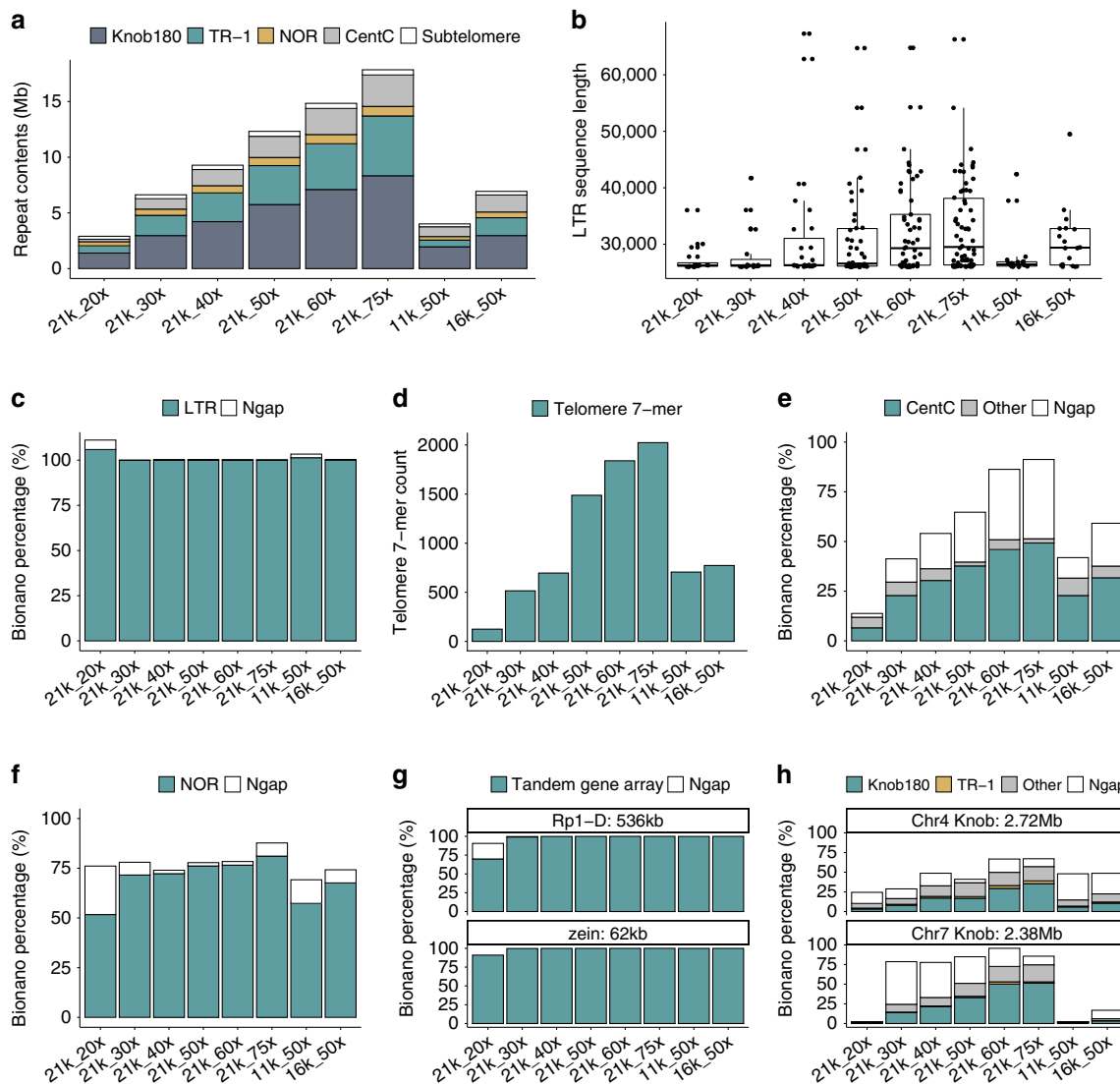
The insertion time of each LTR retrotransposon can be dated based on sequence divergence between terminal repeats<sup>8</sup>. We identified 36% fewer intact LTR retrotransposons in the highly fragmented 21k\_20x assembly (Supplementary Fig. 9) and significantly older LTR elements in the 11k\_50x assembly ( $p < 10^{-5}$ , one-way analysis of variance followed by Tukey's test), suggesting fragmentation of assemblies could bias conclusions of transposon studies. LTR retrotransposons shorter than 26 kb were assembled well across the assemblies (Supplementary Figs. 10 and 11). However, a substantial effect of longer reads and higher depth was observed in the assembly of LTR sequences longer than 26 kb (Fig. 2b). We examined the assemblies of the



**Fig. 1** Assembly of NC358 using various read lengths and coverage. **a** Hybrid scaffolding using the Bionano optical map. A 199 Mb scaffold from chromosome 5 is shown. Gray areas on the chromosome cartoon represent the 199 Mb scaffold; the white area is the remaining 23 Mb scaffold in chromosome 5; the red dot is the centromere. Green tracts represent scaffolded sequences and blue tracts show the contigs that comprise this scaffold with contigs jittered across three levels. **b** Contig NG(x). **c** Scaffold NG(x). **d** Benchmarking Universal Single-Copy Orthologs (BUSCO) of Pilon-polished assemblies. **e** The number of assembly misjoins revealed by DLE-1 conflicts and the number of contigs of each assembly. **f** Regional LTR Assembly Index (LAI) values estimated based on 3 Mb windows with 300 kb steps. The box shows the median, upper, and lower quartiles. Whiskers indicate values  $\leq 1.5\times$  interquartile range. Outliers are plotted as dots. **g** Unique mapping rate of RNA-seq libraries. A total of ten tissues with two biological replicates were sequenced. Each dot represents an RNA-seq library. The box shows the median, upper, and lower quartiles. Whiskers indicate values  $\leq 1.5\times$  interquartile range. **h** Central processing unit (CPU) core hours required for Falcon correction and Canu assembly. **i** Bionano optical map inconsistency. Deletions and insertions are cases where sequences are shorter or longer than the size estimated by the optical map, respectively. Structural variations (SVs) are plotted as dots. The box shows the median, upper and lower quartiles. Whiskers indicate values  $\leq 1.5\times$  interquartile range. Source data underlying Fig. 1b, c, f, g are provided as a Source Data file.

longest LTR sequence clusters using the Bionano optical map and found most assemblies contained no gaps and were virtually complete (Fig. 2c), with the notable exception that the 11k\_50x, 16k\_50x, and 21k\_20x assemblies, which contained large gaps in one of the LTR clusters (Supplementary Table 11). We also inspected the *bz* locus<sup>32</sup>, which has highly nested transposon insertions and an estimated size of 303.5 kb in NC358. The *bz*

locus was well assembled in all but the 21k\_20x assembly, in which only 56.3% of the sequence was included (Supplementary Table 12). In summary, with  $\geq 40\times$  of sequence coverage, long-read sequencing and assembly can traverse most transposon-rich genomic regions including relatively long LTR sequences, although with shorter reads (i.e., read N50 of 11–16 kb) this sequencing depth may not be sufficient.



