

1 *Supplementary Information*

2

3

4

**Improved maize reference genome with single molecule technologies**

5

6

7

8

9

**1. Genome assembly**

10

11 ***De novo* assembly of the long reads:** Two assembly tools, PBcR-MHAP and FALCON,  
12 were independently evaluated for *de novo* assembly of PacBio SMRT Sequencing reads.  
13 For PBcR, following the recommended parameters for large genome assembly<sup>1</sup>, *k*-mer  
14 lengths of 16 and 14 were used to test the performance of assembler. The assembly  
15 redundancy in the unitigs were filtered according to sequencing coverage, according to  
16 the following criteria: coverage  $\geq 2$  reads, and a single read must not cross more 50% of a  
17 unitig. FALCON v.0.4 ([https://github.com/PacificBiosciences/FALCON-](https://github.com/PacificBiosciences/FALCON-integrate/tree/0.4.0)  
18 [integrate/tree/0.4.0](https://github.com/PacificBiosciences/FALCON-integrate/tree/0.4.0)) was also used for *de novo* assembly. The overall design of FALCON  
19 follows the hierarchical genome-assembly process<sup>2</sup>. Instead of BLASR, daligner was  
20 used to overlap reads. To lay out contigs from the assembly graph, the error-correction  
21 module was updated, and the Celera Assembler was replaced by a string graph-based  
22 module. Due to the highly repetitive nature of the maize genome, we adopted more  
23 aggressive parameters to reduce computation time. For the full data set, only reads longer  
24 than 12 kb were corrected. To identify overlaps between raw sequences, we used “-M24 -  
25 l4800 -k18 -h480 -w8” for Daligner. Using these parameters, only overlaps longer than  
26 4,800 bp were considered for error correction with seed matches  $> 480$  bp.

27

28

29

30

31

32

33

34

To ascertain the quality of the three independent assemblies (FALCON, PBcR with  $k=16$  and  $14$ ), the BioNano scaffolding pipeline NGM Hybrid Scaffold (NGM-HS) (version 4304) was used to generate an *in silico* map for sequence contigs from each assembly. The maps were aligned against the genome assemblies using RefAligner<sup>3,4</sup> to identify and resolve potential conflicts in sequence contigs or optical genome maps. The result showed that the PBcR-MHAP assembly ( $k=14$ ) had the fewest conflicts (Extended Data Figure 2a); consequently, it was adopted as the new B73 genome reference.

35 **Curation of the assembly:** Comparison between the contigs and optical map identified  
36 36 conflicts. Next Generation Mapping (NGM-HS) from BioNano Genomics' Irys®  
37 System was used to resolve conflicts between the sequence and optical map assemblies  
38 by cutting either assembly (option: -N 2 -B 2); cut decisions were based on chimeric  
39 scores of labels near the conflict junctions on the optical genome map. The chimeric  
40 score of a label represents the percentage of BioNano molecules that can fully align to the  
41 optical map 55 kb to the left or right of that label. If the chimeric scores of all labels  
42 within 10 kb of the conflict junction were  $\geq 35$ , the scaffolding pipeline suggested a cut  
43 in the sequence contig. If any label in the region had a chimeric score  $< 35$ , a cut was  
44 suggested in the BioNano optical map. All proposed cuts were manually evaluated using  
45 BioNano molecule-to-genome map alignments, molecule-to-sequence contig alignments,  
46 and the BAC-based fingerprint map. Of these 36 conflicts, 18 were chimeric in the long  
47 reads assembly, and 13 were chimeric in the optical map; five were left unresolved.

48 Using alignments of the optical genome map, a total of 1,369 overlaps were  
49 detected among the tails of the contigs. There are two possible reasons for this: the  
50 overlaps could be repeat boundaries between contigs from the Celera assembler<sup>5</sup>, or  
51 alternatively, nearly identical repeats could be over-collapsed in the optical map. The  
52 redundancy at the edges of nearby contigs generated by the Celera assembler was  
53 resolved as follows: if two contigs were detected to have overlap from 0.5-10 kb (based  
54 on the size of PacBio reads) by optical genome map and had sequence identity over 95%  
55 in the overlapped region, they were merged by Minimus2<sup>6</sup>. A total of 670 contigs were  
56 merged into 310 larger contigs.

57  
58 **Pseudomolecule construction:** The curated 2,958 contigs were scaffolded into 625 large  
59 hybrid scaffolds and 269 contigs that are relatively small were not covered by the optical  
60 maps. Using unique BAC sequences as markers, we could order and orient 315 hybrid-  
61 scaffolds and 25 non-scaffolded contigs. In addition, we also incorporated a genetic map  
62 built from an intermated maize recombinant inbred line population (Mo17  $\times$  B73)<sup>7</sup> to  
63 complement pseudomolecule construction and validation. In this new AGP (A Golden  
64 Path) of the maize reference genome, a total of 331 hybrid scaffolds and 45 non-  
65 scaffolded contigs were ordered and oriented. Of 1,907 markers on the genetic map,

66 1,868 could be mapped to the new pseudomolecules, with only one disagreement,  
67 demonstrating the high accuracy of the AGP. During the following gap-filling procedure,  
68 170 gaps were filled by SMRT long reads.

