



**HAL**  
open science

## **ZW sex chromosome structure in *Amborella trichopoda***

Sarah B Carey, Laramie Aközbek, John T Lovell, Jerry Jenkins, Adam L Healey, Shengqiang Shu, Paul Grabowski, Alan Yocca, Ada Stewart, Teresa Jones, et al.

► **To cite this version:**

Sarah B Carey, Laramie Aközbek, John T Lovell, Jerry Jenkins, Adam L Healey, et al.. ZW sex chromosome structure in *Amborella trichopoda*. *Nature Plants*, 2024, 10, pp.1944-1954. 10.1038/s41477-024-01858-x . hal-04768060

**HAL Id: hal-04768060**

**<https://hal.science/hal-04768060v1>**

Submitted on 5 Nov 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## ZW sex chromosome structure in *Amborella trichopoda*

Sarah B. Carey<sup>1</sup>, Laramie Aközbek<sup>1,2</sup>, John T. Lovell<sup>1,3</sup>, Jerry Jenkins<sup>1</sup>, Adam L. Healey<sup>1</sup>, Shengqiang Shu<sup>3</sup>, Paul Grabowski<sup>1,3</sup>, Alan Yocca<sup>1</sup>, Ada Stewart<sup>1</sup>, Teresa Jones<sup>1</sup>, Kerrie Barry<sup>3</sup>, Shanmugam Rajasekar<sup>4</sup>, Jayson Talag<sup>4</sup>, Charlie Scutt<sup>5</sup>, Porter P. Lowry II<sup>6,7</sup>, Jérôme Munzinger<sup>8</sup>, Eric B. Knox<sup>9</sup>, Douglas E. Soltis<sup>10</sup>, Pamela S. Soltis<sup>10</sup>, Jane Grimwood<sup>1,3</sup>, Jeremy Schmutz<sup>1,3</sup>, James Leebens-Mack<sup>11,\*</sup>, Alex Harkess<sup>1,\*</sup>

### AFFILIATIONS

<sup>1</sup>HudsonAlpha Institute for Biotechnology, Huntsville, AL, USA

<sup>2</sup>Department of Crop, Soil, and Environmental Sciences, Auburn University, Auburn, AL, USA

<sup>3</sup>Department of Energy Joint Genome Institute, Lawrence Berkeley National Laboratory, Berkeley, CA, USA

<sup>4</sup>Arizona Genomics Institute, University of Arizona, Tucson, AZ, USA

<sup>5</sup>Laboratoire Reproduction et Développement des Plantes, Univ. Lyon, ENS de Lyon, UCB Lyon-1, CNRS, INRA, F-69342 Lyon, France

<sup>6</sup>Missouri Botanical Garden, 4344 Shaw Blvd., St. Louis, MO, USA

<sup>7</sup>Institut de Systématique, Évolution, et Biodiversité (ISYEB), Muséum National d'Histoire Naturelle, Centre National de la Recherche Scientifique, Sorbonne Université, École Pratique des Hautes Études, Université des Antilles, C.P. 39, 57 rue Cuvier, 75005 Paris, France

<sup>8</sup>AMAP, Univ. Montpellier, IRD, CIRAD, CNRS, INRAE, F-34398 Montpellier, France

<sup>9</sup>Department of Biology, Indiana University, Bloomington, IN, USA

<sup>10</sup>Florida Museum of Natural History, University of Florida, Gainesville, FL, USA

<sup>11</sup>Department of Plant Biology, University of Georgia, Athens, GA, USA

\*Authors for correspondence: jleebsmack@uga.edu (J.L.-M.), aharkess@hudsonalpha.org (A.H.)

## ABSTRACT

Sex chromosomes have evolved hundreds of times across the flowering plant tree of life; their recent origins in some members of this clade can shed light on the early consequences of suppressed recombination, a crucial step in sex chromosome evolution. *Amborella trichopoda*, the sole species on a lineage that is sister to all other extant flowering plants, is dioecious with a young ZW sex determination system. Here we present a haplotype-resolved genome assembly, including highly contiguous assemblies of the Z and W chromosomes. We identify a ~3-Megabase sex-determination region (SDR) captured in two strata that includes a ~300-Kilobase inversion that is enriched with repetitive sequence and contains a homolog of the *Arabidopsis* METHYLTHIOADENOSINE NUCLEOSIDASE (*MTN1-2*) genes, which are known to be involved in fertility. However, the remainder of the SDR does not show patterns typically found in non-recombining SDRs, such as repeat accumulation and gene loss. These findings are consistent with the hypothesis that dioecy is derived in *Amborella* and the sex chromosome pair has not significantly degenerated.

