# The Streptochaeta genome and the evolution of the

## 2 grasses

1

- 3 Arun Seetharam<sup>1†</sup>, Yunqing Yu<sup>2†</sup>, Sébastien Belanger<sup>2</sup>, Lynn G. Clark<sup>1</sup>, Blake C.
- 4 Meyers<sup>2,3</sup>, Elizabeth A. Kellogg<sup>2\*</sup>, Matthew B. Hufford<sup>1\*</sup>
- <sup>1</sup>Department of Ecology, Evolution, and Organismal Biology, Iowa State University,
- 6 Ames, IA, United States.
- <sup>2</sup>Donald Danforth Plant Science Center, St. Louis, MO, United States.
- 8 <sup>3</sup>Division of Plant Sciences, University of Missouri-Columbia, Columbia, MO, United
- 9 States
- 10 <sup>†</sup>These authors have contributed equally to this work and share first authorship

#### 11 \*Correspondence

- 12 Elizabeth A. Kellogg, ekellogg@danforthcenter.org
- 13 Matthew B. Hufford, mhufford@iastate.edu
- 14 Keywords: Streptochaeta angustifolia, grass, evolution, spikelet, APETALA2-like, R2R3
- 15 MYB, small RNA
- Word Count: 9,566; 7 Main Figures, 3 Supplemental Figures, 12 Supplemental Tables

#### 17 **Abstract**

- 18 In this work, we sequenced and annotated the genome of *Streptochaeta angustifolia*,
- one of two genera in the grass subfamily Anomochlooideae, a lineage sister to all other
- 20 grasses. The final assembly size is over 99% of the estimated genome size, capturing
- 21 most of the gene space. Streptochaeta is similar to other grasses in the structure of its
- fruit (a caryopsis or grain) but has peculiar flowers and inflorescences that are distinct
- from those in the outgroups and in other grasses. To provide tools for investigations of
- 24 floral structure, we analyzed two large families of transcription factors, AP2-like and
- 25 R2R3 MYBs, that are known to control floral and spikelet development in rice and maize
- among other grasses. Many of these are also regulated by small RNAs. Structure of the
- 27 gene trees showed that the well documented whole genome duplication at the origin of
- the grasses (p) occurred before the divergence of the Anomochlooideae lineage from
- the lineage leading to the rest of the grasses (the spikelet clade) and thus that the
- common ancestor of all grasses probably had two copies of the developmental genes.
- However, Streptochaeta (and by inference other members of Anomochlooideae) has
- 32 lost one copy of many genes. The peculiar floral morphology of *Streptochaeta* may thus
- have derived from an ancestral plant that was morphologically similar to the spikelet-
- bearing grasses. We further identify 114 loci producing microRNAs and 89 loci
- generating phased, secondary siRNAs, classes of small RNAs known to be influential in
- transcriptional and post-transcriptional regulation of several plant functions.

## Introduction

37

38

39

40

41 42

43

44 45

46

47 48

49

50

51

52

53

54

77

2014).

The grasses (Poaceae) are arguably the most important plant family to humankind due to their agricultural and ecological significance. The diversity of grasses may not be immediately evident given their apparent morphological simplicity. However, the total number of described species in the family is 11,500+ (Soreng et al., 2017), and more continue to be discovered and described. Grasses are cosmopolitan in distribution, occurring on every continent. Estimates vary based on the definition of grassland, but, conservatively, grasses cover 30% of the Earth's land surface (White et al., 2000; Gibson, 2009). Grasses are obviously the major component of grasslands, but grass species also occur in deserts, savannas, forests (both temperate and tropical), sand dunes, salt marshes and freshwater systems, where they are often ecologically dominant (Lehmann et al., 2019). The traits that have contributed to the long-term ecological success of the grasses have also allowed them to be opportunistic colonizers in disturbed areas and agricultural systems (Linder et al., 2018), where grasses are often the main crops, providing humanity with greater than 50% of its daily caloric intake (Sarwar, 2013). The adaptations and morphologies of the grasses that have led to ecological and agronomic dominance represent major innovations relative to ancestral species.

55 Monophyly of the grass family is unequivocally supported by molecular evidence, but grasses also exhibit several uniquely derived morphological or anatomical traits (Grass 56 57 Phylogeny Working Group et al., 2001; Kellogg, 2015; Leandro et al., 2018). These 58 include the presence of arm cells and fusoid cells (or cavities) in the leaf mesophyll; the 59 pollen wall with channels in the outer wall (intraexinous channels); the caryopsis fruit 60 type; and a laterally positioned, highly differentiated embryo. The 30 or so species of the 61 grass lineages represented by subfamilies Anomochlooideae, Pharoideae and 62 Puelioideae, which are successive sisters to the remainder of the family, all inhabit 63 tropical forest understories, and also share a combination of ancestral features including 64 a herbaceous, perennial, rhizomatous habit; leaves with relatively broad. pseudopetiolate leaf blades; a highly bracteate inflorescence; six stamens in two whorls; 65 pollen with a single pore surrounded by an annulus; a uniovulate gynoecium with three 66 67 stigmas; compound starch granules in the endosperm; and the C<sub>3</sub> photosynthetic pathway (GPWG 2001). The BOP (Bambusoideae, Oryzoideae, Pooideae) + PACMAD 68 (Panicoideae, Aristidoideae, Chloridoideae, Micrairoideae, Arundinoideae, 69 70 Danthonioideae) clade encompasses the remaining diversity of the family ((Kellogg, 71 2015); Figure 1A). The majority of these lineages adapted to and diversified in open 72 habitats, evolving relatively narrow leaves lacking both pseudopetioles and fusoid cells 73 in the mesophyll, spikelets with an array of adaptations for dispersal, and flowers with 74 three stamens and two stigmas. The annual habit evolved repeatedly in both the BOP 75 and PACMAD clades, and the 24+ origins of C<sub>4</sub> photosynthesis occurred exclusively 76 within the PACMAD clade (Grass Phylogeny Working Group II, 2012; Spriggs et al.,

- Anomochlooideae, a tiny clade of four species classified in two genera (*Anomochloa*
- and Streptochaeta), is sister to all other grasses (Figure 1A; (Kellogg, 2015)). Its
- 80 phylogenetic position makes it of particular interest for studies of grass evolution and

- 81 biology, particularly genome evolution. All grasses studied to date share a whole 82 genome duplication (WGD), sometimes referred to as p, which is inferred to have occurred just before the origin of the grasses (Paterson et al., 2004; Wang et al., 2005; 83 84 McKain et al., 2016). Not only are ancient duplicated regions found in the grass 85 genomes studied to date, but the phylogenies of individual gene families often exhibit a 86 doubly labeled pattern consistent with WGD (Rothfels, 2021). In this pattern we see, for 87 example, a tree with the topology shown in Figure 1B, which points to a WGD before 88 the divergence of all sequenced grasses, whereas a WGD after divergence of 89 Streptochaeta, would result in the topology shown in Figure 1C. While there is some 90 evidence from individual gene trees that the duplication precedes the divergence of 91 Streptochaeta+Anomochloa (Preston and Kellogg, 2006; Preston et al., 2009;
- 92 Christensen and Malcomber, 2012; Bartlett et al., 2016; McKain et al., 2016), data are 93 sparse. Thus, defining the position of the grass WGD requires a whole genome
- 94 sequence of a species of Anomochlooideae.
- 95 Anomochlooideae is also in a key position for understanding the origins of the 96 morphological innovations of the grass family. All grasses except Anomochlooideae 97 bear their flowers in tiny clusters known as spikelets (little spikes) (Judziewicz et al., 98 1999; Grass Phylogeny Working Group et al., 2001; Kellogg, 2015). Because the 99 number, position, and structure of spikelets affect the total number of seeds produced 100 by a plant, the genes controlling their development are a subject of continual research 101 (e.g., (Whipple, 2017; Huang et al., 2018; Li et al., 2019a, 2019b), to cite just a few). In 102 contrast to the rest of the family, the flowers in Anomochlooideae are borne in complex 103 bracteate structures sometimes called "spikelet equivalents" ((Soderstrom and Ellis, 104 1987; Judziewicz and Soderstrom, 1989; Judziewicz et al., 1999); Figures 2 and 3). 105 These differ from both the conventional monocot flowers of the outgroups and the 106 spikelets of the remainder of the grasses (i.e., the "spikelet clade"; (Sajo et al., 2008, 107 2012; Preston et al., 2009; Kellogg et al., 2013)). The structure of the phylogeny 108 suggests potential interpretations of the origin of the spikelet. One possibility is a "stepwise" model, in which a set of changes to the genetic architecture of floral 109 development occurred before the divergence of Anomochlooideae, leading to the 110 111 formation of spikelet equivalents; these changes were then followed by a second set of 112 changes that led to formation of spikelets in the rest of what would become the spikelet 113 clade. An alternative, which is also consistent with the phylogeny, is a "loss model", in 114 which all the genes and regulatory architecture needed for making spikelets originated 115 before the origin of Anomochlooideae, but portions of that architecture were 116 subsequently lost. Thus, the stepwise model implies that the spikelet equivalents are 117 somehow intermediate between a standard monocot flower and a grass spikelet. 118 whereas the loss model implies that the spikelet equivalents are highly modified or 119 rearranged spikelets. Resolving these hypothetical models will help reveal both how the 120 unique spikelet structure and the overall floral bauplan in grasses evolved.
- 121 Of the handful of species in the Anomochlooideae, Streptochaeta angustifolia (Figures 2 and 3) is the most easily grown from seed and an obvious candidate for ongoing 122 123 functional genomic investigation. Hereafter in this paper, we will refer to S. angustifolia
- 124 simply as Streptochaeta, and use it as a placeholder for the rest of the subfamily. We
- 125 present a draft genome sequence for Streptochaeta that captures the gene-space of

- this species at high contiguity, and we use this genome to assess the position of the
- grass WGD. Genes and small RNAs (sRNAs) are annotated. Because of the distinct
- floral morphology of *Streptochaeta*, we also investigate the molecular evolution of two
- major transcription factor families, APETALA2-like and R2R3 MYB, which are known to
- control floral and spikelet structure in other grasses and are regulated by sRNAs.

## **Materials and Methods**

### Input data

131

132

155

- 133 Streptochaeta leaf tissue was harvested and used to estimate genome size at the Flow
- 134 Cytometry Facility at Iowa State University. DNA was then isolated using Qiagen
- DNeasy plant kits. Three Illumina libraries (paired end and 9- and 11-kb mate pair) were
- generated from these isolations at the Iowa State University (ISU) DNA Facility. One
- lane of 150 bp paired-end HiSeq sequencing (insert size of 180 bp) and one lane of 150
- bp mate-pair HiSeq sequencing (9- and 11-kb libraries pooled) were generated, also at
- the ISU DNA Facility (**Table S1**). Additionally, for the purpose of contig scaffolding,
- 140 Bionano libraries were prepared by first isolating high molecular weight DNA using the
- 141 Bionano Prep™ Plant DNA Isolation Kit followed by sequencing using the Irys system.

### 142 Genome assembly

- We used MaSuRCA v2.21 (Zimin et al., 2013) to generate a draft genome of
- 144 Streptochaeta. The MaSuRCA assembler includes error correction and quality filtering,
- generation of super reads, super read assembly, and gap closing to generate more
- complete and larger scaffolds. Briefly, the config file was edited to include both paired-
- end and mate-pair library data for *Streptochaeta*. The JF\_SIZE parameter was adjusted
- to 20,000,000,000 to accommodate the large input file size, and NUM THREADS was
- set to 128. All other parameters in the config file were left as default. The assembly was
- executed by first generating the assemble.sh script using the config file and submitting
- to a high-memory node using the PBS job scheduler. We then used Bionano
- technology to generate an optical map for the genome and to perform hybrid
- scaffolding. All scripts for assembly and downstream analysis are available at:
- 154 https://github.com/HuffordLab/streptochaeta.

## Assembly evaluation and post-processing

- The Bionano assembly was screened for haplotigs, and additional gaps were filled
- using Redundans v0.13a (Pryszcz and Gabaldón, 2016). Briefly, the scaffolds were
- mapped to themselves using the LAST v719 alignment program (Kielbasa et al., 2011)
- and any scaffold that completely overlapped a longer scaffold with more than 80%
- identity was considered redundant and excluded from the final assembly. Additionally,
- short read data were aligned back to the hybrid assembly and GapCloser v1.12 from
- SOAPdenovo2 (Luo et al., 2012) and SSPACE v3.0 (Boetzer et al., 2011) were run in
- multiple iterations to fill gaps. The final reduced, gap-filled assembly was screened for
- 164 contamination, using Blobtools v0.9.19 (Laetsch and Blaxter, 2017), and any scaffolds
- that matched bacterial genomes were removed. The assembly completeness was then

- evaluated using BUSCO v3.0.2 (Simão et al., 2015) with the plant profile and standard
- assemblathon metrics.
- To annotate the repeats in the genome, we used EDTA v1.8.3 (Ou et al., 2019) with
- default options except for --species, which was set to "others". The obtained TE library
- was then used for masking the genome for synteny analyses. Assembly quality of the
- repeat space was assessed based on the LTR Assembly Index (LAI; (Ou et al., 2018)),
- which was computed using ltr\_retriever v2.9.0 (Ou and Jiang, 2018) and the EDTA-
- 173 generated LTR list.

### Gene prediction and annotation

- Gene prediction was carried out using a comprehensive method combining ab initio
- predictions (from BRAKER; (Hoff et al., 2019)) with direct evidence (inferred from
- transcript assemblies) using the BIND strategy (Seetharam et al., 2019 and citations
- therein). Briefly, RNA-Seg data were mapped to the genome using a STAR (v2.5.3a)-
- indexed genome and an iterative two-pass approach under default options in order to
- generate BAM files. BAM files were used as input for multiple transcript assembly
- programs (Class2 v2.1.7, Cufflinks v2.2.1, Stringtie v2.1.4 and Strawberry v1.1.2) to
- assemble transcripts. Redundant assemblies were collapsed and the best transcript for
- each locus was picked using Mikado (2.0rc2) by filling in the missing portions of the
- ORF using TransDecoder (v5.5.0) and homology as informed by the BLASTX
- 185 (v2.10.1+) results to the SwissProtDB. Splice junctions were also refined using
- Portcullis (v1.2.1) in order to identify isoforms and to correct misassembled transcripts.
- Both ab initio and the direct evidence predictions were analyzed with TESorter (Zhang
- et al., 2019) to identify and remove any TE-containing genes and with phylostratr
- 189 (v0.20; (Arendsee et al., 2019)) to identify orphan genes (i.e., species-specific genes).
- 190 As ab initio predictions of young genes can be unreliable (Seetharam et al., 2019),
- these were excluded. Finally, redundant copies of genes between direct evidence and
- 192 ab initio predictions were identified and removed using Mikado compare (2.0rc2;
- 193 (Venturini et al., 2018)) and merging was performed locus by locus, incorporating
- additional isoforms when necessary. The complete decision table for merging is
- provided in **Table S2**. After the final merge, phylostratr was run again on the
- annotations to classify genes based on their age.
- 197 Functional annotation was performed based on homology of the predicted peptides to
- the curated SwissProt/UniProt set (UniProt Consortium, 2021) as determined by BLAST
- v2.10.1+ (Edgar, 2010). InterProScan v5.48-83 was further used to find sequence
- 200 matches against multiple protein signature databases.