69 To ensure base-pairing accuracy and further polish the pseudomolecules, we  
70 generated ~2300Gb Illumina pair-end reads. To increase the size of reads, the paired-end  
71 library was constructed to be overlapping (~450bp library size, read length: 250bp). After  
72 merging the two reads in a pair, the average size of Illumina reads reached 400bp. These  
73 longer reads also decreased the difficulty in the alignment. About 89.7% of the assembly  
74 had good coverage for the correction (mapping and sequencing quality >20, read depth  
75 >=5). A total of 80k bases, including SNP and small Indel, were corrected, of which 91%  
76 were small indels.

77

78 **Centromere identification by ChIP-seq:** Peaks of CENH3 enrichment were defined by  
79 CENH3 ChIP-seq as described previously<sup>8</sup> using the HOMER findPeaks software<sup>9</sup>. Input  
80 reads from the CENH3 ChIP sample were used as controls. All reads were mapped to the  
81 genome using BWA-MEM<sup>10</sup>. As a first step, all reads, including potential repetitively  
82 mapping reads, were used to identify a set of putative CENH3-enriched regions; the  
83 parameters of HOMER findPeaks were set as follows: -region -size 5000 -minDist 50000  
84 -F 8 -L 0 -C -0. A set of high-confidence peaks was then independently identified using  
85 only uniquely mapping reads (as defined by MAPQ values  $\geq 20$ ) with the following  
86 parameters: -region -size 5000 -F 16 -L 0 -C -0. Putative CENH3-enriched regions that  
87 were either shorter than 100 kb, or that did not overlap with at least one high-confidence  
88 peak, were discarded. To generate the final set of centromeric loci, the remaining putative  
89 CENH3-enriched regions were merged if they were less than 500 kb apart.

90

## 91 **2. Comparison of genome assembly quality in Maize B73 RefGen\_v3 and v4**

92 The Maize Genome Sequencing Pilot Project randomly selected 100 BAC clones for  
93 high-quality sequencing, resulting in 98 curated BACs of finished quality<sup>11</sup>. These BACs  
94 were used for the detailed evaluation of the assembly quality of maize v4 genome. In  
95 total, 25 of the 28 fully completed BACs were spanned by a single contig in RefGen\_v4,  
96 with identity above 99.9%. In addition, the maize pilot sequence contains 57 BACs with

97 ordered contigs and gaps. The gaps of 46 BACs could be closed by a single contig in  
98 RefGen\_v4.

99 Several gene models with assembly errors in the maize B73 RefGen\_v3 have been  
100 corrected in the current maize genome. For example, the *rgh3* locus (JN692485.1) was  
101 involved in an assembly error that arose due to incorrect ordering and orientation of  
102 contigs in the BAC sequence, resulting in mis-annotation of this gene as two distinct gene  
103 models<sup>12</sup>. This problem was successfully fixed in the v4 assembly. Due to correction of  
104 such errors and the increase in contiguity described above, the RefGen\_v4 assembly is  
105 much more robust as a reference genome than the old BAC sequences.

106

### 107 **3. Gene annotation**

108 **Generation of a working gene set:** MAKER-P version 3.1<sup>13</sup> was used to annotate the  
109 maize RefGen\_v4 genome. As evidence, we used all annotated proteins from *Sorghum*  
110 *bicolor*, *Oryza sativa*, *Setaria italica*, *Brachypodium distachyon*, and *Arabidopsis*  
111 *thaliana*, downloaded from Gramene.org release 48<sup>14</sup>. For transcript evidence, the  
112 111,151 high quality transcripts from Iso-seq were further polished by illumina RNA-seq  
113 reads generated from same tissues<sup>15</sup> using Ectools  
114 (<https://github.com/jgurtowski/ectools>). Another set of 69,163 publicly available full-  
115 length cDNAs deposited in Genbank<sup>16</sup>, a total of 1,574,442 Trinity-assembled transcripts  
116 from 94 B73 RNA-Seq experiments<sup>17</sup>, and 112,963 transcripts assembled from deep  
117 sequencing of a B73 seedling<sup>18</sup> were also included as transcript evidence. For gene  
118 prediction, we used Augustus<sup>19</sup> and FGENESH (<http://www.softberry.com/berry.phtml>)  
119 trained on maize and monocots, respectively. For repeat masking, we used  
120 RepeatMasker and the B73-specific TE exemplars<sup>20</sup>. Helitron elements and captured  
121 exons within pack-MULES were removed from this library to prevent the masking of  
122 non-TE-related protein-coding genes. Additional masking was performed using a set of  
123 known TE-derived proteins distributed with the MAKER software package<sup>13</sup>.

124 The final annotation set was built iteratively. The first step, which included all of  
125 the protein evidence, the full-length cDNAs from GenBank, and the Iso-Seq data,  
126 generated 34,088 genes with 56,671 transcripts. For the second step, the gene models  
127 from the first pass were given back to MAKER as models, allowing them to persist

128 unchanged in the annotation set. Next, the additional transcript evidence derived from  
129 short reads was included. This step generated an additional 9,548 genes with 11,475  
130 transcripts. To retain as many genes as possible from the v3 annotations, the third pass  
131 added the previously annotated B73 transcripts and protein translations from the v3  
132 assembly as evidence. This step added 5,449 genes with 5,947 transcripts. MAKER-P is  
133 conservative in annotating alternate transcripts. Additionally, transcripts that contain  
134 large intron retentions, non-canonical splicing, or are expressed at low levels are also  
135 difficult to annotate confidently by computational methods. However, the single-  
136 molecule Iso-Seq transcript sequencing method can unambiguously identify these hard-  
137 to-annotate transcripts. By including the additional unique Iso-Seq transcripts into the  
138 gene models from step 3, we generated a protein-coding gene annotation set of 49,085  
139 genes and 161,680 transcripts (referred to as the working set).

