

1 **Highly expressed maize pollen genes display coordinated**
2 **expression with neighboring transposable elements**
3 **and contribute to pollen fitness**

4
5 Cedar Warman¹, Kaushik Panda², Zuzana Vejlupkova¹, Sam Hokin³, Erica Unger-Wallace⁴, Rex
6 A Cole¹, Antony M Chettoor³, Duo Jiang⁵, Erik Vollbrecht^{4,6,7}, Matthew MS Evans³, R Keith
7 Slotkin², John E Fowler^{1,8*}

8
9 ¹Department of Botany and Plant Pathology, Oregon State University, Corvallis, Oregon, United States of
10 America

11
12 ²Donald Danforth Plant Science Center, St. Louis, Missouri, United States of America

13
14 ³Department of Plant Biology, Carnegie Institution for Science, Stanford, California, United States of
15 America

16
17 ⁴Department of Genetics Development and Cell Biology, Iowa State University, Ames, Iowa, United
18 States of America

19
20 ⁵Department of Statistics, Oregon State University, Corvallis, Oregon, United States of America

21
22 ⁶Bioinformatics and Computational Biology, Iowa State University, Ames, Iowa, United States of America

23
24 ⁷Interdepartmental Genetics, Iowa State University, Ames, Iowa, United States of America

25
26 ⁸Center for Genome Research and Biocomputing, Oregon State University, Corvallis, Oregon, United
27 States of America

28
29
30 * Corresponding author

31 E-mail: fowlerj@science.oregonstate.edu

32 Abstract

33 In flowering plants, the haploid male gametophyte (pollen) is essential for sperm delivery,
34 double fertilization, and subsequent initiation of seed development. Pollen also undergoes
35 dynamic epigenetic regulation of expression from transposable elements (TEs), but how this
36 process interacts with gene regulation and function is not clearly understood. To identify
37 components of these processes, we quantified transcript levels in four male reproductive stages
38 of maize (tassel primordia, microspores, mature pollen, and isolated sperm cells) via RNA-seq.
39 We found that, in contrast to Arabidopsis TE expression in pollen, TE transcripts in maize
40 accumulate as early as the microspore stage and are also present in sperm cells. Intriguingly,
41 coordinated expression was observed between the most highly expressed protein-coding genes
42 and neighboring TEs, specifically in both mature pollen and sperm cells. To test the hypothesis
43 that such elevated expression correlates with functional relevance, we measured the fitness cost
44 (male-specific transmission defect) of GFP-tagged exon insertion mutations in over 50 genes
45 highly expressed in pollen vegetative cell, sperm cell, or seedling (as a sporophytic control).
46 Insertions in genes highly expressed only in seedling or primarily in sperm cells (with one
47 exception) exhibited no difference from the expected 1:1 transmission ratio. In contrast, insertions
48 in over 20% of vegetative cell genes were associated with significant reductions in fitness,
49 showing a positive correlation of transcript level with non-Mendelian segregation. The *gamete*
50 *expressed2* (*gex2*) gene was the single sperm cell gene associated with reduced transmission
51 when mutant (<35% for two independent insertions), and also triggered seed defects when
52 crossed as a male, supporting a role for *gex2* in double fertilization. Overall, our study
53 demonstrates a developmentally programmed and coordinated transcriptional activation of TEs
54 and genes, and further identifies maize pollen as a model in which transcriptomic data have
55 predictive value for quantitative phenotypes.

56
57

58 Author Summary

59 In flowering plants, pollen is essential for delivering sperm cells to the egg and central cell
60 for double fertilization, initiating the process of seed development. In plants with abundant pollen
61 like maize, this process can be highly competitive. In an added layer of complexity, growing
62 evidence indicates expression of transposable elements (TEs) is more dynamic in pollen than in
63 other plant tissues. How these elements impact pollen function and gene regulation is not well
64 understood. We used transcriptional profiling to generate a framework for both detailed analysis
65 of TE expression and quantitative assessment of gene function during maize pollen development.
66 TEs are expressed early and persist, many showing coordinate activation with highly-expressed
67 neighboring genes in the pollen vegetative cell and sperm cells. Measuring fitness costs for a set
68 of over 50 mutations indicates a correlation between elevated transcript level and gene function
69 in the vegetative cell. Finally, we establish a role in fertilization for the *gamete expressed2* (*gex2*)
70 gene, identified based on its specific expression in sperm cells. These results highlight maize
71 pollen as a powerful model for investigating the developmental interplay of TEs and genes, as
72 well as for measuring fitness contributions of specific genes.

73

74 Introduction

75 Sexual reproduction enables the segregation and recombination of genetic material, which
76 increases genetic diversity in populations and contributes to the vast diversity of eukaryotes. In
77 flowering plants, sexual reproduction requires the development of reduced, haploid gametophytes
78 from sporophytic, diploid parents. The mature female gametophyte, the embryo sac, includes the
79 binucleate central cell and the egg cell (reviewed in [1,2]), each of which is fertilized by a sperm
80 cell to generate the triploid endosperm and diploid embryo, respectively. The mature male
81 gametophyte, pollen, consists of a vegetative cell harboring two sperm cells (reviewed in [3,4]).
82 In maize, male gametophytes arise from microspore mother cells in the tassel primordium. The
83 transition from diploid sporophyte to haploid gametophyte occurs when these cells undergo
84 meiosis, each resulting in four haploid microspores. Each microspore then undergoes two rounds
85 of mitosis to produce the pollen grain, first generating the large vegetative cell and a smaller
86 generative cell via asymmetric division, and then producing the two sperm from the generative
87 cell. After the arrival of the pollen grain on the floral stigma, the vegetative cell transports the two
88 sperm cells to the female gametophyte via pollen tube growth (reviewed in [5,6]). Accurate
89 navigation of the pollen tube as it grows down the style is dependent on the architecture of the
90 style's transmitting tract [7] and possible additional signaling and recognition mechanisms that
91 are poorly understood [8]. The final stages of pollen tube growth depend on a complex interplay
92 of signals to guide the pollen tube to the ovule (reviewed in [9]).

93 In maize, a pollen tube must grow up to 30 cm through the silk to reach the female
94 gametophyte, often competing with multiple pollen tubes to eventually enter the embryo sac and
95 release its sperm cells for fertilization (reviewed in [2,5]). Across the plant kingdom, this
96 competitive context for pollen tube development differs, depending on the pollen population as
97 well as sporophytic characters (reviewed in [10]). In a highly competitive environment, successful
98 fertilization is likely enhanced by pollen tubes functioning at full capacity [11–13], as generally
99 only the first tube to reach the micropyle is permitted to enter the female gametophyte. The
100 mechanisms preventing entry of multiple pollen tubes, known as the polytubey block, are not well-
101 understood, but presumably act to reduce polyspermy, which typically leads to sterile offspring
102 [6]. In Arabidopsis, mutant-associated fertilization of only the egg cell or the central cell does not
103 fully activate the block, allowing for the attraction of an additional pollen tube and the completion
104 of double fertilization with sperm cells from two different pollen tubes [14], a process known as
105 heterofertilization. Heterofertilization occurs at low frequencies (0.5-5%, depending on the genetic
106 background) in maize [15,16] and is thought to be associated with abnormal fertilization, but the
107 details remain unclear. In maize, mutations in the genes *MATRILINEAL/NLD/ZmPLA1* and
108 *ZmDMP* have been linked to pollen-induced production of haploid embryos and other seed
109 defects, which are likely associated with aberrant events at fertilization [17–20] or soon after [21].
110 Thus, many mechanisms associated with both pollen tube growth and fertilization remain
111 enigmatic.

112 Given their specialized biological functions and well-defined developmental stages,
113 gametophytes are prime targets for transcriptome analysis. Initial studies of plant gametophytic
114 transcriptomes in Arabidopsis pollen [22,23] and embryo sacs [24,25] described a limited and
115 specialized set of transcripts and identified numerous candidate genes for gametophytic function.
116 In maize, the first RNA-seq study of male and female gametophyte transcriptomes (mature pollen
117 and embryo sacs) similarly identified subsets of developmentally specific genes, with pollen

118 showing the most specialized transcriptome, relative to other tissues assessed [26]. More
119 recently, RNA-seq has been carried out on additional stages of maize reproductive development,
120 including pre-meiotic and meiotic anther cells [27–29], as well as sperm cells, egg cells, and early
121 stages of zygotic development [30].

122 Gametophytic tissues are known to show dynamic expression of transposable elements
123 (TEs). In *Arabidopsis*, global TE expression is derepressed at the late stages of pollen
124 development, occurring in the pollen vegetative nucleus only after pollen mitosis II [31]. The pollen
125 vegetative nucleus undergoes a programmed loss of heterochromatin, resulting in TE activation,
126 TE transposition, TE small interfering RNA production, and subsequent increased RNA-directed
127 DNA methylation [31–35]. A variety of functions have been ascribed to this male gametophytic
128 "developmental relaxation of TE silencing" (DRTS) event [36], including the generation of TE
129 small interfering RNAs that are mobilized to the sperm cells [37], and control of imprinted gene
130 expression after fertilization [38]. However, the dynamics of TE expression during gametophytic
131 development in a transposable element-rich species such as maize have not been investigated.

132 To provide a more full description of transcriptome dynamics across maize male
133 reproductive development, including TE transcriptional activity, we generated RNA-seq datasets
134 from tassel primordia, microspores, mature pollen, and isolated sperm cells. Using these data,
135 we describe differential expression patterns of genes and TEs across these stages, uncovering a
136 coordinated regulation of TEs and their neighboring genes in pollen grains. We then conducted a
137 functional validation of such highly expressed genes by testing over fifty insertional mutations for
138 male-specific fitness effects. Finally, these transcriptome data guided the discovery of mutant
139 alleles in the sperm cell-enriched *gex2*, which induces seed development defects when present
140 in the pollen parent, implying a role in fertilization.

141

142 **Results**

143 **Experimental design and gene expression during maize male reproductive** 144 **development**

145 RNA-seq was performed on four tissues representing integral stages in maize male
146 gametophyte development: immature tassel primordia (TP), isolated unicellular microspores
147 (MS), mature pollen (MP), and isolated sperm cells (SC) (Fig 1A). Techniques were developed to
148 efficiently isolate RNA from TP, MS, and SC (see Methods). RNA was extracted from the inbred
149 maize line B73, with four biological replicates for each tissue. In addition, a single RNA replicate
150 was isolated for the bicellular stage of pollen development (MS-B). Libraries were sequenced
151 using Illumina sequencing (100 bp paired-end reads) and mapped to the B73 AGPv4 reference
152 genome [39]. Principal Component Analysis (PCA) showed samples from each tissue clustering
153 together along PC1 and PC2, which together explained 49.8% of the variance between samples
154 (Fig 1B). One sample, SC1, had significant levels of ribosomal RNA (rRNA) contamination, as
155 well as the fewest number of mapped reads (approximately 1 million). However, to maintain a
156 balanced experimental design with a consistent false discovery rate (FDR), we chose to include
157 SC1 in our analysis of gene expression patterns.

158 Differential gene expression was defined in two ways: in the first, gene expression in later
159 developmental stages was compared to the premeiotic, diploid tassel primordia (TP vs MS, TP
160 vs MP, and TP vs SC); in the second, gene expression was compared between all adjacent
161 developmental stages (TP vs MS, MS vs MP, MS vs SC, MP vs SC) (S2 and S11 Tables).

162 Enriched GO terms highlighted the differences in gene expression among developmental stages
163 and suggested consistency with the established functions of each tissue [22,26,30]. GO terms in
164 MS were consistent with a post-meiotic tissue still at an early stage of development, with terms
165 related to protein synthesis and transport, morphogenesis, and reproduction showing enrichment.
166 MP showed more specific enriched GO terms, including those related to pollen tube growth,
167 signaling, and actin filament-based movement. SC shared many GO terms with MP when
168 compared to MS, but was uniquely enriched for GO terms related to epigenetic regulation of gene
169 expression, such as gene silencing by RNA and histone H3-K9 demethylation.

170

171 **A subset of transposable elements in the maize genome show developmentally** 172 **dynamic expression**

173 To obtain a broad view of TE expression throughout maize development, the RNA-seq
174 data for maize male reproductive development generated by our sequencing (samples with
175 asterisks, S1 Fig) was combined with publicly available RNA-seq of nine-day old above-ground
176 seedlings, juvenile leaves, ovules, another set of independently isolated sperm cells, and three
177 independent studies of pollen RNA-seq [26,30,40,41] & SRP067853. The complete list of
178 samples, their sequencing statistics, references and data availability can be found in S3 Table.
179 All of the raw data were remapped using the same parameters (see Methods). Principal
180 component analysis demonstrates that replicates of the same tissue and growth state typically
181 group together (S1 Fig).

182 We aimed to identify the set of dynamically expressed TEs within the tissues sampled.
183 Thus, the RNA-seq samples were used to calculate expression levels for each individual TE in
184 the genome located more than 2 kb away from annotated non-TE genes. Our rationale was to
185 avoid false positive signals of TE expression due to a TE residing within a gene, and to minimize
186 the influence of read-through transcription from a nearby gene, which could not be distinguished
187 from TE-initiated transcription. To relate TE expression comparatively across development, we
188 used seedling tissue as a baseline against which other tissues were measured. Seedling was
189 chosen for several reasons: it is not a reproductive tissue, it has low to average levels of TE
190 expression, and a large number of TEs show no evidence of expression in this tissue (S2 Fig).

191 Apart from 18.3% of the annotated TEs that are near genes and analyzed separately (see
192 below), we calculated the number of TEs with statistically significant expression differences in
193 each tissue compared to the seedling reference. This identified the subset of TEs that are
194 developmentally dynamic, meaning that they show differential expression in at least one tissue in
195 our dataset compared to the seedling reference. Only 4.4% of all maize annotated TEs are
196 developmentally dynamic, whereas 22.2% of TEs have detectable expression, but do not change
197 in our dataset and therefore are developmentally static (Fig 2A). Finally, the majority of annotated
198 TEs (55%) were not assessed in this analysis, either because no expression was detected in any
199 dataset, or because their sequence lacks the polymorphisms necessary for mapping to a specific
200 TE.

201 Each TE category (Fig 2A) was interrogated for feature overrepresentation. Both dynamic
202 and static TEs are longer than the genome average, and longer than the sets of TEs 'not covered'
203 or 'near genes' (Fig 2B). The finding that expressed TEs as a group (dynamic + static) are longer
204 correlates with Arabidopsis data where longer TE transcripts are overrepresented and
205 differentially regulated when epigenetic repression is lost [42]. Expressed TEs show an under-

206 representation for DNA transposon and SINE families, which are mainly within the ‘near genes’
207 set (Fig 2C). In contrast, the ‘LTR unknown’ TE annotation is over-represented in the dynamic TE
208 set (Fig 2C). Since some LTR retrotransposons are enriched in the pericentromere [43], we tested
209 if the dynamic TE set is enriched in the pericentromere compared to the genome average, but did
210 not detect any correlation (Fig 2D). Therefore, we conclude that expressed TEs are generally
211 longer elements, and the subset of developmentally dynamic TEs are enriched for
212 uncharacterized LTR retrotransposons located throughout the genome.

213

214 **Transposable element transcript levels are up-regulated in the post-meiotic male** 215 **reproductive lineage**

216 From the developmentally dynamic TE set, we calculated the number of differentially
217 expressed TEs in each tissue/stage compared to the seedling reference. In some tissues, such
218 as tassel primordia and ovules, we observed a similar number of TEs up-regulated and down-
219 regulated (Fig 3A), demonstrating that while there are shifts in which TEs are expressed, a
220 genome-scale change in TE expression does not occur. In other tissues, such as juvenile leaves,
221 there is a skew towards increased TE expression. The largest TE up-regulation occurs in the
222 tissues of the male reproductive lineage, including unicellular and bicellular microspores, mature
223 pollen and isolated sperm cells (Fig 3A). The number of up-regulated TEs compared to down-
224 regulated TEs in these tissues suggest that there is a genome-wide activation of TE expression,
225 similar to the DRTS event that occurs in Arabidopsis pollen [31,36]. One important distinction is
226 that TE expression is present in maize sperm cells (Fig 3A), whereas it is not detected in
227 Arabidopsis sperm cells [31]. To verify this finding, we compared our sperm cell RNA-seq data to
228 an independent maize sperm cell dataset [30]. We found that TEs are also significantly expressed
229 in this independent dataset, and 70% of those expressed TEs are also detected in our dataset
230 ($p < 0.001$) (Fig 3B). This shared set of 810 sperm cell-expressed TEs (38% of those detected in
231 our dataset), supports the conclusion that significant expression of TEs occurs in maize sperm
232 cells. Of the sperm cell-expressed TEs, 36% were not observable in total pollen, but rather
233 required the isolation and enrichment of sperm cells for detection (Fig 3B). Overall, we detect 157
234 TEs expressed in both sperm cell datasets that are not expressed throughout development, but
235 specifically in the sperm cells (sperm-cell exclusive).

236 A second notable difference between maize and Arabidopsis is the activation of TE
237 expression early in the male gametophytic phase of maize. A genome-wide increase in TE
238 transcript levels is detected at the earliest post-meiotic stage tested, the microspore, in contrast
239 to low TE expression in the sporophytic tassel primordia (Fig 3A). Arabidopsis TE expression
240 occurs only late in pollen development, after pollen mitosis I when the somatic vegetative cell is
241 generated [31]. To determine if TEs were indeed activated early in maize male reproductive
242 development, we asked if the same TEs that we identified as expressed in the unicellular
243 microspore remain active throughout the male reproductive lineage. We used the set of
244 differentially expressed up-regulated TEs in unicellular microspores (3,335) and found that 62%
245 are still expressed in bicellular microspores and 54% in mature pollen (Fig 3C), demonstrating
246 that once TEs are activated early in development, expression and/or steady-state mRNA
247 frequently remains through pollen maturation. Only some of these male-lineage expressed TEs
248 continue to be expressed in sperm cells (32%), raising the possibility that many TEs with active
249 expression in the early gametophytic stages are under negative/repressive regulation in the

250 gametes. This large-scale developmental activation is potentially limited to the male lineage, as
251 ovules express relatively few TEs (Fig 3A) and only 14% of the male lineage-expressed elements
252 (Fig 3C). Together, our data demonstrate conserved activation of TE expression in the male
253 gametophytes of maize and Arabidopsis, with key differences such as the developmental timing
254 and localization of TE expression in the gamete cells.

255 We determined what types of TEs activate in the male reproductive lineage and sperm
256 cells and compared these to the whole-genome distribution of TEs analyzed. Overall, both male
257 lineage-expressed TEs and sperm cell-expressed TEs reflect the genome-wide TE distribution
258 (Fig 3D). This suggests that TE family type does not have a determining role in the developmental
259 regulation of TE expression. One notable exception is the enrichment of *Mutator* family TE
260 expression in sperm cells (Fig 3D). When normalized for genome-wide TE distribution, *Mutator*
261 element expression is highly enriched across the male lineage, including in sperm cells (Fig 3E).
262 The expression of some *Mutator* TEs in sperm cells is both high confidence (present in both sperm
263 cell datasets) and specific to only that tissue (high confidence sperm cell specific, Fig 3E). LINE
264 L1 elements are also expressed throughout the male lineage and sperm cells, but their expression
265 is general and not specific to these cell types (Fig 3E). Our data demonstrate that there is a
266 general (TE family-independent) activation of TE expression in the male reproductive lineage,
267 with one observable bias towards *Mutator* family expression in both the male lineage and sperm
268 cells.

269

270 **Mature pollen and sperm cells display coexpression of highly expressed genes** 271 **and their neighboring TEs**

272 To determine if TEs have an effect on neighboring gene expression, or vice versa, we
273 analyzed the 36,945 assayable TEs within 2kb of genes from Fig 2A. We calculated the absolute
274 expression level of each genic isoform and categorized them into 100 bins of expression levels
275 for each developmental stage (Fig 4). We find no relationship between gene expression level and
276 the number of up- or down-regulated TEs in tassel primordia or microspores (top row, Fig 4A). In
277 contrast, in both mature pollen and isolated sperm cells there is a positive association between
278 highly expressed genes and the number of up-regulated TEs within 2kb of those genes (bottom
279 row, Fig 4A). Similarly, there is a negative correlation between high gene expression and the
280 number of down-regulated TEs in pollen and sperm cells (Fig 4A). This relationship is not due to
281 the fact that pollen or sperm cell-expressed genes are more likely to be located nearby a TE (S3A
282 Fig). We confirmed that this association between TE and gene regulation in sperm and mature
283 pollen is not due to sample contamination between these two datasets (S3B Fig). It is unclear
284 from these data whether gene expression is influencing TE expression, or TE expression is
285 affecting gene regulation. However, we conclude that specifically in the mature male gametophyte
286 the highest expressed genes are near actively expressing TEs.

287 Comparison of the most highly expressed genes from mature pollen and sperm cells with
288 other tissues assessed in this study (tassel primordia and microspores) shows that such
289 transcripts were generally associated with high tissue-specificity (Fig 4B). Among the top 200
290 highly expressed genes by FPKM value, two-thirds of the genes in each tissue were highly
291 expressed only in that tissue. No single gene was highly expressed in all four tissues. These data
292 are consistent with the idea that genes highly expressed at a particular developmental stage
293 contribute a genetic function specifically required at that stage. We next aimed to determine if

294 these highly-expressed genes, many of which are adjacent to expressed TEs, have measurable
295 functional roles in the male gametophyte.