## **Synteny**

- 202 Synteny of CDS sequences for Strepotchaeta was determined using CoGe (Lyons and
- Freeling, 2008), against the genomes Brachypodium (International Brachypodium)
- 204 Initiative, 2010), Oryza sativa (Ouyang et al., 2007), and Setaria viridis (Mamidi et al.,
- 205 2020). SynMap2 (Haug-Baltzell et al., 2017) was employed to identify syntenic regions
- across these genomes. Dot plots and chain files generated by SynMap2 under default
- 207 options were used for presence-absence analysis. We also performed repeat-masked

- whole genome alignments using minimap2 (Li, 2018) following the Bioinformatics
- 209 Workbook methods (https://bioinformaticsworkbook.org/dataWrangling/genome-
- 210 dotplots.html).
- 211 Identification of APETALA2 (AP2)-like and R2R3 MYB proteins in
- 212 selected monocots
- 213 A BLAST database was built using seven grass species including Streptochaeta and
- 214 two outgroup monocots. Protein and CDS sequences of the following species were
- retrieved from Phytozome 13.0: *Ananas comosus* (Acomosus 321 v3), *Brachypodium*
- 216 distachyon (Bdistachyon\_556\_v3.2), Oryza sativa (Osativa\_323\_v7.0), Spirodela
- 217 polyrhiza (Spolyrhiza 290 v2), Setaria viridis (Sviridis 500 v2.1) and Zea mays
- 218 (Zmays 493 APGv4). Sequences of *Eragrostis tef* were retrieved from CoGe (id50954)
- (VanBuren et al., 2020). Sequences of *Triticum aestivum* were retrieved from Ensemble
- 220 Plant r46 (Triticum aestivum.IWGSCv1) (Table S3).
- AP2 and MYB proteins were identified using BLASTP and hmmscan (HMMER 3.1b2;
- 222 http://hmmer.org/) in an iterative manner. Specifically, 18 Arabidopsis AP2-like proteins
- (Kim et al., 2006) were used as an initial query in a blastp search with an E-value
- threshold of 1e-10. The resulting protein sequences were filtered based on the
- presence of an AP2 domain using hmmscan with an E-value threshold of 1e-3 and
- domain E-value threshold of 0.1. The filtered sequences were used as the query for the
- 227 next round of blastp and hmmscan until the maximal number of sequences was
- retrieved. For MYB proteins, Interpro MYB domain (IPR017930) was used to retrieve
- rice MYBs using Oryza sativa Japonica Group genes (IRGSP-1.0) as the database on
- Gramene Biomart (http://ensembl.gramene.org/biomart/martview/). The number of MYB
- domains was counted by searching for "Myb DNA-bind" in the output of hmmscan, and
- 82 proteins with two MYB domains were used as the initial guery. Iterative blastp and
- 233 hmmscan were performed in the same manner as for AP2 except using a domain E-
- value threshold of 1e-3.
- The number of AP2 or MYB domains was again counted in the final set of sequences in
- the hmmscan output, and proteins with more than one AP2 domain or two MYB
- domains were treated as AP2-like or R2R3 MYB, respectively. To ensure that no
- orthologous proteins were missed due to poor annotation in the AP2 or MYB domain.
- we performed another round of BLASTP searches, and kept only the best hits. These
- sequences were also included in the construction of the phylogenetic trees.
- 241 Construction of phylogenetic trees
- 242 Protein sequences were aligned using MAFFT v7.245 (Katoh and Standley, 2013) with
- 243 default parameters. The corresponding coding sequence alignment was converted
- using PAL2NAL v14 (Suyama et al., 2006) and used for subsequent tree construction.
- For AP2-like genes, the full length coding sequence alignment was used. For MYB, due
- 246 to poor alignment outside of the MYB domain, trimAl v1.2 (Capella-Gutiérrez et al.,
- 247 2009) was used to remove gaps and non-conserved nucleotides with a gap threshold (-
- 248 gt) of 0.75 and percentage alignment conservation threshold (-con) of 30. A maximum
- 249 likelihood tree was constructed using IQ-TREE v1.6.12 (Minh et al., 2020) with default
- settings. Sequences that resulted in long branches in the tree were manually removed,

- and the remaining sequences were used for the final tree construction. Visual formatting
- of the tree was performed using Interactive Tree Of Life (iTOL) v4 (Letunic and Bork,
- 253 2019).

266

## RNA isolation, library construction and sequencing

- We collected tissues from leaf and pistil as well as 1.5 mm, 3 mm and 4 mm anthers.
- 256 Samples were immediately frozen in liquid nitrogen and kept at -80°C prior to RNA
- isolation. Total RNA was isolated using the PureLink Plant RNA Reagent (Thermo
- 258 Fisher Scientific, Waltham, MA, USA). sRNA libraries were published previously (Patel
- et al., 2021). RNA sequencing libraries were prepared from the same material using the
- 260 Illumina TruSeg stranded RNA-seg preparation kit (Illumina Inc., United States)
- following manufacturer's instructions. Parallel analysis of RNA ends (PARE) libraries
- were prepared from a total of 20 µg of total RNA following the method described by Zhai
- et al. (2014). For all types of libraries, single-end sequencing was performed on an
- 264 Illumina HiSeq 2000 instrument (Illumina Inc., United States) at the University of
- 265 Delaware DNA Sequencing and Genotyping Center.

### **Bioinformatic analysis of small RNA data**

- Using cutadapt v2.9 (Martin, 2011), sRNA-seg reads were pre-processed to remove
- adapters (Table S4), and we discarded reads shorter than 15 nt. The resulting 'clean'
- reads were mapped to the Streptochaeta genome using ShortStack v3.8.5 (Johnson et
- al., 2016) with the following parameters: -mismatches 0, -bowtie m 50, -mmap u, -
- 271 dicermin 19, -dicermax 25 and -mincov 0.5 transcripts per million (TPM). Results
- 272 generated by ShortStack were filtered to keep only clusters having a predominant RNA
- size between 20 and 24 nucleotides, inclusively. We then annotated categories of
- microRNAs (miRNAs) and phased small interfering RNAs (phasiRNAs).
- 275 First, sRNA reads representative of each cluster were aligned to the monocot-related
- 276 miRNAs listed in miRBase release 22 (Kozomara and Griffiths-Jones, 2014; Kozomara
- et al., 2019) using NCBI BLASTN v2.9.0+ (Camacho et al., 2009) with the following
- parameters: -strand both, -task blastn-short, -perc identity 75, -no greedy and -
- 279 ungapped. Homology hits were filtered and sRNA reads were considered as known
- 280 miRNA based on the following criteria: (i) no more than four mismatches and (ii) no
- more than 2-nt extension or reduction at the 5' end or 3' end. Known miRNAs were
- and the state of t
- summarized by family. Small RNA reads with no homology to known miRNAs were
- annotated as novel miRNAs using the *de novo* miRNA annotation performed by
- ShortStack. The secondary structure of new miRNA precursor sequences was drawn
- using the RNAfold v2.1.9 program (Lorenz et al., 2011). Candidate novel miRNAs were
- 286 manually inspected, and only those meeting published criteria for plant miRNA
- annotations (Axtell and Meyers, 2018) were retained for subsequent analyses. Then,
- the remaining sRNA clusters were analyzed to identify phasiRNAs based on ShortStack
- 289 analysis reports. sRNA clusters having a "Phase Score" >30 were considered as true
- 290 positive phasiRNAs. Genomic regions corresponding to these phasiRNAs were
- considered as PHAS loci and grouped in categories of 21- and 24-PHAS loci referring to
- the length of phasiRNAs derived from these loci. Other sRNA without miRNA or
- 293 phasiRNA signatures were not considered for analysis or interpretation in this study.

- To compare sRNAs accumulating in *Streptochaeta* anthers with other monocots, we
- analyzed sRNA samples of Asparagus officinalis, Oryza sativa and Zea mays anthers.
- The GEO accession numbers for those datasets are detailed in **Table S3**. We analyzed
- these data as described for the *Streptochaeta* sRNA-seq data.
- We used the upSetR package (UpSetR; Lex et al., 2014; Conway et al., 2017) to
- visualize the overlap of miRNA loci annotated in Streptochaeta, compared to other
- 300 species.

317

### Bioinformatic analysis of PARE data

- We analyzed the PARE data to identify and validate miRNA-target pairs in anther, pistil,
- and leaf of Streptochaeta tissues. Using cutadapt v2.9, PARE reads were pre-
- processed to remove adapters (**Table S4**) and reads shorter than 15 nt were discarded.
- Then, we used PAREsnip2 (Thody et al., 2018) to predict all miRNA-target pairs and to
- validate the effective miRNA-guided cleavage site using PARE reads. We ran
- 307 PAREsnip2 with default parameters using Fahlgren & Carrington targeting rules
- 308 (Fahlgren and Carrington, 2010). We considered only targets in categories 0, 1 and 2
- for downstream analysis. We used the EMBL-EBI HMMER program v3.3 (Potter et al.,
- 2018) to annotate the function of miRNA target genes using the phmmer function with
- 311 the SwissProt database.

### 312 Prediction of miRNA binding sites

- Mature miR172 and miR159 sequences from all available monocots were obtained from
- 314 miRBase (Kozomara et al., 2019). miRNA target sites in AP2-like and R2R3 MYB
- transcripts were predicted on a web server TAPIR (Bonnet et al., 2010) with their default
- settings (score = 4 and free energy ratio = 0.7).

## Results

## 318 Flow Cytometry

- Two replicates of flow cytometry estimated the 1C DNA content for *Streptochaeta* to be
- 1.80 pg and 1.83 pg, which, when converted to base pairs, yields a genome size of
- 321 approximately 1.77 Gb.

## 322 Genome Assembly and post-processing

- Two lanes of short reads (Illumina HiSeg 2500), generated a total of 259 million reads.
- Paired-end reads with a fragment size of 250bp were generated at approximately 25.7x
- genomic coverage, while the mate-pair libraries with 9- and 11-kb insert size collectively
- provided 22.6x coverage. Based on k-mer analysis of these data with the program
- 327 Jellyfish (Marçais and Kingsford, 2011), we estimated the repeat content for the
- 328 Streptochaeta genome to be approximately 51%. Implementation of the MaSuRCA
- assembly algorithm generated an assembly size at 99.8% of the estimated genome
- size, suggesting that a large portion of the genome, including repetitive regions were
- 331 successfully assembled. The MaSuRCA assembler generated a total of 22,591
- scaffolds, with an N50 of 2.4Mb and an L50 of 170.

- The Bionano data produced an optical map near the expected genome size (1.74 Gb)
- with an N50 of 824kb. Through scaffolding with the optical map and collapsing with
- Redundans software, the total number of scaffolds dropped to 17,040, improving the
- N50 to 2.6Mb and the L50 to 161. A total of 79,165 contigs were provided as input for
- Redundans for scaffold reduction (total size 1,898 Mbp). With eight iterations of
- haplotype collapsing, the number of scaffolds was reduced to 17,040 (total size 1,796
- Mbp). Additional rounds of gap-filling using GapCloser reduced the total number of gaps
- (Ns) from 210.13 Mbp to 76.33 Mbp. The improvement in the N50/N90 values with
- each iteration is provided in **Table S5**.
- The final assembly included a total of 3,010 out of 3,278 possible complete Liliopsida
- BUSCOs (91.8%). Of these 2,767 (84.4% of the total) were present as a complete
- 344 single copy. Only 158 BUSCOs were missing entirely with another 110 present as
- fragmented genes. The LAI (LTR Assembly Index) score, which assesses the contiguity
- of the assembled LTR retrotransposons, was 9.02, which is somewhat higher than most
- short-read-based assemblies (Ou et al., 2018), perhaps due to the relatively low repeat
- content of the Streptochaeta genome and the use of mate-pair sequencing libraries. Dot
- 349 plots of *Streptochaeta* contigs aligned to rice revealed substantial colinearity (**Figure**
- 350 **S1**).

359

### **Contamination Detection**

- 352 BlobTools (v0.9.19) (Laetsch and Blaxter, 2017) detected over 95% of the scaffolds
- 353 (1742 Mbp) belonging to the Streptophyta clade out of the 1,797 Mbp of assigned
- scaffolds (GC mean: 0.54). Approximately 2% of the scaffolds mapped to the
- 355 Actinobacteria (36.3Mbp, GC mean: 0.72) and ~0.5% of scaffolds to Chordata (9Mbp,
- 356 GC mean: 0.48). Scaffolds assigned to additional clades by BlobTools collectively
- 357 comprise ~1.46 Mbp and the remaining 8.47 Mbp of scaffolds lacked any hits to the
- database. All bacterial, fungal and vertebrate scaffolds were purged from the assembly.

## Gene prediction and annotation

- 360 *Direct Evidence predictions:* More than 79% of the total RNAseq reads mapped
- uniquely to the Streptochaeta genome with <7% multi-mapped reads. Paired-end reads
- mapped (uniquely) at a higher rate (88.59%) than the single-end RNAseq (70.38%)
- reads. Genome-guided transcript assemblers produced varying numbers of transcripts
- across single-end (SE) and paired-end (PE) data as well as various assemblers.
- 365 Cufflinks produced the highest number of transcripts (SE: 65,552; PE:66,069), followed
- 366 by StringTie (SE: 65,495, PE: 48,111), and Strawberry (SE:68,812; PE:43,882). Class2
- generated fewer transcripts overall (PE: 43,966; SE: 13,173). The best transcript for
- 368 each locus was picked by Mikado from the transcript assemblies based on its
- completeness, homology, and accuracy of splice sites. Mikado also removed any non-
- coding (due to lack of ORFs) or redundant transcripts to generate 28,063 gene models
- 371 (41,857 transcripts). Mikado also identified 19,135 non-coding genes within the provided
- 372 transcript assemblies. Further filtering for transposable-element-containing genes and
- genes with low expression reduced the total number of evidence-based predictions to
- 374 27,082 genes (40,865 transcripts).