**Fig. 2 Assembly of repetitive components in the NC358 genome.** **a** The assembled size of the 180-bp knob repeat, the knob TR-1 element, the chromosome 6 nucleolus organizer region (NOR) region, CentC arrays, and subtelomere arrays in each of the NC358 assemblies. **b** Length distribution of long terminal repeat (LTR) retrotransposons longer than 26 kb. Each dot represents an annotated sequence. The box shows the median, upper and lower quartiles. Whiskers indicate values  $\leq 1.5 \times$  interquartile range. **d** Telomere 7-mer counts in telomere regions of NC358 assemblies. Assembly of **(c)** LTR retrotransposons, **(e)** CentC arrays, **(f)** the chromosome 6 NOR region, **(g)** the *Rp1-D* and *zein* tandem gene arrays, and **(h)** two example knobs in each of the NC358 assemblies. The NC358 Bionano optical map was used to estimate the size of these components. Ngap, estimated gap size. Source data underlying Fig. 2a, b, i are provided as a Source Data file.

**Evaluation of the non-TE repeat assembly space.** The assembly of non-TE tandem repeat space was also evaluated, including telomeres (7 bp repeats), subtelomeres (300–1300 bp repeats), CentC arrays (156 bp repeats), nucleolus organizer region (NOR, ~11 kb repeats), and the two major knob repeats (mixture of 180 and 350 bp repeats) (Fig. 2a and Supplementary Table 13). The effects of sequence read depth and sequence read length were far more pronounced across many of these tandemly duplicated portions of the genome (Fig. 2a).

Telomeres are characterized by 7 bp tandem repeats at the end of each chromosome. Our results showed a substantial increase in the assembled length of telomere sequence with the increase of both read length and sequence coverage (Fig. 2d and Supplementary Table 14). However, a precise estimate of telomere length was not possible with our optical map due to the lack of Bionano DLE-1 sites in these highly repetitive regions. Using the full dataset (21k\_75x), only 10 of 20 telomere–subtelomere combined regions were assembled to >90% of the Bionano estimated size

(Supplementary Table 15), suggesting even longer reads and higher coverage are required for the full assembly of these regions.

The centromere is one of the most repetitive regions of many species’ genomes including maize. We characterized NC358 centromeres based on CentC arrays<sup>33</sup>, which are abundant in functional centromeric regions<sup>34</sup>. Even with the full dataset (21k\_75x), only half of CentC arrays were assembled (Fig. 2e, Supplementary Fig. 12, and Supplementary Table 16). Hybrid scaffolded assemblies with sequence coverage  $\geq 60 \times$  yielded a better approximation to the Bionano estimated size, even though these regions largely consisted of gaps (Fig. 2e). Although assembled sequences were not significantly increased, higher sequence depth resulted in better anchoring of sequences with the Bionano optical map. Only three centromeres, which contained a mixture of CentC arrays, transposons, and intergenic sequences, could be traversed by Bionano DLE-1 labeling due to having a comparatively higher content of low-copy sequence<sup>34</sup>. The size of the remaining centromeres was likely underestimated

(Supplementary Fig. 13) and further improvements in scaffolding technology are required to traversing these regions.

The NOR is enriched with ribosomal DNA (rDNA) and spans approximately 9 Mb on chromosome 6 of NC358 (Supplementary Table 17). Longer read length improved the assembly of this region, but substantial differences were not observed with coverage  $\geq 30\times$  (Fig. 2f). Approximately 72% of the NOR was included in the 21k\_30x assembly and this improved by just 9% to 81% in the 21k\_75x assembly (Supplementary Table 17 and Fig. 2f).

Finally, maize knobs are heterochromatic regions consisting of 180 bp (knob180) and 350 bp (TR-1) repeats<sup>35</sup>. We used the Bionano optical map to assess the assembly of two knobs that together spanned a total of 5 Mb. With longer reads and higher coverage, more knob sequences were assembled, with 6.5% of the two knobs present in the 21k\_20x assembly and up to 65% in the 21k\_75x assembly (Supplementary Table 18 and Fig. 2h).

## Discussion

Recent innovations in long-read and scaffolding technology have made highly contiguous assembly possible across a wide range of species. We have documented how both the completeness and contiguity of assemblies improve with increasing depth and read length. The biological aims of an investigation must be considered when determining the level of investment in depth of sequence. With long-read sequencing, the low-copy gene space (including tandem gene arrays) can be well assembled with as low as  $30\times$  genomic coverage across a range of read lengths. Complete characterization of TEs in complex genomes such as maize will require a greater depth of sequence ( $\sim 40\times$ ) and should employ library preparation protocols that maximize read-length N50. Finally, complete assembly of highly repetitive genomic features such as heterochromatic knobs, telomeres, and centromeres will require substantially more data<sup>36</sup>. In fact, complete assembly of these latter highly repetitive sequences will likely require innovations beyond current sequencing technology.

## Methods

**Sample preparation.** Seeds for the maize NC358 inbred line were obtained from GRIN Global (seed stock ID Ames 27175), grown, and self-pollinated at Iowa State University in 2017. A total of 144 seedlings derived from a single selfed ear were grown in the greenhouse. Leaf tissues from the seedlings at the Vegetative 2 (V2) growth stage were sampled after a 48-hour dark treatment to reduce carbohydrates. A total of 35 g of tissue was collected and flash-frozen. Tissue was sent to the Arizona Genomics Institute (AGI) for high-molecular-weight DNA isolation using a CTAB protocol<sup>37</sup>.

**Illumina and PacBio sequencing.** Pacific BioSciences long-read data for NC358 were generated at AGI using the Sequel platform. Libraries were prepared using the manufacturer's suggested protocol (<https://www.pacb.com/>). The subreads that were generated covered the genome at an estimated 75-fold depth ( $75\times$ ) with a subread N50 of 21,166 bp. Reads from each SMRT cell were inspected and quality metrics were calculated using SequelQC<sup>38</sup>. After validating the PSR (polymerase to subread ratio) and ZOR (ZMW occupancy ratio) were satisfactory, all subreads were used for subsequent steps.