## MAIN TEXT

The evolution of separate sexes, or dioecy, is a rare trait in angiosperms, having been identified in just 5-10% of species<sup>1</sup>. At the same time, dioecy has evolved hundreds of times independently across the flowering plant tree of life,<sup>2</sup> making flowering plants ideal for examining the evolution of sex chromosomes over both deep and shallow time scales. Comparative investigations of sex chromosomes rely on high-quality genome assemblies<sup>2</sup>, and while the availability of genomes for dioecious species has increased, there are only a few where the structure of the sex chromosome pair has been well characterized. While divergence between X and Y sex chromosomes has been described in a growing number of angiosperm species<sup>2,3</sup>, investigations of what some consider less common ZW systems can shed new light on the dynamics and consequences of sex chromosome evolution.

Since its discovery as the likely sister lineage to all other living angiosperms, *Amborella trichopoda* (Amborellaceae; hereafter, *Amborella*)<sup>4-7</sup> has served as a pivotal taxon for investigating the origin and early diversification of flowering plants<sup>8,9</sup>. *Amborella* is an understory shrub or small tree endemic to New Caledonia and the sole extant species in the Amborellales. The flowers of *Amborella* are actinomorphic and have a perianth of undifferentiated tepals, which are characteristics shared with the reconstructed ancestral flower (Fig. 1)<sup>9</sup>. Importantly, however, *Amborella* is dioecious<sup>10</sup> with ZW sex chromosomes that evolved after the lineage diverged from other flowering plants<sup>11</sup>. This implies that dioecy in *Amborella* is derived from a hermaphroditic mating system and that the ancestral angiosperm had perfect flowers, in agreement with ancestral state reconstructions<sup>9</sup>. Significant progress has been made in several angiosperm species to identify the genes involved in the evolution of dioecy<sup>12-17</sup>, but the molecular basis in *Amborella* remains unknown. Here we present a haplotype-resolved assembly of the *Amborella* genome and compare highly contiguous Z and W sex chromosome assemblies to address outstanding questions about their structure and gene content, including putative sex-determining genes.

## RESULTS

### Improved genome assembly and annotation of *Amborella*

The *Amborella* reference genome has been a central anchor for comparative investigations of gene family and gene structure evolution across angiosperms. Despite its demonstrated utility, the 2013 *Amborella* genome used primarily short sequencing reads, which cannot fully resolve repetitive regions<sup>18</sup>. The repeat-derived gaps were filled in a long-read assembly<sup>11</sup>, but both biological haplotypes were collapsed into a single sequence representation. Despite the higher contiguity, the 2022 genome offers limited information regarding sex-determination regions (SDRs) because in this assembly the Z and W chromosomes are a chimeric mix represented as a single chromosome<sup>11</sup>.

To build a haplotype-resolved genome assembly for *Amborella* cv. Santa Cruz 75, we used a combination of PacBio HiFi (mean coverage = 58.81x per haplotype; mean read length = 22,900 bp) and Phase Genomics Hi-C (coverage = 42.31x; Supplementary Table 1) sequencing

technologies. The final haplotype 1 (HAP1) and 2 (HAP2) assemblies include 708.1 Mb in 59 contigs (contig  $N_{50}$  = 36.3 Mb;  $L_{50}$  = 7) and 700.5 Mb in 45 contigs (contig  $N_{50}$  = 44.5 Mb;  $L_{50}$  = 7), respectively; 99.69% and 99.87% of the assembled sequence is contained in the 13 largest scaffolds for HAP1 and HAP2, respectively, corresponding to the expected chromosome number<sup>19</sup> (Supplementary Figure 1). We found the Merqury  $k$ -mer completeness<sup>20</sup> of HAP1 was 95.4% (QV 63) and HAP2 was 95.3% (QV 55), and the combined assemblies exhibit 98.8% completeness (QV 57). Consistent with earlier assemblies, we annotated repeats and found they represent ~56% of the sequence for both haplotypes (Fig. 1; Supplementary Table 2)<sup>18</sup>. To annotate gene models, we used a combination of RNAseq and Iso-Seq (~757 million 2x150 read pairs, ~825K full-length transcripts). We annotated 21,800 gene models in HAP1 and 21,721 in HAP2, with Embryophyte BUSCOs of 98.6% and 98.8%, respectively, an increase from 85.5% in the 2013 release<sup>18</sup>. Overall, the new assemblies represent a great improvement in the *Amborella* genome reference, resolving most of the previous gaps (Supplementary Figure 2, Supplementary Table 2).