- 375 **Ab initio predictions:** BRAKER, with inputs including predicted proteins from the direct
- evidence method (as a gff3 file produced by aligning proteins to a hard-masked
- 377 Streptochaeta genome) and the mapped RNA-Seq reads (as a hints file using the bam
- file), produced a total of 611,013 transcripts on a soft-masked genome. This was then
- subjected to filtering to remove any TE containing genes (244,706 gene models) as well
- as genes only found in *Streptochaeta* (466,839 gene models). After removing both of
- these classes of genes, which overlapped to an extent, the total number of ab initio
- predictions dropped to 40,921 genes (44,013 transcripts).
- 383 BIND (merging BRAKER predictions with directly inferred genes): After comparing
- 384 BRAKER and direct evidence predictions with Mikado compare: 9,617 transcripts were
- exactly identical and direct evidence predictions were retained; 3,263 transcripts from
- 386 Mikado were considered incomplete and were replaced with BRAKER models; 13,360
- 387 BRAKER models were considered incomplete and replaced with direct evidence
- transcripts; 1,884 predictions were adjacent but non-overlapping, and 17,894
- predictions were BRAKER-specific and were retained in the final merged predictions.
- The final gene set included a total of 44,980 genes (58,917 transcripts).
- 391 Functional Annotation: Functional annotation was informed by homology to the
- curated proteins in SwissProt and resulted in the assignment of putative functions for
- 393 38,955 transcripts (10,556 BRAKER predictions, and 28,399 direct evidence
- predictions). Of the unassigned transcripts, 41 predictions had pfam domain matches,
- and 16,918 transcripts had an interproscan hit. Only 3,068 transcripts contained no
- 396 additional information in the final GFF3 file.
- 397 **Phylostrata:** All gene models predicted by the BIND strategy were examined by
- 398 classifying the genes based on their presumed age. More than 8% of the total genes
- 399 (3,742) were specific to the Streptochaeta genus and more than 15% (6,930) of genes
- were Poaceae specific. 19% (8,494) of genes' origins could be traced back to cellular
- organisms and 15% (6,708) to Eukaryotic genes. The distribution of genes based on
- strata and annotation method is provided in **Table S6**.
- 403 *Transposable Element Annotation:* The repeat annotation performed by the EDTA
- package comprised 66.82% of the genome, the bulk of which were LTR class elements
- 405 (42.9% in total; Gypsy: 28.16%, Copia: 8.9%, rest: 5.84%), followed by DNA repeats
- 406 (23.39% in total; DTC-type: 13.65, DTM-type: 5.78%, rest: 3.96%), and MITE class
- 407 repeats (all types 0.54%).

## 408 Molecular evolution of APETALA2-like and R2R3 MYB

- 409 transcription factors
- 410 Our highly contiguous assembly in genic regions combined with gene model and
- functional annotations allowed: 1) an investigation of gene families known to play a role
- in floral development that have potential relevance to the origin of the grass spikelet.
- and 2) evaluation of patterns of orthology between genes in *Streptochaeta* and
- BOP/PACMAD grasses to clarify the timing of the p WGD. Many transcription factor
- families are known to affect spikelet development in the grasses (Hirano et al., 2014;
- 416 Whipple, 2017). Of these, APETALA2 (AP2)-like genes control meristem identity and

- 417 floral morphology, including the number of florets per spikelet (Chuck et al., 1998; Lee
- 418 and An, 2012; Zhou et al., 2012; Debernardi et al., 2020). Several *R2R3 MYB* genes
- are also known to function in floral organ development, especially in anthers (Zhu et al.,
- 420 2008; Aya et al., 2009; Zhang et al., 2010; Schmidt et al., 2013). We explored patterns
- of duplication and loss in these gene families between the origin of the grasses and the
- origin of the spikelet clade, i.e. before and after the divergence of *Streptochaeta*.

#### APETALA2-like

- 424 Previous work on molecular evolution of AP2-like proteins found that the gene family
- was divided into two distinct lineages, euAP2 and AINTEGUMENTA (ANT) (Kim et al.,
- 426 2006). A Maximum Likelihood tree of AP2-like genes was constructed and rooted at the
- branch that separates euAP2 and ANT genes. We found that the euAP2 lineage has
- conserved microRNA172 binding sequences except for a few genes in outgroups, one
- gene in *Eragrostis tef* and one in *Zea mays* (**Figure 4**, **Figure S2**).
- To facilitate the analysis, we name each subclade either by a previously assigned gene
- name within the subclade, or the gene sub-family name with a specific number.
- 432 Streptochaeta orthologs are present in most of the subclades, except IDS1/Q, ANT5,
- 433 BBM4, WRI3 and basalANT1, in which the Streptochaeta copy is lost (Figure 4, Figure
- 434 **S2**). The two most common patterns within each subclade are  $(O_1(S,G))$   $(O_1,O_2(S,G))$
- 435 S, Streptochaeta; G, other grasses) including SHAT1, ANT1, ANT3, ANT4, BBM1,
- 436 ANT7, ANT8 and ANT9, and (S,G) (inferring that outgroup sequence is lost or was not
- retrieved by our search) including *BBM3*, *WRI2* and *WRI4* (**Table S7**). These patterns
- imply that most grass-duplicated AP2-like genes were lost (i.e., the individual subclades
- were returned to single copy) soon after the grass duplication. Some subclades contain
- 440 two Streptochaeta sequences and one copy in other grasses. These Streptochaeta
- 441 sequences are either sisters to each other with the Streptochaeta clade sister to the
- other grasses (O<sub>1</sub>((S1,S2),G)) (RSR1) (**Figure 4**, **Figure S2**, **Table S7**), or successive
- sisters to a clade of grass sequences (O<sub>1</sub>(S1,(S2,G))) (WRI1) (Figure 4, Figure S2,
- 444 **Table S7**).
- In the paired subclades of IDS1/Q-SNB/SID1, ANT5-ANT6, BBM4-BBM2 and
- basalANT1-basalANT2, the grass-duplicated gene pairs were retained, and were also
- found to be syntenic pairs based on a syntelog search of the *Brachypodium distachyon*,
- 448 Oryza sativa or Setaria viridis genomes (**Figure 5**). Interestingly, in these subclade
- pairs, the *Streptochaeta* orthologs are always sister to one member of the syntenic gene
- pair but not the other. Two subclade pairs support a  $\rho$  position before the divergence of
- Streptochaeta, including BBM4-BBM2 with a pattern of (G1.(S.G2)) (Figure 5B) and
- 452 ANT5-ANT6 with a pattern of (G1,((S1,S2),G2)) (Figure 5E). In subclade pairs of
- 453 IDS1/Q-SNB/SID1 and basalANT1-basalANT2, two Streptochaeta sequences are
- successive sisters to one of the grass subclade pairs, forming tree topologies of
- 455 (G1,(S1,(S2,G2))) and (O,(G1,(S1,(S2,G2)))), respectively (**Figure 4**, **Figure S2**, **Table**
- 456 **S7**). These two cases do not fit with a simple history involving  $\rho$  either before or after
- 457 the divergence of *Streptochaeta*, and thus indicate a more complex evolutionary history.

R2R3 MYB

- The maximum likelihood tree of *R2R3 MYBs* was rooted with the CDC5 clade (Jiang and Rao, 2020). Only subclades with bootstrap values larger than 80 at the node of *Streptochaeta* were considered for subsequent analysis. Similar to the *AP2*-like tree, the
- most common tree topology within each subclade is (O,(S,G)), found in 16 individual subclades, followed by (S,G), consisting of 10 subclades. We also found 16 subclades
- with other tree topologies either without or with one or two *Streptochaeta* sequences
- and one copy of the other grass sequences, including (O,G) (MYB48), (O,((S1,S2),G))
- 466 (MYB17, MYB21, GAMYBL2, MYB29 and GAMYBL1), ((S1,S2),G) (MYB78 and
- 467 MYB92), (O,(S1,(S2, G))), (S1,(S2,G)) (MYB56) and ((O,S),G) (MYB47 and MYB83)
- 468 (**Table S7**). Conversely, we also found that 20 subclade pairs retained the grass
- duplicated gene pairs, although their tree topologies vary based on the position of
- 470 Streptochaeta and outgroups. Among these, 15 subclade pairs are also found to be
- 471 syntenic, including MYB1-MYB2, MYB6-MYB7, MYB35-MYB36, MYB42-MYB43,
- 472 *MYB*49-*MYB*50, *MYB*51-*MYB*52, *MYB*53-*MYB*54, *MYB*62-*MYB*63, *MYB*65-*MYB*66,
- 473 SWAM1-SWAM2, MYB75-MYB76, MYB86-MYB87, MYB93-MYB94, MYB103-MYB104
- 474 and MYB105-FDL1 (Figure 5 and Figure 6, Figure S3, Table S7). Together, these
- 475 results indicate that a subset of grass MYB clades have expanded due to the grass
- 476 WGD.
- 477 Among the above subclade pairs that retain both grass sequences, we found that one
- subclade pair, MYB53-MYB54 with tree topology of (O,(S1,S2),(G1,G2)), supports ρ
- having occurred after the divergence of *Streptochaeta* (**Figure 5F**). Conversely, we
- found 10 subclades supporting a  $\rho$  position before the divergence of *Streptochaeta*. The
- subclade MYB93-MYB94 includes three Streptochaeta sequences, one sister to one of
- 482 the grass clades and the other two sister to each other and sister to the other grass
- clade, forming a tree topology of (O<sub>1</sub>((S1,G1),((S2,S3),G2))) (**Figure 5A**). In the other 9
- subclade pairs, one or two Streptochaeta sequences are sister to one of the grass
- syntenic gene pairs but not the other (**Figure 5B-5E**). In subclade pairs *MYB86-MYB87*
- and MYB34-MYB36, one Streptochaeta sequence is sister to one of the grass clades,
- showing (G1,(S,G2)) and (O,(G1,(S,G2))), respectively (**Figure 5B and 5C**). We
- observed more subclades with two sequences of Streptochaeta, either showing
- 489 (O,(G1,((S1,S2),G2))) in *MYB6-MYB7* and *SWAM1* and *SWAM2*, or (G1,((S1,S2),G2))
- 490 in MYB42-MYB43, MYB51-MYB52, MYB65-MYB66, MYB75-MYB76 and MYB105-
- 491 *FDL1*.
- A few subclade pairs have tree topologies that do not support a  $\rho$  position either before
- or after the divergence of Streptochaeta, including (O,(S1,(S2,(G1,G2)))) (MYB1-MYB2
- 494 and MYB62-MYB63), (S1,(G1,(S2,G2))) (MYB22-MYB23) and ((O,S),(G1,G2)) (MYB11-
- 495 MYB12) (**Table S7**). In other cases, the Streptochaeta ortholog is either lost, or
- 496 positioned within the grass clades (**Table S7**). This may indicate a complex evolutionary
- 497 history of *Streptochaeta*. Alternatively, it may be an artifact due to the distant outgroups
- used in this study and poor annotation of some sequences.

- 499 Taken together, both the AP2-like and R2R3 MYB trees support the inference of ρ
- before the divergence of *Streptochaeta* (12 subclades) over  $\rho$  after the divergence of
- 501 Streptochaeta (1 subclade) (Figure 5), consistent with previous findings (McKain et al.,
- 502 2016). In addition, our study suggests that Streptochaeta has often lost one of the
- syntenic paralogs and sometimes has its own duplicated gene pairs.

## Annotation of miRNAs and validation of their targets

- sRNAs are important transcriptional and post-transcriptional regulators that play a role
- in plant development, reproduction, stress tolerance, etc. Identification of the
- 507 complement of these molecules in Streptochaeta can inform our understanding of
- distinguishing features of grass and monocot genomes. To annotate miRNAs present
- in the Streptochaeta genome, we (i) sequenced sRNAs from leaf, anther and pistil
- 510 tissues, (ii) compared miRNAs present in anthers to those of three other representative
- monocots (rice, maize and asparagus), and (iii) validated gene targets of these
- 512 miRNAs. In total, 185.3 million (M) sRNA reads were generated (115.6 M, 33.0 M, and
- 36.7 M reads for anther, pistil, and leaf tissues, respectively) from five sRNA libraries.
- Overall, we annotated 114 miRNA loci, of which 98 were homologous to 32 known
- miRNA families and 16 met strict annotation criteria for novel miRNAs (**Table S8**; **Table**
- 516 **S9**; **Table S10**). Most miRNAs from these loci (85; 90.4%) accumulated in all three
- tissues (**Figure 7**). We found a sub-group (8 miRNAs; 7.0%) of miRNAs abundant in
- anthers but not in the pistil or leaf tissues. Among these miRNAs, we found one copy
- each of miR2118 and miR2275, miRNAs known to function in the biogenesis of
- reproductive phasiRNAs (Johnson et al., 2009; Zhai et al., 2015). Comparing known
- 521 miRNA families expressed in anthers of Streptochaeta with three other monocots, we
- observed that only 25.4% of families overlapped between species. The large number of
- 523 miRNA families detected exclusively in anthers of asparagus (29.9%) and rice (17.9%)
- 524 perhaps explains the small overlap between species.
- We generated parallel analysis of RNA ends (PARE) libraries to identify and validate the
- 526 cleavage of miRNA-target pairs in anther, pistil and leaf of *Streptochaeta* tissues (**Table**
- 11; **Table S12**). Overall, we validated 58, 55 and 66 gene targets in anther, pistil and
- leaf of Streptochaeta tissues, respectively. Half of these targets were detected in all
- 529 tissues (51.9%) while 7 (8.6%), 4 (4.9%) and 14 (17.3%) targets were validated
- exclusively in anther, pistil, and leaf tissues, respectively, and remaining set of targets
- were found in combinations of two tissues. Among the validated targets, we found
- targets for three novel miRNAs, supporting their annotation. As an example, 184 reads
- validated the cleavage site of one novel miRNA target gene (strangu 031733), which is
- homologous to the *GPX6* gene (At4g11600) known to function in the protection of cells
- from oxidative damage in Arabidopsis (Rodriguez Milla et al., 2003). Among targets of
- ion extractive damage in reading to the configuration of any property of
- known miRNAs, we validated the cleavage site of 6 and 4 genes encoding members of
- AP2 and MYB transcription factor families, respectively (Figure S2; Figure S3). We
- observed that miR172 triggered the cleavage of AP2 genes in all tissues, consistent
- with the well-described function of this miRNA (Aukerman and Sakai, 2003; Lauter et
- al., 2005; Chuck et al., 2007, 2008). We also showed that miR159 triggered the

- cleavage of transcripts of four MYB genes, homologous to rice GAMYB genes, in leaf
- and pistil tissues but not in anther.

## **Expression of phasiRNAs is not limited to male reproductive**

- 544 **tissues**
- We used the same sRNA libraries and annotated phasiRNAs expressed in the
- 546 Streptochaeta genome, and compared the abundances of these loci to asparagus.
- maize, and rice. Overall, we detected a total of 89 phasiRNA loci (called *PHAS* loci)
- including 71 21-PHAS and 18 24-PHAS loci (**Table S8**). We made three observations of
- note: First, we observed a switch in the ratio of 21-PHAS to 24-PHAS locus number
- comparing asparagus (< 1), a member of Asparagaceae, to grass species (> 1;
- Poaceae). Second, the number of genomic *PHAS* loci increased, in Poaceae species,
- from *Streptochaeta* to both maize and rice. Third, several *PHAS* loci were also
- expressed in the pistil and leaf tissues -- female reproductive and vegetative tissues,
- respectively. Overall, a total of 23 (32%) 21-PHAS loci and 11 (61%) 24-PHAS loci were
- expressed in the pistil with a median abundance of 32.9% and 12.3% respectively
- compared to phasiRNAs detected in anther tissue. Similarly, 22 (31%) 21-PHAS loci
- and 10 (56%) 24-PHAS loci were detected in leaf tissue with a median abundance of
- 558 53.3% and 13.2% respectively compared to phasiRNAs detected in anthers. This
- expression of 24-nt phasiRNAs in vegetative tissues is unusual.