296

297 **Large-scale insertional mutagenesis supports a relationship between transcript** 298 **level and fitness contribution for vegetative cell-expressed genes**

299 The developmental gene expression dataset generated by this study provided a
300 quantitative framework in which to assess gene function, testing the hypothesis that highly
301 expressed genes contribute significantly to reproductive success – i.e, fitness. The functional
302 validation approach we used relied on a large, sequence-indexed collection of green fluorescent
303 protein (GFP)-marked transposable element (*Ds-GFP*) insertion mutants [44], enabling
304 assessment of the effects of mutations in select genes (Fig 5). We focused on expression data
305 from the MP and SC stages, as these display coordinated expression with TEs (Fig 4). In addition,
306 these stages have distinctive cell fates and roles in reproduction: the vegetative cell generates
307 the pollen tube for competitive delivery of gametes, and the sperm cells accomplish double
308 fertilization. Expression data from seedlings [26] was used to design a comparator sporophytic
309 control. Highly expressed genes, operationally defined as in the top 20% for a tissue by FPKM,
310 were grouped into three mutually exclusive classes: Seedling, Sperm Cell, and Vegetative Cell.
311 The Seedling group also excluded any gene highly expressed in either MP or SC. Due to the
312 significant overlap among genes highly expressed in both MP and SC, we compared expression
313 values to assign each of these genes to a single class. Vegetative Cell genes were not only highly
314 expressed in MP, but were also associated with an FPKM greater in MP than in SC, and vice
315 versa for Sperm Cell genes (S5 Table). All genes in these classes were then cross-referenced
316 with *Ds-GFP* insertion locations to identify potential mutant alleles for study, restricting the search
317 to insertions in coding sequence (CDS), as these were rationalized as most likely to generate
318 loss-of-function effects. Finally, to insure our results were as generalizable as possible, each class
319 list was randomized to identify the specific subset of *Ds-GFP* lines for study. Insertion locations
320 were verified by PCR for 64 of 83 alleles obtained (S6 Table) (see Methods), of which 56,
321 representing mutations in 52 genes, generated sufficient transmission data to include in our final
322 analysis.

323 Mendelian inheritance predicts 50% transmission of mutant and wild-type alleles when a
324 heterozygous mutant is outcrossed to a wild-type plant. However, a mutation that alters the
325 function of a gene expressed during the haploid gametophytic phase can result in a reduced
326 transmission rate if that gene contributes to the fitness of the male gametophyte – i.e., to its ability
327 to succeed in the highly competitive process of pollen tube growth, given that 50% of the pollen
328 population will be wild-type for the same gene. Thus, reduced transmission of a mutant through
329 the male (a male transmission defect) provides not only evidence for gene function in the
330 gametophyte, but also a measure of the mutated gene's contribution to fitness. Transmission
331 rates through the female serve as a control, as 50% transmission through the female would
332 confirm both a single *Ds-GFP* insertion in the genome and male-specificity for any defect
333 identified. To measure the fitness cost of each *Ds-GFP* insertion, heterozygous mutant plants
334 were reciprocally outcrossed with a heavy pollen load to a wild-type plant, maximizing pollen
335 competition within each silk. Transmission rates were then quantified by assessing the ratio of
336 the non-mutant to mutant progeny using a novel scanning system and image analysis pipeline
337 (Fig 5) (see Methods) [45]. Mutant alleles were tracked using linked endosperm markers: either

338 the GFP encoded by the inserted transposable element (Fig 5A-B), or, in ~10% of the lines, a
339 tightly linked *C1*⁺ anthocyanin transgene (present due to the initial *Ds-GFP* generation protocol)
340 (Fig 5C, S7 Table). For an allele to be included in the final dataset, we required a minimum of
341 three independent male outcrosses from two different plants. The number of seeds categorized
342 for each allele ranged from 1,522 to 5,219, with an average of 2,807.

343 Transmission rates for all groups were tested through quasi-likelihood tests on generalized
344 linear models with a logit link function for binomial counts (see Methods, S8 Table). When crossed
345 through the female, no genes showed significant differences from Mendelian inheritance (Fig 6A).
346 When crossed through the male, no genes with insertion alleles in the Seedling category (n=10)
347 showed evidence of abnormal transmission rates (Fig 6B). Most Sperm Cell genes (n=10, 90%)
348 showed no statistically significant transmission defects, with one notable exception (two
349 independent alleles of the *gex2* gene, described in detail below) (Fig 6C). However, among
350 Vegetative Cell genes tested (n=32), a larger proportion of insertion alleles (7 out of 32 or 21.9%)
351 showed significant male transmission defects (quasi-likelihood test, adjusted p-value threshold <
352 0.05) (Fig 6D). The proportions of genes with transmission defects in the three classes were not
353 significantly different by Fisher's exact test (Seedling vs Sperm Cell p-value = 0.500, Seedling vs
354 Vegetative Cell p-value = 0.125, Vegetative Cell vs Sperm Cell p-value = 0.374), likely due to the
355 small number of mutations assessed in the Seedling and Sperm Cell classes.

356 The majority of transmission defects in the Vegetative Cell class genes (six of the seven
357 with significant effects) were modest, at approximately 45% transmission, with only one reducing
358 transmission strongly, to ~30%. Two of the seven Vegetative Cell genes with transmission defects
359 were adjacent to TEs (S7 Table). Notably, all but one of the genes associated with significant
360 defects were measured at a $\log_2(\text{FPKM}) > 8$ (i.e., in the top 5% of Vegetative Cell genes by
361 FPKM). Above this expression level, the percentage of genes associated with a significant defect
362 rises to 50% (six out of twelve). Thus, the most highly expressed Vegetative Cell genes are
363 significantly more likely to be associated with non-Mendelian transmission than the group of
364 Vegetative Cell genes below this expression threshold (1 out of 20) (Fisher's exact test, p-value
365 = 0.00572). Consistent with this observation, an increase in $\log_2(\text{FPKM})$ was associated with both
366 reduced transmission rate and an increase in $-\log_{10}(\text{p-value})$ (linear regression, p-value = 0.0120,
367 0.0255, respectively; adjusted $R^2 = 0.151, 0.116$, respectively). Thus, our data suggest that higher
368 transcript level in the Vegetative Cell predicts a gene-specific contribution to male gametophytic
369 fitness. Vegetative Cell genes showing non-Mendelian inheritance had a range of predicted
370 cellular functions, including cell wall modification, cell signaling, protein folding, vesicle trafficking,
371 and actin binding (Table 1).

372

373 **Table 1. Characteristics of genes showing non-Mendelian inheritance.**

Category	Gene designation (v4)	Ds-GFP allele	Trans. rate	Adjusted p-value	Best BLAST Hit, <i>A. thaliana</i>	Predicted Function (B73v4 Gramene)	Cellular process (inferred)
Vegetative Cell	Zm00001d028437	tdsgR04A02	43.84%	3.29E-04	AT3G61050	Calcium-dependent lipid-binding (CaLB domain) family protein	Cell signaling
Vegetative Cell	Zm00001d037695	tdsgR102H01	45.50%	3.17E-02	AT1G52080	Actin binding protein family	Cytoskeleton
Vegetative Cell	Zm00001d022250	tdsgR33F03	43.95%	1.38E-04	AT2G02370	SNARE associated Golgi protein family	Vesicle trafficking
Vegetative Cell	Zm00001d003431	tdsgR49F11	43.87%	1.38E-04	AT3G05610	Pectinesterase 5	Cell wall modification
Vegetative Cell	Zm00001d014731	tdsgR67C09	44.09%	1.06E-03	AT2G29960	Peptidyl-prolyl cis-trans isomerase CYP20-1	Protein folding
Vegetative Cell	Zm00001d014782	tdsgR92F08	45.06%	1.41E-02	AT1G19940	Endoglucanase 2	Cell wall modification
Vegetative Cell	Zm00001d015901	tdsgR96C12	29.51%	0.00E+00	AT2G24450	Fasciclin-like arabinogalactan protein 3	Cell wall modification
Sperm Cell	Zm00001d005781	tdsgR82A03	33.43%	4.15E-14	AT5G49150	Protein GAMETE EXPRESSED 2	Fertilization
Sperm Cell	Zm00001d005781	tdsgR84A12	23.14%	0.00E+00	AT5G49150	Protein GAMETE EXPRESSED 2	Fertilization

374

375 To ensure the experimental design was robust, we examined two potential confounding
 376 variables: the presence of the *wx1-m7::Ac* allele in a subset of lines tested and the potential for
 377 epigenetic silencing of GFP transgenes (see S1 Methods). We found no evidence that the
 378 presence of *wx1-m7::Ac* significantly impacted the overall conclusions drawn from the dataset,
 379 as well as no evidence of epigenetic silencing of GFP transgenes.

380

381 **Insertional mutants in sperm cell-expressed *gex2* cause paternally triggered** 382 **aberrant seed development**

383 The male-specific transmission defect for the sole affected gene in the Sperm Cell class,
 384 Zm00001d005781 (GRMZM2G036832), was significantly higher than the average defect across
 385 all genes identified with decreased transmission through the male. Sequencing confirmed that the
 386 *Ds-GFP* elements in the two independent alleles, tdsgR82A03 and tdsgR84A12, associated with
 387 average transmission rates of 33.4% and 23.1%, respectively, were inserted into their predicted
 388 exonic locations (Fig 7A). In addition, both alleles showed an unusual phenotype of
 389 underdeveloped or aborted seeds and ovules with no apparent seed development when crossed
 390 through the male, despite heavy pollination (Fig 7B). These features motivated further
 391 investigation of this gene.

392 Across maize tissues, Zm00001d005781 is highly and specifically expressed in sperm
 393 cells [41] (Fig 7C). Zm00001d005781 was previously identified in maize sperm cells via EST
 394 sequencing and named *gamete expressed 2*, or *gex2*, but no mutant has been described [46].
 395 The gene is among the top 200 expressed in sperm cells, and like many highly tissue-specific
 396 genes expressed in pollen and sperm, it is within 2kb of a transcriptionally active TE, both a
 397 downstream RLG retrotransposon that displays sperm cell-specific activation and an upstream

398 DHH family TE that is not detectably transcribed (Fig 7C). Predicted *gex2*-like genes are widely
399 distributed throughout the currently sequenced Embryophyta taxa. The Arabidopsis ortholog,
400 *GEX2*, has been described as necessary for gamete attachment and effective double fertilization
401 [47]. Maize *gex2*, a single copy gene, encodes a protein with predicted domain structure similar
402 to that of Arabidopsis *GEX2* (44.2% protein similarity overall). Both proteins harbor a large N-
403 terminal non-cytoplasmic region including filimin repeat-like domains, predicted transmembrane
404 domains near their C-termini, and a small C-terminal cytoplasmic region (Fig 7D).

405 Small and aborted seeds were quantified for both *gex2* insertion alleles when outcrossed
406 to wild-type plants as heterozygotes, as well as for outcrosses from *gex2-tdsgR84A12*
407 homozygotes (S10 Table, S4 Fig). As controls, the same assessment was made for pollinations
408 from two different heterozygous *Ds-GFP* insertion lines that were not associated with transmission
409 defects (*tdsgR12H07*, *tdsgR46C04*), as well as heterozygous plants carrying the *Ds-GFP*
410 associated with the strongest male transmission defect (29.5% transmission) in the Vegetative
411 Cell class (*tdsgR96C12*). These three *Ds-GFP* insertions showed similar percentages of aborted
412 seeds (Fig 7E). In contrast, pollination with both *gex2::Ds-GFP* insertion alleles was associated
413 with increased percentages of small or aborted seeds, significantly so in *gex2-tdsgR84A12*
414 (pairwise t-test against *Ds-GFP* controls separately, all p-values < 0.05). Pollination from *gex2-*
415 *tdsgR84A12* homozygotes approximately doubled the percentage of small or aborted seed
416 percentages from *gex2::Ds-GFP* heterozygous plants (pairwise t-test against *Ds-GFP* controls
417 separately, all p-values < 0.01). From the heterozygous *gex2::Ds-GFP* crosses, small seeds with
418 endosperm large enough for DNA preparation were genotyped, and 79.2% were found to harbor
419 the *gex2* mutation, whereas in crosses from the *tdsgR46C04* control, the *Ds-GFP* insertion
420 showed Mendelian segregation in small seeds. These data support the hypothesis that aberrant
421 seed development is associated with fertilization by *gex2::Ds-GFP* sperm.

422 If *GEX2* acts to promote double fertilization, the arrival of a *gex2::Ds-GFP* pollen tube at
423 the embryo sac could lead to failure of one or both fertilization events. Given an active polytubey
424 block, this could produce the observed gaps between seeds on the ear, resulting from ovules
425 associated with completely failed fertilization, or with very early seed abortion due to single
426 fertilization. To explore this possibility, seedless area was measured in ear projection images for
427 ears pollinated by *gex2::Ds-GFP* heterozygous and homozygous mutants, as well as for the *Ds-*
428 *GFP* control ears already described. As heterozygotes, all three *Ds-GFP* controls as well as the
429 *gex2-tdsgR82A03* allele are associated with low levels of seedless area (each at 4%), whereas
430 the *gex2-tdsgR84A12* shows a non-significant increase in seedless area (8.70%) (S5 Fig).
431 However, ears pollinated by homozygous *gex2-tdsgR84A12* pollen had 31.48% seedless area, a
432 significant increase (pairwise t-test against *Ds-GFP* controls separately, all p-values < 0.0001).
433 To test for aberrant fertilization more directly, seed development was assessed at 4 days post-
434 pollination with either wild-type or *gex2-tdsgR84A12* homozygous pollen (Fig 8, Table 2). Typical
435 embryo and endosperm development, as well as indication of the polytubey block (i.e., arrival of
436 only single pollen tubes at the embryo sac), was observed in all ovules assessed from wild-type
437 pollination. In contrast, half of the ovules assessed following pollination with *gex2::Ds-GFP*
438 showed significant evidence of abnormal double fertilization, demonstrating single fertilization of
439 either embryo or endosperm or indication of arrival of more than one pollen tube at the embryo
440 sac (Fisher's exact test, p-value = 0.000241). We conclude that maize *GEX2* is part of the sperm
441 cell machinery that helps insure proper double fertilization.

442
443
444

Table 2. Seed development at 4 days after pollination by wild-type or *gex2::Ds-GFP* pollen

Pollen parent	One synergid penetrated by a pollen tube			Both synergids penetrated by a pollen tube		
	Both embryo and endosperm	Endosperm without embryo	Embryo without endosperm	Both embryo and endosperm	Endosperm without embryo	Embryo without endosperm
<i>gex2-tdsgR84A12/gex2-tdsgR84A12</i>	6	2	2	0	0	2
Wild-type	28	0	0	0	0	0

445
446

447 Discussion

448 Transposable element dynamics in the maize male gametophyte

449 Our analysis of TE expression during maize male reproductive development provides an
450 informative comparison to similar analyses in Arabidopsis, an evolutionarily distant plant with a
451 genome landscape that is quite distinct from maize. Although maize has a higher number and
452 percentage of its genome occupied by TEs compared to Arabidopsis, we found that only a fraction
453 of maize TEs are developmentally dynamic with regards to transcript accumulation. These
454 ‘dynamic’ TEs tend to be longer elements than average, and are enriched for *Mutator* family DNA
455 transposons and ‘unknown’ classification LTR retrotransposons. From this dynamic TE set, we
456 were able to identify individual elements that are expressed in a number of specific tissues.
457 However, more globally, there is a trend towards activation of TE transcription over the course of
458 the development of the mature male gametophyte. This finding confirms that both monocots and
459 eudicots have developmental activation of TE expression in pollen [31]. This conservation
460 suggests that the roles of TE and TE-induced small RNAs during reproductive development may
461 also be conserved between monocots and eudicots [38,48,49]. Consistent with our findings, a
462 recent study found that spontaneous retrotransposon mutations are much more frequent through
463 the male than the female in certain maize lines [34].

464 Although TE activation is conserved in maize and Arabidopsis pollen, we have identified
465 key differences in the timing and location. Maize TE activation is detected earlier (in the unicellular
466 microspore) compared to when it is thought to occur in Arabidopsis [31]. Transcripts from these
467 early-activated TEs in the microspores typically remain detectable through pollen development
468 and in the mature pollen grain, which may be due to continued expression or transcript stability.
469 A second distinction is the location of activation, which in Arabidopsis is confined to the pollen
470 vegetative cell nucleus [31,35], whereas in maize also occurs in sperm cells. *Mutator* family TE
471 transcripts are overrepresented in the pool of sperm-cell TE transcripts, suggesting that this family
472 of TEs may have evolved (or co-opted) specific regulatory mechanism(s) such as an enhancer
473 element that confers expression in this cell type.

474 We also examined the correlation between gene and TE expression. We found that
475 generally there is no association between highly expressed genes and TE activation. However,
476 in mature pollen and sperm cells there is a positive correlation: the more highly expressed a gene
477 is, the more likely it is to have expressed TEs within 2kb. This tissue-specific correlation is not
478 due to there being more TEs near pollen- or sperm-expressed genes. This demonstrates a
479 developmentally specific co-regulation of gene and TE expression. Several potential mechanisms
480 account for this observation. First, the programmed activation of TE expression may influence
481 chromatin, enhancer, or other regulatory function that influences the neighboring genes. Second,
482 the genes and TEs may be directly controlled by the same mechanism of large-scale epigenetic
483 activation. The same mechanisms that repress TE activation may be responsible for limiting the
484 expression of some genes to a single specific tissue or developmental time point. Third, the gene
485 activation may control the expression of the neighboring TE. Ongoing studies aim to decipher
486 these possibilities.

487

488 **The *gex2* gene has a conserved role in promoting double fertilization**

489 One highly and specifically expressed gene in maize sperm cells with a nearby activated
490 TE is *gex2*. The gene *gex2* is located between two transposable elements, one of which (an RLG
491 retrotransposon) was also specifically transcriptionally activated in sperm cells. Mutations in *gex2*
492 led not only to severely reduced transmission through the male, but also, in contrast to other
493 mutations analyzed in this study, to paternally triggered post-fertilization defects, such as an
494 increase in underdeveloped and aborted seeds and an increase in the seedless surface of the
495 ear. *gex2* was first identified in maize by sperm cell EST sequencing [46], which led to
496 identification of the orthologous gene in Arabidopsis, *GEX2*, and its sperm cell-specific promoter
497 [50]. In Arabidopsis, seed abortion and single fertilization events were observed at increased
498 frequency in *gex2* mutant-pollinated plants, both of the egg cell (leading to seeds that contained
499 only an embryo) and of the central cell (leading to seeds that contained only endosperm) [47].
500 Our results are similar, with maize *gex2* mutant pollen causing single fertilization events in embryo
501 sacs, small and aborted seeds, and leading to aberrant early seed phenotypes consistent with
502 single fertilization. Mechanistically, Arabidopsis *GEX2* appears to contribute to gamete
503 attachment through interactions between plasma membrane-localized *GEX2* in the sperm cell
504 and either the female egg or central cells. The two orthologues share a predicted domain
505 structure, including a large N-terminal non-cytoplasmic region containing filimin repeat domains
506 potentially acting in gamete attachment [47], raising the possibility that maize *GEX2* acts similarly
507 during double fertilization. However, we cannot rule out additional functions for *gex2* in the earlier
508 phase of pollen tube growth.

509 Few genes with clear fertilization-associated functions have been identified in maize.
510 Three mutations which cause paternally-triggered aberrant ('rough') endosperm development
511 have been isolated, but their corresponding genes have not been molecularly identified [51]. Two
512 genes that are more clearly linked to double fertilization, *MTL/NLD/ZmPLA1* [17–19] and *ZmDMP*
513 [20], have been identified based on their contributions to pollen-induced haploid induction
514 associated with the 'Stock 6' line of maize. In our RNA-seq dataset, *MTL/NLD/ZmPLA1* is
515 detected at low levels in SC, but is 7-fold higher in MP, suggesting enriched expression in the
516 vegetative cell; *ZmDMP* is not detected in either SC or MP. Similar to *gex2*, mutants in both
517 *MTL/NLD/ZmPLA1* and *ZmDMP* cause small/aborted seed phenotypes when used as a pollen

518 parent, raising the possibility that the three genes could act in related fertilization mechanisms.
519 However, the exact mechanisms by which mutations in *MTL/NLD/ZmPLA1* and *ZmDMP* lead to
520 haploid induction are unclear. Single fertilization of the central cell has been suggested as a
521 possible mechanism [52], but more recent work points toward degradation of the paternal
522 chromosomes in the zygote following double fertilization [53–55], while some data support the
523 idea that both mechanisms contribute [56]. Some of the abnormal and early aborted seeds in
524 *gex2* mutant pollinations are not explained simply by single fertilization events, which are
525 expected to produce germless or endospermless seeds. The latter of these abort early and
526 resemble unfertilized ovules [57], and thus would likely contribute to the observed gaps on the
527 ears. Abnormal seeds with some endosperm development may indicate other post-fertilization
528 defects caused by *gex2* loss-of-function, similar to *mtl/nld/zmpla1* or *zmdmp* loss-of-function.
529 Direct, detailed comparison of seed phenotypes induced by pollination with these mutants could
530 prove useful to understanding mechanisms both of pollen-induced haploid induction, and of
531 fertilization in general.