140

141 **Compara gene tree construction:** The Ensembl Compara gene tree pipeline<sup>21</sup> was used  
142 to define gene families, construct phylogenetic gene trees, and infer orthologs and  
143 paralogs. Updated protocols used in the Ensembl version 81 software are detailed  
144 elsewhere

145 ([http://jul2015.archive.ensembl.org/info/genome/compara/homology\\_method.html](http://jul2015.archive.ensembl.org/info/genome/compara/homology_method.html)). The  
146 analysis included annotated protein-coding genes from both the v3 and v4 gene sets of  
147 maize B73, as well as 17 additional species (five monocots, four dicots, one basal  
148 angiosperm, three lower plants, and four non-plants), which were downloaded from the  
149 Ensembl core databases within Gramene Release-41. Tree reconciliation to classify  
150 duplication and speciation nodes, and the assignment of taxon levels to nodes, used the  
151 following input species tree derived from the NCBI Taxonomy database<sup>21</sup>:

152 ((((((((((sorghum\_bicolor,(zea\_mays\_v3,zea\_mays\_v4)N)Andropogoneae,setaria\_italica)  
153 Panicoideae,(brachypodium\_distachyon,oryza\_sativa)BEP\_clade)Poaceae,musa\_acumina  
154 ta)commelinids,(((arabidopsis\_thaliana,glycine\_max),vitis\_vinifera))rosids,solanum\_lyc  
155opersicum)Eudicot)Mesangiospermae,amborella\_trichopoda)Magnoliophyta,selaginella\_  
156moellendorffii)Tracheophyta,physcomitrella\_patens)Embryophyta,chlamydomonas\_reinh  
157ardtii)Viridiplantae,(((caenorhabditis\_elegans,drosophila\_melanogaster)Ecdysozoa,homo  
158\_sapiens)Bilateria,saccharomyces\_cerevisiae)Opisthokonta)Eukaryota;

159 Synteny maps, which relate collinear chains of orthologous genes between two genomes,  
160 were built using DAGchainer<sup>22</sup> in combination with other previously described  
161 methods<sup>20,23</sup>.

162 **Generation of the filtered gene set:** The working set of protein-coding gene annotations  
163 is expected to contain TEs that were not masked prior to annotation, long noncoding  
164 RNAs annotated as protein-coding genes, and annotations with little supporting evidence.  
165 We filtered the working set based on evidentiary support, transposon screening, long non-  
166 coding RNA screening, homology support, and valid CDSs. The approach is  
167 schematically represented in Extended data Figure 4a.

168 **tRNA annotation:** tRNAs were identified using tRNAscan-SE<sup>24</sup> within the MAKER-P  
169 framework<sup>25</sup>. A total of 2,305 tRNAs were identified: 1,451 decode standard amino acids,  
170 four decode seleno-Cys, seven are putative suppressors, 13 contain an undeterminable  
171 anti-codon sequence, and 830 are apparent pseudogenes. Compared to the v3 assembly,  
172 v4 contains 59 additional complete tRNAs and 54 additional putative tRNA pseudogenes.  
173 This increase in identifiable tRNAs provides further evidence that v4 is a more complete  
174 genome assembly than v3.

175

#### 176 **4. Comparison of gene annotation between RefGen\_v3 and v4**

177 **Alignment of v3 genes to the v4 genome:** We used two pipelines to map the v3 genes to  
178 the v4 genome, Genome Assembly Converter and Mummer pipeline<sup>26</sup>. In Genome  
179 Assembly Converter, the ATAC pipeline<sup>27</sup> was used to create the alignment chain file  
180 between two assemblies, and then CrossMap<sup>28</sup> was used to convert the coordinates of the  
181 v3 gene annotation. Due to the complexity of repeats in maize genome, only the one-to-  
182 one alignment blocks were saved to build the chain. In the chain file, the v3 genome  
183 covered 89.7% of v4 genome, whereas v4 covered 92.5% of v3 genome. A chromosome-  
184 to-chromosome alignment was first performed using Mummer to map the v3 genes to the  
185 v4 genome. Genes from v3 that could not be mapped to the same chromosome in v4 were  
186 then aligned to the whole v4 assembly. Only unique hits with identity above 98% and  
187 100% coverage were retained for merging with the Genome Assembly Converter  
188 pipeline. Disagreements between the two pipelines were resolved as follows: if the

189 Genome Assembly Converter pipeline had 100% coverage for a given gene, then those  
190 coordinates were kept; otherwise, the result from the Mummer alignment was used.