## **Discussion**

560

561

## Genome assembly, contiguity, structure.

- The Streptochaeta genome presented here provides a resource for comparative
- 563 genomics, genetics, and phylogenetics of the grass family. It represents the subfamily
- Anomochlooideae, which is sister to all other grasses and thus is equally
- 565 phylogenetically distant to the better-known species rice, Brachypodium, sorghum, and
- maize (Clark et al., 1995; Grass Phylogeny Working Group et al., 2001; Saarela et al.,
- 567 2018). The genome assembly captures nearly all of the predicted gene space at high
- contiguity (complete BUSCOs 91.8%, liliopsida odb10 profile, n = 3278), with the
- genome size matching predictions based on flow cytometry. The genome-wide LTR
- Assembly Index (LAI), for measuring the completeness of intact LTR elements, was
- 571 9.02. This score classifies the current genome as "draft" in quality, and is on par with
- other assemblies using similar sequencing technology (Apple (v1.0) (Velasco et al.,
- 573 2010), Cacao (v1.0) (Argout et al., 2011)).
- Our comprehensive annotation strategy identified a high proportion of genes specific to
- 575 the genus *Streptochaeta*, also known as orphan genes (3,742). Many previous studies
- 576 have indicated that orphan genes may comprise 3-10% of the total genes in plants and
- 577 can, in certain species, range up to 30% of the total (Arendsee et al., 2014). Overall the
- average gene length (3,956bp), average mRNA length (3,931bp) and average CDS
- length (1,060bp) are similar to other grass species gueried in Ensembl (Howe et al.,
- 580 2021).

- Previous phylogenetic work based on transcriptomes (McKain et al., 2016) or individual
- gene tree analyses (Preston and Kellogg, 2006; Whipple et al., 2007; Christensen and
- 583 Malcomber, 2012; McKain et al., 2016)) suggested that *Streptochaeta* shared the same
- WGD (p) as the rest of the grasses but that it might also have its own duplication.
- 585 Among the large sample (200) of clades in the transcriptome gene trees from McKain et
- al. (2016), 44% of these showed topologies consistent with ρ before the divergence of
- 587 Streptochaeta (e.g., topologies shown in Figure 2 Ai, Aii, and Aiv), with 39% being
- ambiguous (Figure 2 Aiii, Bii). Fewer than 20% of the clades identified by (McKain et
- al., 2016) had topologies consistent with the ρ duplication occurring after the divergence
- of Streptochaeta (Figure 2 Bi).
- 591 Streptochaeta contigs show good collinearity with the rice genome, a finding that is also
- 592 consistent with the hypothesis that ρ preceded the divergence of *Streptochaeta* as
- suggested by most of our gene trees. Mapping the *Streptochaeta* contigs against
- themselves also hints at another *Streptochaeta*-specific duplication, although the timing
- of this duplication cannot be inferred purely from the dot plot. Analysis of individual
- clades within large gene families (see below) support the same conclusion.
- 597 Analyzing the AP2-like and MYB subclades through the lens of grass WGD events, we
- 598 found 12 and 1 cases supporting ρ before and after the divergence of *Streptochaeta*,
- thus confirming previous transcriptomic data (Preston and Kellogg, 2006; Whipple et al.,
- 2007; Christensen and Malcomber, 2012; McKain et al., 2016). We also found that
- Streptochaeta often lost one copy of the syntenic paralogs, not only in MADS-box genes
- 602 (Preston and Kellogg, 2006; Christensen and Malcomber, 2012) but also in AP2-like
- and R2R3 MYB families. In addition, there are often two Streptochaeta sequences sister
- to a grass clade (**Figure 5**, **Table S7**), underscoring the fact that *Streptochaeta* does
- not simply represent an ancestral state for polarization of grass evolution, but has its
- own unique evolutionary history.
- 607 Genome structure and phylogenetic trees of *Streptochaeta* genes and their orthologs
- support the "loss model" shown in **Figure 1B iv**, in which many of the genes known to
- 609 control the structure of the grass spikelet were found in an ancestor of both
- Streptochaeta and the spikelet clade, but have then been lost in Streptochaeta. This
- provides circumstantial evidence that the common ancestor of all grasses including
- Streptochaeta (and Anomochloa) might have borne its flowers in spikelets, and the
- truly peculiar "spikelet equivalents" of Anomochlooideae are indeed highly modified.
- 614 Complex evolutionary history of *Streptochaeta* may contribute to its
- 615 unique characteristics
- 616 Previous studies have focused on the evolution of MADS-box genes in shaping grass
- spikelet development. For example, the A-class gene in flower development
- 618 FRUITFULL (FUL) duplicated at the base of Poaceae before the divergence of
- 619 Streptochaeta, but FUL1/VRN1 in Streptochaeta was subsequently lost (Preston and
- Kellogg, 2006). Similarly, paralogous LEAFY HULL STERILE1 (LHS1) and Oryza sativa
- 621 MADS5 duplicated at the base of Poaceae, but Streptochaeta has only one gene sister

- to the *LHS1* clade (Christensen and Malcomber, 2012). However, in another study on
- the B-class MADS-box gene PISTILLATA (PI), Streptochaeta has orthologs in both the
- 624 *PI1* and *PI2* clades (Whipple et al., 2007).
- Here we focused on AP2-like and R2R3 MYB transcription factor families, both of which
- 626 include members regulating inflorescence and spikelet development. The *euAP2*
- 627 lineage of the AP2-like genes determines the transition from spikelet meristem to floral
- meristem (Hirano et al., 2014). In the maize mutant *indeterminate spikelet1* (ids1), extra
- florets are formed within the spikelets in both male and female flowers (Chuck et al.,
- 630 1998). The double mutant of *ids1* and its syntenic paralog *sister of indeterminate*
- spikelet1 (sid1) produce repetitive glumes (Chuck et al., 2008). Consistently, the rice
- 632 mutants of SUPERNUMERARY BRACT (SNB), which is an ortholog of SID1, also
- exhibit multiple rudimentary glumes, due to the delay of transition from spikelet
- 634 meristem to floral meristem. Such mutant phenotypes are somewhat analogous to the
- 635 Streptochaeta "spikelet equivalents", which possess 11 or 12 bracts. In situ
- 636 hybridization studies on *FUL* and *LHS1* showed that the outer bracts 1-5 resemble the
- expression pattern of glumes in other grass spikelets, while inner bracts 6-8 resemble
- the expression pattern of lemma and palea (Preston et al., 2009). Our phylogenetic
- analysis suggests that the ortholog of *IDS1* in *Streptochaeta* is lost (**Figure 4**, **Figure**
- 640 **S2**). Instead, *Streptochaeta* has two sequences orthologous to *SID1/SNB*, and these
- two sequences are successively sister to each other with a tree pattern of
- 642 (G1,(S1,(S2,G2)) in *IDS1/Q-SID1/SNB* subclade pairs, leaving the evolutionary history
- of Streptochaeta ambiguous (Figure 4, Figure S2, Table S7). Both IDS1 and SID1 are
- targets of miRNA172 in maize (Chuck et al., 2007, 2008). Our PARE analyses did
- validate the cleavage of all six *Streptochaeta euAP2* by miRNA172 (**Table S12**),
- demonstrating that the miRNA172 post-transcriptional regulation of *euAP2* is functional
- in Streptochaeta. Detailed spatial gene expression analysis may further reveal whether
- and how these *euAP2* genes contribute to floral structure in *Streptochaeta*.
- 649 BABY BOOM genes (BBMs) belong to the euANT lineage of the AP2-like genes, and
- are well known for their function in induction of somatic embryogenesis (Boutilier et al.,
- 2002) and application for in vitro tissue culture (Lowe et al., 2016). Ectopic expression
- of BBM in Arabidopsis and Brassica results in pleiotropic defects in plant development
- 653 including changes in floral morphology (Boutilier et al., 2002). The grasses have four
- annotated *BBMs*, although it is not known whether other *ANT* members share similar
- functions. BBM4 and BBM2 subclades appeared to be duplicated paralog pairs due to
- the grass WGD. Similar to the cases in previous studies (Preston and Kellogg, 2006;
- 657 Christensen and Malcomber, 2012), Streptochaeta has apparently lost its BBM4 copy
- and contains one copy in the *BBM2* subclade (**Figure 4**, **Figure 5**, and **Figure S2**).
- 659 R2R3 MYB is a large transcription factor family, some of which are crucial for anther
- development. The rice carbon starved anther (csa) mutants show decreased sugar
- content in floral organs including anthers, resulting in a male sterile phenotype (Zhang
- et al., 2010). DEFECTIVE in TAPETAL DEVELOPMENT and FUNCTION1 (TDF1) is
- required for tapetum programmed cell death (Zhu et al., 2008; Cai et al., 2015). GAMYB
- positively regulates GA signaling by directly binding to the promoter of GA-responsive
- genes in both *Arabidopsis* and grasses (Tsuji et al., 2006; Aya et al., 2009; Alonso-Peral

- 666 et al., 2010). OsGAMYB is highly expressed in stamen primordia, tapetum cells of the 667 anther and aleurone cells, and its expression is regulated by miR159. Nonfunctional mutants of OsGAMYB are defective in tapetum development and are male sterile 668 669 (Kaneko et al., 2004; Tsuji et al., 2006). We found conserved miRNA159 binding sites in 670 GAMYBs and its closely related subclades, including MYB27, MYB28, GAMYBL2, 671 MYB29, GAMYBL1, MYB30 and GAMYB (Figure 4). Our PARE analyses also validated 672 the cleavage of Streptochaeta GAMYB and GAMYBL1 in leaf and pistil tissues but not 673 in anthers, suggesting the expression of *Streptochaeta GAMYB* and *GAMYBL1* may be 674 suppressed by miR159 in tissues other than anthers, at least at the developmental 675 stages we investigated (Table S12). Streptochaeta has two sequences in each of the GAMYBL2, MYB29, GAMYBL1 and GAMYB clades, either with a tree topology of 676
- 677 (O,(S1,S2),G) in GAMYBL2, MYB29 and GAMYBL1, or a tree topology of
- (O,(S1,(S2,G)) in GAMYB (Figure 6, Figure 4, Table S7). This again indicates that 678
- Streptochaeta has a complex duplication history. 679

## A survey of small RNAs in the Streptochaeta genome

- 681 miRNAs are major regulators of mRNA levels, active in pathways important to plant 682 developmental transitions, biotic and abiotic stresses, and others, miRNAs generally act 683 as post-transcriptional regulators by homology-dependent cleavage of target gene transcripts, when loaded to the RNA-induced silencing complex (RISC). Plant genomes 684 685 encode a variety of sRNAs that can act in a transcriptional or post-transcriptional 686 regulation mode. In this paper, we focused on miRNA and phasiRNA. The list of miRNA 687 annotated in this study is likely incomplete because the Streptochaeta sRNA-seg data 688 were limited to anther, pistil and leaf tissues, and would miss miRNAs expressed 689 specifically in other tissues/cell types or at growth conditions not sampled. Thus, 690 miRNAs missed in our data may well be encoded in the Streptochaeta genome. That
- 691 being said, our miRNA characterization provides a starting point with which to describe 692 Streptochaeta miRNAs, and our sequencing depth and tissue diversity was likely
- 693 sufficient to identify many if not the majority of miRNAs encoded in the genome.
- 694 Phased short interfering RNAs (phasiRNAs) are 21-nt or 24-nt sRNAs generated from 695 the recursive cleavage of a double-stranded RNA from a well-defined terminus; these 696 transcripts define their precursor PHAS loci (Axtell and Meyers, 2018). Reproductive 697 phasiRNAs are a subset abundant in anthers and in some cases essential to male 698 fertility. Genomes of grass species are particularly rich in reproductive PHAS loci (Patel 699 et al., 2021), expressed in anthers but not in female reproductive tissues or vegetative 700 tissues. Previous species studies identified hundreds of PHAS loci in anthers of maize 701 (Zhai et al., 2015) to thousands of PHAS loci in rice (Fei et al., 2016), barley (Bélanger 702 et al., 2020) and bread wheat (Bélanger et al., 2020; Zhang et al., 2020). Additionally, work in maize (Teng et al., 2020) and rice (Fan et al., 2016) showed that 21-nt and 24-
- 703
- 704 nt phasiRNAs are essential to ensure proper development of meiocytes and to
- 705 guarantee male fertility under normal growth conditions. However, Streptochaeta has a
- 706 different internal anatomy than the rest of the grasses. Specifically, anthers in
- 707 Streptochaeta are missing the "middle layer" between the endothecium and the tapetum
- (Sajo et al., 2009, 2012) such that the microsporangium has only three cell lavers. 708

- Given that most of our data (> 100 M reads) were collected from anthers, we have good
- resolution for annotation of phasiRNAs in this tissue. We characterized their
- absence/presence in the three-layer anthers of *Streptochaeta*. We annotated tens of
- 712 PHAS loci in Streptochaeta showing that anthers express phasiRNAs even in the
- absence of the middle layer. Likewise, in maize, Zhai et al. (2015) showed that the
- miRNA and phasiRNA precursors are dependent on the epidermis, endothecium, and
- tapetum, and the phasiRNAs accumulate in the tapetum and meiocytes, so the middle
- layer is apparently not involved. We observed a shift in the ratio of 21-PHAS to 24-
- 717 PHAS loci from asparagus (< 1), an Asparagaceae, to grass species (> 1), although the
- 718 implications of this shift are as yet unclear.
- We also observed that several 21-nt and 24-nt phasiRNAs accumulate in either pistil or
- leaf tissues, inconsistent with prior results. A small number of 21-nt *PHAS* loci are likely
- 721 trans-acting-siRNA-generating (TAS) loci, important in vegetative tissues, but typically
- there are only a few *TAS* loci per genome (Xia et al., 2017), not the 20 loci that we
- observed. Additionally, we found no previous reports of 24-nt phasiRNAs accumulating
- in vegetative tissues or female reproductive tissues.