*Amborella*'s ancient divergence ~140 MYA<sup>21</sup> from all other living angiosperms provides an opportunity to examine conserved features that were likely present in the ancestral genome of all flowering plants. For example, the repeat-dense pericentromeric region and gene-dense chromosome arms of *Amborella* (Fig. 1) mirror those of most angiosperm genomes, in stark contrast to the more uniform gene and repeat density of most conifers, ferns, and mosses<sup>22-24</sup>. The pericentromeric regions are enriched in Long Terminal Repeats (LTRs), specifically *Ty3* and *Ty1* elements, as is often seen in other monocentric angiosperms<sup>25,26</sup>. Interestingly, unlike many previously examined sex chromosomes, the *Amborella* Z/W do not stand out as notable exceptions in terms of gene or repeat density (Fig. 1).

### Identification of the phased *Amborella* sex chromosomes

Sex chromosomes have unique inheritance patterns relative to autosomes. In a ZW system, the non-recombining SDR of the W chromosome is only inherited by females, while the remaining pseudoautosomal region (PAR) recombines freely and is expected to show a similar lack of divergence between the sexes as the autosomes. Identification of the boundary between the SDR and PAR of sex chromosomes is nontrivial, and PAR/SDR boundaries have been shown to vary among populations in some species<sup>27,28</sup>. Standard approaches for boundary identification employ combinations of methodologies like sex-biased read coverage and population genomic analyses<sup>29</sup>.

To delimit the PAR/SDR boundary we performed a  $k$ -mer analysis<sup>12,30</sup> to identify sequences that are unique to the *Amborella* SDR (henceforth, W-mers), using four different sampling strategies (Supplementary Methods). We found the W-mers densely mapped to Chr09 at ~44.32-47.26 Mb of HAP1 (Fig. 1-2, S3-6), supporting its identity as the W chromosome. This location is consistent with previous analyses<sup>11</sup>, although we find assessing W-mers to a haplotype-resolved assembly narrows the estimated size of the SDR from ~4 Mb to 2.94 Mb (Fig. 2, S7). Importantly, the W-mers show consistent coverage on Chr09 in HAP1, with low and

sporadic coverage along any other chromosome or unincorporated scaffold in the assembly (e.g., when using the Island-wide sampling, 97.73% of the mapped W-mers are within the SDR; Supplementary Figure 3-6; Supplementary Table 3). In contrast to the chimeric Z/W in the previous assembly, the resulting sex chromosome assemblies are nearly complete with only four unresolved gaps in the SDR (zero gaps in the homologous region on the Z; HZR) and are fully phased (Supplementary Figure 7).

A key characteristic of sex chromosomes is suppressed recombination of the SDR, and in many species, structural variants have been identified as the causal mechanism. To examine this in *Amborella*, we first used genome alignments to identify the HZR. The HZR is located on Chr09 of HAP2 at 44.52-47.12 (~2.60 Mb; Supplementary Figure 8), suggesting the SDR is only 340 Kb larger than the HZR, which is consistent with the observed cytological homomorphy of the ZW pair<sup>19</sup>. In the SDR, we found evidence for a ~292-Kb inversion located ~20 Kb within the beginning of the boundary and containing the majority of the W-specific sequence (Figs. 1B, S9). We could not, however, find evidence for inversions or other large structural variants surrounding the remaining portion of the SDR. Instead, the Z and W chromosomes are highly syntenic with one another, similar to the autosomes (Figs. 1, S8). We investigated other potential mechanisms for suppressed recombination, such as proximity to centromeres, where the existing low recombination has been shown to facilitate SDR evolution in some species<sup>31</sup>. In *Amborella*, the SDR is not located near the centromere; rather, it is approximately 1.82 Mb away from the *Ty3*-retrotransposon-rich pericentromeric region (Fig. 2). In the absence of obvious structural variants encompassing the SDR, it suggests that *Amborella* has a non-canonical mechanism to enforce non-recombination between the Z and the W.