191 Alignment of the RefGen\_v3 and v4 genome assemblies indicated that the two  
192 versions are highly consistent with each other in gene space. A total of 36,725 (94%) v3  
193 gene models could be mapped to the new RefGen\_v4 genome without sequence changes.  
194 Most of the remaining v3 genes (1,356) could be mapped, but crossed multiple contigs in  
195 RefGen\_v3, with gaps; consequently, it is very likely that they were incorrectly  
196 assembled in v3. In RefGen\_v4, most of these genes were contained within continuous  
197 sequences, indicating the improvement of the genomic sequences of these genes. In  
198 addition, 92 of the 146 genes previously unanchored in RefGen\_v3 were anchored to  
199 chromosomes in the RefGen\_v4 assembly.

200 **Core promoter elements:** Core promoter elements were analyzed in both RefGen\_v3  
201 and v4 with a published pipeline<sup>29,30</sup>. Comparison of core promoter elements, especially  
202 the TATA-box, CCAAT-box, and Y patch in the new assembly to those in the previously  
203 published assembly revealed 17.5% of genes in the new assembly contained a TATA-  
204 box, whereas in the previous assembly only 12.8% genes contained this element.  
205 Similarly, 7.2% genes contained a CCAAT-box and 58.17% contained Y patch in maize  
206 B73 RefGen\_v4, versus 2.4% and 41.5%, respectively in v3.

207 **Gene orientation:** Of 30,926 genes that could be mapped between the v3 and v4  
208 annotations, 2,151 genes were switched to a different strand. To evaluate this, we  
209 compared gene orientation to sorghum orthologs within syntenic blocks. Among 652  
210 genes that could be tracked in this manner, the orientation of 589 (90.3%) was conserved  
211 with sorghum. Thus, in the vast majority of cases, the re-orientation of a gene in v4  
212 brought the configuration into closer agreement with sorghum, further lending confidence  
213 to strand reassignments of v4 genes.

214 **Identification of missing genes in maize genome:** We identified 22,048 orthologous  
215 gene sets that originated prior to, or within, the grass common ancestor, and cataloged  
216 deficiencies in gene content among annotations of the five grass species (maize, rice,  
217 sorghum, *Setaria*, and *Brachypodium*). Of these sets, ~69% were found in all five species,  
218 and of individual species, rice, *Setaria*, and sorghum had the most complete  
219 representation, possessing from 91% to 92% of ortholog sets. By contrast, despite the

220 fact that maize is a product of whole-genome duplication, maize genes were found in  
221 only 86% of ortholog sets, representing a deficit of over 3,000 genes. To minimize  
222 artifacts, we restricted analysis to 592 ortholog sets containing 668 sorghum genes that 1)  
223 are syntenic with an outgroup species (either rice, *Brachypodium*, or *Setaria*), 2) are  
224 flanked by genes contained within a synteny block that maps to a single maize contig in  
225 both the A and B subgenomes, and 3) lack alignment of CDS features to the v4 reference.

226

227

228

## 5. Structural identification of transposable elements

229

230

231

232

233

234

235

236

237

238

**LTR retrotransposons:** LTR retrotransposons were identified using LTRharvest<sup>31</sup> and LTRdigest<sup>32</sup>. LTRharvest searches sequence data for structural characteristics of LTR retrotransposons; in an analysis of the *Drosophila X* chromosome, it was shown to be the most sensitive among available structural search tools<sup>33</sup>. To be consistent with known LTR retrotransposons in maize, we adjusted default parameters including LTR length (100–7000 bp) and element length (1000–20000 bp). All searches required target site duplications (TSDs) of 4–6 bp (allowing one mismatch) and a 2-nt inverted motif at the terminal ends of each LTR (5' TG..CA 3', allowing one mismatch). If multiple overlapping elements were found, the one with the highest percent identity between LTRs was chosen with the '-overlaps best' option.

239

240

241

242

243

244

245

246

247

The resultant TE models were further annotated with LTRdigest<sup>32</sup>, which identifies sequence features such as primer binding site, polypurine tract, and protein domains associated with previously identified retrotransposons from any organism. We used all eukaryotic tRNA entries from the UCSC gtRNA database to predict primer binding sites, and amino-acid HMM profiles of retrotransposon-associated proteins as deposited in GyDB (<http://gydb.org>)<sup>34</sup>. If RNase H, reverse transcriptase, and integrase domains were present, gene order was used to classify elements into the Ty1/Copia (integrase upstream of RNase H) and Ty3/gypsy (RNase H upstream of integrase) superfamilies.

248

249

250

251

LTR retrotransposons dominate the intergenic space of the maize genome. To capture the nested structure of these elements, generated when a newly arriving TE inserts into a TE already present at that genomic location, we computationally excised each LTR retrotransposon copy and repeated the structural search on this subtracted



252 pseudo-genome. We repeated this computational subtraction for 80 rounds, increasing the  
253 element length by 1000 bp for each round to accommodate sequence contributed by TE  
254 fragments and TEs of other orders.

255 **SINE and LINE:** Because SINES are transcribed by RNA polymerase III, they are often  
256 derived from one of three classes of Pol III–transcribed molecules (tRNA, 7SL, 5s  
257 rRNA). Animal SINES of all three classes are known, whereas plant SINES are  
258 exclusively tRNA-derived<sup>35</sup>. We used SINE-finder<sup>35</sup> to search for tRNA-derived SINES  
259 containing RNA polymerase III A and B boxes near the polyA tail. The default A and B  
260 box consensus (RVTGG; GTTCRA), a 25–50 bp spacer between the A and B boxes,  
261 and a spacer of 20–500 bp between the B box and polyA tail were applied. Structural  
262 SINES were predicted only on the forward strand of the genome. LINES were identified  
263 using TARGeT and mTEA as below for TIR elements, using LINE exemplars and 15 bp  
264 target site duplications.