## 725 Utility of Streptochaeta for understanding grass evolution and

## 726 **genetics**

- 727 The four species of Anomochlooideae are central to understanding the evolution of the
- grasses and the many traits that make them unique. We have highlighted the unusual
- floral and inflorescence morphology of *Streptochaeta* and have compared it to grass
- 530 spikelets, but Streptochaeta can also illuminate the evolution and genetic basis of other
- 731 important traits. It is common to compare traits between members of the BOP clade
- 732 (e.g. Oryza, Brachypodium, or Triticum) and the PACMAD clade (e.g. Zea, Sorghum,
- 733 Panicum, Eragrostis), but, because these comparisons involve two sister clades, it is
- impossible to determine whether the BOP or the PACMAD clade character state is
- ancestral. Streptochaeta functions as an outgroup in such comparisons and can help
- establish the direction of change. Here, we highlight just a few of the traits whose
- analysis may be helped in future studies by reference to Streptochaeta and its genome
- 738 sequence.
- 739 Drought intolerance, shade tolerance. The grasses, including not only
- Anomochlooideae, but also Pharoideae and Puelioideae, the three subfamilies that are
- successive sister groups of the rest of the family, appear to have originated in
- environments with low light and high humidity (Edwards and Smith, 2010; Gallaher et
- al., 2019). The shift from shady, moist habitats to open, dry habitats where most grass
- species are now found promises insights into photosynthesis and water use efficiency,
- among other physiological traits.
- 746 Streptochaeta, like other forest grasses, has broad, spreading leaf blades and a
- 747 pseudopetiole that results in higher leaf angle and increased light interception (Gallaher
- et al., 2019). Leaf angle is an important agronomic trait, with selection during modern
- 5749 breeding often favoring reduced leaf angle to maximize plant density and yield (Liu et
- al., 2019; Mantilla-Perez et al., 2020). A close examination of *Streptochaeta* may
- provide insight into how leaf angle is controlled in diverse grasses. Leaf width in maize

- 752 is controlled particularly by the WOX3-like homeodomain proteins NARROWSHEATH1
- 753 (*NS1*) and *NS2*, which function in cells at the margins of leaves (Scanlon et al., 1996;
- Conklin et al., 2020). Duplication patterns and expression of NS1 and NS2 genes in the
- 755 Streptochaeta genome could test whether the models developed for maize were
- present in the earliest of the grasses.
- Leaf anatomy. The grass outgroup *Joinvillea* develops colorless cells in the mesophyll
- 758 (Leandro et al., 2018). These appear to form from the same ground tissue that is
- responsible for the cavity-like "fusoid" cells in Anomochlooideae, Pharoideae, and
- Puelioideae as well as the bambusoid grasses. These cells, which appear to be a
- shared derived character for the grasses, form from the collapse of mesophyll cells and
- may play a role in the synthesis and storage of starch granules early in plant
- development (Leandro et al., 2018). While the genetic basis of leaf anatomy is, at the
- moment, poorly understood, Streptochaeta will be a useful system for understanding the
- development of fusoid cells in early diverging and other grasses.
- Grass leaves also contain silica bodies in the epidermis; the vacuoles of these cells are
- 767 filled with amorphous silica (SiO<sub>2</sub>). In *Streptochaeta* the silica bodies are a distinctive
- shape, being elongated transverse to the long axis of the blade (Judziewicz and
- 769 Soderstrom, 1989). The genetic basis of silica deposition has been studied in rice (Yu et
- al., 2020) and the availability of the Streptochaeta genome now permits examination of
- the evolution of these genes in the grasses.
- 772 Anther and pollen development. Streptochaeta differs from most other grasses (and
- indeed some Poales as well) in details of its anthers and pollen development, and the
- current genome provides tools for comparative analyses. The sRNAs described above
- are produced in the epidermis, endothecium and tapetum of most grasses and we
- presume they are also produced in those tissues in *Streptochaeta*. In all grasses except
- 777 Anomochlooideae and Pharoideae, the microsporangium has four concentric layers of
- cells the epidermis, the endothecium, the middle layer, and the tapetum which
- surround the archesporial cells (Walbot and Egger, 2016). Cells in the middle layer and
- the tapetum are sisters, derived from division of a secondary parietal cell. The inner
- walls of the endothecial cells also mature to become fibrous (Artschwager and McGuire.
- 782 1949; Furness and Rudall, 1998). In *Streptochaeta* and *Pharus*, however, the middle
- layer is absent (Sajo et al., 2007, 2009, 2012) and the endothecial cells lack fibrous
- thickenings. It is tempting to speculate that the middle layer may have a role in
- coordinating maturation of the endothecium. Lack of the middle layer is apparently
- derived within Streptochaeta and Pharus. In known mutants of maize and rice, loss of
- the middle layer leads to male sterility (Walbot and Egger, 2016) so the functional
- 788 implications of its absence in *Streptochaeta* are unclear.
- 789 Development of microsporangium layers may also be related to the position of
- microspores inside the locule. In most grasses, the microspores and mature pollen
- grains form a single layer adjacent to the tapetum, with the pore of the pollen grain
- facing the tapetum, unlike many non-grasses in which the microsporocytes fill the locule
- and have a haphazard arrangement. The condition in *Streptochaeta* is unclear, with
- contradictory reports in the literature (Kirpes et al., 1996; Sajo et al., 2009, 2012).

- The exine, or outer layer, of grass pollen is distinct from that of its close relatives due to
- the presence of channels that pass through the exine. While controls of this particular
- aspect of the pollen wall are unknown in the grasses, we find that Streptochaeta and its
- 798 grass sisters have several GAMYB genes, which are known to be involved in exine
- formation in rice (Aya et al., 2009) and to have played a role more broadly in
- reproductive processes, including microspore development in early vascular plants (Aya
- 801 et al., 2011).
- 802 Chromosome number in the early grasses. Estimates of the ancestral grass
- chromosome number and karyotype have reached different conclusions (e.g., (Salse et
- al., 2008; Murat et al., 2010; Wang et al., 2016)). Genomes of Streptochaeta and other
- early diverging grasses will be useful for resolving this open question, but will require
- 806 psuedomolecule-quality assemblies. Two other species of Streptochaeta have been
- reported to have n=11 chromosomes (Valencia, 1962; Pohl and Davidse, 1971;
- Hunziker et al., 1982), well below the number reported for the sister species
- 809 Anomochloa marantoidea, n=18 (Judziewicz and Soderstrom, 1989). The outgroups
- 310 Joinvillea plicata and Ecdeiocolea monostachya have n=18 (Newell, 1969) and n=19
- 811 (Hanson et al., 2005), respectively. However, without high quality genomes and good
- 812 cytogenetic data for these species, the ancestral chromosome number and structure of
- the genomes of ancestral grasses remains a matter of speculation.
- Finally, these are but a few of the opportunities for understanding trait evolution in the
- grasses based on investigation of *Streptochaeta*, with additional insights possible in, for
- example, the study of embryo development, caryopsis modifications, endosperm/starch
- evolution and branching/tillering. We have demonstrated that genomes of targeted, non-
- 818 model species, particularly those that are sister to large, better-studied groups, can
- provide out-sized insight about the nature of evolutionary transitions and should be an
- increased focus now that genome assembly is a broadly accessible component of the
- 821 biologist's toolkit.

## Data Availability

- The sRNA-seg data were reported in a previous study (Patel et al., 2021). Also, one
- library of RNA-Seg (SRR3233339) used for annotation was previously published
- (Givnish et al., 2010). Otherwise, all data utilized in this study are original. The complete
- set of raw WGS, RNA-seq, sRNA-seq and PARE-seq reads were deposited in the
- 827 Sequence Read Archive under the BioProject ID PRJNA343128. Alignments and
- phylogenies for AP2-like and MYB R2R3 genes have been deposited at datadryad.org,
- accession #XXX (to be added after acceptance). The scripts and commands used for
- generating assembly, annotations, small RNA analyses and phylogenetic analyses are
- documented in the GitHub repository accessible here:
- 832 https://github.com/HuffordLab/streptochaeta

**Acknowledgments** 

833

847

852

855

- We thank Sandra Mathioni for construction of the RNA-seq and PARE libraries. Y.Y.
- was supported by NSF grant IOS-1938086 to E.A.K. and by an Enterprise-Rent-a-Car
- Foundation award through the Donald Danforth Plant Science Center, also to E.A.K.
- 837 S.B. was supported by USDA | National Institute of Food and Agriculture "BTT EAGER"
- award no. 2018–09058 to B.C.M., as well as resources from the Donald Danforth Plant
- 839 Science Center and the University of Missouri-Columbia. A.S. was supported by NSF
- grant IOS-1822330 to M.B.H. This work used 1) Extreme Science and Engineering
- Discovery Environment (XSEDE)(National Science Foundation Grant No. ACI-1548562)
- via Blacklight HPC environment allocation TG-MCB140103 and 2) HPC equipment at
- lowa State University, some of which has been purchased through funding provided by
- NSF under MRI grant number 1726447. We thank Dr. Philip Blood for his assistance
- with MaSuRCA optimization, which was made possible through the XSEDE Extended
- 846 Collaborative Support Service (ECSS) program.

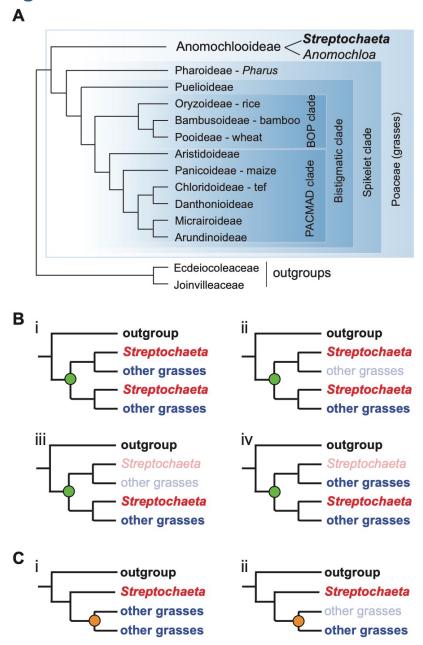
#### **Author contributions statement**

- M.B.H., A.S., E.A.K, and L.G.C. designed the project. L.G.C. and E.A.K. provided plant
- material. M.B.H. and A.S. generated sequence data and assembled the genome. S.B.
- and B.C.M. analyzed data on small RNAs. Y.Y. analyzed AP2 and MYB sequence data.
- All authors drafted and edited the manuscript, and produced figures and tables.

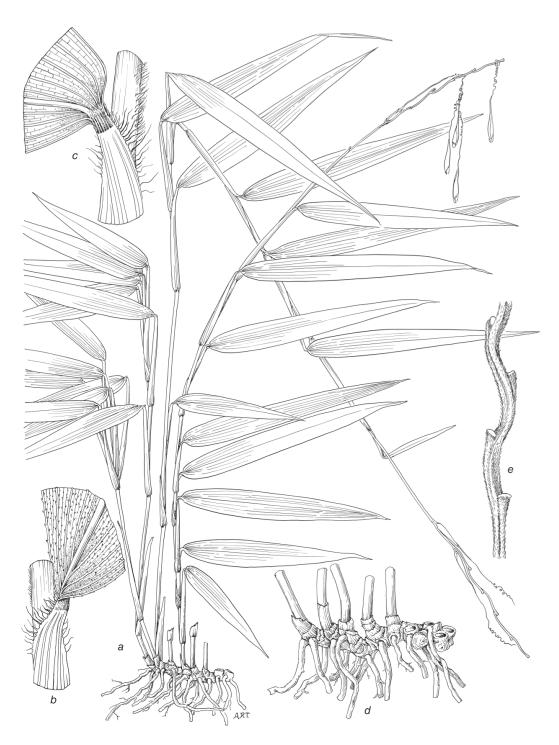
#### Conflict of Interest

- The authors declare that the research was conducted in the absence of any commercial
- or financial relationships that could be construed as a potential conflict of interest.

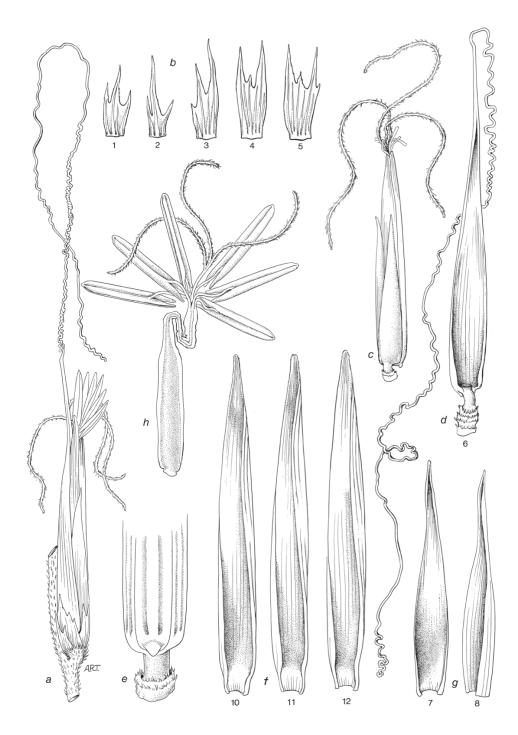
## **Figures**



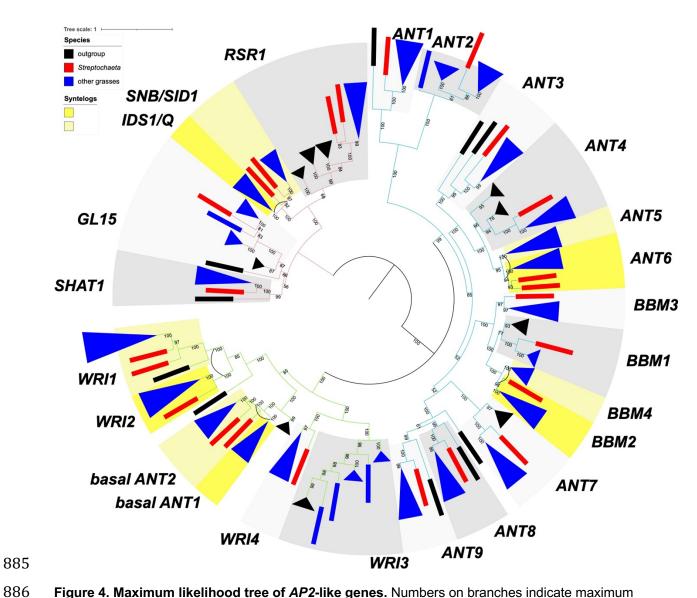
**Figure 1.** The phylogenetic placement of *Streptochaeta*. (A) Phylogenetic tree depicting the BOP (Bambusoideae, Oryzoideae, Pooideae) + PACMAD (Panicoideae, Aristidoideae, Chloridoideae, Micrairoideae, Arundinoideae, Danthonioideae) clade and the basal placement of focal organism Streptochaeta. (B) and (C) Possible patterns of whole genome duplication (WGD) and gene loss. (B) WGD before the divergence of *Streptochaeta* assuming (i) no gene loss; (ii) loss of one clade of non-*Streptochaeta* paralogs soon after WGD; (iii) loss of all grass paralogs soon after WGD; (iv) loss of one *Streptochaeta* paralog soon after WGD. (C) WGD after divergence of *Streptochaeta*. (i) no gene loss; (ii) loss of one clade of non-*Streptochaeta* grass paralogs soon after WGD. Note that patterns (Biii) and (Cii) are indistinguishable.



**Figure 2.** *Streptochaeta angustifolia.* **(A)** Habit (× 0.5). **(B)** Mid-region of leaf showing summit of sheath and upper surface of blade (× 4.5). **(C)** Mid-region of leaf showing summit of sheath and lower surface of blade (× 5). **(D)** Rhizome system with culm base (× 1). **(E)** Portion of rachis enlarged (× 1.5) All drawings based on *Soderstrom & Sucre 1969* (US). Illustration by Alice R. Tangerini. Reprinted from Soderstrom (1981, Some evolutionary trends in the Bambusoideae (Poaceae), *Annals of the Missouri Botanical Garden 68*: 15-47, originally Figure 5, p. 31), with permission from the Missouri Botanical Garden Press.