Paired-end Illumina data for NC358 were generated at the Georgia Genomics and Bioinformatics Core from the same DNA extraction as was used for the long-read sequencing. Quality control of DNA was conducted using Qubit and Fragment Analyzer to determine the concentration and size distribution of the DNA. The library was constructed using the KAPA Hyper Prep Kit (catalog number KK8504). During library preparation, DNA was fragmented by acoustic shearing with Covaris before end repair and A-tailing. Barcoded adaptors were ligated to DNA fragments to form the final sequencing library. Libraries were purified and cleaned with solid phase reversible immobilization (SPRI) beads before being amplified with PCR. Final libraries underwent another bead cleanup before being evaluated by Qubit, quantitative PCR (qPCR) (KAPA Library Quantification Kit catalog number KK4854), and Fragment Analyzer. The final pool undergoing Illumina's Dilute and Denature Libraries protocol was diluted to 2.2 pM for loading onto the sequencer and then sequenced with 1% PhiX by volume. Libraries were sequenced on the NextSeq 500 instrument using PE150 cycles. The demultiplexing was done on Illumina's BaseSpace.

PacBio SMRT subreads for the maize inbred line B73 (sequenced to  $68\times$  depth) were retrieved from the NCBI SRA database with accession ID SRX1472849

(ref. 13). PacBio SMRT subreads for the human HG002 sample (sequenced to  $147\times$  depth) were retrieved with accession IDs SRX1033793 and SRX1033794 (ref. 29).

**Downsampling full dataset.** The  $75\times$  SMRT Sequel full data from maize NC358 was downsampled to  $60\times$ ,  $50\times$ ,  $40\times$ ,  $30\times$ , and  $20\times$  data using seqtk (v1.2) (<https://github.com/lh3/seqtk>). Downsampling was performed as serial titration, in which each dataset was the superset of the next smaller dataset, and was sampled to have similar length distributions (Supplementary Fig. 3). The N50 of the downsampled datasets were almost identical to the N50 of the full  $75\times$  data (Table 1). There is a chance that the subsampling process could pick a disproportionate number of poor reads and result in a poor assembly. There is also a chance that the subsampling process could pick, on average, better quality reads and result in a higher quality assembly. However, evidence suggests it is more likely that the small sample size (i.e.,  $20\times$ ) could not provide enough sequence information to reconstruct the original genome, as has been observed in a study based on simulated and empirical long-read data to test the effect of different depths of Prokaryotic genome assemblies<sup>19</sup>.

**Shifting read-length distribution of subreads.** Two more NC358 datasets were downsampled and trimmed from the original  $75\times$  SMRT dataset to match the read-length distribution of the maize B73 data<sup>13</sup> and the human HG002 data<sup>29</sup>, which had read N50 lengths of  $\sim 16$  kb and  $\sim 11$  kb, respectively (Supplementary Fig. 3). To do this, first, the read lengths of the maize B73 and human HG002 data were each sorted in descending order. For each read-length value, all subreads from NC358 that were longer than said value were randomly sampled without replacement and clipped to have matched read length. The unused clipped part of the read was put back in the pool for further use with short read length. This distribution-shifting approach was chosen to achieve a realistic distribution of read length rather than trimming all reads by fixed lengths. These datasets were labeled as 16k and 11k based on their N50 of subread data of 16,765, and 11,092, respectively.

**RNA tissue sampling and sequencing.** Samples from ten tissues throughout development were collected to generate expression evidence for gene annotation. Two biological replicates were collected for each tissue type and each replicate consisted of three individual plants. The tissues that were sampled were as follows: (1) primary root at 6 days after planting; (2) shoot and coleoptile at 6 days after planting; (3) base of the tenth leaf at the Vegetative 11 (V11) growth stage; (4) middle of the tenth leaf at the V11 growth stage; (5) tip of the tenth leaf at the V11 growth stage; (6) meiotic tassel at the Vegetative 18 (V18) growth stage; (7) immature ear at the V18 growth stage; (8) anthers at the Reproductive 1 (R1) growth stage; (9) endosperm at 16 days after pollination; and (10) embryo at 16 days after pollination. Tissue from developmental stage V11 and older were taken from field-grown plants, while all younger tissue samples were taken from greenhouse-grown plants. For the endosperm and embryo samples, tissue from 50 kernels per plant (150 total per biological replicate) were sampled. Greenhouse-grown plants were planted in Metro-Mix300 (Sun Gro Horticulture) with no additional fertilizer and grown under greenhouse conditions ( $27^\circ\text{C}/24^\circ\text{C}$  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 inch row spacing at  $\sim 52,000$  plants per hectare. RNA was extracted using the Qiagen RNeasy plant mini kit following the manufacturer's suggested protocol.

The quality of the total RNA was assessed by Bioanalyzer or Fragment analyzer to determine RNA concentration and integrity. The sample concentration was normalized in 25  $\mu\text{L}$  of nuclease-free  $\text{H}_2\text{O}$  before library preparation. Libraries were prepared using KAPA's stranded mRNA-seq kit with halved reaction volumes. During library preparations, mRNA was selected using oligo-dT beads, the RNA was fragmented, and cDNA was generated using random hexamer priming. Single or dual indices were ligated depending on the desired level of multiplexing. The number of cycles for library PCR was determined based on kit recommendations for the amount of total RNA used during library preparation. Libraries were quality control checked using Qubit or plate reader, depending on the number of samples in the batch for library concentration, and fragment analyzer for the size distribution of the library. The pooling of samples was based on qPCR. The pooled libraries were then checked by Qubit, Fragment Analyzer, and qPCR.

RNA libraries were prepared for sequencing on Illumina instruments using Illumina's Dilute and Denature protocol. Pooled libraries were diluted to 4 nM, then denatured using NaOH. The denatured library was further diluted to 2.2 pM and PhiX was added at 1% of the library volume. RNA pools were sequenced on a NextSeq 550 to generate 75 bp pair-end reads. On average, 24.5 million pair-end reads were generated per replicate per tissue type, for a total of 489 million reads across all samples. Data were demultiplexed and trimmed of adapter and barcode sequences on BaseSpace (Supplementary Fig. 14).

**Bionano data generation.** The DNA extraction was performed using the Bionano Prep™ Plant Tissue DNA Isolation Kit according to a modified version of the Plant Tissue DNA Isolation Base Protocol. Approximately 0.5 g leaf tissue was collected from young etiolated seedlings germinated in soil-free conditions and grown in the dark for  $\sim 2$  weeks after germination. Freshly cut leaves were treated with a 2%

formaldehyde fixing solution and then washed, cut into small pieces and homogenized using a Qiagen TissueRuptor probe. Free nuclei were concentrated by centrifugation at 2000 × g, washed, isolated by gradient centrifugation, and embedded into a low-melting-point agarose plug. After proteinase K and RNase A treatments, the agarose plug was washed four times in Wash Buffer and five times in TE (Tris and EDTA) buffer. Finally, purified ultra-high-molecular weight nuclear DNA (uHMW nDNA) was recovered by melting the plug, digesting it with agarase and subjecting the resulting sample to drop dialysis against TE.