265 **TIR:** Exemplar elements from the maize TE consortium (MTEC) annotation<sup>20</sup> were used  
266 as nucleotide queries in TARGeT<sup>36</sup>, a pipeline designed to recover high-copy transposon  
267 and gene families. The number of elements clustered in the PHI step was increased to  
268 10000 copies, and 200 bp of flanking sequence on either edge of genomic matches was  
269 extracted (-p\_f 200). This approach recovered candidate TE sequences, but the TE  
270 boundaries and flanking sequence were unknown. To identify the boundaries of each  
271 element, we scanned each candidate and verified the presence of terminal inverted repeat  
272 (TIR) and TSD sequences indicative of the TE superfamily (see the table below), using  
273 mTEA ([https://github.com/hyphaltip/mTEA/blob/master/scripts/id\\_TIR\\_in\\_FASTA.pl](https://github.com/hyphaltip/mTEA/blob/master/scripts/id_TIR_in_FASTA.pl);  
274 modified to use mafft for alignment), Although TSDs and TIRs should be identical for  
275 most superfamilies upon insertion into the genome, mutations arising at the background  
276 genomic mutation rate can generate differences. Thus, we allowed mismatches of 80% of  
277 the length of a TSD or TIR to accommodate identification of these older, degraded  
278 copies.

279

280

281

282

283

**DNA TIR TE Superfamily TSD & TIR Classification**

Superfamily	TSD Length (sequence restrictions, if any)	TIR Length (sequence restrictions, if any)
DTT <i>Tc1/Mariner</i>	2 bp (TA)	13 bp
DTA <i>hAT</i>	8 bp	11 bp
DTM <i>Mutator</i>	9 bp	40 bp
DTH <i>Pif/Harbinger</i>	3 bp (TNN)	14 bp
DTC <i>CACTA</i>	3 bp	13 bp (CACTNNNNNNNNNN)

284

285 In addition, MiteHunter<sup>37</sup> and detectMITE<sup>38</sup> were used to identify *de novo* structural  
 286 MITEs, searching for TIR and TSDs in genomic sequences. We filtered MITE output by  
 287 TSD and TIR length, and all exemplars with TIRs and TSDs of anticipated length for the  
 288 superfamily were used to search using mTEA, as described above.

289 **Helitron:** HelitronScanner<sup>39</sup> with default parameters was deployed to identify upstream  
 290 and downstream termini of helitrons, and to join upstream and downstream termini  
 291 within 200–20,000 bp of each other into helitron TE copies. We predicted helitrons in  
 292 both the direct and reverse complement orientations.

293 **Family clustering:** Families were identified within each superfamily of TIR TE and  
 294 order of retrotransposon using the 80–80–80 rule<sup>40</sup>, which requires that elements within a  
 295 family must share 80% homology over at least 80 base pairs of 80% of the element's  
 296 functional or internal domains. For LTR retrotransposons, the 5' LTR was used to cluster  
 297 families, consistent with previous annotations in maize<sup>41</sup>. The entire element sequence  
 298 was used to group TIRs, LINEs, and SINEs, because functional domains are short, and  
 299 because a large proportion of non-autonomous elements lack protein-coding domains.  
 300 Because the internal regions of maize helitrons are diverse and clustering methods  
 301 applied to the entire element yield almost exclusively singletons<sup>42</sup>, we used a family  
 302 classification previously applied to maize helitrons that relies on 80% identity of the 30  
 303 bp at the 3' end of each copy<sup>43</sup>, a region of hairpin-forming sequence important for  
 304 rolling circle replication. All family definitions are consistent with those used previously  
 305 in the maize genome sequencing project<sup>20,41</sup>, although we implemented clustering of  
 306 families in SiLiX<sup>44</sup>. Additionally, for each structurally defined TE in the genome, we  
 307 assigned a unique identifier that indicates its superfamily and family.

308 **Calculating genomic composition and resolving TE overlaps:** As structural searches  
309 were run independently for each TE order, we filtered overlapping insertions in order to  
310 count each genomic position as derived from only one transposable element and generate  
311 a filtered set of TE annotations. As subsequent transposition into existing TEs causes  
312 them to occupy larger ranges along the genome, larger TEs are expected to be older.  
313 Since the chance of false homology increases as requirements of sequence identity are  
314 reduced, we filtered out LTR retrotransposons that occupy over 100kb along the genome,  
315 as these old large elements are more likely to be false positives. As nested insertions from  
316 most orders of TEs are known<sup>45-48</sup> (LTR into helitron, helitron into LTR; TIR into LTR,  
317 LTR into TIR), we retain TE copies entirely nested within another copy, but remove  
318 insertions that overlap boundaries of other copies. When copies overlap, we retain first  
319 LTR retrotransposons, next TIR, next SINE and LINE, and finally helitrons. This  
320 removal order was chosen to favor TE orders with stronger structural signatures.