### **The *Amborella* sex chromosomes are evolutionarily young**

*Amborella*'s sex chromosomes have previously been shown to have evolved after the lineage split from other living flowering plants<sup>11</sup>. With our phased Z/W pairs, we can better determine Z- and W-linked genes, providing a more confident estimate of the age of the SDR, and examine gene gain events. A classic signature of multiple recombination suppression events is a stepwise pattern of synonymous substitutions (Ks) of neighboring genes on the sex chromosomes<sup>32</sup>. Genes captured into the SDR in the same event are expected to have similar levels of Ks (i.e., evolutionary strata), whereas the older strata will have higher divergence between the Z and W compared to younger strata.<sup>32</sup> Understanding this timing of gene gain is essential to understanding the genetic mechanism for sex determination, because the candidate sex-determining genes are likely to have ceased recombining first (barring turnovers<sup>29</sup>).

To examine gene gain in the *Amborella* SDR, we calculated Ks of one-to-one orthologs on the W and Z chromosomes (i.e., gametologs). We compared the Ks values of 45 identifiable gametologs to 1,397 one-to-one orthologs in the PARs. We found that Ks varies across the SDR-HZR portion of the sex chromosomes (0.002-0.20; mean Ks=0.0298, SD=0.032) and is significantly higher than Ks in the PARs (mean Ks=0.004, SD=0.019; Kruskal-Wallis  $p < 0.00001$ ); Supplementary Figure 10), consistent with the expectation that the SDR is diverging

from the HZR on the Z chromosome. Interestingly, the gametolog pair with the highest Ks within the SDR is a homolog of *Arabidopsis* METHYLTHIOADENOSINE NUCLEOSIDASE *MTNI-2*, a gene involved in fertility, suggesting it resides in the oldest portion of the SDR; notably, the location of the W-linked *MTNI-2* homolog is within the SDR inversion.

We found the Ks values have two distinct steps, with the higher Ks values in the region corresponding to the inversion, suggesting two strata of gene capture into the SDR (Fig. 3). Defining the precise boundary between strata without obvious structural variants can be a challenge. To delineate stratum one (S1) from two (S2), we used a change point analysis on Ks and the average nucleotide differences between sampled females and males (Nei's dXY), which suggested that S1 ends at ~46.08Mb (Supplementary Figure 11). We found Ks to be significantly different between the strata (S1 mean Ks=0.037, SD=0.037, n=25; S2 mean Ks=0.021, SD=0.023, n=20; Mann-Whitney U, p=0.0014) as was the extent of nonsynonymous changes in proteins (Ka; Mann-Whitney U, p=0.008; Fig. 3), supporting inference of two strata. We also found dXY of genes to be significantly different (Mann Whitney U, p<2.6e-6), higher in S1 (mean = 0.0169, SD = 0.007, n=57) than S2 (mean = 0.0089, SD = 0.006, n=40). Using Ks, we also estimated the age of the SDR in *Amborella*. Following the previously applied approach<sup>11</sup>, we found S1 to be ~4.97 MYA while S2 is nearly half as old at ~2.41 MYA. These analyses indicate that the *Amborella* sex chromosomes are evolutionarily young, similar to several well-characterized XY systems<sup>3</sup>, and further suggest that the sex chromosomes evolved well after the lineage split from the rest of all living angiosperms.

### **The *Amborella* W shows little degeneration**

The recent origin of the *Amborella* sex chromosomes provides an opportunity to examine the early stages of their evolution. The lack of recombination in an SDR reduces the efficacy of natural selection and drives the accumulation of slightly deleterious mutations<sup>33,34</sup>. Two parallel signatures of deleterious mutations seen across independent evolutions of sex chromosomes is the accumulation of repeats and the loss of genes<sup>35-38</sup>. However, the tempo of this process of degeneration is not well understood.