The Bionano Saphyr platform, in combination with the Direct Label and Stain (DLS) process, was used to generate chromosome-level sequence scaffolds and validate PacBio sequence contigs. Direct labeling was performed using the Direct Labeling and Staining Kit (Bionano Genomics Catalog 80005) according to the manufacturer's recommendations, with some modifications<sup>39</sup>. In total, 1 µg uHMW nDNA was incubated for 2:20 h at 37 °C, followed by 20 min at 70 °C in the presence of DLE-1 Enzyme, DL-Green, and DLE-1 Buffer. Following proteinase K digestion and cleanup of the unincorporated DL-Green label, the labeled DNA was combined with Flow Buffer, DL-Dithiothreitol (DTT), and incubated overnight at 4 °C. DNA was quantified and stained by adding Bionano DNA Stain to a final concentration of 1 µL per 0.1 µg of final DNA. The labeled sample was loaded onto a Bionano chip flow cell and molecules separated, imaged and digitized in a Bionano Genomics Saphyr System and server according to the manufacturer's recommendations (<https://bionanogenomics.com/support-page/saphyr-system/>).

Data visualization, processing, DLS map assembly, and hybrid scaffold construction were all performed using the Bionano Genomics software Access, Solve, and Tools. A filtered subset of 1,282,746 molecules (353,596 Mb total length) with a minimum size of 150 kb and a maximum size of 3 Mb were assembled without pre-assembly using the non-haplotype parameters with no CMPR cut and without extend-split.

**Genome assembly.** To determine the assembly approach to apply to each of the datasets, six different methods were tested on the complete dataset, including Falcon-only, Canu-only, WTDBG2-only, a Falcon–Canu hybrid approach, a Falcon–Peregrine hybrid approach (the longest 23× corrected reads were used), and a Falcon–Flye hybrid approach. We also downloaded PacBio sequencing data for the B73 v4 genome (68× subreads) for comparison of the different approaches.

The Falcon genome assemblies were performed using the `falcon_kit` pipeline v0.7 (ref. 21) with some modifications. TANmask and REPmask were not used due to their extensive masking for the maize genome. Error correction for subreads was performed on the longest 50× coverage, with the average read correction rate set to 75% (-e 0.75) and local alignments for at least 3000 bp (-l 3000). The usage of -l 3000 instead of -l 2500 was done because of the omitted repeat masking, which works better for highly repetitive genome species like maize. A minimum of two reads and a maximum of 200 reads were used for error corrections (-min\_cov 2 --max\_n\_read 200). For sequence assembly, the exact matching k-mers between two reads was set to 24 bp (-k 24) with read correction rate as 95% (-e 0.95) and local alignments at least 1000 bp (-l 1000). The longest 20× coverage reads were used for assembly with a minimum coverage of two and maximum coverage of 80 (-min\_cov 2 --max\_cov 80). Full parameter sets are included in the Supplementary Note 1.

For Canu read correction and assembly, Canu v1.7 (ref. 22) was used. K-mers more frequent than 500 were not used to seed overlaps (ovlMerThreshold = 500). The genome size of 2,272,400,000 bp and 2,500,000,000 bp for NC358 and B73, respectively, were used in this study<sup>20</sup>. Other parameters were used as default. Due to a bug in the Canu v1.7 program, truncations of large contigs would occur during the consensus process (<https://github.com/marbl/canu/releases/tag/v1.8>). As the program was not expecting the superlong contigs that were being generated for our NC358 assemblies, we found a total of nine large contigs that suffered from consensus truncations. To fix these truncation gaps, consensus-free contigs were generated using Canu v1.7 (cnsConsensus = quick), then blastn was used to search for 5 kb boundaries of truncation gaps in consensus-free assemblies. Truncated sequences were retrieved and patched to the truncated contigs.

For the Falcon–Canu hybrid approach, the error correction was performed by Falcon, and the trimming and assembly were performed by Canu using the versions and parameters described above. All the assemblies were performed on the DNAnexus cloud platform. CPU core hour and maximum memory usage were recorded every 10 minutes for each Falcon error correction and Canu assembly job. For Falcon error correction of the 21k datasets, the CPU core hour ( $y$ ) could be predicted by subread depth ( $m$ ) with

$$y = 20603100000 + (3136.685 - 20603100000)/(1 + (m/1932.377)^4.148144) \quad (1)$$

For Canu assembly of the 21k datasets, the CPU core hour ( $y$ ) could be predicted by corrected read depth ( $n$ ) with

$$y = 6438752000 + (1284.689 - 6438752000)/(1 + (n/56334.74)^1.872455) \quad (2)$$

These curves were fit using the <https://mycurvfit.com/> website and plotted in R.

For the Falcon–Flye hybrid approach, the Falcon-corrected reads (44×) were assembled by Flye (v2.6)<sup>24</sup> with the genome size parameter set to 2,272,000,000. For the Falcon–Peregrine hybrid approach, the longest 23× Falcon-corrected reads were assembled by Peregrine (pg0.1.6.1)<sup>25</sup> with parameters 24 24 24 12 24 24 24 24

24 --with-consensus --shimmer-w 80 --shimmer-r 3 --best\_n\_ovlp 24 --mc\_upper 640. The WTDBG2 genome assembly was performed using the WTDBG2 pipeline (v2.5)<sup>23</sup> with default parameters using the 75× uncorrected subreads. The estimated genome size was set to 2,272,400,000.

Three of the assembly approaches were tested using both maize NC358 and B73 reads. For both inbred lines, a similar assembly size was generated by the Falcon-only, Canu-only, and Falcon–Canu hybrid approaches. However, the Falcon–Canu hybrid approach yielded the longest contig length (78.4 Mb and 19.7 Mb, respectively), the highest contig NG50 (23.0 Mb and 3.0 Mb, respectively), and the lowest number of assembly errors based on Bionano DLE-1 conflict (21 and 64, respectively; Supplementary Table 2). The gene-space completeness evaluated using BUSCOs v3.0.2<sup>28</sup> and the repeat space continuity evaluated using the LAI (vbeta3.2)<sup>8</sup> were similar between the Canu and the hybrid approach and higher than those assemblies that were created using the Falcon assembler (Supplementary Table 2). This was likely due to the consensus approach used at the end of the Canu program, which was missing in the Falcon program.