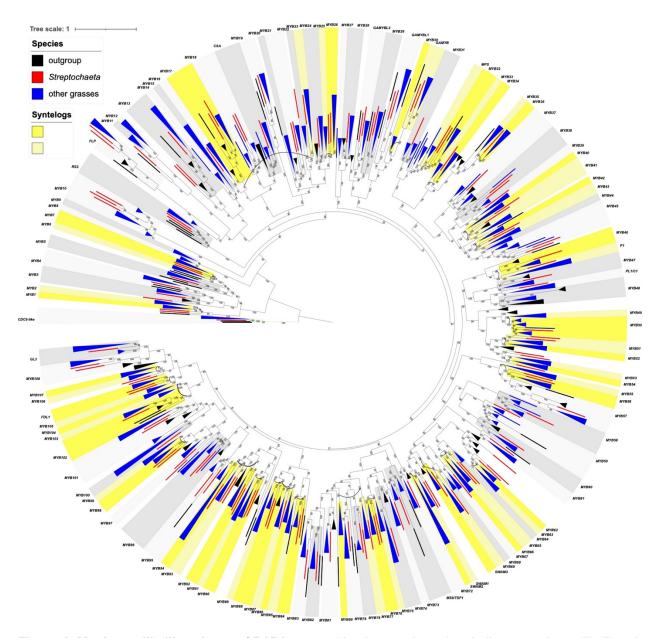
**Figure 3.** *Streptochaeta angustifolia.* **(A)** Pseudospikelet (× 4.5). **(B)** Series of bracts 1-5 from the base of the pseudospikelet (× 6). **(C)** Pseudospikelet with basal bracts 1-5 removed and showing bracts 7 and 8, whose bases are overlapping (× 4.5). **(D)** Bract 6 with long coiled awn (× 4.5). **(E)** Back portion of the base of bract 6 showing region where embryo exits at germination. **(F)** Bracts 10-12 (× 6). **(G)** Bracts 7 and 8 (× 6). Bract 9, which exists in other species, has not been found here. **(H)** Ovary with long style and three stigmas, surrounded by the thin, fused filaments of the 6 stamens (× 4.5). All drawings based on *Soderstrom & Sucre 1969* (US). Illustration by Alice R. Tangerini. Reprinted from Soderstrom (1981, Some evolutionary trends in the Bambusoideae (Poaceae), *Annals of the Missouri Botanical Garden* 68: 15-47, originally Figure 6, p. 33), with permission from the Missouri Botanical Garden Press.



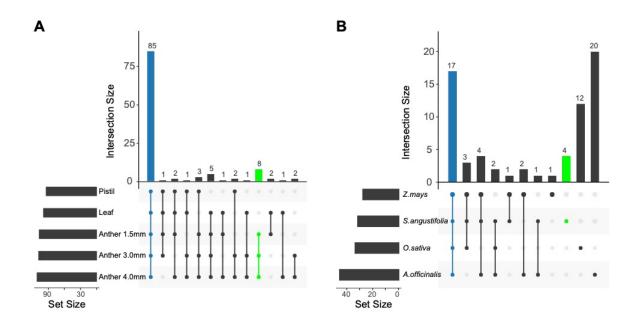
**Figure 4. Maximum likelihood tree of** *AP2***-like genes.** Numbers on branches indicate maximum likelihood bootstrap values. A single gene is denoted by a rectangle, and collapsed branches are denoted by triangles. Each subclade is shaded in two grey colors and named either by known genes within the subclade or subfamily name with a number. Subclades with syntenic genes in *Brachypodium*, *Oryza* or *Setaria* are shaded in two colors of yellow, and syntenic pairs are connected by an arc. Outgroup, *Streptochaeta* and other grasses are shown in black, red and blue colors.



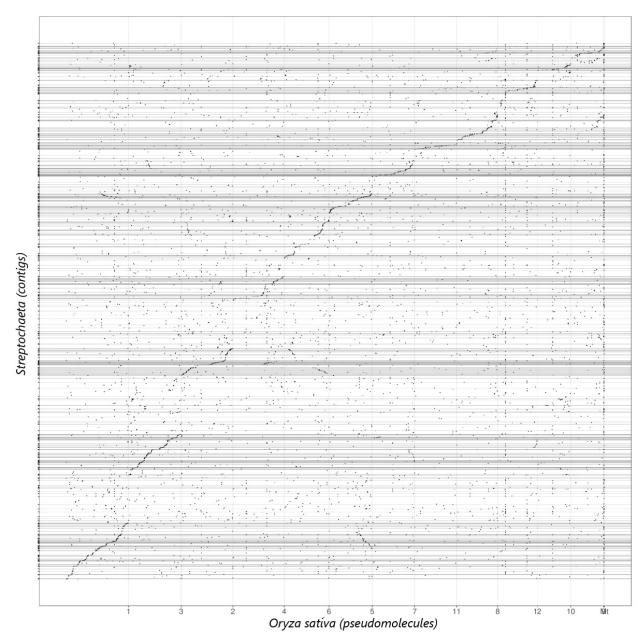
Figure 5. Tree topologies of paired syntenic subclades that support grass whole genome duplication (WGD) before or after the divergence of *Streptochaeta*. (A-E) Grass WGD before the divergence of *Streptochaeta*. Tree topologies: (A) (O,(S1,G1),((S2,S3),G2)). (B) (G1,(S2,G2)). (C) (O,(G1,(S2,G2))). (D) (O,(G1,((S1,S2),G2))). (E) (G1,((S1,S2),G2)). (F) Grass WGD after the divergence of *Streptochaeta* with tree pattern of (O,(S1,S2),(G1,G2)).



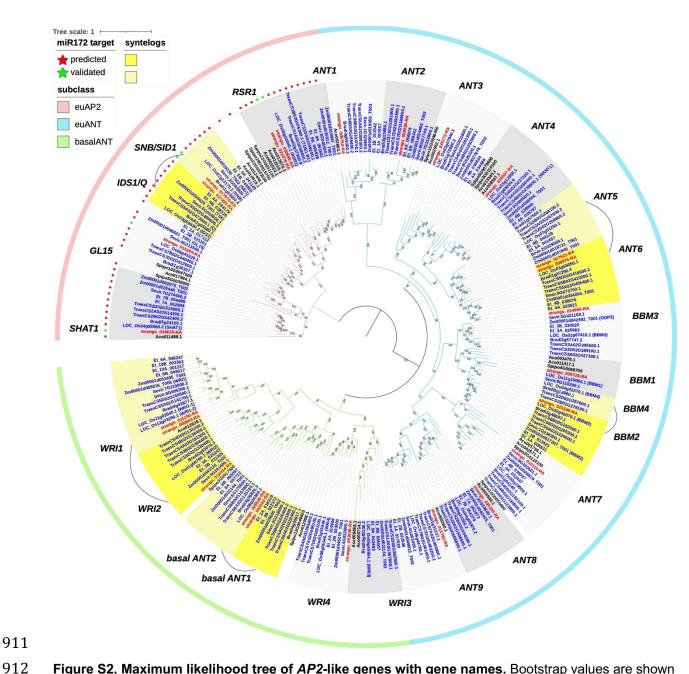
**Figure 6. Maximum likelihood tree of R2R3 genes.** Numbers on branches indicate maximum likelihood bootstrap values. A single gene is denoted by a rectangle, and collapsed branches are denoted in triangles. Each subclade is shaded in two grey colors and named either by known genes within the subclade or subfamily name with a number. Subclades with syntenic genes in *Brachypodium*, *Oryza* or *Setaria* are shaded in two colors of yellow, and syntenic pairs are connected by an arc. Outgroup, *Streptochaeta* and other grasses are shown in black, red and blue colors.



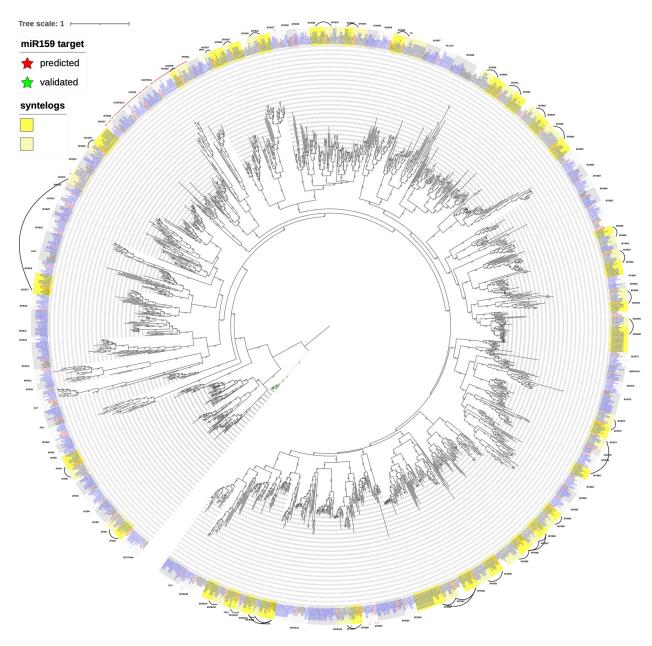
**Figure 7:** Overlap of miRNA loci annotated in *Streptochaeta* tissues **(A)** and miRNA families annotated in *Streptochaeta* anthers compared to three other monocots **(B)**.



**Figure S1. Dot plot of** *Streptochaeta* **versus** *Oryza sativa*. The contigs of the draft *Streptochaeta* assembly plotted against *Oryza sativa* (Nipponbare; (Ouyang et al., 2007)) pseudomolecules.



**Figure S2. Maximum likelihood tree of** *AP2***-like genes with gene names.** Bootstrap values are shown on the branches. Each subclade is shaded in two grey colors and named either by known genes within the subclade or subfamily name with a number. Subclades with syntenic genes in *Brachypodium*, *Oryza* or *Setaria* are shaded in two colors of yellow, and syntenic pairs are connected by an arc. Predicted and experimentally validated miR172 binding sites are denoted by red and green stars, respectively.



**Figure S3. Maximum likelihood tree of** *R2R3* **genes with gene names.** Bootstrap values are shown on the branches. Each subclade is shaded in two grey colors and named either by known genes within the subclade or subfamily name with a number. Subclades with syntenic genes in *Brachypodium*, *Oryza* or *Setaria* are shaded in two colors of yellow, and syntenic pairs are connected by an arc. Predicted and experimental validated miR159 binding sites are denoted by red and green stars, respectively.

## Tables (see supplemental excel file)

Supplemental Table 1: Short reads (raw data) used for the assembly and their estimated coverage based on a genome size of 1.8 Gbp

Supplemental Table 2: Criteria for merging ab initio gene models with the direct evidence models. The codes are as described in the Mikado compare manual. For the

928 comaprision, BRAKER gene models were used as prediction and evidence models 929 were used as reference. 930 Supplemental Table 3: Source of genome, annotation version, and sRNA-seq data used 931 in this study 932 Supplemental Table 4: 5' and 3' adapters used to construct RNA-seq libraries. 933 Supplemental Table 5: Summary statistics of the genome assembly after each iteration 934 of Redundans. 935 Supplemental Table 6: Phylostrata distribution of the genes predicted by BIND strategy 936 Supplemental Table 7: Tree topologies of the subclades in the AP2-like and R2R3 MYB 937 trees. O: outgroup; S: Streptochaeta; G: grasses other than Streptochaeta. If 938 Streptochaeta and/or outgroup genes are inside of a grass clade, it is labeled as S-G or 939 O-G. 940 Supplemental Table 8: Summary of miRNA and phasiRNA annotated in anthers of 941 Streptocheata angustifolia and other monocots. 942 Supplemental Table 9: Coordinates and abundance of the 114 annotated miRNAs in 943 Streptochaeta angustifolia. 944 Supplemental Table 10: Candidate novel miRNA annotated in Streptochaeta. This table 945 details the sequence and abundance of each mature miRNA and miRNA-star plus the 946 sequence of the locus and the predicted RNA secondary structure in dot-bracket 947 notation. 948 Supplemental Table 11: Summary of miRNA targets validated via PARE-Seq. The 949 described miRNAs were captured in Streptochaeta angustifolia anthers. 950 Supplemental Table 12: Details of PARE-validated miRNA cleavage sites detected in 951 anther, pistil and leaf tissues in Streptochaeta.

References

952

- Alonso-Peral, M. M., Li, J., Li, Y., Allen, R. S., Schnippenkoetter, W., Ohms, S., et al.
- 954 (2010). The microRNA159-regulated GAMYB-like genes inhibit growth and promote
- programmed cell death in Arabidopsis. *Plant Physiol.* 154, 757–771.
- 956 Arendsee, Z., Li, J., Singh, U., Seetharam, A., Dorman, K., and Wurtele, E. S. (2019).
- 957 phylostratr: a framework for phylostratigraphy. *Bioinformatics* 35, 3617–3627.
- Arendsee, Z. W., Li, L., and Wurtele, E. S. (2014). Coming of age: orphan genes in plants. *Trends Plant Sci.* 19, 698–708.
- Argout, X., Salse, J., Aury, J.-M., Guiltinan, M. J., Droc, G., Gouzy, J., et al. (2011). The genome of Theobroma cacao. *Nat. Genet.* 43, 101–108.
- Artschwager, E., and McGuire, R. C. (1949). Cytology of reproduction in Sorghum vulgare. *J. Agric. Res.* 78, 659–673.
- Aukerman, M. J., and Sakai, H. (2003). Regulation of flowering time and floral organ identity by a MicroRNA and its APETALA2-like target genes. *Plant Cell* 15, 2730–

966 2741.

- 967 Axtell, M. J., and Meyers, B. C. (2018). Revisiting Criteria for Plant MicroRNA 968 Annotation in the Era of Big Data. *Plant Cell* 30, 272–284.
- Aya, K., Hiwatashi, Y., Kojima, M., Sakakibara, H., Ueguchi-Tanaka, M., Hasebe, M., et al. (2011). The Gibberellin perception system evolved to regulate a pre-existing
- 971 GAMYB-mediated system during land plant evolution. *Nat. Commun.* 2, 544.
- 972 Aya, K., Ueguchi-Tanaka, M., Kondo, M., Hamada, K., Yano, K., Nishimura, M., et al.
- 973 (2009). Gibberellin modulates anther development in rice via the transcriptional
- 974 regulation of GAMYB. *Plant Cell* 21, 1453–1472.
- 975 Bartlett, M., Thompson, B., Brabazon, H., Del Gizzi, R., Zhang, T., and Whipple, C.
- 976 (2016). Evolutionary Dynamics of Floral Homeotic Transcription Factor Protein-
- 977 Protein Interactions. *Mol. Biol. Evol.* 33, 1486–1501.
- 978 Bélanger, S., Pokhrel, S., Czymmek, K., and Meyers, B. C. (2020). Premeiotic, 24-
- Nucleotide Reproductive PhasiRNAs Are Abundant in Anthers of Wheat and Barley
- 980 but Not Rice and Maize. *Plant Physiol.* 184, 1407–1423.
- 981 Boetzer, M., Henkel, C. V., Jansen, H. J., Butler, D., and Pirovano, W. (2011).
- 982 Scaffolding pre-assembled contigs using SSPACE. *Bioinformatics* 27, 578–579.
- 983 Bonnet, E., He, Y., Billiau, K., and Van de Peer, Y. (2010). TAPIR, a web server for the
- prediction of plant microRNA targets, including target mimics. *Bioinformatics* 26,
- 985 1566–1568.