532

533 **Maize pollen provides a powerful model for quantifying gene-specific** 534 **contributions to fitness**

535 Despite the explosion of omic-scale methods to characterize genomes and to measure
536 molecular characters (e.g., transcript levels), our ability to predict phenotypic relevance for
537 specific genes is limited, particularly in multicellular organisms with complex genomes. A simple
538 assumption is that a high transcript level at a particular developmental stage implies an important
539 function for the associated gene at that stage, thus pointing toward a potential phenotypic role.
540 However, this has been difficult to test at larger scales, as generating gene-specific quantitative
541 phenotypic data is laborious and standardization can be challenging. This study begins to address
542 this assumption with an initial systematic assessment of the functional relevance of highly
543 expressed genes in maize pollen, taking advantage of the ease of reciprocal outcross pollination
544 in maize, the availability of a sufficient number of marked and likely null mutations, and the
545 development of an imaging technique that enables sensitive quantitation.

546 The progamic (i.e., post-pollination) phase of male gametophyte development can be
547 thought of as comprising two stages. In the first stage, pollen tube growth through the silk, male
548 gametophytes compete to be the first to reach the female gametophyte. Each pollen grain is an
549 independent multicellular organism, often genetically-distinct from others in the population due to
550 meiotic recombination, with its sole purpose to deliver the sperm cells for double fertilization.
551 Thus, in an outcrossing plant with an extensive stigma and style like maize, there is likely a
552 heightened context for competition among individuals in pollen populations [10]. We reasoned
553 that genes involved in this stage could be identified by measuring fitness costs for mutations in
554 genes highly expressed in the vegetative cell, which is responsible for pollen tube growth. We
555 found that CDS-insertion alleles for 7 out of 32 (21.9%) tested genes in this class are associated
556 with mild to moderate male-specific transmission defects. These results support the idea that
557 competitive pollen tube growth requires a broad array of genes, as genes associated with
558 demonstrated fitness contributions were assigned to a variety of predicted cellular functions.
559 However, it is worth noting that three of the seven are predicted to directly influence modification
560 of the cell wall (Table 1). Extrapolating from this dataset to the entire genome gives an estimate

561 of 600 vegetative cell genes in maize leading to non-Mendelian segregation when inactivated, at
562 least under competitive conditions.

563 In the second stage, following the arrival of the pollen tube at the embryo sac, the pollen
564 tube must properly penetrate the embryo sac and release the sperm cells it contains, which then
565 must fuse with the egg cell and central cell for double fertilization. Due to the presence of the
566 polytubey block, the second stage contrasts with the first in that competition is likely minimal.
567 Upon delivery, sperm cells must perform the tasks of adhesion, communication, fusion with the
568 egg cell and central cell, and karyogamy. However, due to the lack of competition, minor problems
569 with these processes may not lead to failure. We aimed to investigate this second stage by testing
570 genes highly expressed in the sperm cell. Our limited dataset identified 1 out of 10 of tested genes
571 in this category as associated with a male transmission defect. This mutant gene, *gex2*, had
572 among the strongest transmission defects observed in this study – and, intriguingly, is also
573 associated with the highest FPKM among the sperm cell genes tested. Overall, our results
574 suggest a scenario in which mutations in a larger proportion of genes operating during pollen tube
575 growth lead to slight reductions in fitness, whereas mutations of relatively fewer genes operating
576 in double fertilization lead to strong reductions in fitness. Testing a larger number of *Ds-GFP*
577 insertions in vegetative cell- and sperm cell-enriched genes, encompassing a wider range of
578 expression levels, would help establish if these trends are borne out.

579 To our knowledge, this study is the largest yet in plants to explore the possibility of a
580 relationship between transcriptomic data and quantitative phenotypic effects of mutations. The
581 largest phenomic studies (hundreds of mutant lines assessed) of leaf or reproduction-related
582 characters using the Arabidopsis T-DNA mutant collection have identified phenotypic effects in
583 ~4% of lines assessed, although these were not guided by transcriptomic data [58,59]. Thus, the
584 higher frequency of phenotypic effects we found could be due to sampling only mutations in genes
585 that are most highly expressed at developmental stages most relevant to the phenotype. Genome
586 scale measurement of the fitness costs of gene knockouts via competitive assays have been
587 carried out in yeast [60][61] and bacteria [62][63]. All of these studies found that, for particular
588 environmental conditions, there was little to no correlation between the expression level of a gene
589 and its impact on fitness in that condition. In contrast, we found a small but significant correlation
590 between high expression and fitness in our smaller dataset of vegetative cell genes (Fig 6D). This
591 could be indicative of biological differences between single-celled organisms and more complex
592 organisms like maize, or could be a feature of the highly specialized pollen tube. Interestingly, our
593 results are consistent with a recent meta-analysis that identified higher mRNA expression levels
594 as a feature distinguishing gene models with known phenotypes from the overall population of
595 gene models defined by sequencing and other molecular approaches [64]. Given that genes –
596 i.e., those gene models with a clear functional role in determining the characteristics of an
597 organism – represent only a subset of total predicted gene models, genome-scale approaches
598 that can associate quantitative phenotypes with specific genes appear important for achieving a
599 more global understanding of genotype-to-phenotype relationships.

600

601 **Methods**

602 **Plant materials**

603 Maize inbred line B73 was used for all RNA isolations. Plants were grown in a controlled
604 greenhouse environment (16 hrs light, 8 hrs dark, 80 F day/70 F night) and in the field at the

605 Botany & Plant Pathology Field Lab (Oregon State University, Corvallis, OR) using standard
606 practices. Lines containing *Ds-GFP* insertion alleles were acquired from the Maize Genetics
607 Cooperation Stock Center.

608

609 **RNA isolation, library preparation and sequencing**

610 Detailed methods are available in S1 Methods. Briefly, tissue was isolated either by
611 dissection (TP), differential density centrifugation (MS, MS-B and SC), or collection at anthesis
612 (MP). Total RNA from TP, MS, and MP was extracted using a modified Trizol Reagent (Life
613 Technologies) protocol; SC total RNA was extracted via a phenol/chloroform protocol. Poly-A
614 RNA (mRNA) was isolated using streptavidin magnetic beads (New England Biolabs, # S1420S)
615 and a biotin-linked poly-T primer. RNA libraries were prepared and sequenced by the Central
616 Services Lab (CSL) at the Center for Genome Research and Biocomputing (CGRB, Oregon State
617 University) using WaferGen robotic strand specific RNA preparation (WaferGen Biosystems) with
618 an Illumina TruSeq RNA LT (single index) prep kit and run on an Illumina HiSeq 3000 with 100
619 bp paired-end reads.

620

621 **Mapping reads to genes, differential expression assessment and GO enrichment** 622 **analysis**

623 Ribosomal reads (rRNA) were removed from all samples using STAR, version 2.5.1b [65]
624 to map reads to a repository of maize rRNA sequences (parameters: --outSAMunmapped Within
625 --outSAMattributes NH HI AS NM MD --outSAMstrandFieldintronMotif --limitBAMsortRAM
626 50000000000 --outReadsUnmapped Fastx). The number of mappable reads generated from
627 each sample after rRNA removal ranged from approximately 1 million to approximately 41 million,
628 with an average mappable reads of approximately 18 million per sample. Total reads, mapped
629 reads, rRNA contamination, and other statistics are summarized in S1 Table.

630 rRNA-filtered sequences were mapped to the maize reference genome, version B73
631 RefGen_v4.33 [39] using STAR, keeping only unique alignments (parameters: --
632 outSAMunmapped Within --outSAMattributes NH HI AS NM MD --outSAMstrandField intronMotif
633 --outFilterMultimapNmax 1 --limitBAMsortRAM 50000000000). Transcript levels of annotated
634 gene isoforms were measured using Cufflinks, version 2.2.1 [66]. FPKM (fragments per kilobase
635 of transcript per million mapped reads) values are shown in S4 Table. Differential expression was
636 calculated between each tissue with Cuffdiff, version 1.0.2, using default parameters. FPKM
637 counts were normalized using the geometric library normalization method. A pooled dispersion
638 method was used by Cuffdiff to model variance. Differential expression results are summarized
639 in S11 Table.

640 Gene ontology (GO) terms were found for enriched genes in each tissue using the AgriGO
641 2: GO Analysis Toolkit [67]. Enriched genes were defined as the top 300 significantly differentially
642 expressed genes (q-value) from Cuffdiff output, with ties broken by \log_2 fold change. Enriched
643 sets were split into up- and down-expressed genes. GO term enrichment was calculated using
644 the singular enrichment analysis method with a Fisher test and Yekutieli multi-test adjustment.
645 GO annotations were based off the maize-GAMER annotation set [68].

646

647 **Mapping reads to transposable elements**

648 The rRNA-filtered reads were quality trimmed (QC30) and adapter sequences were
649 removed using BBDUK2 [69]. The remaining sequences were mapped to the whole genome using
650 STAR, allowing mapping to at most 100 'best' matching loci. (parameters: --outMultimapperOrder
651 Random --outSAMmultNmax -1 --outFilterMultimapNmax 100). For paired-end reads, the
652 unmapped reads were re-mapped using single-end approach to maximize the number of
653 mappable reads. The mapping percentage is reported in S3 Table. Because 19% of the total
654 reads in the dataset mapped to more than one location, such reads were mapped to only their
655 best match in the genome, and when multiple best matches existed, they were mapped to all of
656 these loci, and then counted fractionally. For example, if one read maps to 4 TE locations equally
657 well, each TE would receive a weighted value of 0.25 mapped reads. Because the TE expression
658 of the aberrant SC1 biological replicate did not cluster with the other three SC replicates (S1 Fig),
659 it was removed from all subsequent analysis of TE expression.

660

661 **Principal component analysis (PCA)**

662 Using the maize gene and TE annotation file available from Ensembl Genomes (v38) [39],
663 a combined annotation file was generated for both genes and TEs to run PCA for all samples.
664 FeatureCounts [70] was used to calculate the accumulation of each gene and TE in all samples
665 following fractional assignment of reads (parameters: -O --largestOverlap -M --fraction -p -C). This
666 counts file was used in DESeq-2 [71] to generate the PCA plot.

667

668 **Analysis of transposable elements**

669 From the featurecount file (described above), counts of TEs (farther than 2kb from genes)
670 were retained for further analysis. First, normalized read counts for all TEs were obtained (data
671 in S2 Fig) and then, after selecting seedling as the reference tissue, pairwise volcano plots were
672 generated for all samples against the reference seedling tissue. The number of TEs statistically
673 significantly up- and down-regulated in each tissue was calculated and plotted (Fig 3).

674 All TEs less than 2kb away from a gene were categorized as 'near genes' TEs. Any TE
675 with low expression that was excluded by DESeq for pairwise comparison, was counted in 'not
676 covered' category. Compared to seedling, if a TE was found to be up-regulated or down-regulated
677 in any tissue, it was categorized as a dynamic TE. Remaining TEs with p-value > 0.05 compared
678 to seedling were categorized as static TEs, as no evidence of TE expression was observed over
679 different developmental time points analyzed. For all categories, the length, family or distance
680 from centromere was calculated based on the published TE annotation file.

681

682 **Validation of *Ds-GFP* insertion sites**

683 A FASTA file containing 2 kb of genomic sequence surrounding each *Ds-GFP* insertion
684 site was used as input to a primer3-based tool to generate a pair of specific primers to genotype
685 individual plants from each line (<https://vollbrechtlab.gdcb.iastate.edu/tools/primer-server/>). The
686 primers used for each *Ds-GFP* line are listed in S6 Table.

687 To genotype the plants, two 7 mm discs of leaf tissue were collected from each plant using
688 a modified paper punch. The samples were collected in 1.2 ml tubes that fit within a labeled 96
689 well plate/rack (<https://vollbrechtlab.gdcb.iastate.edu/tools/tissue-sample-plate-mapper/>) (Phenix
690 Research Products, Candler, NC; M845 and M845BR or equivalent). Genomic DNA was isolated

691 from the leaf punches [72] with the following modifications. An additional centrifugation (3,000 g
692 for 10 min.) was added to clear the leaf extracts prior to loading onto a 96-well glass fiber filter
693 plate (Pall, 8032). DNA was eluted from filter plates in 125 μ L water, and 2 μ L was used as
694 template for PCR. Amplification followed standard PCR conditions using GoTaq Green Master
695 Mix (Promega) with 4% DMSO (v/v) and amplicons were resolved using agarose gel
696 electrophoresis. Lines were genotyped using the pair of *Ds-GFP* line gene-specific primers plus
697 one *Ds*-specific primer (JSR01 GTTCGAAATCGATCGGGATA or JGP3
698 ACCCGACCGGATCGTATCGG). All lines were also screened by PCR for the presence of *wx1-*
699 *m7::Ac* using primers for *wx1* (CACAGCACGTTGCGGATTTC) and *Ac*
700 (CCGGATCGTATCGGTTTTTCG). Followup PCR to test for co-segregation of GFP fluorescence
701 with the presence of the insertion used the appropriate set of three PCR primers (two gene-
702 specific and one *Ds*-specific) and DNA prepared either from endosperm or seedling leaves [73].
703

704 **Insertional mutagenesis transmission quantification and statistics**

705 Heterozygous lines with PCR-validated *Ds-GFP* insertion alleles were planted in the
706 Botany & Plant Pathology Field Lab (Oregon State University, Corvallis, OR). All insertions were
707 in coding sequence (CDS) sites. Heterozygous *Ds-GFP* plants were outcrossed to tester plants
708 (*c1/c1 wx1/wx1* or *c1/c1* genetic background) through both the female and the male, with male
709 pollinations made with a heavy pollen load on extended silks (silks that had been allowed to grow
710 for at least two days following cutback). Following harvest, resulting ears were imaged using a
711 custom rotational scanner in the presence of a blue light source and orange filter for GFP seed
712 illumination [45]. Briefly, videos were captured of rotating ears, which were then processed to
713 generate flat cylindrical projections covering the surface of the ear (for examples, see Figs 5 and
714 7). Seeds were manually counted using the Cell Counter plugin of the Fiji distribution of ImageJ
715 [74]. Ears showing evidence of more than a single *Ds-GFP* insertion (~75% GFP transmission)
716 were excluded from further analysis. Seed transmission rates of remaining ears were quantified
717 using a generalized linear model with a logit link function for binomial counts and a quasi-binomial
718 family to correct for overdispersion between parent lines. Non-Mendelian inheritance was
719 assessed with a quasi-likelihood test with p-values corrected for multiple testing using the
720 Benjamini-Hochberg procedure to control the false discovery rate at 0.05. Significant non-
721 Mendelian segregation was defined with an adjusted p-value < 0.05. Proportions of genes with
722 male-specific transmission defects in the Seedling, Sperm Cell, and Vegetative Cell categories
723 were compared using a two-sided Fisher's exact test, with significance defined as a p-value <
724 0.05. A two-sided Fisher's exact test was also used to compare the proportions of male-specific
725 transmission defects in the most highly expressed genes and the less highly expressed in the
726 vegetative cell category. A two-sided test for equality of proportions with continuity correction was
727 used to compare transmission rates in families with partial *Ac* presence. A Git repository
728 containing statistical tests and plotting information for this portion of the study can be found at
729 https://github.com/fowler-lab-osu/maize_gametophyte_transcriptome.
730

731 **gex2 sequence analysis and phenotype characterization**

732 Maize *gex2* protein sequence (Zm00001d005781_T002) was retrieved from the Maize
733 Genetics and Genomics Database (MaizeGDB) hosting of the B73 v4 genome [39,75].
734 Arabidopsis *GEX2* protein sequence (AT5G49150.3) was retrieved from the Arabidopsis

735 Information Portal (ARAPORT) Col-0 Araport11 release [76,77]. Protein sequences were aligned
736 using EMBOSS Needle [78]. Maize and Arabidopsis *gex2* protein domains were predicted by
737 InterPro [79], with transmembrane helix predictions by TMHMM [80]. Prediction of land plant
738 species *gex2* conservation was retrieved from PLAZA, gene family HOM04M006791 [81]. Maize
739 *gex2* gene duplication searches were performed using BLAST [82] and the B73 v4 genome. To
740 confirm the predicted insertion sites for the two *gex2::Ds-GFP* alleles, flanking insertion site
741 fragments were PCR-amplified with a gene-specific primer and a *Ds-GFP*-specific primer
742 (*DsGFP_3UTR* – TGCAAGCTCGAGTTTCTCCA) and sequenced via Sanger sequencing.

743 To quantify small seed phenotype, mature, dried down maize ears were imaged prior to
744 seed removal from the ear. For small seeds selection, the ear was first visually scanned row by
745 row from the top to the bottom of the ear. Seeds that were noticeably smaller than their
746 surrounding (regular-sized) seeds are carefully removed from the ear using a pin tool. This
747 sometimes required the removal of regular-sized surrounding seeds, which were saved for later
748 counting. A second visual inspection of the ear often resulted in additional small seeds and is
749 recommended. All remaining seeds were then removed from the ear by hand or using a hand
750 corn sheller tool (Seedburo Equip. Co., Chicago, IL). The ear was screened again for any small
751 (flat/tiny) seeds that could have been missed previously. The cob was inspected prior to
752 discarding, and if any small seed was left behind it was removed and accounted for. Small/smaller
753 seeds and regular-sized seeds were counted and counts were recorded (S10 Table). To measure
754 seedless area, ears were scanned as previously described to create flat surface projections.
755 "Seedless area" was defined as ear surface area that lacked mature or partially developed seeds.
756 Seedless area was quantified as a percentage of total area, as measured with the "Freehand
757 selection" tool of the Fiji distribution of ImageJ [74]. A Git repository containing statistical tests
758 and plotting information for this portion of the study can be found at [https://github.com/fowler-lab-
759 osu/maize_gametophyte_transcriptome](https://github.com/fowler-lab-osu/maize_gametophyte_transcriptome).

760 For analysis of embryo sacs by confocal microscopy, tissues were stained with acriflavine,
761 followed by propidium iodide staining [83,84]. After staining, samples were dehydrated in an
762 ethanol series and cleared in methyl salicylate. Samples were visualized on a Leica SP8 point-
763 scanning confocal microscope using excitations of 436 nm and 536 nm and emissions of $540 \pm$
764 20 nm and 640 ± 20 nm.

765
766

767 **Acknowledgements**

768 We thank O. Childress, H. Fowler, B. Galardi, B. Hamilton, R. Hartman, and C. Lambert
769 for their tireless seed counting, genotyping, field work, and other technical assistance; and Dr.
770 Lian Zhou for her contributions to maize field genetics. We also thank K. Wimalanathan and T.
771 Shibamoto for computational support at ISU, and M. Dasenko and the Center for Genome
772 Research and Biocomputing for library preparation, sequencing and computational support at
773 OSU. We thank D. Auger for reading the manuscript.

774
775

776 **References**

777 1. Yang W-C, Shi D-Q, Chen Y-H. Female gametophyte development in flowering plants.