In the SDR of *Amborella*, curiously we do not find the expected patterns of repeat expansions found in other SDRs. At 51.66% repeat elements, the SDR is lower than the genome average (56%) and 0.05% lower than the HZR, even when considering S1 and S2 separately (S1=52.13%; S2=50.98%; Supplementary Table 4). The only observed enrichment in repeats is within the inversion in S1, where we find more *Ty3* LTRs (4.32% increase relative to the HZR; Fig. 2). Otherwise, only a slight distinction between the SDR and its HZR is evident: the SDR exhibits a marginal increase ranging between 0.01-0.13% in the density of some superfamily elements (Fig. 2; Supplementary Table 4). We examined the distribution of the divergence values for intact LTRs as a proxy for their age<sup>39</sup> but found no patterns of distinctly younger or older LTRs within the W or Z (Supplementary Figure 12). Moreover, to assess genome-wide repeat expansion across the major Transposable Element (TE) superfamilies<sup>40</sup>, we used repeat landscapes, which showed a comparable pattern within the Z/W (Fig. 3, S13). These

observations support previous characterization of TE insertions in the *Amborella* genome as being quite old with little proliferation over the last 5 MY<sup>18</sup>. It has been proposed that a loss of active transposases or silencing may be playing a role in reducing TE activity across the *Amborella* genome<sup>18</sup> including the SDR.

Gene loss in an SDR has been hypothesized to contribute to the evolution of heteromorphy seen in many sex chromosome pairs<sup>41,42</sup>. In *Amborella*, of the 97 annotated models in the SDR and 84 in the HZR, 37 were W-specific and 24 Z-specific. To examine whether these models were missing from the other haplotype for technical or biological reasons, we also used dXY and presence-absence variation (PAV; Tables S5-6) between the sexes to evaluate gene content. For most of the W-specific models, males showed presence, and dXY within females was comparable to that of identifiable gametologs (mean dXY = 0.0136; Supplementary Table 7). Only seven models showed absence in coverage in males (dXY = 0 in females), suggesting conservatively that these represent W-specific genes, four of which are in the SDR inversion. Similarly, we identified only six Z-specific gene models. These analyses suggest that the Z and W have similar numbers of haplotype-specific genes and that the SDR has experienced similar levels of gene loss as the HZR.

Together, these results provide little evidence that degenerative processes, associated with cessation of recombination, have occurred in the *Amborella* SDR. This region is younger than that of *Rumex* (5-10 MYA<sup>43</sup>) and *Silene* (10 MYA<sup>44</sup>), which both show signatures of degeneration<sup>38,45</sup>. However, in *Spinacia oleracea*, a younger SDR (2-3 MYA) does show signs of degeneration<sup>46,47</sup>. The tempo of degeneration is apparently slower in *Amborella*, and there has not been sufficient time for gene loss or an accumulation of repeats as a consequence of the loss of recombination. One possible reason for the slower relative tempo is that most analyses of degeneration have focused on Y chromosomes, which are expected to degenerate faster than Ws due to male-biased mutation rates and stronger sexual selection<sup>48</sup>. Comparisons to other W chromosomes across independent origins are necessary to see if this holds true.

### **Candidate sex-determining genes in *Amborella***

ZW sex chromosomes have been less well characterized in plants than in animals; thus, *Amborella* can provide unique insights regarding the genetic mechanisms associated with their evolution. The two-gene model for sex chromosome evolution associated with a transition from hermaphroditism to dioecy posits that distinct genes with antagonistic impacts on female and male function experience strong selection for tight linkage (i.e., loss of recombination)<sup>49</sup>. Under this model, evolution of a ZW sex chromosome pair requires a dominant mutation causing male sterility arising on a proto-W chromosome, followed by a recessive loss-of-female-function mutation on the proto-Z (assuming a gynodioecious intermediate)<sup>49</sup>. As more sex chromosome pairs have been assembled, new models<sup>50,51</sup> have emerged that could be congruent with the data presented here, including the possibility that recombination suppression around a sterility locus could expand due to the sheltering of deleterious mutations.