The remaining three approaches were tested using only NC358 reads. For assemblies generated by the Falcon–Flye hybrid, Falcon–Peregrine hybrid, and WTDBG2-only approaches, the assembled sizes were 16–51% larger than the estimated genome size with 28–202 times more contigs compared to the assembly generated by the Falcon–Canu hybrid approach (Supplementary Table 2). Other quality metrics of these assemblies, such as the longest contig, contig N50, Bionano DLE-1 conflict, and LAI were all inferior compared to those of the Falcon–Canu assembly (Supplementary Table 2). BUSCO of the Falcon–Peregrine assembly was higher than that of the Falcon–Canu raw assembly, but lower than that of the Pilon or Arrow-polished assembly. The correction-free WTDBG2 assembly was the most fragmented probably due to the lack of error correction that hindered the assembly of repetitive sequences, which is demonstrated in the low LAI value (LAI = 2.5).

Due to the consistently high quality of the assemblies generated from the Falcon–Canu hybrid approach, we used this approach to assemble each of the NC358 datasets with varying sequence depth and read length. Full parameter sets of all assembly approaches are included in the Supplementary Note 1.

**Genome polishing.** Two polishing approaches were tested on the 21k\_75x assembly. The first was done using Arrow with PacBio subreads (75× coverage). Read mapping to the assembly was done using BLASR<sup>40</sup> with default parameters (--minMatch 12 --bestn 10 --minPctSimilarity 70.0 --refineConcordantAlignments). The Arrow tool in the SMRT Link (v5.1.0) software package was then applied to correct for sequencing errors with default parameters. A second approach for polishing was done using Pilon with Illumina pair-end reads (30.7× coverage). Read mapping to the assembly was done using Minimap2 (v2.16)<sup>41</sup> with the short read option (-ax sr). Pilon (v1.23.0)<sup>42</sup> was then applied to correct for sequencing errors including SNPs and small indels (--fix bases) on sites with a minimum depth of 10 and a minimum mapping quality of 30 (--mindepth 10 --minmq 30).

With both approaches, minimal differences were observed in the contiguity statistics (Supplementary Table 2) or the repeat content for the 21k\_75x assembly (Supplementary Fig. 15), and it is expected that this minimal impact would be observed across all of the NC358 assemblies. A more substantial difference in BUSCO scores was observed with both the Arrow-polished and the Pilon-polished 21k\_75x assemblies (Supplementary Table 2). As the polishing had a substantial impact on this metric, the other NC358 assemblies were also polished using Pilon with the same parameter settings and similar improvement of BUSCO scores was observed (Table 1; Supplementary Table 4).

**Generation of pseudomolecules.** Hybrid scaffolds for the assemblies were generated with Bionano Direct Label and Stain data using Bionano Solve (v3.2.1\_04122018). Overlaps of contigs within Bionano map space were resolved by placing 13 bp of Ns (13 N gaps) at the overlap site. In addition to arranging contigs into scaffolds, the hybrid scaffold was also used to detect misassembly and to assess the completeness of the assembled genome and repeat elements.

The pseudomolecules were constructed from the hybrid scaffolds using ALLMAPS (v0.8.12)<sup>43</sup>. Both pan-genome anchor markers<sup>27</sup> and GoldenGate markers<sup>26</sup> were used with equal weights for ordering and orientating the scaffolds. For pan-genome anchor markers, data were downloaded from the CyVerse Data Commons ([http://datacommons.cyverse.org/browse/iplant/home/shared/panzea/genotypes/GBS/v27/Lu\\_2015\\_NatCommun\\_panGenomeAnchors20150219.txt.gz](http://datacommons.cyverse.org/browse/iplant/home/shared/panzea/genotypes/GBS/v27/Lu_2015_NatCommun_panGenomeAnchors20150219.txt.gz)) and a bed file with 50 bp upstream and downstream of the B73 v3 coordinates were generated. A text file with marker name and predicted distance was also constructed from the same file. The extracted markers were mapped to HiSat2 (v2.1.0)<sup>44</sup> indexed assemblies of NC358 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 of 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 GoldenGate markers were downloaded from MaizeGDB ([https://www.maizegdb.org/data\\_center/map?id=1203673](https://www.maizegdb.org/data_center/map?id=1203673)). 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 using it in ALLMAPS. ALLMAPS was run with default options, and the pseudomolecules were finalized after inspecting the marker placement plot and the scaffold directions. Synteny dotplots were generated using the scaffolds as well as pseudomolecule assemblies against the B73 genome by following the ISUgenomics Bioinformatics Workbook (<https://bioinformaticsworkbook.org/dataWrangling/genome-dotplots.html>)<sup>45</sup>. Briefly, the repeats were masked using RepeatMasker (v4.0.9)<sup>46</sup> and the Maize TE Consortium (MTEC) curated library<sup>3</sup>. RepeatMasker was configured to use the NCBI engine (rmbblastn) with a quick search option (-q) and GFF as a preferred output. The repeat-masked genomes were then aligned using Minimap2 (v2.2)<sup>41</sup> and set to break at 5% divergence (-x asm5). The paf files were filtered to eliminate alignments less than 1 kb and dotplots were generated using the R package dotPlotly (<https://github.com/tpoorten/dotPlotly>).

**Gene annotation and RNA-seq mapping.** The MAKER-P pipeline<sup>47</sup> was used to annotate protein-coding genes for Pilon-polished NC358 21k\_20x and 21k\_75x genome assemblies. The baseline evidence used in annotating the B73 v4 genome<sup>13</sup> was applied. Before gene annotation, the MTEC curated TE library<sup>3</sup> and RepeatMasker was used to mask repetitive sequences. For gene prediction, we used Augustus<sup>48</sup> and FGENESH<sup>49</sup> (<http://www.softberry.com/berry.phtml>) with training sets based on maize and monocots, respectively. To identify genes that were missing in the 21k\_20x assembly, total coding sequences (CDS) from the 21k\_75x annotation was masked by total CDS from the 21k\_20x annotation using RepeatMasker (-div 2 -cutoff 1000 -q -no\_is -norma -nolow). The 21k\_75x CDS that were masked less than 20% were determined missing in the 21k\_20x annotation. These missing CDS were blast against the 21k\_20x assembly and those that had less than 20% similarity were also determined to be missing in the 21k\_20x assembly.

A total of 20 RNA-seq libraries were sequenced from NC358 tissue samples. Each library was sequenced to  $21.9 \times \pm 0.7 \times$  coverage with a mapping rate of  $86.4\% \pm 1.0\%$  to the B73 v4 using STAR (v2.5.2b)<sup>50</sup> (Supplementary Fig. 16 and Supplementary Table 19). To benchmark the gene-space assembly, STAR (v2.5.2b)<sup>50</sup> was used to map the RNA-seq reads against the Pilon-polished NC358 assemblies. Unmapped reads from the 21k\_20x assembly were extracted using SAMtools<sup>51</sup> and remapped to the 21k\_75x assembly with STAR. Genes with read coverage  $\geq 20\%$  were extracted using BEDtools<sup>52</sup>, and blast against the 21k\_20x assembly for the identification of full-length copies. The NC358 TE library (see next section for details on library generation) was used to identify TE fragments in genes with aligned reads (Supplementary Table 9). In addition, TESorter (v1.1.4)<sup>53</sup> (<https://github.com/zhangrengang/TEsorter>) was used to identify TE-related protein domains in genes with default parameters (Supplementary Table 9).