- 778 Annu Rev Plant Biol. 2010;61: 89–108. doi:10.1146/annurev-arplant-042809-112203
- 779 2. Zhou L-Z, Juranić M, Dresselhaus T. Germline Development and Fertilization Mechanisms
780 in Maize. *Mol Plant*. 2017;10: 389–401. doi:10.1016/j.molp.2017.01.012
- 781 3. McCormick S. Male Gametophyte Development. *Plant Cell*. 1993;5: 1265–1275.
782 doi:10.1105/tpc.5.10.1265
- 783 4. Hafidh S, Fila J, Honys D. Male gametophyte development and function in angiosperms: a
784 general concept. *Plant Reprod*. 2016;29: 31–51. doi:10.1007/s00497-015-0272-4
- 785 5. Dresselhaus T, Sprunck S, Wessel GM. Fertilization Mechanisms in Flowering Plants. *Curr*
786 *Biol*. 2016;26: R125–39. doi:10.1016/j.cub.2015.12.032
- 787 6. Zhou L-Z, Dresselhaus T. Chapter Seventeen - Friend or foe: Signaling mechanisms during
788 double fertilization in flowering seed plants. In: Grossniklaus U, editor. *Current Topics in*
789 *Developmental Biology*. Academic Press; 2019. pp. 453–496.
790 doi:10.1016/bs.ctdb.2018.11.013
- 791 7. Lausser A, Kliwer I, Srilunchang K-O, Dresselhaus T. Sporophytic control of pollen tube
792 growth and guidance in maize. *J Exp Bot*. 2010;61: 673–682. doi:10.1093/jxb/erp330
- 793 8. Mizukami AG, Inatsugi R, Jiao J, Kotake T, Kuwata K, Ootani K, et al. The AMOR
794 Arabinogalactan Sugar Chain Induces Pollen-Tube Competency to Respond to Ovular
795 Guidance. *Curr Biol*. 2016;26: 1091–1097. doi:10.1016/j.cub.2016.02.040
- 796 9. Higashiyama T, Takeuchi H. The mechanism and key molecules involved in pollen tube
797 guidance. *Annu Rev Plant Biol*. 2015;66: 393–413. doi:10.1146/annurev-arplant-043014-
798 115635
- 799 10. Williams JH, Reese JB. Evolution of development of pollen performance. *Curr Top Dev*
800 *Biol*. 2019;131: 299–336. doi:10.1016/bs.ctdb.2018.11.012
- 801 11. Arthur KM, Vejlupekova Z, Meeley RB, Fowler JE. Maize ROP2 GTPase provides a
802 competitive advantage to the male gametophyte. *Genetics*. 2003;165: 2137–2151.
803 Available: <https://www.ncbi.nlm.nih.gov/pubmed/14704193>
- 804 12. Cole RA, Synek L, Zarsky V, Fowler JE. SEC8, a subunit of the putative Arabidopsis
805 exocyst complex, facilitates pollen germination and competitive pollen tube growth. *Plant*
806 *Physiol*. 2005;138: 2005–2018. doi:10.1104/pp.105.062273
- 807 13. Huang JT, Wang Q, Park W, Feng Y, Kumar D, Meeley R, et al. Competitive Ability of
808 Maize Pollen Grains Requires Paralogous Serine Threonine Protein Kinases STK1 and
809 STK2. *Genetics*. 2017;207: 1361–1370. doi:10.1534/genetics.117.300358
- 810 14. Maruyama D, Hamamura Y, Takeuchi H, Susaki D, Nishimaki M, Kurihara D, et al.
811 Independent control by each female gamete prevents the attraction of multiple pollen tubes.
812 *Dev Cell*. 2013;25: 317–323. doi:10.1016/j.devcel.2013.03.013
- 813 15. Robertson DS. A study of heterofertilization in diverse lines of maize. *J Hered*. 1984;75:
814 457–462. doi:10.1093/oxfordjournals.jhered.a109985
- 815 16. Wu C-C, Diggle PK, Friedman WE. Kin recognition within a seed and the effect of genetic
816 relatedness of an endosperm to its compatriot embryo on maize seed development. *Proc*
817 *Natl Acad Sci U S A*. 2013;110: 2217–2222. doi:10.1073/pnas.1220885110
- 818 17. Kelliher T, Starr D, Richbourg L, Chintamanani S, Delzer B, Nuccio ML, et al.
819 MATRILINEAL, a sperm-specific phospholipase, triggers maize haploid induction. *Nature*.
820 2017;542: 105–109. doi:10.1038/nature20827
- 821 18. Gilles LM, Khaled A, Laffaire J-B, Chaignon S, Gendrot G, Laplaige J, et al. Loss of pollen-

- 822 specific phospholipase NOT LIKE DAD triggers gynogenesis in maize. *EMBO J.* 2017;36:
823 707–717. doi:10.15252/embj.201796603
- 824 19. Liu C, Li X, Meng D, Zhong Y, Chen C, Dong X, et al. A 4-bp Insertion at ZmPLA1
825 Encoding a Putative Phospholipase A Generates Haploid Induction in Maize. *Mol Plant.*
826 2017;10: 520–522. doi:10.1016/j.molp.2017.01.011
- 827 20. Zhong Y, Liu C, Qi X, Jiao Y, Wang D, Wang Y, et al. Mutation of ZmDMP enhances
828 haploid induction in maize. *Nat Plants.* 2019;5: 575–580. doi:10.1038/s41477-019-0443-7
- 829 21. Kelliher T, Starr D, Su X, Tang G, Chen Z, Carter J, et al. One-step genome editing of elite
830 crop germplasm during haploid induction. *Nat Biotechnol.* 2019;37: 287–292.
831 doi:10.1038/s41587-019-0038-x
- 832 22. Honys D, Twell D. Comparative analysis of the Arabidopsis pollen transcriptome. *Plant*
833 *Physiol.* 2003;132: 640–652. doi:10.1104/pp.103.020925
- 834 23. Becker JD, Boavida LC, Carneiro J, Haury M, Feijó JA. Transcriptional profiling of
835 Arabidopsis tissues reveals the unique characteristics of the pollen transcriptome. *Plant*
836 *Physiol.* 2003;133: 713–725. doi:10.1104/pp.103.028241
- 837 24. Steffen JG, Kang I-H, Macfarlane J, Drews GN. Identification of genes expressed in the
838 Arabidopsis female gametophyte: Female gametophyte-expressed genes. *Plant J.* 2007;51:
839 281–292. doi:10.1111/j.1365-313X.2007.03137.x
- 840 25. Jones-Rhoades MW, Borevitz JO, Preuss D. Genome-wide expression profiling of the
841 Arabidopsis female gametophyte identifies families of small, secreted proteins. *PLoS*
842 *Genet.* 2007;3: 1848–1861. doi:10.1371/journal.pgen.0030171
- 843 26. Chetoor AM, Givan SA, Cole RA, Coker CT, Unger-Wallace E, Vejrupkova Z, et al.
844 Discovery of novel transcripts and gametophytic functions via RNA-seq analysis of maize
845 gametophytic transcriptomes. *Genome Biol.* 2014;15: 414. doi:10.1186/s13059-014-0414-2
- 846 27. Zhai J, Zhang H, Arikait S, Huang K, Nan G-L, Walbot V, et al. Spatiotemporally dynamic,
847 cell-type-dependent premeiotic and meiotic phasiRNAs in maize anthers. *Proc Natl Acad*
848 *Sci U S A.* 2015;112: 3146–3151. doi:10.1073/pnas.1418918112
- 849 28. Nelms B, Walbot V. Defining the developmental program leading to meiosis in maize.
850 *Science.* 2019;364: 52–56. doi:10.1126/science.aav6428
- 851 29. Begcy K, Nosenko T, Zhou L-Z, Fragner L, Weckwerth W, Dresselhaus T. Male Sterility in
852 Maize after Transient Heat Stress during the Tetrad Stage of Pollen Development. *Plant*
853 *Physiol.* 2019; doi:10.1104/pp.19.00707
- 854 30. Chen J, Strieder N, Krohn NG, Cyprys P, Sprunck S, Engelmann JC, et al. Zygotic Genome
855 Activation Occurs Shortly after Fertilization in Maize. *Plant Cell.* 2017;29: 2106–2125.
856 doi:10.1105/tpc.17.00099
- 857 31. Slotkin RK, Vaughn M, Borges F, Tanurdzić M, Becker JD, Feijó JA, et al. Epigenetic
858 reprogramming and small RNA silencing of transposable elements in pollen. *Cell.*
859 2009;136: 461–472. doi:10.1016/j.cell.2008.12.038
- 860 32. Schoft VK, Chumak N, Mosiolek M, Slusarz L, Komnenovic V, Brownfield L, et al. Induction
861 of RNA-directed DNA methylation upon decondensation of constitutive heterochromatin.
862 *EMBO Rep.* 2009;10: 1015–1021. doi:10.1038/embor.2009.152
- 863 33. Calarco JP, Borges F, Donoghue MTA, Van Ex F, Jullien PE, Lopes T, et al.
864 Reprogramming of DNA methylation in pollen guides epigenetic inheritance via small RNA.
865 *Cell.* 2012;151: 194–205. doi:10.1016/j.cell.2012.09.001

- 866 34. Dooner HK, Wang Q, Huang JT, Li Y, He L, Xiong W, et al. Spontaneous mutations in
867 maize pollen are frequent in some lines and arise mainly from retrotranspositions and
868 deletions. *Proc Natl Acad Sci U S A*. 2019; doi:10.1073/pnas.1903809116
- 869 35. He S, Vickers M, Zhang J, Feng X. Natural depletion of histone H1 in sex cells causes DNA
870 demethylation, heterochromatin decondensation and transposon activation. *Elife*. 2019;8.
871 doi:10.7554/eLife.42530
- 872 36. Martínez G, Slotkin RK. Developmental relaxation of transposable element silencing in
873 plants: functional or byproduct? *Curr Opin Plant Biol*. 2012;15: 496–502.
874 doi:10.1016/j.pbi.2012.09.001
- 875 37. Martínez G, Panda K, Köhler C, Slotkin RK. Silencing in sperm cells is directed by RNA
876 movement from the surrounding nurse cell. *Nat Plants*. 2016;2: 16030.
877 doi:10.1038/nplants.2016.30
- 878 38. Martinez G, Wolff P, Wang Z, Moreno-Romero J, Santos-González J, Conze LL, et al.
879 Paternal easiRNAs regulate parental genome dosage in Arabidopsis. *Nat Genet*. 2018;50:
880 193–198. doi:10.1038/s41588-017-0033-4
- 881 39. Jiao Y, Peluso P, Shi J, Liang T, Stitzer MC, Wang B, et al. Improved maize reference
882 genome with single-molecule technologies. *Nature*. 2017;546: 524–527.
883 doi:10.1038/nature22971
- 884 40. Lunardon A, Forestan C, Farinati S, Axtell MJ, Varotto S. Genome-Wide Characterization of
885 Maize Small RNA Loci and Their Regulation in the required to maintain repression6-1
886 (*rmr6-1*) Mutant and Long-Term Abiotic Stresses. *Plant Physiol*. 2016;170: 1535–1548.
887 doi:10.1104/pp.15.01205
- 888 41. Walley JW, Sartor RC, Shen Z, Schmitz RJ, Wu KJ, Urich MA, et al. Integration of omic
889 networks in a developmental atlas of maize. *Science*. 2016;353: 814–818.
890 doi:10.1126/science.aag1125
- 891 42. Panda K, Ji L, Neumann DA, Daron J, Schmitz RJ, Slotkin RK. Full-length autonomous
892 transposable elements are preferentially targeted by expression-dependent forms of RNA-
893 directed DNA methylation. *Genome Biol*. 2016;17: 170. doi:10.1186/s13059-016-1032-y
- 894 43. Wolfgruber TK, Sharma A, Schneider KL, Albert PS, Koo D-H, Shi J, et al. Maize
895 centromere structure and evolution: sequence analysis of centromeres 2 and 5 reveals
896 dynamic Loci shaped primarily by retrotransposons. *PLoS Genet*. 2009;5: e1000743.
897 doi:10.1371/journal.pgen.1000743
- 898 44. Li Y, Segal G, Wang Q, Dooner HK. Gene Tagging with Engineered Ds Elements in Maize.
899 In: Peterson T, editor. *Plant Transposable Elements: Methods and Protocols*. Totowa, NJ:
900 Humana Press; 2013. pp. 83–99. doi:10.1007/978-1-62703-568-2_6
- 901 45. Warman C, Fowler JE. Custom built scanner and simple image processing pipeline enables
902 low-cost, high-throughput phenotyping of maize ears [Internet]. *bioRxiv*. 2019. p. 780650.
903 doi:10.1101/780650
- 904 46. Engel ML, Chaboud A, Dumas C, McCormick S. Sperm cells of *Zea mays* have a complex
905 complement of mRNAs. *Plant J*. 2003;34: 697–707. doi:10.1046/j.1365-313X.2003.01761.x
- 906 47. Mori T, Igawa T, Tamiya G, Miyagishima S-Y, Berger F. Gamete attachment requires GEX2
907 for successful fertilization in Arabidopsis. *Curr Biol*. 2014;24: 170–175.
908 doi:10.1016/j.cub.2013.11.030
- 909 48. Wang G, Jiang H, Del Toro de León G, Martinez G, Köhler C. Sequestration of a

- 910 Transposon-Derived siRNA by a Target Mimic Imprinted Gene Induces Postzygotic
911 Reproductive Isolation in Arabidopsis. *Dev Cell*. 2018;46: 696–705.e4.
912 doi:10.1016/j.devcel.2018.07.014
- 913 49. Borges F, Parent J-S, van Ex F, Wolff P, Martínez G, Köhler C, et al. Transposon-derived
914 small RNAs triggered by miR845 mediate genome dosage response in Arabidopsis. *Nat*
915 *Genet*. 2018;50: 186–192. doi:10.1038/s41588-017-0032-5
- 916 50. Engel ML, Holmes-Davis R, McCormick S. Green sperm. Identification of male gamete
917 promoters in Arabidopsis. *Plant Physiol*. 2005;138: 2124–2133. doi:10.1104/pp.104.054213
- 918 51. Bai F, Daliberti M, Bagadion A, Xu M, Li Y, Baier J, et al. Parent-of-Origin-Effect rough
919 endosperm Mutants in Maize. *Genetics*. 2016;204: 221–231.
920 doi:10.1534/genetics.116.191775
- 921 52. Sarkar KR, Coe EH. A genetic analysis of the origin of maternal haploids in maize.
922 *Genetics*. 1966;54: 453–464. Available: <https://www.ncbi.nlm.nih.gov/pubmed/17248321>
- 923 53. Zhao X, Xu X, Xie H, Chen S, Jin W. Fertilization and uniparental chromosome elimination
924 during crosses with maize haploid inducers. *Plant Physiol*. 2013;163: 721–731.
925 doi:10.1104/pp.113.223982
- 926 54. Qiu F, Liang Y, Li Y, Liu Y, Wang L, Zheng Y. Morphological, cellular and molecular
927 evidences of chromosome random elimination in vivo upon haploid induction in maize.
928 *Current Plant Biology*. 2014;1: 83–90. doi:10.1016/j.cpb.2014.04.001
- 929 55. Li X, Meng D, Chen S, Luo H, Zhang Q, Jin W, et al. Single nucleus sequencing reveals
930 spermatid chromosome fragmentation as a possible cause of maize haploid induction. *Nat*
931 *Commun*. 2017;8: 991. doi:10.1038/s41467-017-00969-8
- 932 56. Tian X, Qin Y, Chen B, Liu C, Wang L, Li X, et al. Hetero-fertilization together with failed
933 egg-sperm cell fusion supports single fertilization involved in in vivo haploid induction in
934 maize. *J Exp Bot*. 2018;69: 4689–4701. doi:10.1093/jxb/ery177
- 935 57. Phillips AR, Evans MMS. Analysis of stunter1, a maize mutant with reduced gametophyte
936 size and maternal effects on seed development. *Genetics*. 2011;187: 1085–1097.
937 doi:10.1534/genetics.110.125286
- 938 58. Wilson-Sánchez D, Rubio-Díaz S, Muñoz-Viana R, Pérez-Pérez JM, Jover-Gil S, Ponce
939 MR, et al. Leaf phenomics: a systematic reverse genetic screen for Arabidopsis leaf
940 mutants. *Plant J*. 2014;79: 878–891. doi:10.1111/tpj.12595
- 941 59. Rutter MT, Murren CJ, Callahan HS, Bisner AM, Leebens-Mack J, Wolyniak MJ, et al.
942 Distributed phenomics with the unPAK project reveals the effects of mutations. *Plant J*.
943 2019; doi:10.1111/tpj.14427
- 944 60. Giaever G, Chu AM, Ni L, Connelly C, Riles L, Véronneau S, et al. Functional profiling of
945 the *Saccharomyces cerevisiae* genome. *Nature*. 2002;418: 387–391.
946 doi:10.1038/nature00935
- 947 61. Berry DB, Guan Q, Hose J, Haroon S, Gebbia M, Heisler LE, et al. Multiple means to the
948 same end: the genetic basis of acquired stress resistance in yeast. *PLoS Genet*. 2011;7:
949 e1002353. doi:10.1371/journal.pgen.1002353
- 950 62. Price MN, Deutschbauer AM, Skerker JM, Wetmore KM, Ruths T, Mar JS, et al. Indirect
951 and suboptimal control of gene expression is widespread in bacteria. *Mol Syst Biol*. 2013;9:
952 660. doi:10.1038/msb.2013.16
- 953 63. Helmann TC, Deutschbauer AM, Lindow SE. Genome-wide identification of *Pseudomonas*

- 954 syringae genes required for fitness during colonization of the leaf surface and apoplast.
955 Proc Natl Acad Sci U S A. 2019; doi:10.1073/pnas.1908858116
- 956 64. Schnable JC. Genes and gene models, an important distinction. *New Phytol.* 2019;
957 doi:10.1111/nph.16011
- 958 65. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
959 universal RNA-seq aligner. *Bioinformatics.* 2013;29: 15–21.
960 doi:10.1093/bioinformatics/bts635
- 961 66. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al. Transcript
962 assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform
963 switching during cell differentiation. *Nat Biotechnol.* 2010;28: 511–515.
964 doi:10.1038/nbt.1621
- 965 67. Tian T, Liu Y, Yan H, You Q, Yi X, Du Z, et al. agriGO v2.0: a GO analysis toolkit for the
966 agricultural community, 2017 update. *Nucleic Acids Res.* 2017;45: W122–W129.
967 doi:10.1093/nar/gkx382
- 968 68. Wimalanathan K, Friedberg I, Andorf CM, Lawrence-Dill CJ. Maize GO Annotation-
969 Methods, Evaluation, and Review (maize-GAMER). *Plant Direct.* 2018;2: e00052.
970 doi:10.1002/pld3.52
- 971 69. Bushnell B. BBTools software package. URL <http://sourceforge.net/projects/bbmap>. 2014;
- 972 70. Liao Y, Smyth GK, Shi W. featureCounts: an efficient general purpose program for
973 assigning sequence reads to genomic features. *Bioinformatics.* 2014;30: 923–930.
974 doi:10.1093/bioinformatics/btt656
- 975 71. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-
976 seq data with DESeq2. *Genome Biol.* 2014;15: 550. doi:10.1186/s13059-014-0550-8
- 977 72. Gao H, Smith J, Yang M, Jones S, Djukanovic V, Nicholson MG, et al. Heritable targeted
978 mutagenesis in maize using a designed endonuclease. *Plant J.* 2010;61: 176–187.
979 doi:10.1111/j.1365-313X.2009.04041.x
- 980 73. Vejlupekova Z, Fowler JE. Maize DNA preps for undergraduate students: a robust method
981 for PCR genotyping. *Maize Genetics Cooperation Newsletter.* 2003;77: 24–25. Available:
982 <https://mnl.maizegdb.org/mnl/77/57vejlupekova.html>
- 983 74. Schindelin J, Arganda-Carreras I, Frise E, Kaynig V, Longair M, Pietzsch T, et al. Fiji: an
984 open-source platform for biological-image analysis. *Nat Methods.* 2012;9: 676–682.
985 doi:10.1038/nmeth.2019
- 986 75. Portwood JL 2nd, Woodhouse MR, Cannon EK, Gardiner JM, Harper LC, Schaeffer ML, et
987 al. MaizeGDB 2018: the maize multi-genome genetics and genomics database. *Nucleic*
988 *Acids Res.* 2019;47: D1146–D1154. doi:10.1093/nar/gky1046
- 989 76. Krishnakumar V, Hanlon MR, Contrino S, Ferlanti ES, Karamycheva S, Kim M, et al.
990 Araport: the Arabidopsis information portal. *Nucleic Acids Res.* 2015;43: D1003–9.
991 doi:10.1093/nar/gku1200
- 992 77. Cheng C-Y, Krishnakumar V, Chan AP, Thibaud-Nissen F, Schobel S, Town CD.
993 Araport11: a complete reannotation of the Arabidopsis thaliana reference genome. *Plant J.*
994 2017;89: 789–804. doi:10.1111/tpj.13415
- 995 78. Li W, Cowley A, Uludag M, Gur T, McWilliam H, Squizzato S, et al. The EMBL-EBI
996 bioinformatics web and programmatic tools framework. *Nucleic Acids Res.* 2015;43: W580–
997 4. doi:10.1093/nar/gkv279

- 998 79. Mitchell AL, Attwood TK, Babbitt PC, Blum M, Bork P, Bridge A, et al. InterPro in 2019:
999 improving coverage, classification and access to protein sequence annotations. *Nucleic*
1000 *Acids Res.* 2019;47: D351–D360. doi:10.1093/nar/gky1100
- 1001 80. Krogh A, Larsson B, von Heijne G, Sonnhammer EL. Predicting transmembrane protein
1002 topology with a hidden Markov model: application to complete genomes. *J Mol Biol.*
1003 2001;305: 567–580. doi:10.1006/jmbi.2000.4315
- 1004 81. Van Bel M, Diels T, Vancaester E, Kreft L, Botzki A, Van de Peer Y, et al. PLAZA 4.0: an
1005 integrative resource for functional, evolutionary and comparative plant genomics. *Nucleic*
1006 *Acids Res.* 2018;46: D1190–D1196. doi:10.1093/nar/gkx1002
- 1007 82. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST
1008 and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids*
1009 *Res.* 1997;25: 3389–3402. doi:10.1093/nar/25.17.3389
- 1010 83. Vollbrecht E, Hake S. Deficiency analysis of female gametogenesis in maize. *Dev Genet.*
1011 1995;16: 44–63. doi:10.1002/dvg.1020160109
- 1012 84. Running MP, Clark SE, Meyerowitz EM. Chapter 15 Confocal Microscopy of the Shoot
1013 Apex. In: Galbraith DW, Bohnert HJ, Bourque DP, editors. *Methods in Cell Biology.*
1014 Academic Press; 1995. pp. 217–229. doi:10.1016/S0091-679X(08)61456-9
- 1015 85. Heuer S, Lörz H, Dresselhaus T. The MADS box gene *ZmMADS2* is specifically expressed
1016 in maize pollen and during maize pollen tube growth. *Sex Plant Reprod.* 2000;13: 21–27.
1017 doi:10.1007/PL00009838
- 1018 86. Townsley BT, Covington MF, Ichihashi Y, Zumstein K, Sinha NR. BrAD-seq: Breath
1019 Adapter Directional sequencing: a streamlined, ultra-simple and fast library preparation
1020 protocol for strand specific mRNA library construction. *Front Plant Sci.* 2015;6: 366.
1021 doi:10.3389/fpls.2015.00366
- 1022 87. McClintock B. Chromosome organization and genic expression. *Cold Spring Harb Symp*
1023 *Quant Biol.* 1951;16: 13–47. doi:10.1101/SQB.1951.016.01.004
- 1024 88. Levy AA, Walbot V. Regulation of the timing of transposable element excision during maize
1025 development. *Science.* 1990;248: 1534–1537. doi:10.1126/science.2163107
- 1026 89. Box MS, Coustham V, Dean C, Mylne JS. Protocol: A simple phenol-based method for 96-
1027 well extraction of high quality RNA from *Arabidopsis*. *Plant Methods.* 2011;7: 7.
1028 doi:10.1186/1746-4811-7-7
- 1029 90. Vollbrecht E, Duvick J, Schares JP, Ahern KR, Deewatthanawong P, Xu L, et al. Genome-
1030 wide distribution of transposed Dissociation elements in maize. *Plant Cell.* 2010;22: 1667–
1031 1685. doi:10.1105/tpc.109.073452
- 1032

1033 Figure legends

1034 Fig 1. Experimental design and transcriptome replicate assessment.

1035 (A) mRNA was isolated from four developmental stages of maize male reproductive development,
1036 with four biological replicates for each: pre-meiotic tassel primordia (TP), post-meiotic, unicellular
1037 microspores (MS), mature pollen (MP), and sperm cells (SC). A single biological replicate of
1038 mRNA from the bicellular stage of pollen development was also isolated and sequenced (MS-B,
1039 not shown). Nuclei were stained with either DAPI or Dyecycle green. (B) Principal component
1040 analysis of genic transcriptomic data generated by this study, showing the 2 major components
1041 (explaining 49.8% of the variance) on x- and y-axis. The four biological replicates of each of the

1042 four sequenced tissues clustered with other replicates from the same tissue. TP and MS were
1043 clearly separated in principal component space, whereas SC and MP samples displayed less
1044 separation from each other.