321 **Homology Search:** After a TE inserts into a position in the genome, it is subject to  
322 subsequent mutations. Because features will erode over time, making identification  
323 difficult, these changes can complicate its ascertainment by structural methods. To  
324 identify these waning TE-derived sequences, we used RepeatMasker  
325 (<http://www.repeatmasker.org>) to mask the B73 RefGen\_v4 pseudomolecules with a  
326 repeat library consisting of structurally defined TEs. These consist of the filtered TE set  
327 described above, but with LTR retrotransposon families containing greater than 10 copies  
328 additionally downsampled to reduce computational runtime. This is necessary due to the  
329 existence of large families with tens of thousands of nearly identical copies. For these  
330 LTR retrotransposon families, we algorithmically selected exemplar elements, based on  
331 the length distribution of the TE family. Briefly, we used a Dirichlet Process Prior to  
332 identify the most likely number of normal distributions needed to generate the observed  
333 length distribution, and identified cluster membership for each element in the family.  
334 Then, we selected the copy with a length closest to the mean of each inferred normal  
335 distribution. These copies were used as exemplars in the homology search.

336

337 **Comparison of transposable element annotations in v3 and v4:** To compare our  
338 annotation approach with existing TE annotations generated based on homology to the

339 MTEC repeat library (www.maizetdb.org), we annotated the AGPv3 assembly using the  
340 structural methods applied to AGPv4. We then assessed the overlap between the  
341 available RepeatMasker annotation of AGPv3 and this new annotation. This analysis  
342 revealed that only 0.6% (11,017 of 1,695,362) of LTR retrotransposons in RepeatMasker  
343 AGPv3 annotation are full-length and contain TSDs. Such striking underrepresentation is  
344 anticipated when homology-based methods are used to identify diverse TEs<sup>49</sup>. In addition  
345 to the improved quality of the annotation, the AGPv4 genome allows more complete  
346 reconstruction of the entire sequence of each TE. For example, we recovered 68% more  
347 Ty1/Copia and Ty3/Gypsy LTR retrotransposons with evidence of all proteins required  
348 for retrotransposition (42,929 in AGPv4 vs. 25,412 in AGPv3); in AGPv3, many of these  
349 internal domains were represented by gaps between contigs.

350

351 **Diversification of maize LTR retrotransposons:** To investigate the evolutionary  
352 dynamics of retrotransposition in maize since divergence from sorghum, we applied our  
353 annotation approach for LTR retrotransposons to the *Sorghum bicolor* genome (Sorbi1).  
354 Sequences matching HMM models of RT\_crm.hmm (Ty3/Gypsy) and RT\_sire.hmm  
355 (Ty1/Copia) were extracted from each non-nested LTR TE they matched. As the  
356 estimated divergence time between maize and sorghum (12 Mya) predicts greater  
357 divergence than the 80% identity used to define families, generated a consensus sequence  
358 for each family using emboss cons<sup>50</sup> to track differences between species. We aligned  
359 these family consensus with MAFFT mafft<sup>51</sup> and built a maximum likelihood  
360 phylogenetic tree with fasttree2<sup>52</sup>. We then collapsed sister tips on the tree if they arose  
361 from the same species, and summed the number of copies belonging to each of these  
362 species-specific lineages. Hence, monophyletic lineages of TEs, with respect to the  
363 genome they were ascertained from, are shown in Figure 2.

364 **Data Availability:** Scripts, parameters, and intermediate files of each TE superfamily are  
365 available at

366 [https://github.com/mcstitzer/agpv4\\_te\\_annotation/tree/master/ncbi\\_pseudomolecule](https://github.com/mcstitzer/agpv4_te_annotation/tree/master/ncbi_pseudomolecule)

367

368

369

370

371

372 **Reference**

373

- 374 1 Berlin, K. *et al.* Assembling large genomes with single-molecule sequencing  
375 and locality-sensitive hashing. *Nature biotechnology* **33**, 623-630,  
376 doi:10.1038/nbt.3238 (2015).
- 377 2 Chin, C. S. *et al.* Nonhybrid, finished microbial genome assemblies from long-  
378 read SMRT sequencing data. *Nat Methods* **10**, 563-569,  
379 doi:10.1038/nmeth.2474 (2013).
- 380 3 Nguyen, J. V. *Genomic mapping: a statistical and algorithmic analysis of the*  
381 *optical mapping system.* (University of Southern California, 2010).
- 382 4 Anantharaman, T. & Mishra, B. in *Algorithms in Bioinformatics, First*  
383 *International Workshop, WABI.* 27-40.
- 384 5 Myers, E. W. *et al.* A whole-genome assembly of *Drosophila*. *Science* **287**,  
385 2196-2204 (2000).
- 386 6 Sommer, D. D., Delcher, A. L., Salzberg, S. L. & Pop, M. Minimus: a fast,  
387 lightweight genome assembler. *BMC bioinformatics* **8**, 64, doi:10.1186/1471-  
388 2105-8-64 (2007).
- 389 7 Ganai, M. W. *et al.* A large maize (*Zea mays* L.) SNP genotyping array:  
390 development and germplasm genotyping, and genetic mapping to compare  
391 with the B73 reference genome. *PLoS One* **6**, e28334,  
392 doi:10.1371/journal.pone.0028334 (2011).
- 393 8 Gent, J. I., Wang, K., Jiang, J. & Dawe, R. K. Stable Patterns of CENH3  
394 Occupancy Through Maize Lineages Containing Genetically Similar  
395 Centromeres. *Genetics* **200**, 1105-1116, doi:10.1534/genetics.115.177360  
396 (2015).
- 397 9 Heinz, S. *et al.* Simple combinations of lineage-determining transcription  
398 factors prime cis-regulatory elements required for macrophage and B cell  
399 identities. *Molecular cell* **38**, 576-589, doi:10.1016/j.molcel.2010.05.004  
400 (2010).
- 401 10 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools.  
402 *Bioinformatics* **25**, 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).
- 403 11 Haberer, G. *et al.* Structure and architecture of the maize genome. *Plant*  
404 *Physiol* **139**, 1612-1624, doi:10.1104/pp.105.068718 (2005).
- 405 12 Fouquet, R. *et al.* Maize rough endosperm3 encodes an RNA splicing factor  
406 required for endosperm cell differentiation and has a nonautonomous effect  
407 on embryo development. *The Plant cell* **23**, 4280-4297,  
408 doi:10.1105/tpc.111.092163 (2011).
- 409 13 Cantarel, B. L. *et al.* MAKER: an easy-to-use annotation pipeline designed for  
410 emerging model organism genomes. *Genome research* **18**, 188-196,  
411 doi:10.1101/gr.6743907 (2008).
- 412 14 Monaco, M. K. *et al.* Gramene 2013: comparative plant genomics resources.  
413 *Nucleic Acids Res* **42**, D1193-1199, doi:10.1093/nar/gkt1110 (2014).
- 414 15 Wang, B. *et al.* Unveiling the complexity of the maize transcriptome by single-  
415 molecule long-read sequencing. *Nature communications* **7**, 11708,  
416 doi:10.1038/ncomms11708 (2016).