**Assessment of genome assembly quality.** The quality of the different NC358 assemblies was assessed on the unpolished assemblies unless noted. For continuity, N50, NG50, NG(x), the number of contigs, and maximum contig length were estimated. NG(x) values were the length of the contig at the top x percent of the estimated genome size (2.2724 Gb) consisting of the longest contigs. NG50 is a commonly used case of NG(x) values. NG(x) values were calculated using GenomeQC (<https://github.com/HuffordLab/GenomeQC>)<sup>54</sup>. The gene-space completeness was estimated using BUSCO (v3.0.2)<sup>28</sup> with the Embryophyta odb9 dataset ( $n = 1440$ ) and BLAST (v2.6)<sup>55</sup>, Augustus (v3.3)<sup>48</sup>, EMBOSS (v6.6.0)<sup>56</sup>, and HMMER (v3.1b2)<sup>57</sup>.

The repeat space contiguity was accessed using the LAI (vbeta3.2)<sup>8</sup>. To annotate LTR retrotransposons, LTR\_retriever (v2.6)<sup>58</sup> was used to identify intact LTR retrotransposons and construct LTR libraries for each NC358 assembly with default parameters. To generate a high-quality LTR library for NC358, assembly-specific LTR libraries were aggregated and masked by the MTEC curated LTR library using RepeatMasker (v4.0.7)<sup>46</sup>. Library sequences masked over 90% were removed and redundant sequences were also removed using utility scripts (cleanup\_tandem.pl and cleanup\_tandem.pl) from the EDTA package<sup>59</sup>. Non-redundant NC358-specific LTR sequences were added to the MTEC curated LTR library to form the final LTR library for NC358. The final library was then used to mask the 21k\_75x assembly for the estimation of total LTR content. The total LTR content of 76.34% and LTR identity of 94.854% was used to estimate LAI values of all NC358 assemblies (-totLTR 76.34 -iden 94.854). The LAI of the other maize line genomes, including PH207 (GeneBank Accession: GCA\_002237485.1)<sup>60</sup>, CML247 (GeneBank Accession: GCA\_002682915.2)<sup>27</sup>, Mo17<sup>61</sup>, and GeneBank Accession: GCA\_003185045.1<sup>62</sup>, W22 (GeneBank Accession: GCA\_001644905.2)<sup>63</sup>, and B73 v4 (GeneBank Accession: GCA\_000005005.6)<sup>13</sup> were also evaluated for context.

Effective assembly size, which is the length of the uniquely mappable sequences of an assembly, was estimated using unique 150-mers in each sequence assembly and quantified using Jellyfish (v2.0)<sup>64</sup> with default parameters.

**Misassemblies identification with optical maps.** The Bionano optical mapping was used as an orthogonal method to identify misassemblies in genomes. Bionano de novo assembled optical maps were aligned to the sequence pseudomolecules to characterize structural inconsistencies using the structural variant calling pipeline of BionanoSolve 3.4. Default parameters were employed from the non-haplotype\_noES\_DLE file. Homozygous calls with a confidence of 0.1, a size of 500 bp, and non-overlaps with gap regions were regarded as insertions and deletions in sequence assemblies.

**Assembly quality evaluation in repeat space.** The coordinates of CentC arrays, knob180, TR-1 knobs, and NOR in the assemblies were identified by blasting CentC, knob180, TR-1 knob consensus sequences<sup>34</sup>, and the rDNA intergenic spacer (AF013103.1) against each assembly. An individual repeat array was defined as clusters of repetitive sequences that had less than 100 kb interspace between repeated elements. The level of repeats and gaps were then quantified in each defined repeat array. Respective sizes of each repeat array in the Bionano maps were estimated using the Bionano labels closest to the start and end coordinates in the assemblies.

To identify the telomere-subtelomere boundaries of the NC358 assemblies, seven maize subtelomere repeat sequences were downloaded from NCBI (EU253568.1, S46927.1, S46926.1, S46925.1, CL569186.1, AF020266.1, and AF020265.1) and used as queries to blast against the NC358 21k\_75x assembly. Subtelomere boundaries were first identified at the start and end of chromosomes where blast hits were clustering then cross-checked with subtelomere-specific fluorescence in-situ hybridization (FISH) data<sup>65</sup>. The blast results were concordant with FISH results, showing the beginning of chromosomes 7, 8, 9, and 10 lack subtelomeres (Supplementary Table 15). Telomeres were defined as the distance between the boundary of subtelomeres to the end of pseudomolecules of the 21k\_75x assembly, which were used as the basis for estimating the telomere size and count of the telomeric repeat sequences (5'-TTTAGGG-3' and 5'-CCCTAAA-3' in reverse complementation) in all other NC358 assemblies.

To identify the *bz* locus in the NC358 assemblies, the sequence of the maize W22 *bz* locus was first downloaded from NCBI (EU338354.1)<sup>32</sup>. The starting and ending 2 kb of the W22 *bz* locus were used to blast against the NC358 21k\_75x assembly and the longest matches on chromosome 9 were used as the location of the *bz* locus in the NC358 21k\_75x assembly. The obtained NC358 *bz* locus is 289,103 bp in length (chr9:11625031.11914133), which is 50 kb longer than that of the W22 *bz* locus (238,141 bp). Similarly, the 2 kb flanking sequences of the NC358 21k\_75x *bz* locus were used to locate the *bz* locus coordinates in the other NC358 assemblies.

The *zein* sequence was downloaded from NCBI (AF031569.1) and the *Rp1-D* from MaizeGDB (AC152495.1\_FG002). The same method as described for the *bz* locus was used to identify coordinates in the NC358 assemblies based on blast results using 2 kb flanking sequences.