1045

1046 **Fig 2. Characterization of developmentally dynamic transcription from transposable**
1047 **elements (TEs).**

1048 **(A)** Distribution of different categories of TEs based on their expression. Number of TEs are in
1049 parentheses. **(B)** Length of TEs in the different TE categories from part A. The box represents
1050 lower and upper quartile, the line is the median, and the whiskers represent 10-90% range. Red
1051 asterisk denotes the mean. **(C)** Observed / expected Log₂ ratios of TE family proportions in the
1052 different TE categories from part A. Grey indicates no data. **(D)** Distance from the annotated
1053 centromere for different TE categories from panel A.

1054

1055 **Fig 3. High TE expression in the maize male gametophyte lineage.**

1056 **(A)** Number of differentially expressed TEs in seven tissues compared to seedlings. The inset
1057 volcano plot shows for mature pollen how differentially expressed TEs were identified. Green and
1058 red numbers within the volcano plot indicate how many TEs were statistically up- or down-
1059 regulated, respectively. **(B)** Number of up-regulated TEs in mature pollen compared to isolated
1060 sperm cells from this study and a previously published distinct isolation and sequencing of sperm
1061 cell mRNA. **(C)** Starting with TEs differentially up-regulated in unicellular microspores (boxed, far
1062 left volcano plot), we determined how many of these same TEs are expressed at other
1063 developmental time points. **(D)** Raw distribution of expressed TE family annotations. 'Male' refers
1064 to the set of TEs expressed in any male lineage dataset (MS, MS-B, MP, SC). 'Male specific' are
1065 TEs expressed in only the male lineage (not other tissues / timepoints). 'Male extensive' TEs are
1066 expressed in all of the male lineage datasets, and 'male extensive + specific' refers to TEs
1067 expressed all male lineage datasets and not other tissues / timepoints. 'High-confidence sperm'
1068 refers to TEs identified as expressed in both analyzed sperm cell datasets from part B. **(E)**
1069 Expressed TE family annotations normalized to the genome-wide TE distribution of TEs >2 kb
1070 from genes. Categories are the same as part D.

1071

1072 **Fig 4. Co-regulation of TE and gene expression in the male gametophyte.**

1073 **(A)** For each tissue type, the top 20,000 most highly expressed genes are distributed along the
1074 X-axis in bins of 200, with the most highly expressed bin on the far left. For each bin the number
1075 of up- and down-regulated TEs near (<2kb) its genes is then summed on the Y-axis. In unicellular
1076 microspores (top right) there is little correlation, while in mature pollen (bottom left) the most highly
1077 expressed genes are near primarily up-regulated TEs. The bin location of *gex2* (see Fig. 7) is
1078 annotated in the sperm cell data (bottom right). **(B)** High expression levels are associated with
1079 developmental specificity: approximately 2/3 of the genes associated with the highest FPKM
1080 values in each of the four sample types are highly expressed in only that sample type.

1081

1082 **Fig 5. Large-scale tracking of seed marker transmission frequencies was accomplished by**
1083 **generating ear projections with a custom built rotational scanner.**

1084 **(A)** When crossed either through the male or the female, *Ds-GFP* mutant allele *tdsgR107C12* (in
1085 gene Zm00001d012382), marked by green fluorescent seeds, shows 1:1 Mendelian inheritance

1086 (50% transmission of the GFP seed marker). Images captured in blue light with an orange filter.
1087 **(B)** Mutant alleles in other genes, such as *tdsgR102H01* (Zm00001d037695), showed non-
1088 Mendelian segregation when crossed through the male (37.5% GFP transmission). Segregation
1089 through the female remained Mendelian, indicating a male-specific transmission defect. **(C)** For
1090 some mutant alleles (~10% of lines in this study), the anthocyanin transgene *C1* was tightly linked
1091 to the insertion mutant. In these cases, seeds carrying a mutant allele of a gene of interest could
1092 be tracked by their purple color. Here, insertion *tdsgR96C12* (Zm00001d015901) shows a strong
1093 male-specific transmission defect (24.8% *C1* transmission through the male). Images captured in
1094 full spectrum visible light.

1095

1096 **Fig 6. Highly expressed pollen genes are more likely to be associated with decreased male**
1097 **transmission rates.**

1098 Alleles with CDS insertions were tested for differences from Mendelian inheritance using a quasi-
1099 likelihood test, with p-values corrected for multiple testing using the Benjamini-Hochberg
1100 procedure; alleles with quasi-likelihood adjusted p-value < 0.05 are represented in pink. Alleles
1101 are plotted by the $\log_2(\text{FPKM})$ of their respective gene according to that gene's expression class
1102 (Seedling, Vegetative Cell, or Sperm Cell). Insertion alleles distributed among the classes as
1103 follows: Vegetative Cell, 35 alleles; Sperm Cell, 11 alleles; Seedling, 10 alleles. Genes
1104 represented by two independent insertion alleles are enclosed by dotted lines. **(A)** Transmission
1105 rates of 56 mutant allele seed markers for heterozygous *Ds-GFP* mutant plants crossed through
1106 the female. **(B)** Transmission rates for alleles in the negative control Seedling class when crossed
1107 through the male. **(C)** For genes belonging to the Sperm Cell class, one out of ten (10%) was
1108 associated with significant non-Mendelian inheritance. The single gene with a male transmission
1109 defect in this group (*gex2*) showed a strong defect for both of the independent alleles tested. **(D)**
1110 An increased proportion of the genes in the Vegetative Cell class were associated with significant
1111 non-Mendelian inheritance when mutant (7/32 genes, 21.9%). In this class, an increase in
1112 $\log_2(\text{FPKM})$ was significantly associated with a decrease in marker transmission (linear
1113 regression, $p = 0.0120$).

1114

1115 **Fig 7. Mutations in the sperm cell-specific *gex2* gene cause aberrant seed development.**

1116 **(A)** The exon/intron structure of *gex2* (Zm00001d005781/GRMZM2G036832), showing the
1117 locations of the two independent *Ds-GFP* insertion mutants. **(B)** Ear projections of *gex2* mutant
1118 outcrosses. Top: heterozygote outcrossed as female, showing 1:1 transmission of the GFP-
1119 tagged allele. Middle: heterozygote outcrossed as a male, with 26.1% transmission of the mutant
1120 allele. Additionally, small seeds and occasional, small gaps between seeds are visible. Bottom:
1121 homozygous mutant outcrossed as a male, with many small seeds and large gaps, despite heavy
1122 pollination. **(C)** Genomic neighborhood of the GEX2 locus, with two nearby TEs, and their RNA-
1123 seq expression levels across male reproductive development. **(D)** Predicted domain structure of
1124 GEX2. **(E)** Quantification of small/aborted seeds resulting from pollination by *gex2* mutant plants
1125 and controls. Controls included two *Ds-GFP* lines that did not show transmission defects
1126 (*tdsgR12H07* and *tdsgR46C04*) and one *Ds-GFP* line that showed a strong transmission defect
1127 in the vegetative cell group (*tdsgR96C12*). A higher percentage of small/aborted seeds was
1128 present following pollination by heterozygous *gex2* plants representing the two mutant alleles

1129 (*tdsgR82A03* and *tdsgR84A12*), and pollination by homozygous *gex2-tdsgR84A12* plants further
1130 increased the percentage of small/aborted seeds.

1131

1132 **Fig 8. Pollination by *gex2-tdsgR84A12* leads to aberrant fertilization events and developing**
1133 **seed phenotypes.**

1134 **(A)** Seed development in a typical ovule pollinated by wild-type pollen, with one synergid
1135 penetrated by a pollen tube, and both embryo and endosperm development initiated. **(B-D)**
1136 Abnormal phenotypes seen following *gex2* pollination. **(B)** Ovule with developing (cellularizing)
1137 endosperm but unfertilized egg cell. **(C)** Ovule with developing embryo but unfertilized central cell.
1138 **(D)** Ovule with both synergids penetrated by a pollen tube, and a developing embryo and
1139 unfertilized central cell. emb=embryo; endo=endosperm; ps=synergid penetrated by a pollen
1140 tube; PN=polar nuclei.

1141

1142

1143 **Supporting information**

1144 **S1 Fig. Principal component analysis of gene and transposable element (TE) expression**
1145 **levels.**

1146 Two major components, on x- and y-axis, explain 89% of the variance in gene and TE expression
1147 levels. Asterisk (*) mark indicates the sample generated as part of this study, whereas other
1148 datasets are publicly available. For the sperm cells isolated in this study (SC), the TE expression
1149 of one biological replicate did not cluster with the other three (SC1), and therefore was removed
1150 from subsequent analyses of expression from TEs. MP-2014, SE, and OV are from [26]; MP-WEB
1151 is from [41]; LF is from [40]; MP-LM is from NCBI BioProject 306885 (2015); SC-TD is from [30].

1152

1153 **S2 Fig. Seedling tissue is the appropriate reference for comparison of TE activity.**

1154 **(A)** Steady-state mRNA accumulation of all TEs in different tissues. Datasets generated in this
1155 study are marked with an asterisk. **(B)** The number of TEs with zero or near-zero expression
1156 levels in different tissues. Seedlings (SE) have the most TEs with low expression levels.

1157

1158 **S3 Fig. Abundance of TEs near genes in each tissue.**

1159 **(A)** For each tissue type, the top 20,000 expressed genes are distributed along the X-axis in bins
1160 of 200, with the highest expressed bin on the far left. The number of TEs near (<2kb) these genes
1161 is then counted on the Y-axis. **(B)** Genes filtered for either higher expression in pollen (MP) over
1162 sperm cells (SC) (left) or SC>MP (right) were used to determine if the association in Figure 4 is
1163 due to sample contamination between SC and MP. Once genes were filtered, the top expressed
1164 genes in that tissue were distributed along the X-axis in bins of 200 based on their expression
1165 values, with the highest expressed bin on the far left. The number of up- and down-regulated TEs
1166 near (<2kb) these genes is then counted on the Y-axis.

1167

1168 **S4 Fig. *gex2* mutant pollen is associated with increased small and aborted seeds in**
1169 **outcross progeny.**

1170 **(A)** Seeds were removed from ears, arranged according to size, and counted. Images of
1171 representative seed populations are shown, with the top two rows in each image showing

1172 representative fully developed seeds. Rows below the top two contain all of the smaller or aborted
1173 seed from that particular ear. **(B)** PCR genotyping of small endosperm seeds from two
1174 independent crosses for the two *gex2* alleles show the majority of small seeds harbor the
1175 *gex2::Ds-GFP* allele, despite overall reduced transmission of the insertion alleles through the
1176 male. Small seeds from control *tsdgR46C04* crosses segregate in a Mendelian fashion.

1177

1178 **S5 Fig. Characterization of *gex2* seedless ear area.**

1179 Seedless area was quantified from scanned ear images for *gex2 Ds-GFP* alleles and *Ds-GFP*
1180 controls. Pollen from heterozygous *gex2* plants did not show significantly increased seedless area
1181 (*gex2-tdsgR82A03* pairwise t-test p-values relative to GFP line 1, GFP line 2, and VC mutant
1182 0.95, 0.96, and 0.74, respectively; *gex2-tdsgR84A12* pairwise t-test p-values 0.19, 0.13, and 0.06,
1183 respectively), whereas pollen from homozygous *gex2-tdsgR84A12* plants had significantly
1184 increased seedless area (pairwise t-test, p-value < 0.0001).

1185

1186 **S1 Methods. Tissue sample preparation, RNA extraction, and analysis of potential** 1187 **confounding variables in insertional mutagenesis lines.**

1188

1189 **S1 Table. Gene sequencing statistics and availability.**

1190 Summary statistics for sequencing data generated in the study.

1191

1192 **S2 Table. GO term enrichment results.**

1193 Differentially expressed genes in developmental categories examined in the study, as well as
1194 significantly enriched GO terms associated with these genes.

1195

1196 **S3 Table. Transposable element sequencing statistics and availability.**

1197 Summary statistics and availability for expression datasets used in the analysis of transposable
1198 element expression.

1199

1200 **S4 Table. Genic isoform abundance (FPKM) across developmental stages.**

1201 Cufflinks output describing isoform expression by developmental stages, separated by biological
1202 replicate.

1203

1204 **S5 Table. Top 20% transcripts by FPKM in Mature Pollen, Sperm Cell and Seedling** 1205 **datasets.**

1206 List of top 20% highly expressed genes assigned to the Vegetative Cell, Sperm Cell or Seedling
1207 Only classes.

1208

1209 **S6 Table. Insertional mutagenesis alleles and primers.**

1210 List of alleles tested for the presence of *Ds-GFP* insertions by PCR, including primers
1211 sequences.

1212

1213 **S7 Table. Insertional mutagenesis results by line.**

1214 Insertional mutagenesis results, separated by line, including marker transmission rates and
1215 expression category information.

1216

1217 **S8 Table. Insertional mutagenesis results by allele.**

1218 Insertional mutagenesis results, separated by allele, including marker transmission rates and
1219 expression category information.

1220

1221 **S9 Table. Concordance of seed phenotype with DsGFP genotype.**

1222 PCR results from testing *Ds-GFP* presence GFP and non-GFP seeds for selected alleles.

1223

1224 **S10 Table. *gex2* small seed phenotyping.**

1225 *gex2* small seed counting and seedless area results.