- 986 Boutilier, K., Offringa, R., Sharma, V. K., Kieft, H., Ouellet, T., Zhang, L., et al. (2002).
- 987 Ectopic expression of BABY BOOM triggers a conversion from vegetative to
- 988 embryonic growth. *Plant Cell* 14, 1737–1749.
- 989 Cai, C.-F., Zhu, J., Lou, Y., Guo, Z.-L., Xiong, S.-X., Wang, K., et al. (2015). The
- functional analysis of OsTDF1 reveals a conserved genetic pathway for tapetal
- development between rice and Arabidopsis. Sci Bull. Fac. Agric. Kyushu Univ. 60,
- 992 1073–1082.
- 993 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., and Bealer, K.
- 994 (2009). BLAST plus: architecture and applications. BMC Bioinformatics. *BioMed*
- 995 *Central* 10, 1.
- 996 Capella-Gutiérrez, S., Silla-Martínez, J. M., and Gabaldón, T. (2009). trimAl: a tool for
- automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*
- 998 25, 1972–1973.
- 999 Christensen, A. R., and Malcomber, S. T. (2012). Duplication and diversification of the
- 1000 LEAFY HULL STERILE1 and Oryza sativa MADS5 SEPALLATA lineages in
- graminoid Poales. *Evodevo* 3, 4.
- 1002 Chuck, G., Meeley, R. B., and Hake, S. (1998). The control of maize spikelet meristem
- fate by the APETALA2-like gene indeterminate spikelet1. *Genes Dev.* 12, 1145–
- 1004 1154.
- 1005 Chuck, G., Meeley, R., and Hake, S. (2008). Floral meristem initiation and meristem cell
- fate are regulated by the maize AP2 genes ids1 and sid1. Development 135, 3013–
- 1007 3019.
- 1008 Chuck, G., Meeley, R., Irish, E., Sakai, H., and Hake, S. (2007). The maize tasselseed4
- microRNA controls sex determination and meristem cell fate by targeting
- Tasselseed6/indeterminate spikelet1. *Nat. Genet.* 39, 1517–1521.
- 1011 Clark, L. G., Zhang, W., and Wendel, J. F. (1995). A Phylogeny of the Grass Family
- 1012 (Poaceae) Based on ndhF Sequence Data. Syst. Bot. 20, 436–460.
- 1013 Conklin, P. A., Johnston, R., Conlon, B. R., Shimizu, R., and Scanlon, M. J. (2020).
- Plant homeodomain proteins provide a mechanism for how leaves grow wide.
- Development 147.
- 1016 Conway, J. R., Lex, A., and Gehlenborg, N. (2017). UpSetR: an R package for the
- visualization of intersecting sets and their properties. *Bioinformatics* 33, 2938–
- 1018 2940.
- Debernardi, J. M., Greenwood, J. R., Jean Finnegan, E., Jernstedt, J., and Dubcovsky,
- J. (2020). APETALA 2-like genes AP2L2 and Q specify lemma identity and axillary
- floral meristem development in wheat. *Plant J.* 101, 171–187.

- Edgar, R. C. (2010). Search and clustering orders of magnitude faster than BLAST.
- 1023 Bioinformatics 26, 2460–2461.
- Edwards, E. J., and Smith, S. A. (2010). Phylogenetic analyses reveal the shady history
- of C4 grasses. *Proc. Natl. Acad. Sci. U. S. A.* 107, 2532–2537.
- Fahlgren, N., and Carrington, J. C. (2010). miRNA Target Prediction in Plants. *Methods*
- 1027 *Mol. Biol.* 592, 51–57.
- 1028 Fan, Y., Yang, J., Mathioni, S. M., Yu, J., Shen, J., Yang, X., et al. (2016). PMS1T,
- producing phased small-interfering RNAs, regulates photoperiod-sensitive male
- 1030 sterility in rice. *Proc. Natl. Acad. Sci. U. S. A.* 113, 15144–15149.
- Fei, Q., Yang, L., Liang, W., Zhang, D., and Meyers, B. C. (2016). Dynamic changes of
- small RNAs in rice spikelet development reveal specialized reproductive phasiRNA
- 1033 pathways. *J. Exp. Bot.* 67, 6037–6049.
- 1034 Furness, C. A., and Rudall, P. J. (1998). The tapetum and systematics in
- 1035 monocotyledons. *Bot. Rev.* 64, 201–239.
- Gallaher, T. J., Adams, D. C., Attigala, L., Burke, S. V., Craine, J. M., Duvall, M. R., et
- al. (2019). Leaf shape and size track habitat transitions across forest-grassland
- boundaries in the grass family (Poaceae). *Evolution* 73, 927–946.
- 1039 Gibson, D. J. (2009). Grasses and Grassland Ecology. Oxford, UK: Oxford University
- 1040 Press.
- 1041 Givnish, T. J., Ames, M., McNeal, J. R., McKain, M. R., Roxanne Steele, P.,
- dePamphilis, C. W., et al. (2010). Assembling the Tree of the Monocotyledons:
- 1043 Plastome Sequence Phylogeny and Evolution of Poales1. Ann. Mo. Bot. Gard 97,
- 1044 584–616.
- 1045 Grass Phylogeny Working Group, Barker, N. P., Clark, L. G., Davis, J. I., Duvall, M. R.,
- Guala, G. F., et al. (2001). Phylogeny and Subfamilial Classification of the Grasses
- 1047 (Poaceae). Ann. Mo. Bot. Gard. 88, 373–457.
- 1048 Grass Phylogeny Working Group II (2012). New grass phylogeny resolves deep
- evolutionary relationships and discovers C4 origins. *New Phytol.* 193, 304–312.
- Hanson, L., Boyd, A., Johnson, M. A. T., and Bennett, M. D. (2005). First nuclear DNA
- 1051 C-values for 18 eudicot families. *Ann. Bot.* 96, 1315–1320.
- Haug-Baltzell, A., Stephens, S. A., Davey, S., Scheidegger, C. E., and Lyons, E. (2017).
- 1053 SynMap2 and SynMap3D: web-based whole-genome synteny browsers.
- 1054 Bioinformatics 33, 2197–2198.
- Hirano, H.-Y., Tanaka, W., and Toriba, T. (2014). "Grass Flower Development," in
- 1056 Flower Development: Methods and Protocols, eds. J. L. Riechmann and F. Wellmer

- 1057 (New York, NY: Springer New York), 57–84.
- Hoff, K. J., Lomsadze, A., Borodovsky, M., and Stanke, M. (2019). "Whole-Genome
- Annotation with BRAKER," in *Gene Prediction: Methods and Protocols*, ed. M.
- 1060 Kollmar (New York, NY: Springer New York), 65–95.
- Howe, K. L., Achuthan, P., Allen, J., Allen, J., Alvarez-Jarreta, J., Amode, M. R., et al.
- 1062 (2021). Ensembl 2021. Nucleic Acids Res. 49, D884–D891.
- 1063 Huang, Y., Zhao, S., Fu, Y., Sun, H., Ma, X., Tan, L., et al. (2018). Variation in the
- regulatory region of FZP causes increases in secondary inflorescence branching
- and grain yield in rice domestication. *Plant J.* 96, 716–733.
- Hunziker, J. H., Wulff, A. F., and Soderstrom, T. R. (1982). Chromosome Studies on the
- Bambusoideae (Gramineae). *Brittonia* 34, 30.
- 1068 International Brachypodium Initiative (2010). Genome sequencing and analysis of the
- model grass Brachypodium distachyon. *Nature* 463, 763–768.
- Jiang, C.-K., and Rao, G.-Y. (2020). Insights into the Diversification and Evolution of
- 1071 R2R3-MYB Transcription Factors in Plants. *Plant Physiol.* 183, 637–655.
- Johnson, C., Kasprzewska, A., Tennessen, K., Fernandes, J., Nan, G.-L., Walbot, V., et
- al. (2009). Clusters and superclusters of phased small RNAs in the developing
- inflorescence of rice. Genome Res. 19, 1429–1440.

Johnson, N. R., Yeoh, J. M., Coruh, C., and Axtell, M. J. (2016). Improved Placement of

- 1076 Multi-mapping Small RNAs. *G*3 6, 2103–2111.
- 1077 Judziewicz, E. J., Clark, L. G., Londoño, X., and Stern M. J. (1999). American
- 1078 Bamboos. Washington, D.C.: Smithsonian Books.
- 1079 Judziewicz, E. J., and Soderstrom, T. R. (1989). Morphological, anatomical, and
- taxonomic studies in Anomochloa and Streptochaeta (Poaceae: Bambusoideae).
- 1081 Smithson. Contrib. Bot. Available at:
- https://repository.si.edu/bitstream/handle/10088/6984/scb-0068.pdf.
- Kaneko, M., Inukai, Y., Ueguchi-Tanaka, M., Itoh, H., Izawa, T., Kobayashi, Y., et al.
- 1084 (2004). Loss-of-function mutations of the rice GAMYB gene impair alpha-amylase
- expression in aleurone and flower development. *Plant Cell* 16, 33–44.
- 1086 Katoh, K., and Standley, D. M. (2013). MAFFT multiple sequence alignment software
- version 7: improvements in performance and usability. *Mol. Biol. Evol.* 30, 772–780.
- Kellogg, E. A. (2015). "Poaceae," in *The Families and Genera of Vascular Plants*, ed. K.
- 1089 Kubitzki (Springer), 1–416.
- Kellogg, E. A., Camara, P. E. A. S., Rudall, P. J., Ladd, P., Malcomber, S. T., Whipple,

- 1091 C. J., et al. (2013). Early inflorescence development in the grasses (Poaceae).
- 1092 Front. Plant Sci. 4, 250.
- Kielbasa, S. M., Wan, R., Sato, K., Horton, P., and Frith, M. C. (2011). Adaptive seeds
- tame genomic sequence comparison. *Genome Research* 21, 487–493.
- Kim, S., Soltis, P. S., Wall, K., and Soltis, D. E. (2006). Phylogeny and domain evolution
- in the APETALA2-like gene family. *Mol. Biol. Evol.* 23, 107–120.
- Kirpes, C. C., Clark, L. G., and Lersten, N. R. (1996). Systematic significance of pollen
- arrangement in microsporangia of Poaceae and Cyperaceae: review and
- observations on representative taxa. *Am. J. Bot.* 83, 1609–1622.
- Kozomara, A., Birgaoanu, M., and Griffiths-Jones, S. (2019). miRBase: from microRNA
- sequences to function. *Nucleic Acids Res.* 47, D155–D162.
- 1102 Kozomara, A., and Griffiths-Jones, S. (2014). miRBase: annotating high confidence
- microRNAs using deep sequencing data. *Nucleic Acids Research* 42, D68–D73.
- Laetsch, D. R., and Blaxter, M. L. (2017). BlobTools: Interrogation of genome
- 1105 assemblies. *F1000Res.* 6, 1287.
- Lauter, N., Kampani, A., Carlson, S., Goebel, M., and Moose, S. P. (2005).
- microRNA172 down-regulates glossy15 to promote vegetative phase change in
- 1108 maize. *Proc. Natl. Acad. Sci. U. S. A.* 102, 9412–9417.
- 1109 Leandro, T. D., Rodrigues, T. M., Clark, L. G., and Scatena, V. L. (2018). Fusoid cells in
- the grass family Poaceae (Poales): a developmental study reveals homologies and
- suggests new insights into their functional role in young leaves. Ann. Bot. 122,
- 1112 833–848.
- Lee, D.-Y., and An, G. (2012). Two AP2 family genes, supernumerary bract (SNB) and
- Osindeterminate spikelet 1 (OsIDS1), synergistically control inflorescence
- architecture and floral meristem establishment in rice. *Plant J.* 69, 445–461.
- Lehmann, C. E. R., Griffith, D. M., Simpson, K. J., Michael Anderson, T., Archibald, S.,
- Beerling, D. J., et al. (2019). Functional diversification enabled grassy biomes to fill
- 1118 global climate space. *biorxiv* . doi:10.1101/583625.
- 1119 Letunic, I., and Bork, P. (2019). Interactive Tree Of Life (iTOL) v4: recent updates and
- new developments. *Nucleic Acids Res.* 47, W256–W259.
- 1121 Lex, A., Gehlenborg, N., Strobelt, H., Vuillemot, R., and Pfister, H. (2014). UpSet:
- 1122 Visualization of Intersecting Sets. *IEEE Trans. Vis. Comput. Graph.* 20, 1983–1992.
- Li, C., Lin, H., Chen, A., Lau, M., Jernstedt, J., and Dubcovsky, J. (2019a). Wheat
- VRN1, FUL2 and FUL3 play critical and redundant roles in spikelet development
- and spike determinacy. *Development* 146.

- Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34, 3094–3100.
- Linder, H. P., Lehmann, C. E. R., Archibald, S., Osborne, C. P., and Richardson, D. M.
- 1129 (2018). Global grass (Poaceae) success underpinned by traits facilitating
- colonization, persistence and habitat transformation. *Biol. Rev. Camb. Philos. Soc.*
- 1131 93, 1125–1144.
- 1132 Liu, K., Cao, J., Yu, K., Liu, X., Gao, Y., Chen, Q., et al. (2019). Wheat TaSPL8
- Modulates Leaf Angle Through Auxin and Brassinosteroid Signaling1. *Plant*
- 1134 Physiol. 181, 179–194.
- Li, Y., Zhu, J., Wu, L., Shao, Y., Wu, Y., and Mao, C. (2019b). Functional Divergence of PIN1 Paralogous Genes in Rice. *Plant Cell Physiol.* 60, 2720–2732.
- Lorenz, R., Bernhart, S. H., Höner Zu Siederdissen, C., Tafer, H., Flamm, C., Stadler, P. F., et al. (2011). ViennaRNA Package 2.0. *Algorithms Mol. Biol.* 6, 26.
- 1139 Lowe, K., Wu, E., Wang, N., Hoerster, G., Hastings, C., Cho, M.-J., et al. (2016).
- Morphogenic Regulators Baby boom and Wuschel Improve Monocot
- 1141 Transformation. *Plant Cell* 28, 1998–2015.
- 1142 Luo, R., Liu, B., Xie, Y., Li, Z., Huang, W., Yuan, J., et al. (2012). SOAPdenovo2: an
- empirically improved memory-efficient short-read de novo assembler. *Gigascience*
- 1144 1, 18.
- Lyons, E., and Freeling, M. (2008). How to usefully compare homologous plant genes
- and chromosomes as DNA sequences. *Plant J.* 53, 661–673.
- 1147 Mamidi, S., Healey, A., Huang, P., Grimwood, J., Jenkins, J., Barry, K., et al. (2020). A
- genome resource for green millet Setaria viridis enables discovery of agronomically
- 1149 valuable loci. *Nat. Biotechnol.* 38, 1203–1210.
- 1150 Mantilla-Perez, M. B., Bao, Y., Tang, L., Schnable, P. S., and Salas-Fernandez, M. G.
- 1151 (2020). Toward "Smart Canopy" Sorghum: Discovery of the Genetic Control of Leaf
- Angle Across Layers. *Plant Physiol.* 184, 1927–1940.
- 1153 Marçais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel
- 1154 counting of occurrences of k-mers. *Bioinformatics* 27, 764–770.
- 1155 Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput
- sequencing reads. *EMBnet.journal* 17, 10–12.
- 1157 McKain, M. R., Tang, H., McNeal, J. R., Ayyampalayam, S., Davis, J. I., dePamphilis, C.
- 1158 W., et al. (2016). A Phylogenomic Assessment of Ancient Polyploidy and Genome
- 1159 Evolution across the Poales. *Genome Biol. Evol.* 8, 1150–1164.
- 1160 Minh, B. Q., Schmidt, H. A., Chernomor, O., Schrempf, D., Woodhams, M. D., von

- Haeseler, A., et al. (2020). IQ-TREE 2: New Models and Efficient Methods for
- Phylogenetic Inference in the Genomic Era. *Mol. Biol. Evol.* 37, 1530–1534.
- 1163 Murat, F., Xu, J.-H., Tannier, E., Abrouk, M., Guilhot, N., Pont, C., et al. (2010).
- Ancestral grass karyotype reconstruction unravels new mechanisms of genome
- shuffling as a source of plant evolution. *Genome Res.* 20, 1545–1557.
- Newell, T. K. (1969). A study of the genus Joinvillea (Flagellariaceae). *J. Arnold Arbor.*
- 1167 50, 527–555.
- Ou, S., Chen, J., and Jiang, N. (2018). Assessing genome assembly quality using the
- 1169 LTR Assembly Index (LAI). *Nucleic Acids Res.* 46, e126.
- Ou, S., and Jiang, N. (2018). LTR retriever: A Highly Accurate and Sensitive Program
- for Identification of Long Terminal Repeat Retrotransposons. *Plant Physiology* 176,
- 1172 1410–1422.
- 1173 Ou, S., Su, W., Liao, Y., Chougule, K., Agda, J. R. A., Hellinga, A. J., et al. (2019).
- Benchmarking transposable element annotation methods for creation of a
- streamlined, comprehensive pipeline. *Genome Biol.* 20, 275.
- Ouyang, S., Zhu, W., Hamilton, J., Lin, H., Campbell, M., Childs, K., et al. (2007). The
- 1177 TIGR Rice Genome Annotation Resource: improvements and new features.
- 1178 *Nucleic Acids Res.* 35, D883–7.
- Patel, P., Mathioni, S. M., Hammond, R., Harkess, A. E., Kakrana, A., Arikit, S., et al.
- 1180 (2021). Reproductive phasiRNA loci and DICER-LIKE5, but not microRNA loci,
- diversified in monocotyledonous plants. *Plant Physiol.* 185, 1764–1782.
- Paterson, A. H., Bowers, J. E., and Chapman, B. A. (2004). Ancient polyploidization
- predating divergence of the cereals, and its consequences for comparative
- 1184 genomics. *Proc. Natl. Acad. Sci. U. S. A.* 101, 9903–9908.

Pohl, R. W., and Davidse, G. (1971). Chromosome Numbers of Costa Rican Grasses.

- 1186 Brittonia 23, 293.
- 1187 Potter, S. C., Luciani, A., Eddy, S. R., Park, Y., Lopez, R., and Finn, R. D. (2018).
- HMMER web server: 2018 update. *Nucleic Acids Res.* 46, W200–W204.
- 1189 Preston, J. C., Christensen, A., Malcomber, S. T., and Kellogg, E. A. (2009). MADS-box
- gene expression and implications for developmental origins of the grass spikelet.
- 1191 Am. J. Bot. 96, 1419–1429.
- 1192 Preston, J. C., and Kellogg, E. A. (2006). Reconstructing the evolutionary history of
- paralogous APETALA1/FRUITFULL-like genes in grasses (Poaceae). *Genetics*
- 1194 174, 421–437.
- 1195 Pryszcz, L. P., and Gabaldón, T. (2016). Redundans: an assembly pipeline for highly

heterozygous genomes. *Nucleic Acids Res.* 44, e113.

- Rodriguez Milla, M. A., Maurer, A., Rodriguez Huete, A., and Gustafson, J. P. (2003).
- Glutathione peroxidase genes in Arabidopsis are ubiquitous and regulated by
- abiotic stresses through diverse signaling pathways. *Plant J.* 36, 602–615.
- 1200 Rothfels, C. J. (2021). Polyploid phylogenetics. New Phytol. 230, 66-72.
- Saarela, J. M., Burke, S. V., Wysocki, W. P., Barrett, M. D., Clark, L. G., Craine, J. M.,
- et al. (2018). A 250 plastome phylogeny of the grass family (Poaceae): topological
- support under different data partitions. *PeerJ* 6, e4299.
- 1204 Sajo, M. D. G., Furness, C. A., and Rudall, P. J. (2009). Microsporogenesis is
- simultaneous in the early-divergent grass Streptochaeta, but successive in the
- 1206 closest grass relative, Ecdeiocolea. *Grana* 48, 27–37.
- 1207 Sajo, M. G., Longhi-Wagner, H. M., and Rudall, P. J. (2008). Reproductive morphology
- of the early-divergent grass Streptochaeta and its bearing on the homologies of the
- grass spikelet. Plant Syst. Evol. 275, 245.
- Sajo, M. G., Longhi-Wagner, H., and Rudall, P. J. (2007). Floral Development and
- 1211 Embryology in the Early-Divergent Grass Pharus. *Int. J. Plant Sci.* 168, 181–191.
- 1212 Sajo, M. G., Pabón-Mora, N., Jardim, J., Stevenson, D. W., and Rudall, P. J. (2012).
- Homologies of the flower and inflorescence in the early-divergent grass
- 1214 Anomochloa (Poaceae). *Am. J. Bot.* 99, 614–628.
- Salse, J., Bolot, S., Throude, M., Jouffe, V., Piegu, B., Quraishi, U. M., et al. (2008).
- 1216 Identification and characterization of shared duplications between rice and wheat
- provide new insight into grass genome evolution. *Plant Cell* 20, 11–24.

1218 Sarwar, H. (2013). The importance of cereals (Poaceae: Gramineae) nutrition in human

- health: A review. *J. Cereals Oilseeds* 4, 32–35.
- 1220 Scanlon, M. J., Schneeberger, R. G., and Freeling, M. (1996). The maize mutant narrow
- sheath fails to establish leaf margin identity in a meristematic domain. *Development*
- 1222 122, 1683–1691.
- 1223 Schmidt, R., Schippers, J. H. M., Mieulet, D., Obata, T., Fernie, A. R., Guiderdoni, E., et
- al. (2013). MULTIPASS, a rice R2R3-type MYB transcription factor, regulates
- adaptive growth by integrating multiple hormonal pathways. *Plant J.* 76, 258–273.
- Seetharam, A., Singh, U., Li, J., Bhandary, P., Arendsee, Z., and Wurtele, E. S. (2019).
- Maximizing prediction of orphan genes in assembled genomes. *biorxiv*
- 1228 doi:10.1101/2019.12.17.880294.
- 1229 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., and Zdobnov, E. M.
- 1230 (2015). BUSCO: assessing genome assembly and annotation completeness with

single-copy orthologs. *Bioinformatics* 31, 3210–3212.

- Soderstrom, T. R., and Ellis, R. P. (1987). "The position of bamboo genera and allies in
- a system of grass classification," in *Grass Systematics and Evolution*, eds. T. R.
- Soderstrom, K. W. Hilu, C. S. Campbell, and M. E. Barkworth (Washington, DC:
- 1235 Smithsonian Institution Press), 225–238.
- Soreng, R. J., Peterson, P. M., Romaschenko, K., Davidse, G., Teisher, J. K., Clark, L.
- 1237 G., et al. (2017). A worldwide phylogenetic classification of the Poaceae
- 1238 (Gramineae) II: An update and a comparison of two 2015 classifications:
- 1239 Phylogenetic classification of the grasses II. *J. Syst. Evol.* 55, 259–290.
- Spriggs, E. L., Christin, P.-A., and Edwards, E. J. (2014). C4 photosynthesis promoted
- species diversification during the Miocene grassland expansion. *PLoS One* 9,
- 1242 e97722.
- Suyama, M., Torrents, D., and Bork, P. (2006). PAL2NAL: robust conversion of protein
- sequence alignments into the corresponding codon alignments. *Nucleic Acids Res.*
- 1245 34, W609–12.
- Teng, C., Zhang, H., Hammond, R., Huang, K., Meyers, B. C., and Walbot, V. (2020).
- Dicer-like 5 deficiency confers temperature-sensitive male sterility in maize. *Nat.*
- 1248 *Commun.* 11, 2912.
- 1249 Thody, J., Folkes, L., Medina-Calzada, Z., Xu, P., Dalmay, T., and Moulton, V. (2018).
- 1250 PAREsnip2: a tool for high-throughput prediction of small RNA targets from
- degradome sequencing data using configurable targeting rules. *Nucleic Acids Res.*
- 1252 46, 8730–8739.
- 1253 Tsuji, H., Aya, K., Ueguchi-Tanaka, M., Shimada, Y., Nakazono, M., Watanabe, R., et
- al. (2006). GAMYB controls different sets of genes and is differentially regulated by
- microRNA in aleurone cells and anthers. *Plant J.* 47, 427–444.
- 1256 UniProt Consortium (2021). UniProt: the universal protein knowledgebase in 2021.
- 1257 Nucleic Acids Res. 49, D480–D489.

1258 UpSetR Github Available at: https://github.com/hms-dbmi/UpSetR [Accessed February

- 1259 21, 2021].
- 1260 Valencia, J. I. (1962). Los cromosomas de Streptochaeta spicata Schrad. (Gramineae).
- 1261 Darwiniana 12, 379–383.
- 1262 VanBuren, R., Man Wai, C., Wang, X., Pardo, J., Yocca, A. E., Wang, H., et al. (2020).
- 1263 Exceptional subgenome stability and functional divergence in the allotetraploid
- 1264 Ethiopian cereal teff. *Nat. Commun.* 11, 884.
- 1265 Velasco, R., Zharkikh, A., Affourtit, J., Dhingra, A., Cestaro, A., Kalyanaraman, A., et al.
- 1266 (2010). The genome of the domesticated apple (Malus × domestica Borkh.). *Nat.*

1267 Genet. 42, 833–839.

- 1268 Venturini, L., Caim, S., Kaithakottil, G. G., Mapleson, D. L., and Swarbreck, D. (2018).
- Leveraging multiple transcriptome assembly methods for improved gene structure
- annotation. *Gigascience* 7, giy093.
- Walbot, V., and Egger, R. L. (2016). Pre-Meiotic Anther Development: Cell Fate
- 1272 Specification and Differentiation. *Annu. Rev. Plant Biol.* 67, 365–395.
- 1273 Wang, J., Yu, J., Sun, P., Li, Y., Xia, R., Liu, Y., et al. (2016). Comparative Genomics
- 1274 Analysis of Rice and Pineapple Contributes to Understand the Chromosome
- Number Reduction and Genomic Changes in Grasses. Front. Genet. 7, 174.
- Wang, X., Shi, X., Hao, B., Ge, S., and Luo, J. (2005). Duplication and DNA segmental loss in the rice genome: implications for diploidization. *New Phytol.* 165, 937–946.
- Whipple, C. J. (2017). Grass inflorescence architecture and evolution: the origin of novel signaling centers. *New Phytol.* 216, 367–372.
- Whipple, C. J., Zanis, M. J., Kellogg, E. A., and Schmidt, R. J. (2007). Conservation of B
- class gene expression in the second whorl of a basal grass and outgroups links the
- origin of lodicules and petals. *Proc. Natl. Acad. Sci. U. S. A.* 104, 1081–1086.
- White, R. P., Murray, S., Rohweder, M., Prince, S. D., Thompson, K. M., and Others (2000). *Grassland ecosystems*. World Resources Institute Washington, DC, USA.
- Xia, R., Xu, J., and Meyers, B. C. (2017). The Emergence, Evolution, and Diversification of the miR390-TAS3-ARF Pathway in Land Plants. *Plant Cell* 29, 1232–1247.
- 1287 Yu, Y., Woo, M.-O., Rihua, P., and Koh, H.-J. (2020). The DROOPING LEAF (DR) gene
- encoding GDSL esterase is involved in silica deposition in rice (Oryza sativa L.).
- 1289 PLoS One 15, e0238887.
- 1290 Zhai, J., Arikit, S., Simon, S. A., Kingham, B. F., and Meyers, B. C. (2014). Rapid
- construction of parallel analysis of RNA end (PARE) libraries for Illumina
- 1292 sequencing. *Methods* 67, 84–90.
- 1293 Zhai, J., Zhang, H., Arikit, S., Huang, K., Nan, G.-L., Walbot, V., et al. (2015).
- Spatiotemporally dynamic, cell-type-dependent premeiotic and meiotic phasiRNAs
- in maize anthers. *Proc. Natl. Acad. Sci. U. S. A.* 112, 3146–3151.
- Zhang, H., Liang, W., Yang, X., Luo, X., Jiang, N., Ma, H., et al. (2010). Carbon starved
- anther encodes a MYB domain protein that regulates sugar partitioning required for
- rice pollen development. *Plant Cell* 22, 672–689.
- Zhang, R.-G., Wang, Z.-X., Ou, S., and Li, G.-Y. (2019). TEsorter: lineage-level
- classification of transposable elements using conserved protein domains. *biorxiv*,
- 1301 doi:10.1101/800177.

1302 Zhang, R., Huang, S., Li, S., Song, G., Li, Y., Li, W., et al. (2020). Evolution of PHAS 1303 loci in the young spike of Allohexaploid wheat. BMC Genomics 21, 200. Zhou, Y., Lu, D., Li, C., Luo, J., Zhu, B.-F., Zhu, J., et al. (2012). Genetic control of seed 1304 1305 shattering in rice by the APETALA2 transcription factor shattering abortion1. Plant 1306 Cell 24, 1034-1048. Zhu, J., Chen, H., Li, H., Gao, J.-F., Jiang, H., Wang, C., et al. (2008). Defective in 1307 Tapetal development and function 1 is essential for anther development and tapetal 1308 1309 function for microspore maturation in Arabidopsis. *Plant J.* 55, 266–277. 1310 Zimin, A. V., Marçais, G., Puiu, D., Roberts, M., Salzberg, S. L., and Yorke, J. A. (2013). 1311 The MaSuRCA genome assembler. *Bioinformatics* 29, 2669–2677. 1312