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

## Data availability

Data supporting the findings of this work are available within the paper and its Supplementary Information files. A reporting summary for this Article is available as a Supplementary Information file. The genome assemblies and gene model annotations generated and analyzed during the current study are available for download through CyVerse ([https://de.cyverse.org/dl/d/42938190-C7AF-4BF5-B953-0BB28F61887D/NC358\\_genome\\_annotation\\_April2020\\_public\\_release.tar](https://de.cyverse.org/dl/d/42938190-C7AF-4BF5-B953-0BB28F61887D/NC358_genome_annotation_April2020_public_release.tar)). PacBio and Illumina sequencing reads for the NC358 line used in this study are available with EBI Biosample ID ERS3120561 (<https://www.ebi.ac.uk/ena/data/view/ERS3120561>). PacBio SMRT subreads for the maize inbred line B73 (sequenced to 68x depth) were retrieved from the NCBI SRA database with accession ID SRX1472849 (<https://www.ncbi.nlm.nih.gov/sra/SRX1472849>). PacBio SMRT subreads for the human HG002 sample (sequenced to 147x depth) were retrieved with accession IDs SRX1033793 (<https://www.ncbi.nlm.nih.gov/sra/SRX1033793>) and SRX1033794 (<https://www.ncbi.nlm.nih.gov/sra/SRX1033794>). The source data underlying Figs. 1b, c, f, g and 2a, b, i, and Supplementary Figs. 2, 4, 6–11, 15, and 16, as well as Supplementary Table 3 are provided as a Source Data file.

## Code availability

All code developed for this study is available at [https://github.com/HuffordLab/Maize\\_NC358](https://github.com/HuffordLab/Maize_NC358).

Received: 23 November 2019; Accepted: 9 April 2020;

Published online: 08 May 2020

## References

- Adams, M. D. et al. The genome sequence of *Drosophila melanogaster*. *Science* **287**, 2185–2195 (2000).



2. Lander, E. S. et al. Initial sequencing and analysis of the human genome. *Nature* **409**, 860–921 (2001).
3. Schnable, P. S. et al. The B73 maize genome: complexity, diversity, and dynamics. *Science* **326**, 1112–1115 (2009).
4. Yu, J. A draft sequence of the rice genome (*Oryza sativa* L. ssp. indica). *Science* **296**, 79–92 (2002).
5. Jaillon, O. et al. The grapevine genome sequence suggests ancestral hexaploidization in major angiosperm phyla. *Nature* **449**, 463–467 (2007).
6. Ming, R. et al. Genome of the long-living sacred lotus (*Nelumbo nucifera* Gaertn.). *Genome Biol.* **14**, R41 (2013).
7. Alkan, C., Sajjadian, S. & Eichler, E. E. Limitations of next-generation genome sequence assembly. *Nat. Methods* **8**, 61–65 (2011).
8. Ou, S., Chen, J. & Jiang, N. Assessing genome assembly quality using the LTR Assembly Index (LAI). *Nucleic Acids Res.* **46**, e126 (2018).
9. Miga, K. H. et al. Telomere-to-telomere assembly of a complete human X chromosome. *bioRxiv* 735928 <https://doi.org/10.1101/735928> (2019).
10. Belser, C. et al. Chromosome-scale assemblies of plant genomes using nanopore long reads and optical maps. *Nat. Plants* **4**, 879–887 (2018).
11. Du, H. et al. Sequencing and de novo assembly of a near complete indica rice genome. *Nat. Commun.* **8**, 15324 (2017).
12. Jain, M. et al. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat. Biotechnol.* **36**, 338–345 (2018).
13. Jiao, Y. et al. Improved maize reference genome with single-molecule technologies. *Nature* **546**, 524–527 (2017).
14. Sun, S. et al. Extensive intraspecific gene order and gene structural variations between Mo17 and other maize genomes. *Nat. Genet.* **50**, 1289–1295 (2018).
15. Yang, N. et al. Genome assembly of a tropical maize inbred line provides insights into structural variation and crop improvement. *Nat. Genet.* **51**, 1052–1059 (2019).
16. Van Bel, M., Bucchini, F. & Vandepoele, K. Gene space completeness in complex plant genomes. *Curr. Opin. Plant Biol.* **48**, 9–17 (2019).
17. Jayakumar, V. & Sakakibara, Y. Comprehensive evaluation of non-hybrid genome assembly tools for third-generation PacBio long-read sequence data. *Brief. Bioinformatics* **20**, 866–876 (2019).
18. Goldstein, S., Bekas, L., Graf, J. & Klassen, J. L. Evaluation of strategies for the assembly of diverse bacterial genomes using MinION long-read sequencing. *BMC Genomics* **20**, 23 (2019).
19. Wick, R. R. & Holt, K. E. Benchmarking of long-read assemblers for prokaryote whole genome sequencing. *F1000Res* **8**, 2138 (2019).
20. Chia, J. -M. et al. Maize HapMap2 identifies extant variation from a genome in flux. *Nat. Genet.* **44**, 803–807 (2012).
21. Chin, C. -S. et al. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat. Methods* **13**, 1050–1054 (2016).
22. Koren, S. et al. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res.* **27**, 722–736 (2017).
23. Ruan, J. & Li, H. Fast and accurate long-read assembly with wtdbg2. *Nat. Methods* **17**, 155–158 (2019).
24. Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of long, error-prone reads using repeat graphs. *Nat. Biotechnol.* **37**, 540–546 (2019).
25. Chin, C. -S. & Khalak, A. Human genome assembly in 100 minutes. *bioRxiv* 705616 <https://doi.org/10.1101/705616> (2019).
26. Yan, J. et al. Genetic characterization and linkage disequilibrium estimation of a global maize collection using SNP markers. *PLoS ONE* **4**, e8451 (2009).
27. Lu, F. et al. High-resolution genetic mapping of maize pan-genome sequence anchors. *Nat. Commun.* **6**, 6914 (2015).
28. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**, 3210–3212 (2015).
29. Zook, J. M. et al. Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci. Data* **3**, 160025 (2016).
30. Collins, N. et al. Molecular characterization of the maize Rp1-D rust resistance haplotype and its mutants. *Plant Cell* **11**, 1365–1376 (1999).
31. Song, R., Ilaca, V., Linton, E. & Messing, J. Sequence, regulation, and evolution of the maize 22-kD alpha zein gene family. *Genome Res.* **11**, 1817–1825 (2001).
32. Dooner, H. K. & He, L. Maize genome structure variation: interplay between retrotransposon polymorphisms and genic recombination. *Plant Cell* **20**, 249–258 (2008).
33. Jin, W. et al. Maize centromeres: organization and functional adaptation in the genetic background of oat. *Plant Cell* **16**, 571–581 (2004).
34. Gent, J. I., Wang, N. & Dawe, R. K. Stable centromere positioning in diverse sequence contexts of complex and satellite centromeres of maize and wild relatives. *Genome Biol.* **18**, 121 (2017).
35. Santos-Serejo, J. A., Gardingo, J. R., Mondin, M. & Aguiar-Perecin, M. L. R. Alterations in heterochromatic knobs in maize callus culture by breakage-fusion-bridge cycle and unequal crossing over. *Cytogenet. Genome Res.* **154**, 107–118 (2018).
36. Liu, J. et al. Gapless assembly of maize chromosomes using long read technologies. *bioRxiv* 906230 <https://doi.org/10.1101/2020.01.14.906230> (2020).
37. Doyle, J. J. & Doyle, J. L. A rapid DNA isolation procedure for small quantities of fresh leaf tissue. *Phytochem. Bull.* <https://worldveg.tind.io/record/33886/> (1987).
38. Hufnagel, D. E., Hufford, M. B. & Seetharam, A. S. SequelQC: analyzing PacBio sequel raw sequence quality. *bioRxiv* 611814 <https://doi.org/10.1101/611814> (2019).
39. Deschamps, S. et al. A chromosome-scale assembly of the sorghum genome using nanopore sequencing and optical mapping. *Nat. Commun.* **9**, 4844 (2018).
40. Chaisson, M. J. & Tesler, G. Mapping single molecule sequencing reads using basic local alignment with successive refinement (BLASR): application and theory. *BMC Bioinformatics* **13**, 238 (2012).
41. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
42. Walker, B. J. et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS ONE* **9**, e112963 (2014).
43. Tang, H. et al. ALLMAPS: robust scaffold ordering based on multiple maps. *Genome Biol.* **16**, 3 (2015).
44. Kim, D., Paggi, J. M., Park, C., Bennett, C. & Salzberg, S. L. Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* **37**, 907–915 (2019).
45. Seetharam, A. et al. ISUGenomics/bioinformatics-workbook: 2019-10-11 Release of the Bioinformatics Workbook <https://doi.org/10.5281/zenodo.3482894> (2019).
46. Smit, A. F. A., Hubley, R. & Green, P. RepeatMasker Open-4.0. 2013–2015 <http://www.repeatmasker.org/> (2015).
47. Campbell, M. S., Holt, C., Moore, B. & Yandell, M. Genome annotation and curation using MAKER and MAKER-P. *Curr. Protoc. Bioinformatics* **48**, 4–11 (2014).
48. Stanke, M. et al. AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Res.* **34**, W435–9 (2006).
49. Salamov, A. & Solovyev, V. Fgenesh multiple gene prediction program. <http://www.softberry.com/berry.phtml> (1998).
50. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
51. Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
52. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
53. Zhang, R. -G., Wang, Z. -X., Ou, S. & Li, G. -Y. TEsorter: lineage-level classification of transposable elements using conserved protein domains. *bioRxiv* 800177 <https://doi.org/10.1101/800177> (2019).
54. Manchanda, N. et al. GenomeQC: a quality assessment tool for genome assemblies and gene structure annotations. *bioRxiv* 795237 <https://doi.org/10.1101/795237> (2019).
55. Johnson, M. et al. NCBI BLAST: a better web interface. *Nucleic Acids Res.* **36**, W5–W9 (2008).
56. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open Software Suite. *Trends Genet.* **16**, 276–277 (2000).
57. Mistry, J., Finn, R. D., Eddy, S. R., Bateman, A. & Punta, M. Challenges in homology search: HMMER3 and convergent evolution of coiled-coil regions. *Nucleic Acids Res.* **41**, e121 (2013).
58. Ou, S. & Jiang, N. LTR\_retriever: a highly accurate and sensitive program for identification of long terminal repeat retrotransposons. *Plant Physiol.* **176**, 1410–1422 (2018).
59. Ou, S. et al. Benchmarking transposable element annotation methods for creation of a streamlined, comprehensive pipeline. *Genome Biol.* **20**, 657890 (2019).
60. Hirsch, C. N. et al. Draft assembly of elite inbred line PH207 provides insights into genomic and transcriptome diversity in maize. *Plant Cell* **28**, 2700–2714 (2016).
61. Xin, M. et al. Dynamic expression of imprinted genes associates with maternally controlled nutrient allocation during maize endosperm development. *Plant Cell* **25**, 3212–3227 (2013).
62. Yang, N. et al. Contributions of Zea mays subspecies mexicana haplotypes to modern maize. *Nat. Commun.* **8**, 1874 (2017).
63. Springer, N. M. et al. The maize W22 genome provides a foundation for functional genomics and transposon biology. *Nat. Genet.* **50**, 1282–1288 (2018).
64. Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* **27**, 764–770 (2011).
65. Albert, P. S., Gao, Z., Danilova, T. V. & Birchler, J. A. Diversity of chromosomal karyotypes in maize and its relatives. *Cytogenet. Genome Res.* **129**, 6–16 (2010).