1226

1227 **S11 Table. Differential expression results.**

1228 Cuffdiff output comparing expression between tissues examined in this study.

1229

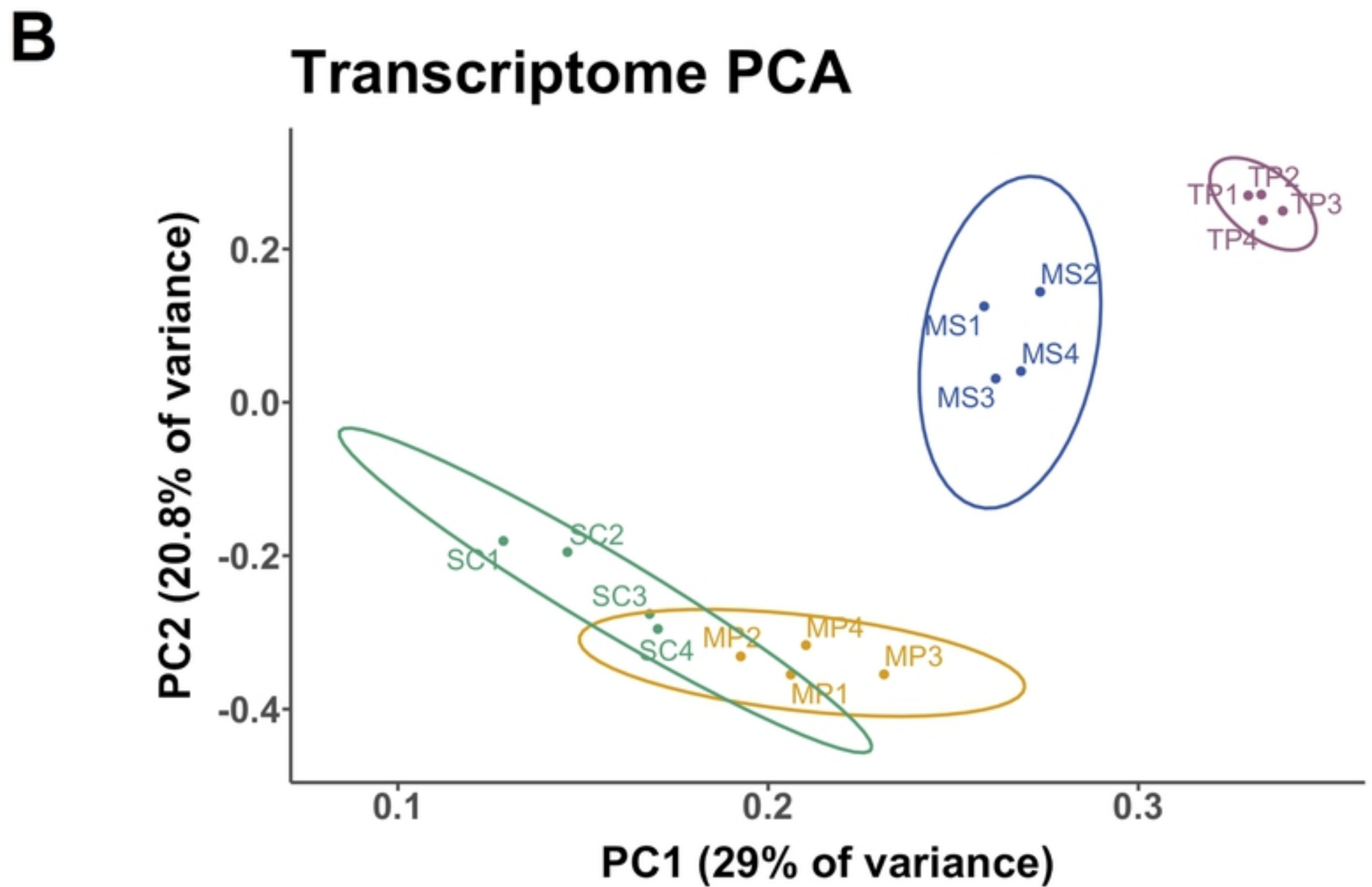
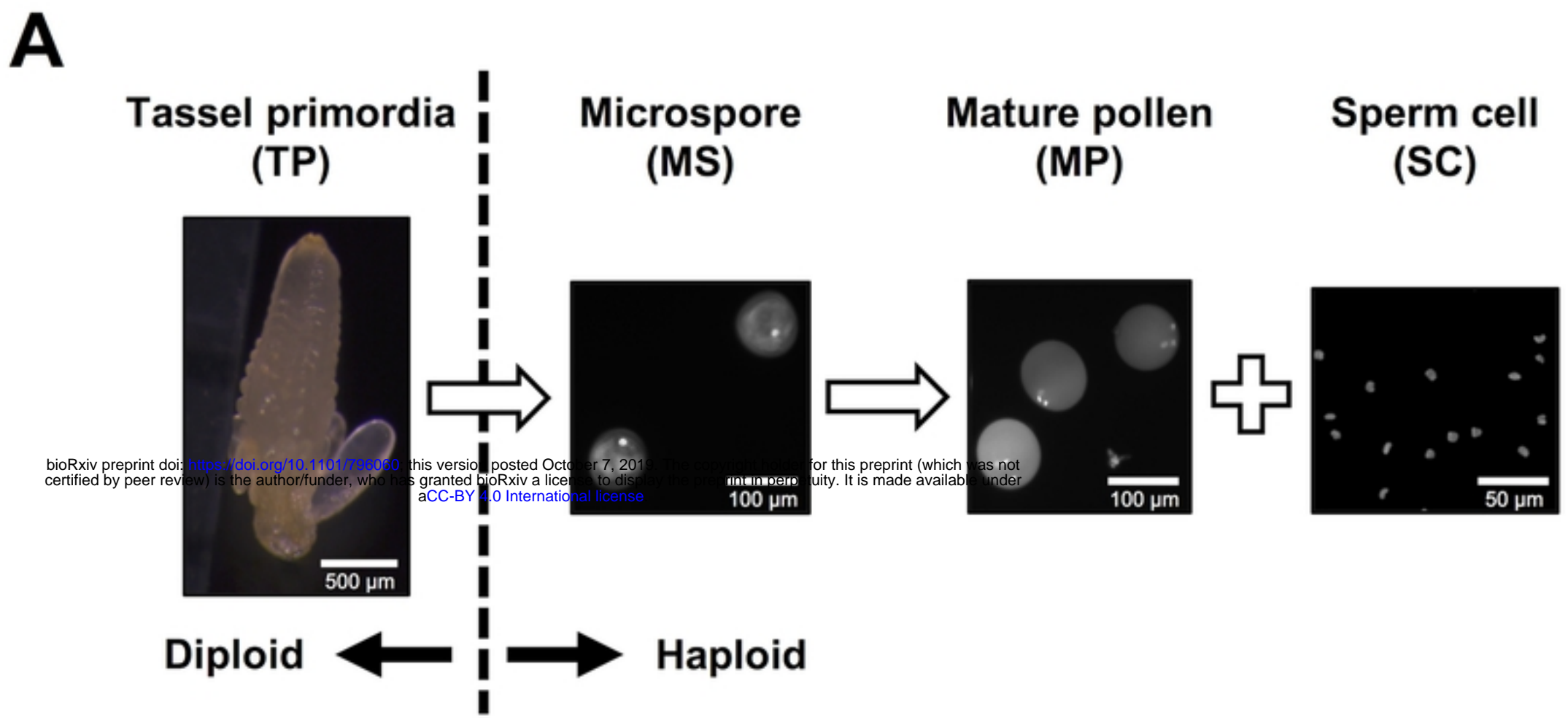
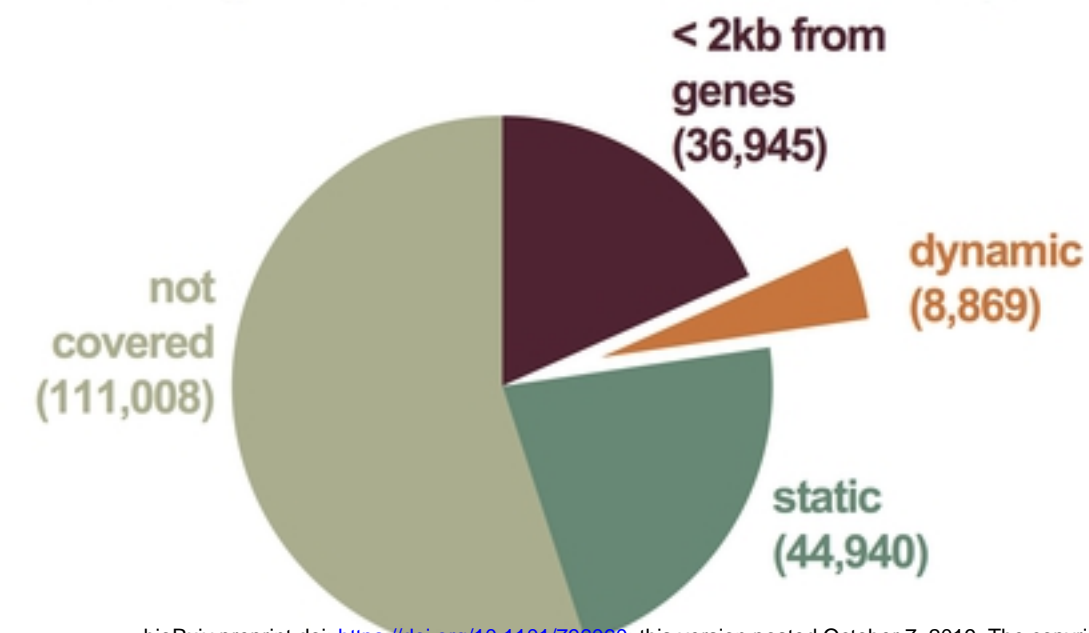
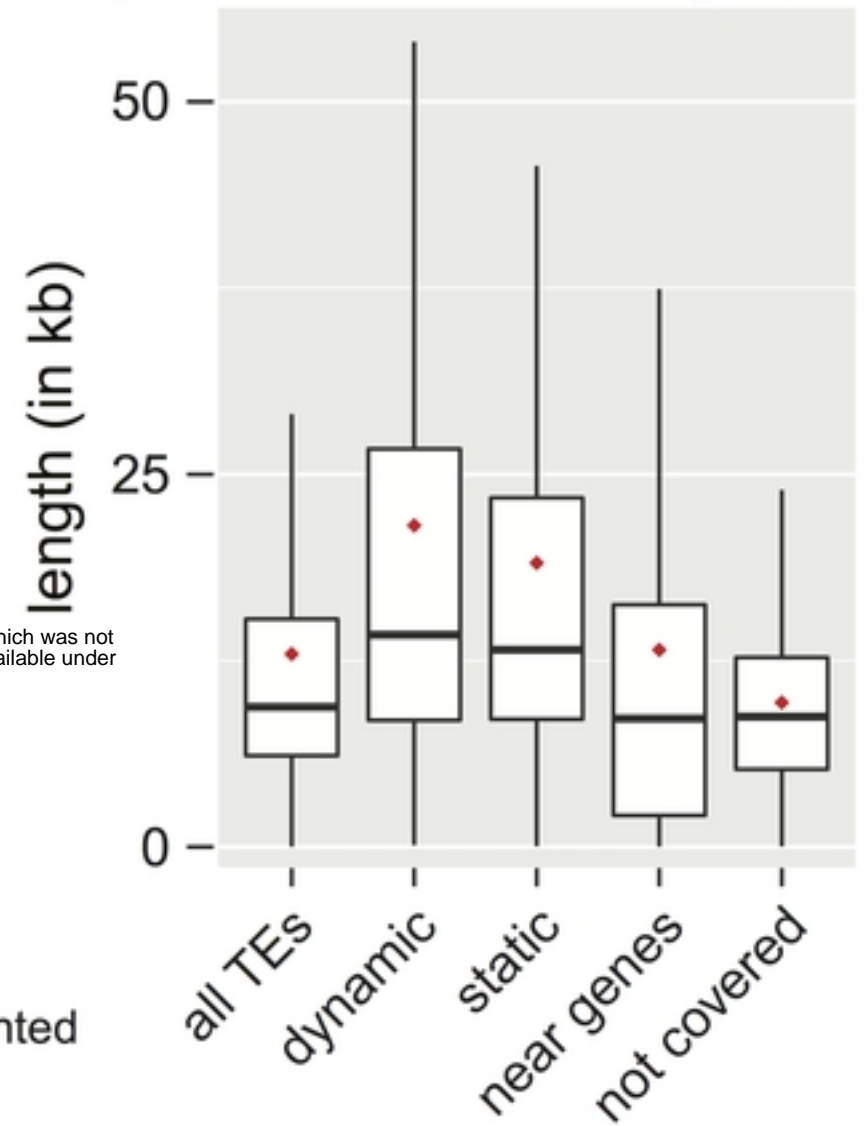


Figure 1

A Categorization of TEs based on expression

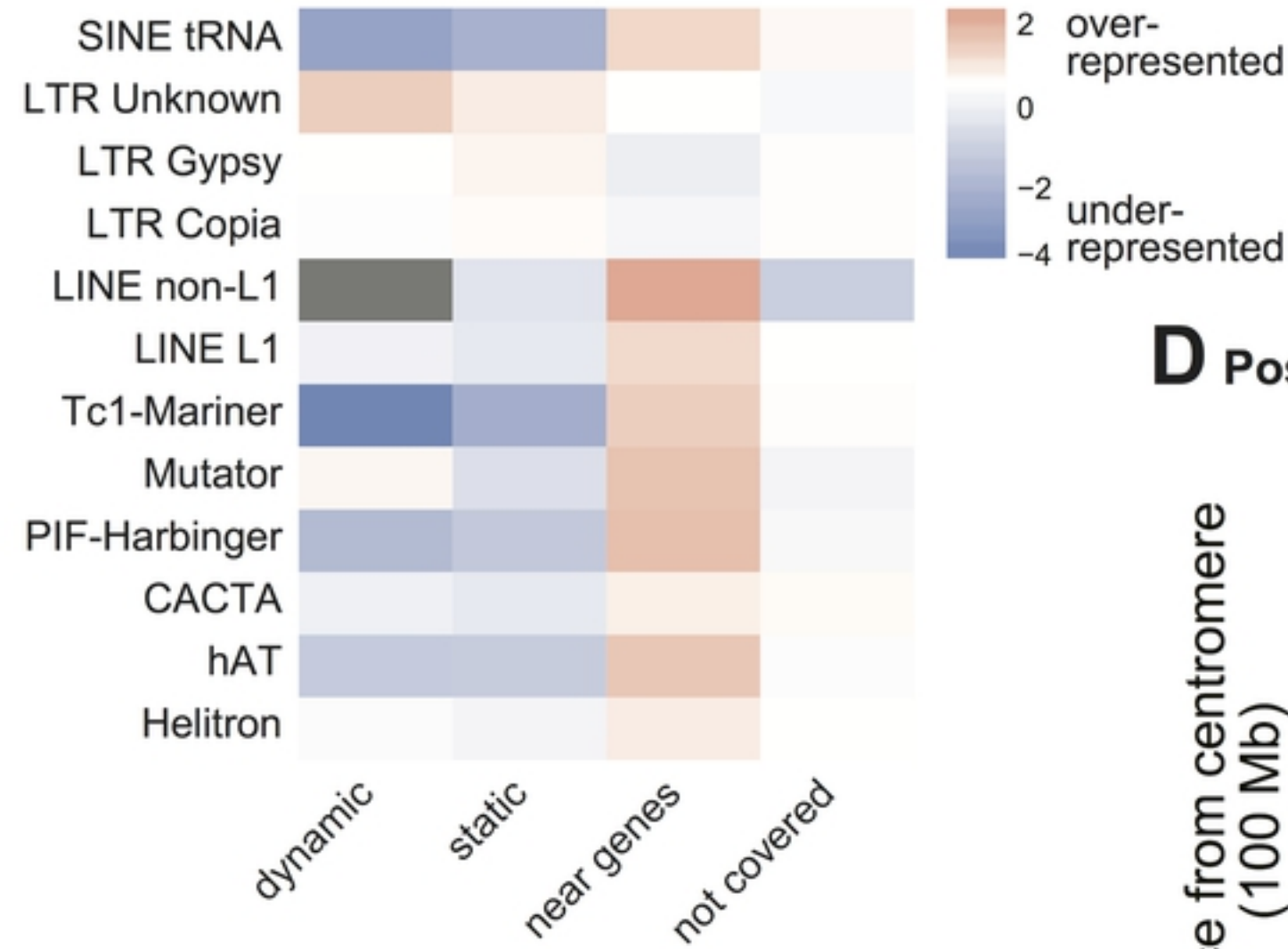


B Length distribution of categorized TEs



bioRxiv preprint doi: <https://doi.org/10.1101/796060>; this version posted October 7, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