417 16 Soderlund, C. *et al.* Sequencing, mapping, and analysis of 27,455 maize full-  
418 length cDNAs. *PLoS Genet* **5**, e1000740, doi:10.1371/journal.pgen.1000740  
419 (2009).

420 17 Law, M. *et al.* Automated update, revision, and quality control of the maize  
421 genome annotations using MAKER-P improves the B73 RefGen\_v3 gene  
422 models and identifies new genes. *Plant Physiol* **167**, 25-39,  
423 doi:10.1104/pp.114.245027 (2015).

424 18 Martin, J. A. *et al.* A near complete snapshot of the *Zea mays* seedling  
425 transcriptome revealed from ultra-deep sequencing. *Scientific reports* **4**,  
426 4519, doi:10.1038/srep04519 (2014).

427 19 Keller, O., Kollmar, M., Stanke, M. & Waack, S. A novel hybrid gene prediction  
428 method employing protein multiple sequence alignments. *Bioinformatics* **27**,  
429 757-763, doi:10.1093/bioinformatics/btr010 (2011).

430 20 Schnable, P. S. *et al.* The B73 maize genome: complexity, diversity, and  
431 dynamics. *Science* **326**, 1112-1115, doi:10.1126/science.1178534 (2009).

432 21 Vilella, A. J. *et al.* EnsemblCompara GeneTrees: Complete, duplication-aware  
433 phylogenetic trees in vertebrates. *Genome research* **19**, 327-335,  
434 doi:10.1101/gr.073585.107 (2009).

435 22 Haas, B. J., Delcher, A. L., Wortman, J. R. & Salzberg, S. L. DAGchainer: a tool  
436 for mining segmental genome duplications and synteny. *Bioinformatics* **20**,  
437 3643-3646, doi:10.1093/bioinformatics/bth397 (2004).

438 23 Youens-Clark, K. *et al.* Gramene database in 2010: updates and extensions.  
439 *Nucleic Acids Res* **39**, D1085-1094, doi:10.1093/nar/gkq1148 (2011).

440 24 Lowe, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection of  
441 transfer RNA genes in genomic sequence. *Nucleic Acids Res* **25**, 955-964  
442 (1997).

443 25 Campbell, M. S. *et al.* MAKER-P: a tool kit for the rapid creation, management,  
444 and quality control of plant genome annotations. *Plant Physiol* **164**, 513-524,  
445 doi:10.1104/pp.113.230144 (2014).

446 26 Kurtz, S. *et al.* Versatile and open software for comparing large genomes.  
447 *Genome biology* **5**, R12, doi:10.1186/gb-2004-5-2-r12 (2004).

448 27 Istrail, S. *et al.* Whole-genome shotgun assembly and comparison of human  
449 genome assemblies. *Proc Natl Acad Sci U S A* **101**, 1916-1921,  
450 doi:10.1073/pnas.0307971100 (2004).

451 28 Zhao, H. *et al.* CrossMap: a versatile tool for coordinate conversion between  
452 genome assemblies. *Bioinformatics* **30**, 1006-1007,  
453 doi:10.1093/bioinformatics/btt730 (2014).

454 29 Kumari, S. & Ware, D. Genome-wide computational prediction and analysis of  
455 core promoter elements across plant monocots and dicots. *PLoS One* **8**,  
456 e79011, doi:10.1371/journal.pone.0079011 (2013).

457 30 Smith, A. D., Sumazin, P., Xuan, Z. & Zhang, M. Q. DNA motifs in human and  
458 mouse proximal promoters predict tissue-specific expression. *Proc Natl Acad  
459 Sci U S A* **103**, 6275-6280, doi:10.1073/pnas.0508169103 (2006).