## Acknowledgements

This work was supported by NSF Plant Genome Research Program grant IOS-1744001 to R.K.D., D.W., and M.B.H., and grant IOS-1546727 to C.N.H., USDA ARS 5030-21000-

068-00D to M.W., and USDA ARS 58-8062-2100-044 to D.W. B.P.W., S.K., and A.M.P. were supported by the Intramural Research Program of the National Human Genome Research Institute. We acknowledge Jonathan Gent for helpful discussion on repeat space analyses. This research was supported in part by the U.S. Department of Agriculture, Agricultural Research Service. Mention of trade names or commercial products in this publication is solely for the purpose of providing specific information and does not imply recommendation or endorsement by the U.S. Department of Agriculture. USDA is an equal opportunity provider and employer.

### Author contributions

R.K.D., C.N.H., M.B.H., and D.W. supervised the research. R.K.D., C.N.H., M.B.H., D.W., S.O., and A.F. designed the research. A.F., A.S., V.L., C.-S.C., K.F., S.K., and A.M.P. performed genome assembly. S.O., J.L., K.M.C., A.F., A.S., J.C.S., V.L., N.M., A.M.G., S.W., D.E.H., S.P., S.J.S., B.P.W., and B.H. performed the research and analyzed data. S.O., A.S., and M.W. carried out data submission. S.O., J.L., K.M.C., A.F., A.S., V.L., N.M., C.N.H., and M.B.H. wrote and revised the manuscript. All authors read and approved the final version of this manuscript.

### Competing interests

C.S.-C. and A.F. are employees of DNAnexus, Inc. B.T.H. was an employee of DNAnexus, Inc. at the time research was completed. DNAnexus, Inc. offers genome assembly services. K.F. and V.L. are employees of Corteva Agriscience, a PacBio and Bionano Certified provider. All the other authors declare no competing interests.

### Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41467-020-16037-7>.

**Correspondence** and requests for materials should be addressed to R.K.D., C.N.H., M.B.H. or D.W.

**Peer review information** *Nature Communications* thanks Roger Barthelson, and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Peer reviewer reports are available.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

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



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

This is a U.S. government work and not under copyright protection in the U.S.; foreign copyright protection may apply 2020