C Family distribution of categorized TEs



D Position of categorized TEs

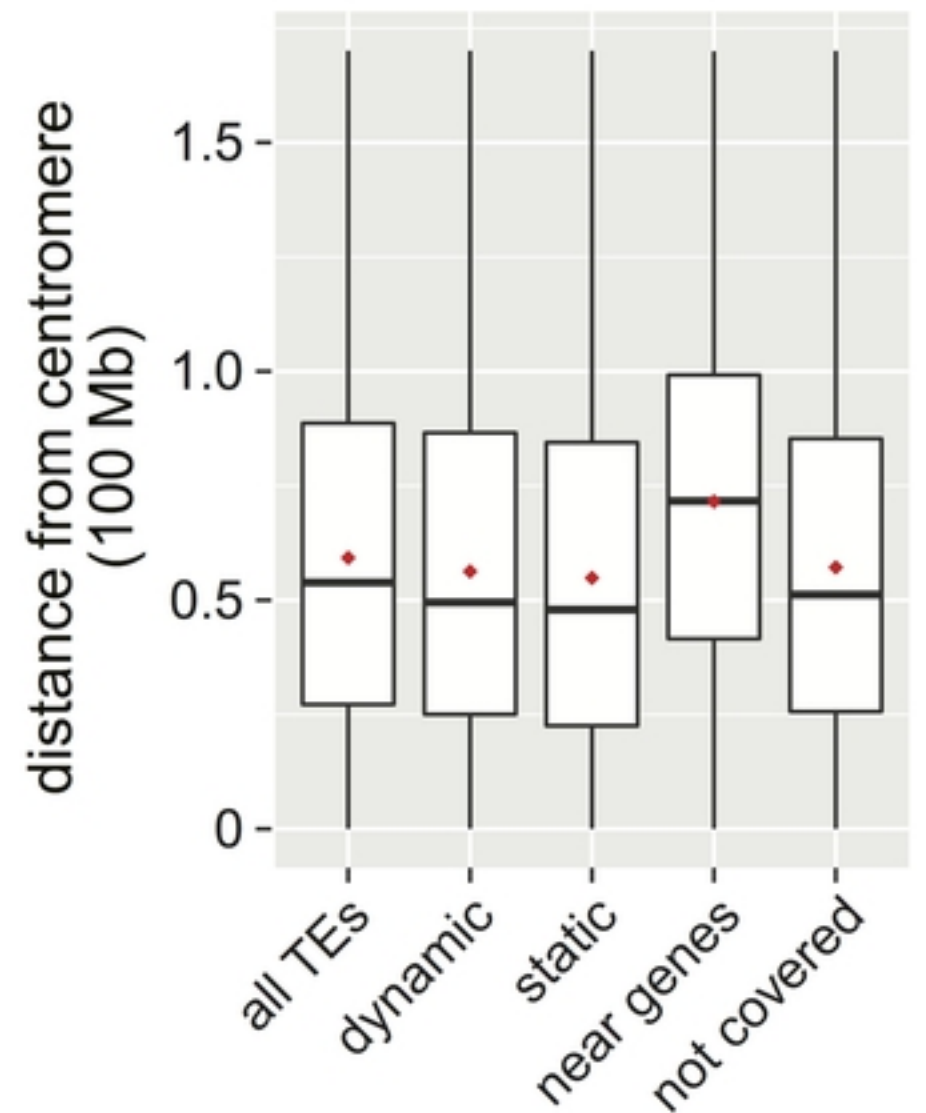


Figure2

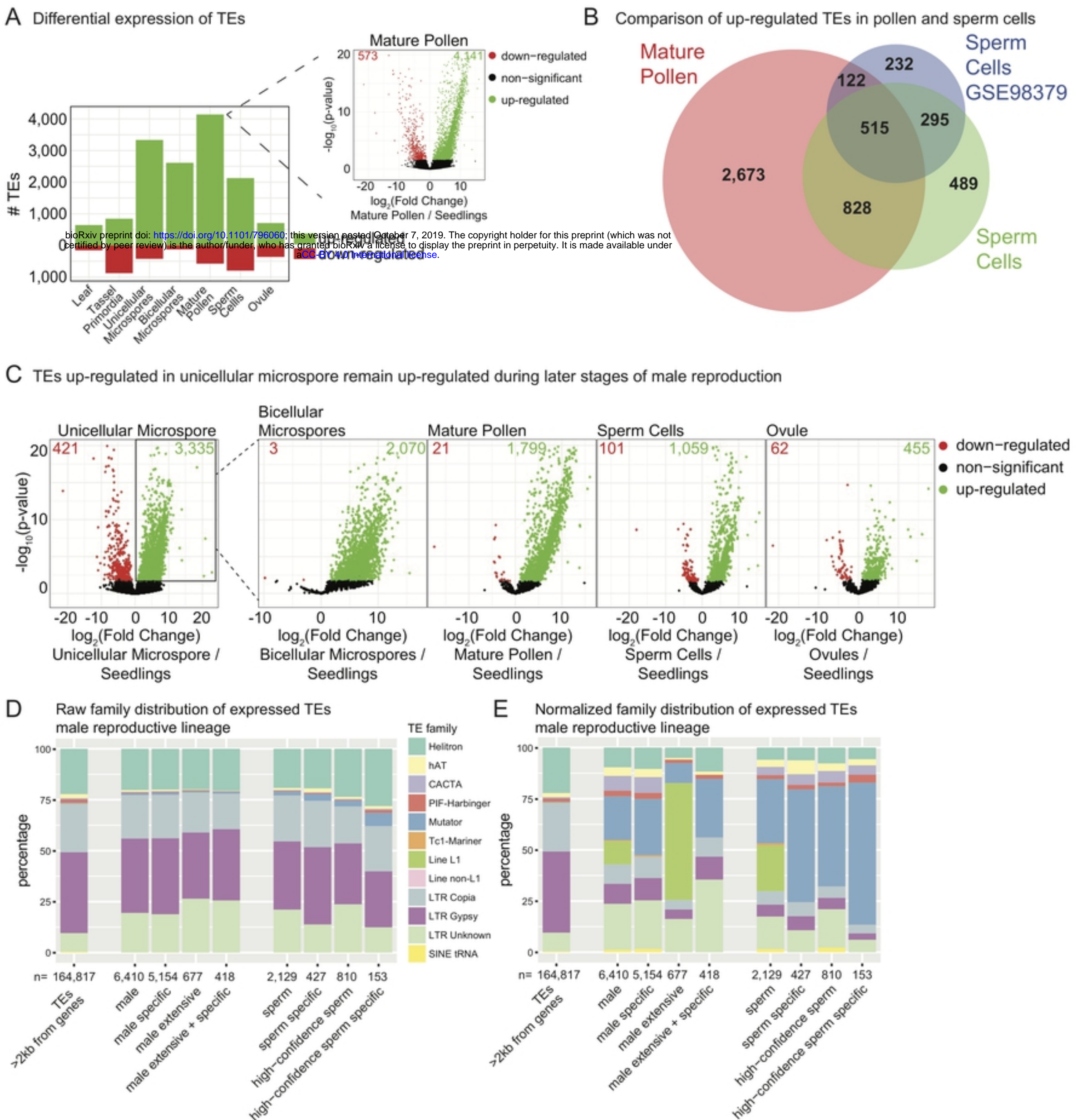


Figure3

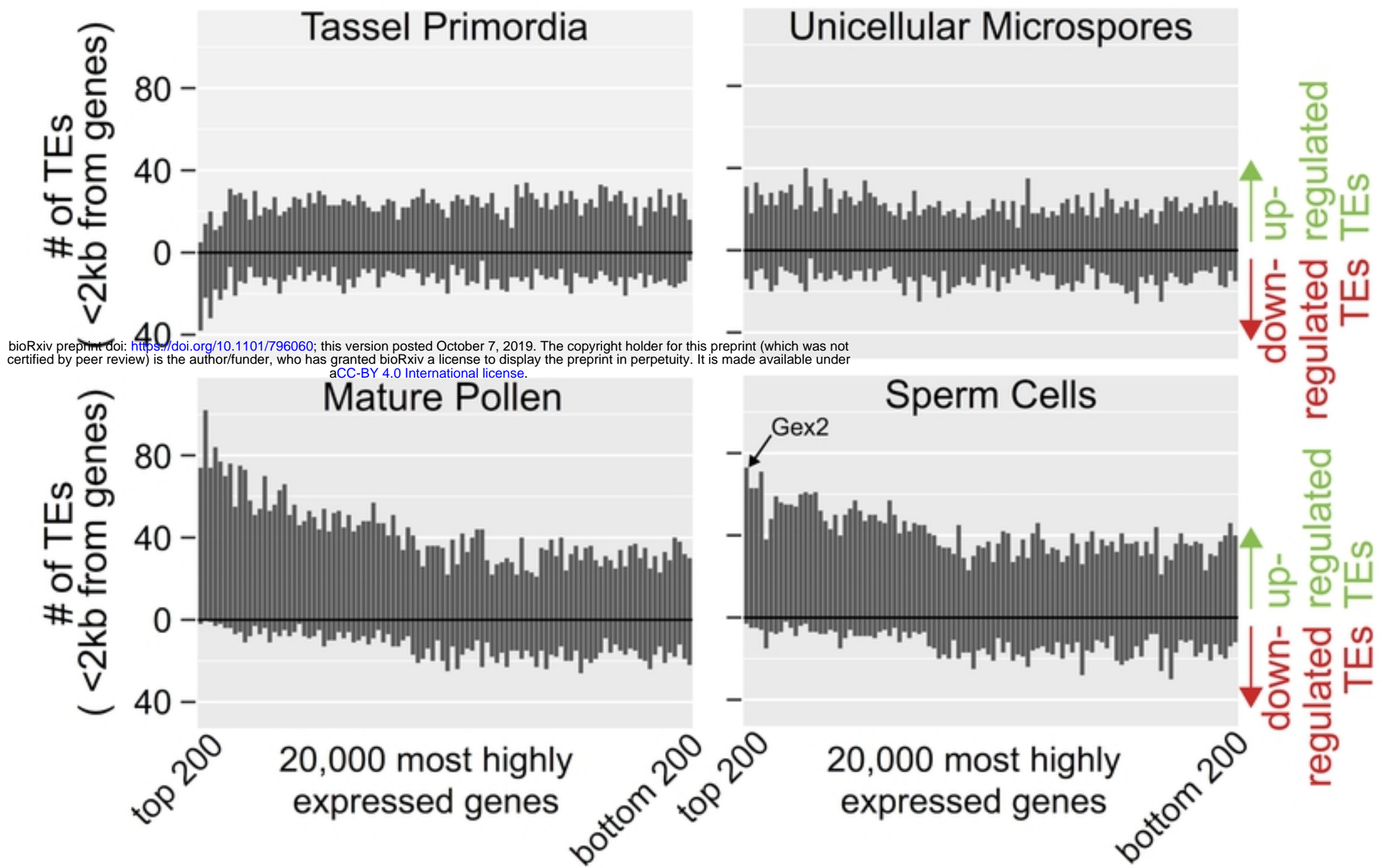
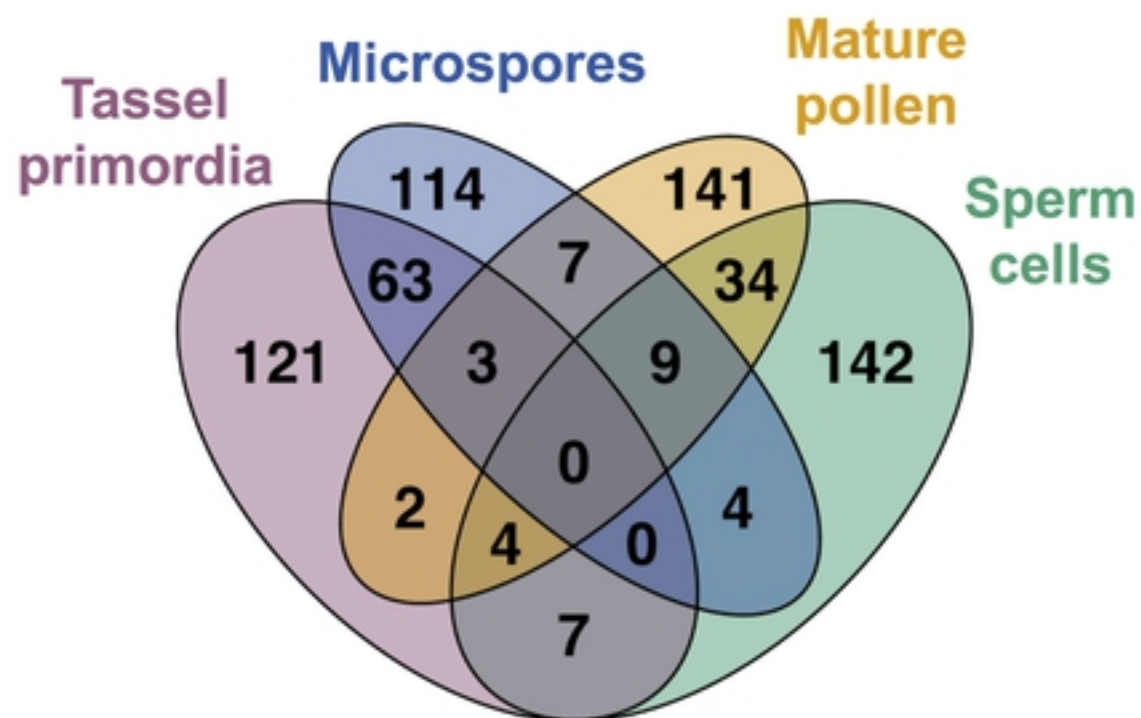
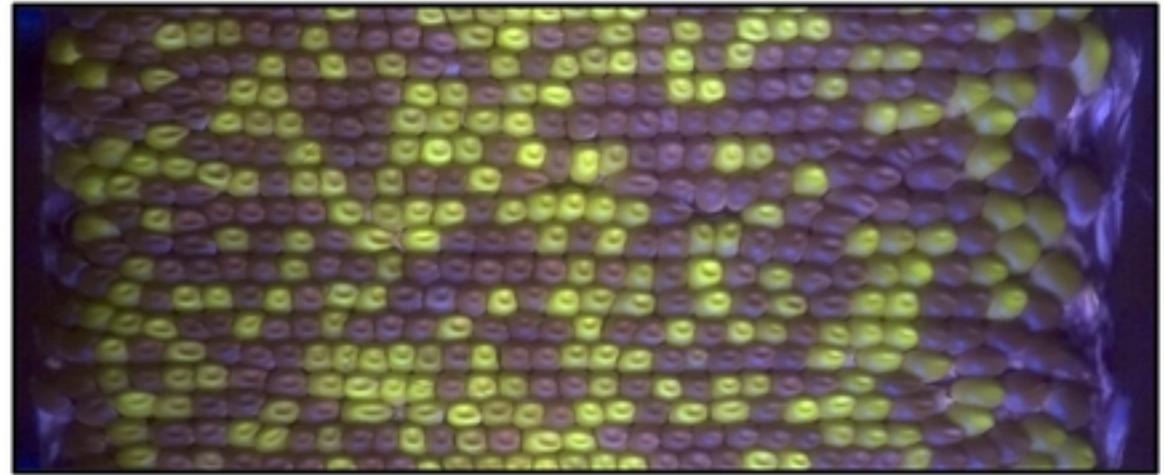
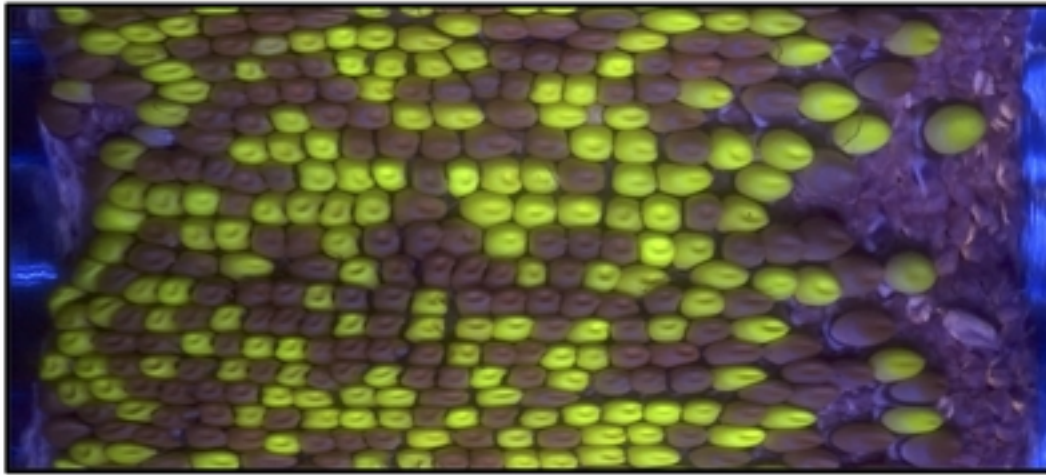
A**B****Top 200 FPKM genes per tissue**

Figure 4

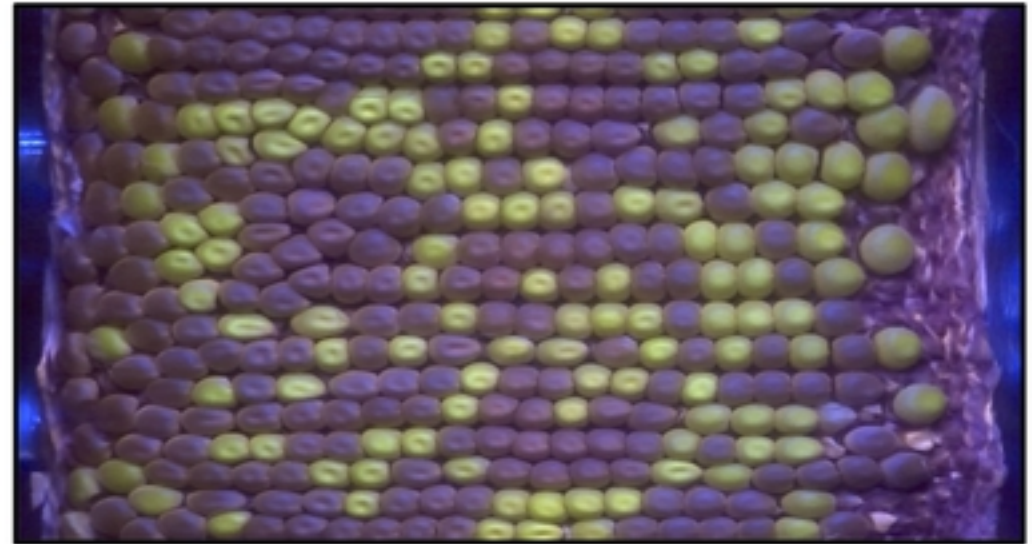
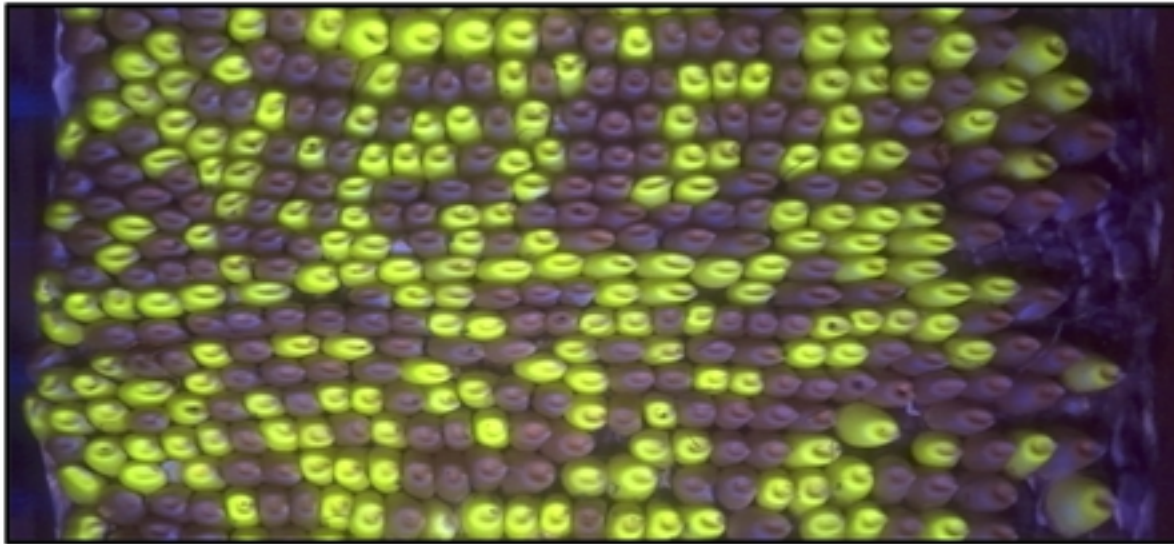
Female cross

Male cross

A



B



C

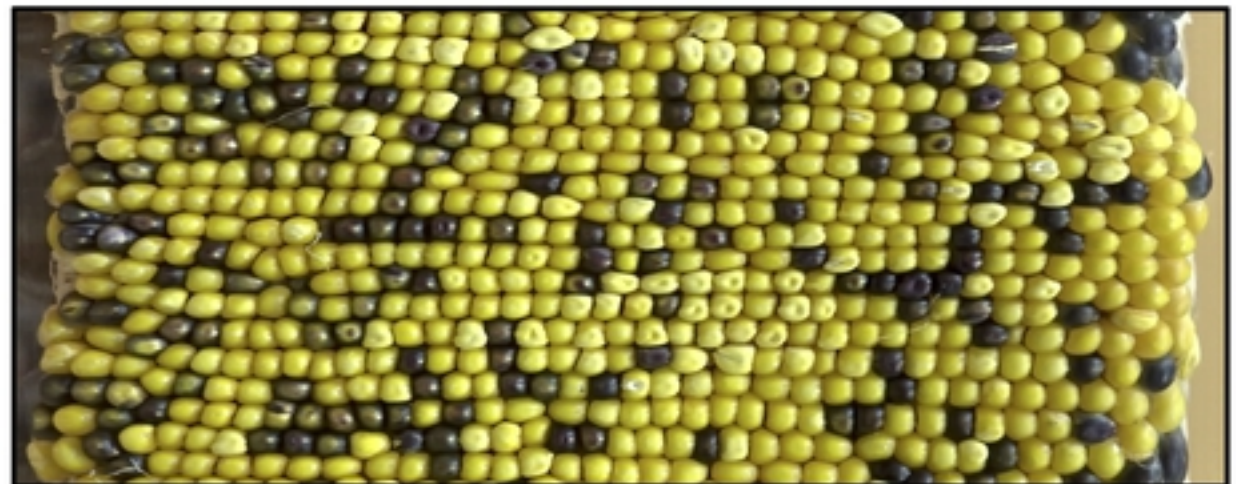
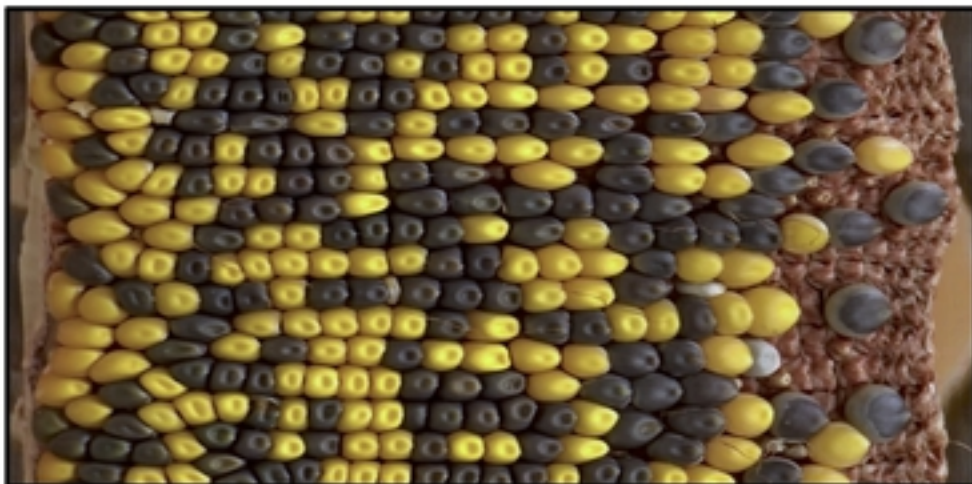


Figure 5

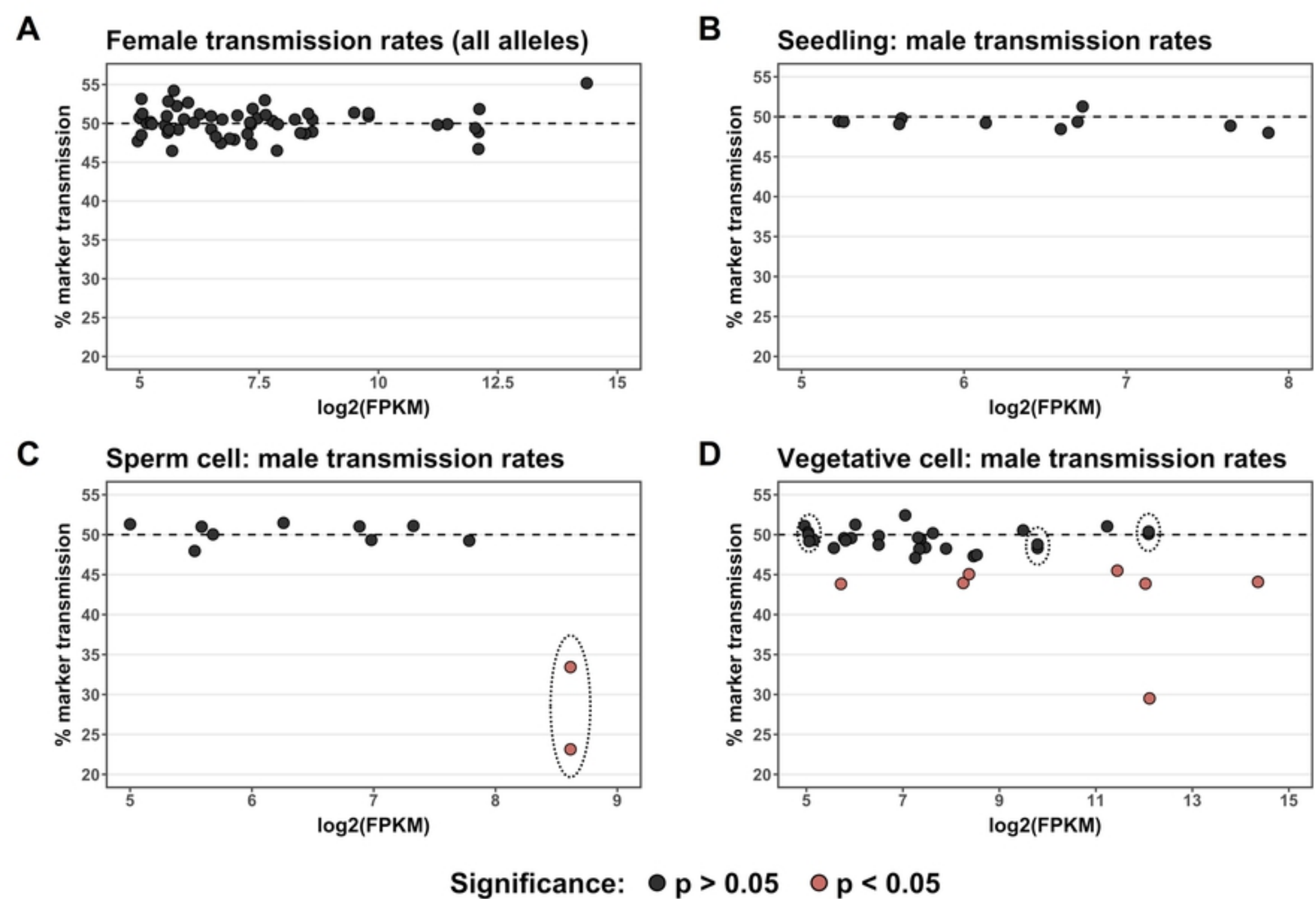
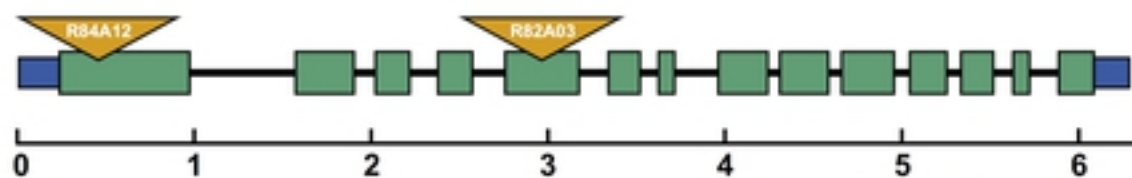


Figure6

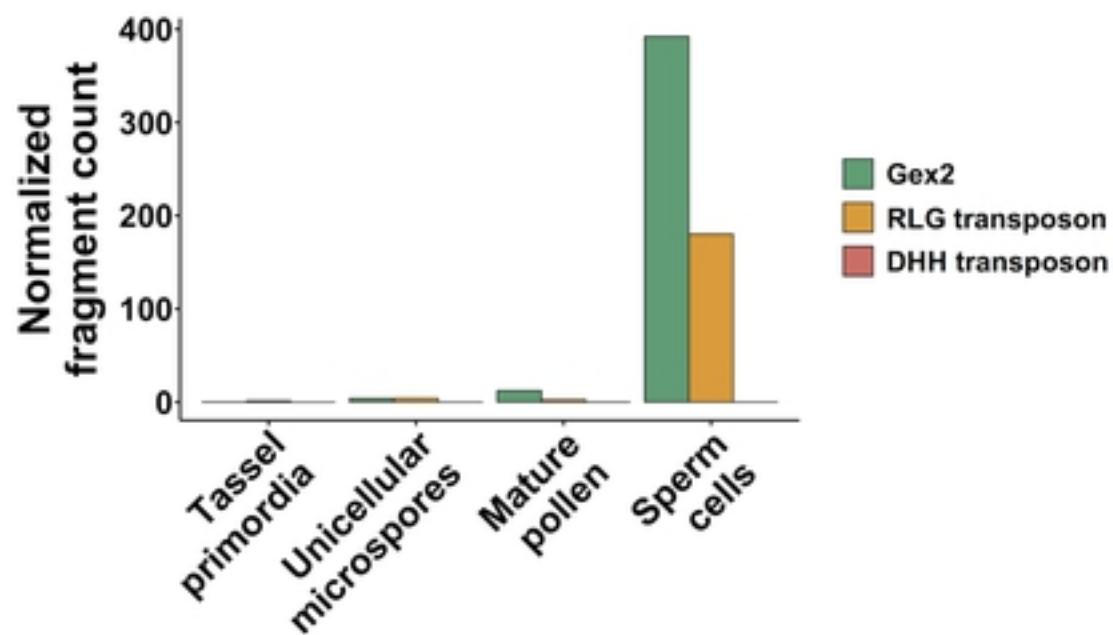
A

Ds-GFP insertions



bioRxiv preprint doi: <https://doi.org/10.1101/796060>; this version posted October 7, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