460 31 Ellinghaus, D., Kurtz, S. & Willhoeft, U. LTRharvest, an efficient and flexible  
461 software for de novo detection of LTR retrotransposons. *BMC bioinformatics*  
462 **9**, 18, doi:10.1186/1471-2105-9-18 (2008).

463 32 Steinbiss, S., Willhoeft, U., Gremme, G. & Kurtz, S. Fine-grained annotation and  
464 classification of de novo predicted LTR retrotransposons. *Nucleic Acids Res*  
465 **37**, 7002-7013, doi:10.1093/nar/gkp759 (2009).

466 33 Lerat, E., Bulet, N., Biemont, C. & Vieira, C. Comparative analysis of  
467 transposable elements in the melanogaster subgroup sequenced genomes.  
468 *Gene* **473**, 100-109, doi:10.1016/j.gene.2010.11.009 (2011).

469 34 Llorens, C. *et al.* The Gypsy Database (GyDB) of mobile genetic elements:  
470 release 2.0. *Nucleic Acids Res* **39**, D70-74, doi:10.1093/nar/gkq1061 (2011).

471 35 Wenke, T. *et al.* Targeted identification of short interspersed nuclear element  
472 families shows their widespread existence and extreme heterogeneity in  
473 plant genomes. *The Plant cell* **23**, 3117-3128, doi:10.1105/tpc.111.088682  
474 (2011).

475 36 Han, Y., Burnette, J. M., 3rd & Wessler, S. R. TARGeT: a web-based pipeline for  
476 retrieving and characterizing gene and transposable element families from  
477 genomic sequences. *Nucleic Acids Res* **37**, e78, doi:10.1093/nar/gkp295  
478 (2009).

479 37 Han, Y. & Wessler, S. R. MITE-Hunter: a program for discovering miniature  
480 inverted-repeat transposable elements from genomic sequences. *Nucleic*  
481 *Acids Res* **38**, e199, doi:10.1093/nar/gkq862 (2010).

482 38 Ye, C., Ji, G. & Liang, C. detectMITE: A novel approach to detect miniature  
483 inverted repeat transposable elements in genomes. *Scientific reports* **6**,  
484 19688, doi:10.1038/srep19688 (2016).

485 39 Xiong, W., He, L., Lai, J., Dooner, H. K. & Du, C. HelitronScanner uncovers a  
486 large overlooked cache of Helitron transposons in many plant genomes. *Proc*  
487 *Natl Acad Sci U S A* **111**, 10263-10268, doi:10.1073/pnas.1410068111  
488 (2014).

489 40 Wicker, T. *et al.* A unified classification system for eukaryotic transposable  
490 elements. *Nature reviews. Genetics* **8**, 973-982, doi:10.1038/nrg2165 (2007).

491 41 Baucom, R. S. *et al.* Exceptional diversity, non-random distribution, and rapid  
492 evolution of retroelements in the B73 maize genome. *PLoS Genet* **5**,  
493 e1000732 (2009).

494 42 Sweredoski, M., DeRose-Wilson, L. & Gaut, B. S. A comparative computational  
495 analysis of nonautonomous helitron elements between maize and rice. *BMC*  
496 *genomics* **9**, 467, doi:10.1186/1471-2164-9-467 (2008).

497 43 Yang, L. & Bennetzen, J. L. Distribution, diversity, evolution, and survival of  
498 Helitrons in the maize genome. *Proc Natl Acad Sci U S A* **106**, 19922-19927,  
499 doi:10.1073/pnas.0908008106 (2009).

500 44 Miele, V., Penel, S. & Duret, L. Ultra-fast sequence clustering from similarity  
501 networks with SiLiX. *BMC bioinformatics* **12**, 116, doi:10.1186/1471-2105-  
502 12-116 (2011).

503 45 Gupta, S., Gallavotti, A., Stryker, G. A., Schmidt, R. J. & Lal, S. K. A novel class of  
504 Helitron-related transposable elements in maize contain portions of multiple  
505 pseudogenes. *Plant molecular biology* **57**, 115-127 (2005).

506 46 Morgante, M. *et al.* Gene duplication and exon shuffling by helitron-like  
507 transposons generate intraspecies diversity in maize. *Nat Genet* **37**, 997-  
508 1002 (2005).

509 47 Jameson, N. *et al.* Helitron mediated amplification of cytochrome P450  
510 monooxygenase gene in maize. *Plant molecular biology* **67**, 295-304 (2008).  
511 48 Jiang, N. & Wessler, S. R. Insertion preference of maize and rice miniature  
512 inverted repeat transposable elements as revealed by the analysis of nested  
513 elements. *The Plant cell* **13**, 2553-2564 (2001).  
514 49 Platt, R. N., Blanco-Berdugo, L. & Ray, D. A. Accurate transposable element  
515 annotation is vital when analyzing new genome assemblies. *Genome Biology*  
516 *and Evolution*, evw009 (2016).  
517 50 Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European molecular biology  
518 open software suite. *Trends in genetics* **16**, 276-277 (2000).  
519 51 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software  
520 version 7: improvements in performance and usability. *Molecular biology and*  
521 *evolution* **30**, 772-780 (2013).  
522 52 Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2—approximately maximum-  
523 likelihood trees for large alignments. *PLoS One* **5**, e9490 (2010).  
524