C

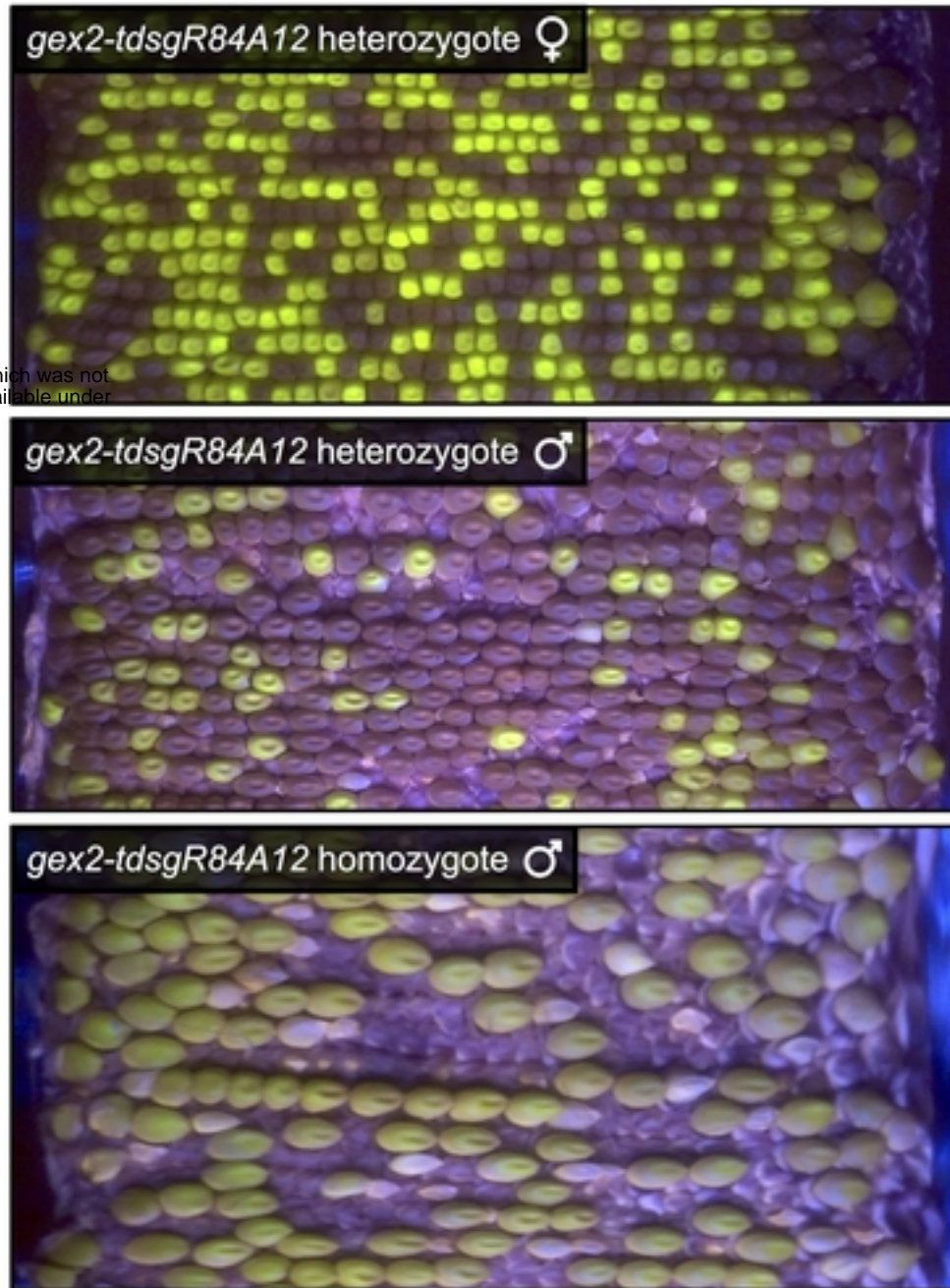


D

Predicted domains



B



E

Frequency of small or aborted seeds

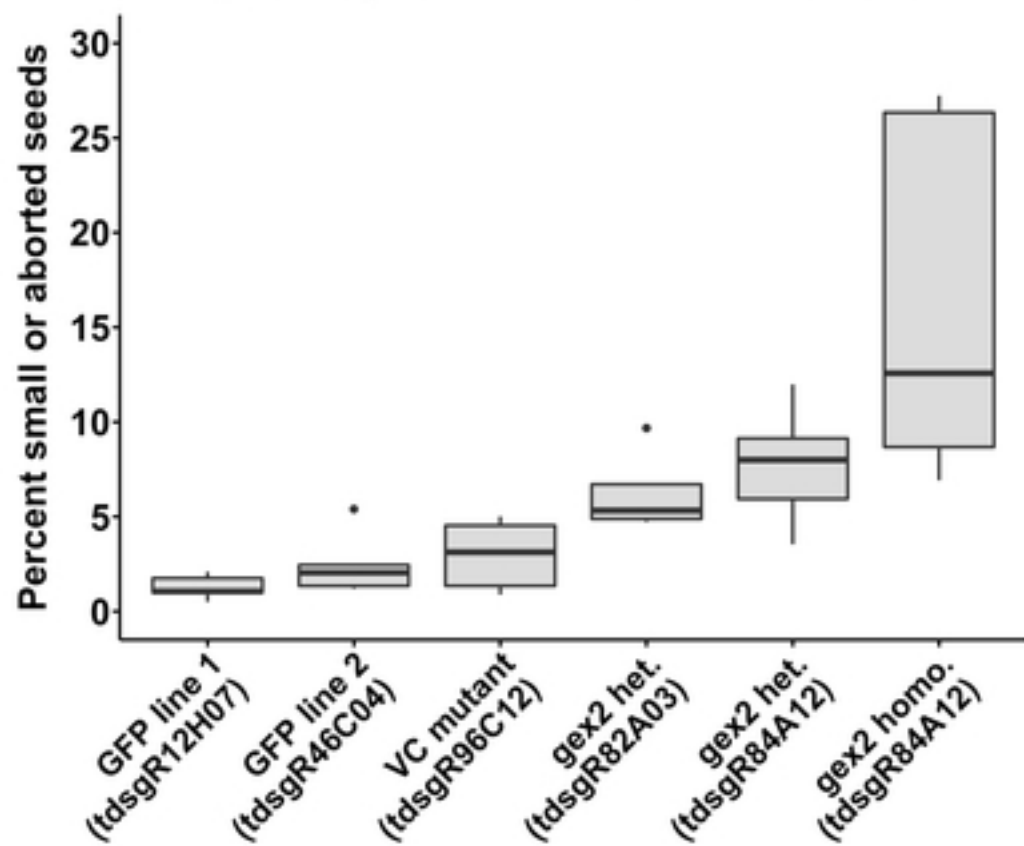


Figure 7

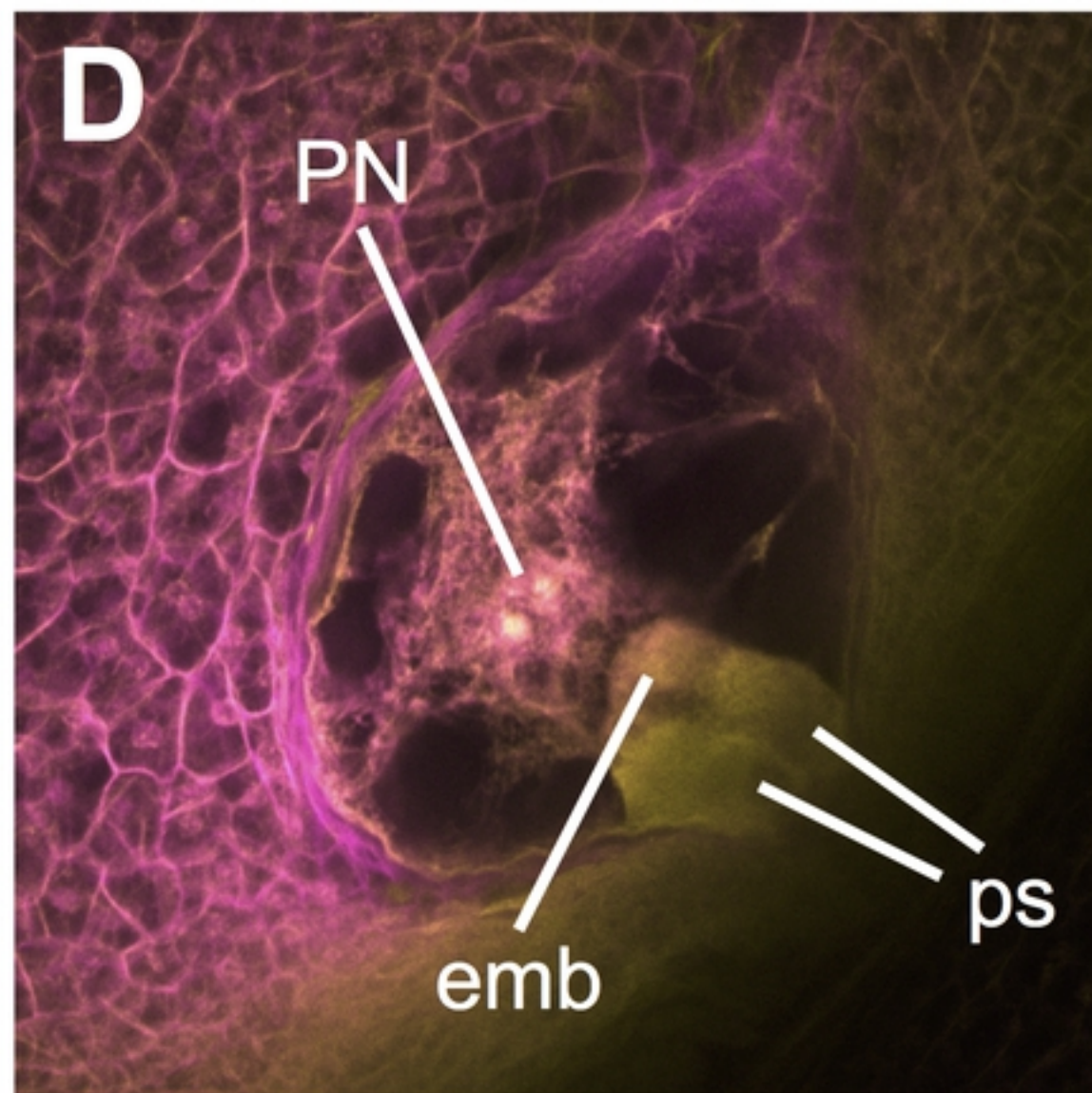
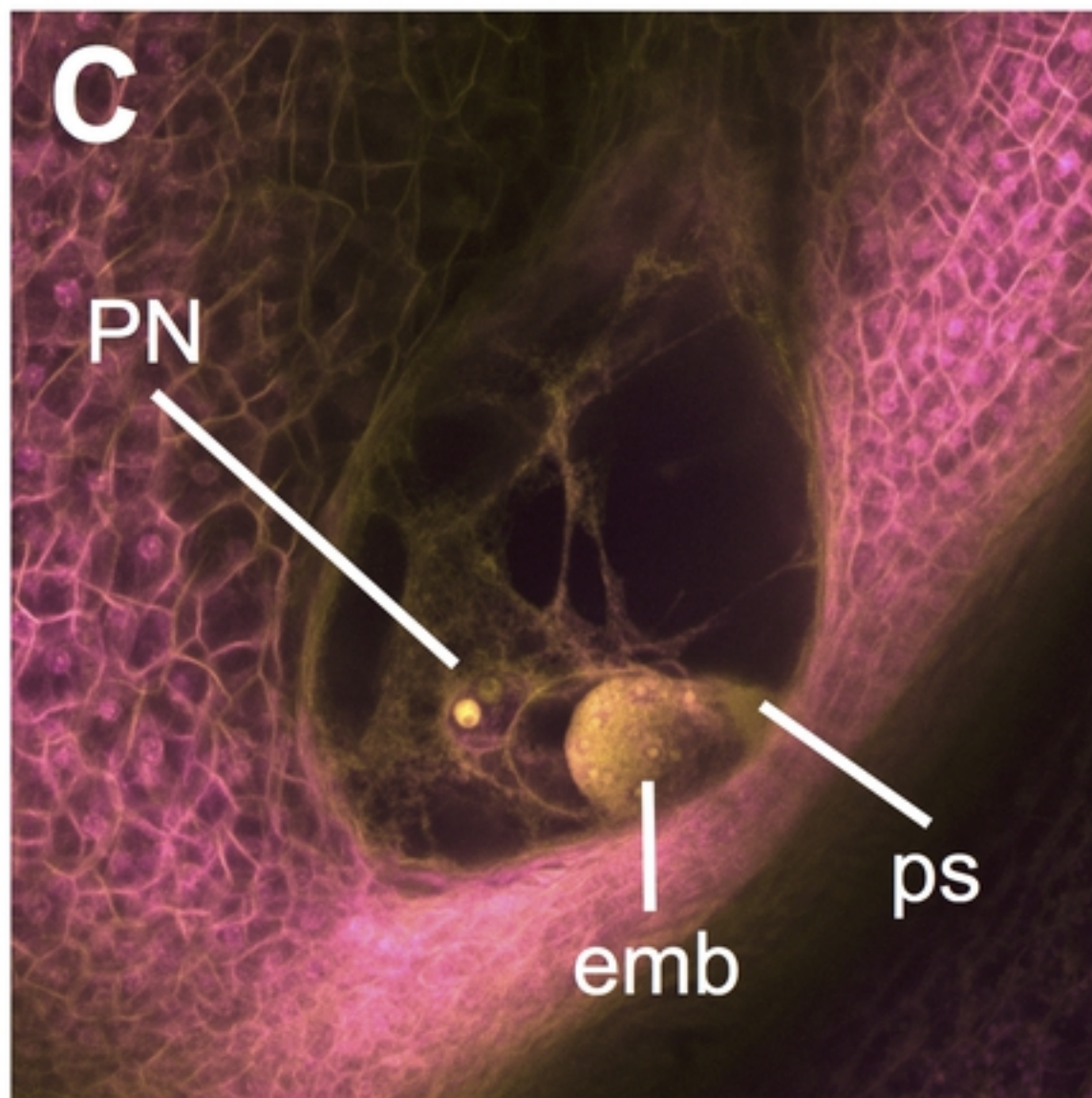
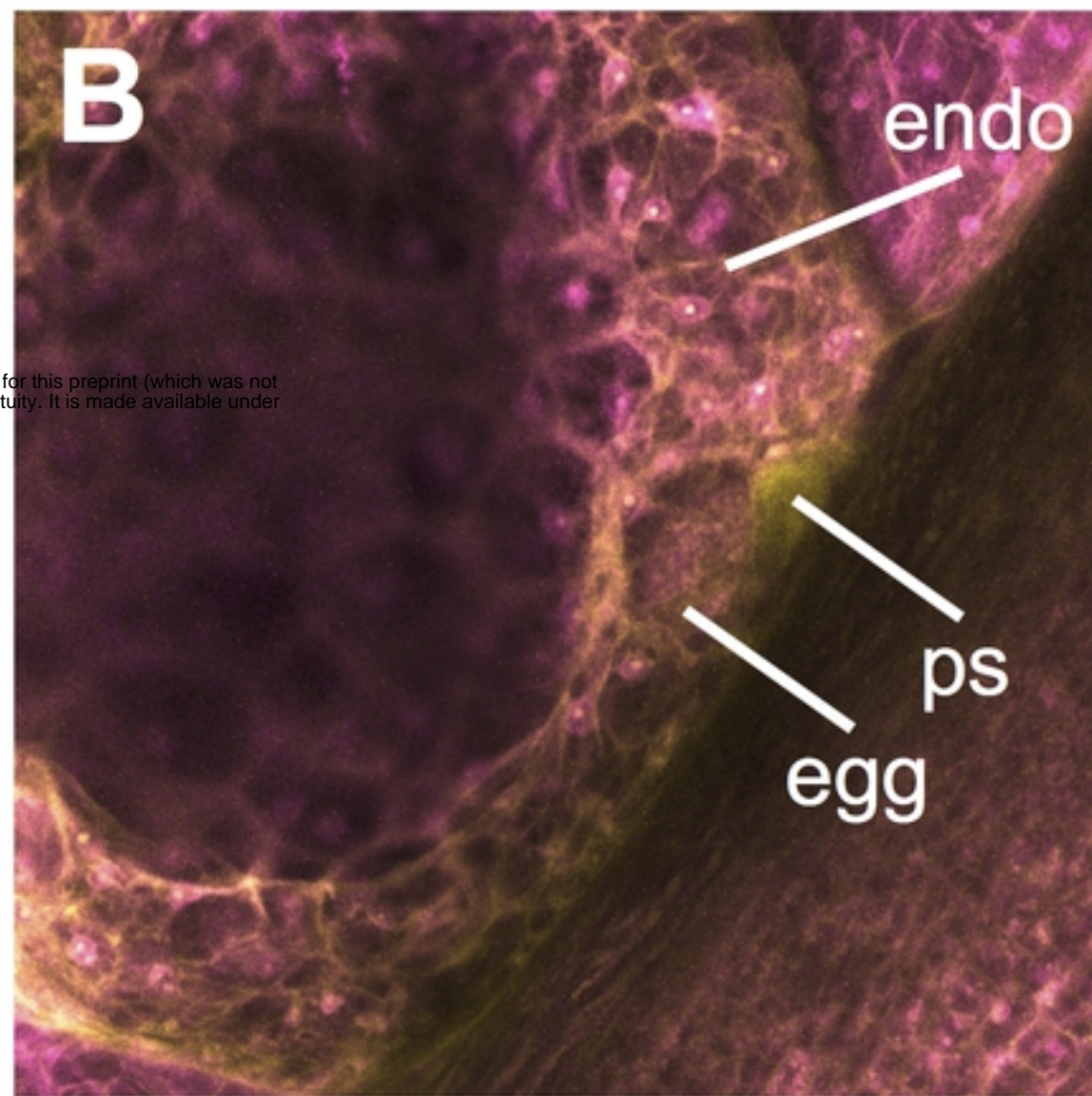
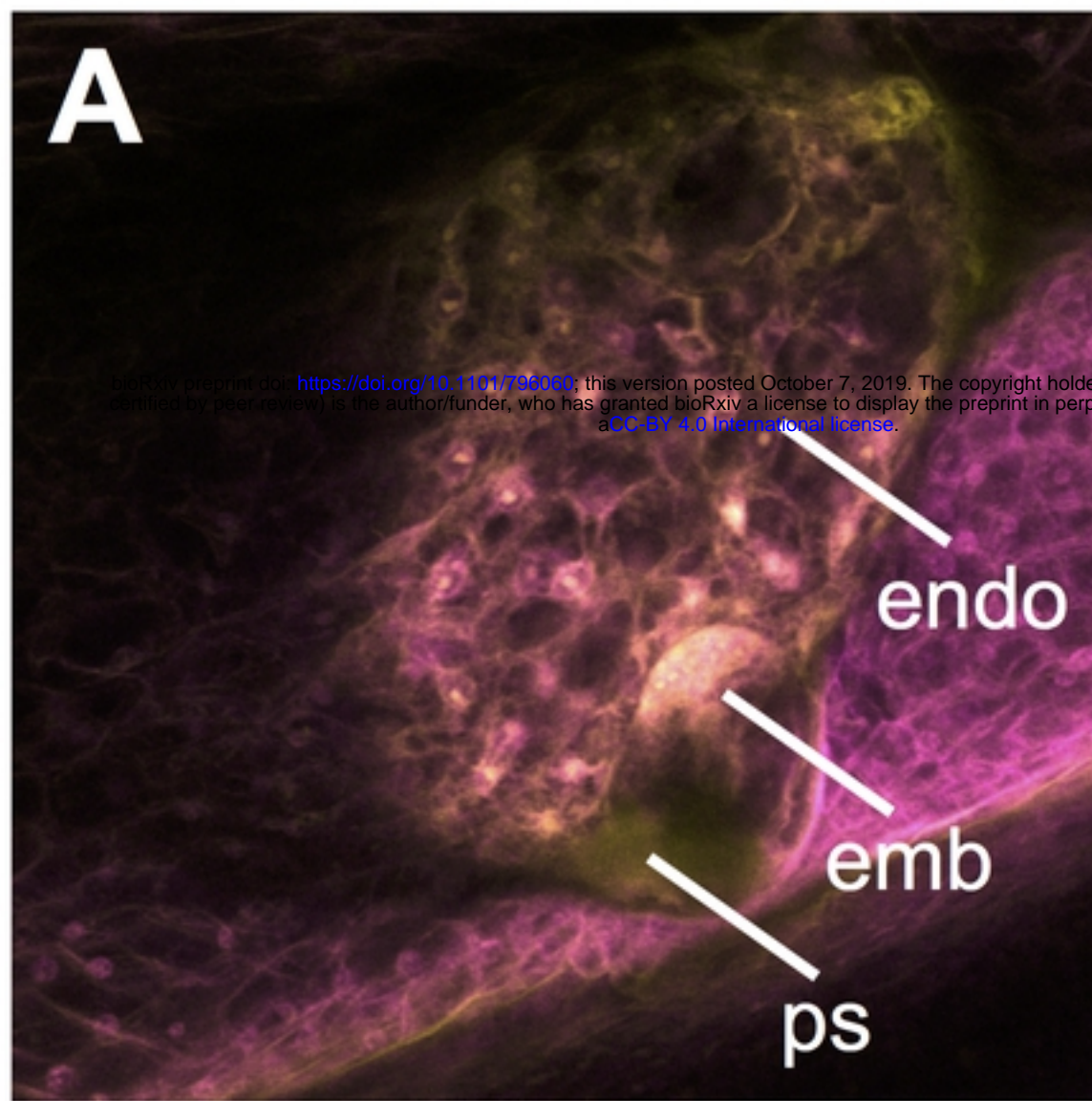


Figure8