Identification of these sex-determining genes requires an understanding of when sterility arises in the carpel and stamen developmental pathways. In *Amborella*, ontogenetic differences between female and male flowers are seen early in development<sup>52</sup>. Whereas male flowers produce an average of 12 stamens spiraling into the center of the flower, female flowers typically initiate a few staminodes just inside the tepals, but carpel initiation replaces staminode initiation as organ development proceeds towards the center of the flower<sup>52</sup> (Fig. 1). To identify candidate sex-determining genes, we examined differential expression between female and male flower buds during stage 5/6 of flower development, when carpels, stamens, and microsporangia develop<sup>11,52,53</sup>. We found 1,777 significantly differentially expressed genes at an adjusted p-value greater than 0.05. Of these, 34 are in the SDR, several of which are well-known flower development genes, including homologs of *MTNI-2*, *WUSCHEL (WUS)*, *LONELY GUY (LOG)*, *MONOPTEROS/Auxin Response Factor 5 (MP/ARF5)*, and *small auxin up-regulated RNA (SAUR)* gene families (Supplementary Figure 14; Supplementary Table 8-9). We found *ambMTN* and *ambLOG* had higher transcript abundance in females, while *ambWUS*, *ambMP*, and *ambSAUR* had greater expression in males. To further examine the sex-specific expression of SDR genes, we used the EvoRepro database (<https://evorepro.sbs.ntu.edu.sg/>), which has transcriptome data for 16 different tissue types for *Amborella*<sup>54</sup>. We contrasted female and male buds and flowers and found three genes with male-biased transcript abundance: *ambWUS* and a *DUF827* gene in buds and *ambLOG* in flowers, the latter differing in which sex has higher abundance from the analyses using stage 5/6 flowers. Given the known functions of these genes in *Arabidopsis* flower development, they are strong candidates for investigation of sex determination in *Amborella*.

While functional analyses are not currently possible in *Amborella*, comparisons to other species implicate the function of candidate genes that may be playing roles in *Amborella* sex determination. *WUS* encodes a homeobox transcription factor that is required for the maintenance of the floral meristem and has been shown to influence gynoecium and anther development<sup>55,56</sup>. In *Arabidopsis*, *WUS* knockouts have sepals, petals, a single stamen, and no carpel<sup>57</sup>. *WUS* has also been implicated in sex determination or shown sex-specific expression in several species that have unisexual flowers. In monoecious castor bean (*Ricinus communis*), *WUS* expression was only found in the shoot apical meristem of male flowers<sup>58</sup>, and in cucumbers (*Cucumis sativus*), *WUS* expression is three times greater in the carpel primordia of male flowers than females<sup>59</sup>. In *Silene*, gynoecium suppression is controlled by the *WUSCHEL-CLAVATA* feedback loop<sup>16</sup>. However, we do not see male-biased expression of the *CLV3* ortholog in *Amborella*, but we do see female-biased transcript abundance of the *Amborella* *CLE40* ortholog. In *Arabidopsis*, *WUS* promotes *CLV3* expression in the central zone of the inflorescence meristem while suppressing *CLE40* expression in the peripheral zone<sup>60</sup>. It is possible that the smaller floral meristem seen in female development relative to male floral meristems is due to reduced *ambWUS* expression driving increased *ambCLE40* expression and encroachment of peripheral zone cells into the central zone of the floral meristem. The role of *WUS* in maintaining meristematic zonation, coupled with its position in S1 in the SDR, makes

*ambWUS* a strong candidate for playing some role in gynoeceium suppression. Another strong candidate is *ambLOG*. *LOG* mutants were originally characterized in rice as producing floral phenotypes with a single stamen and no carpels<sup>61</sup>; in date palms (*Phoenix dactylifera*), a *LOG*-like gene was identified as a candidate Y-chromosome-linked female suppression gene<sup>13</sup>. In *Amborella*, *ambLOG* showed greater expression in females in the stage 5/6 data but was male-biased when considering all 16 tissues in the EvoRepro dataset. This switch in sex bias, and the fact *ambLOG* is located in the younger stratum of SDR (S2), suggest that differential *ambWUS* (and *ambCLE40*) expression may have been a first step in the divergence of male and female flower development. Like *ambLOG*, the *ambMP* and *ambSAUR* genes were captured in S2, and their functions in *Arabidopsis* suggest other roles in sex-specific development. *MP* has been shown to be involved with apical patterning of the embryo axis<sup>62,63</sup>. *SAURs* are a large gene family and in general play a role in cell elongation<sup>64</sup>, including in pollen tube growth<sup>65</sup>, stamen filament elongation<sup>66</sup>, and pistil growth<sup>67</sup>. Without functional validation in *Amborella*, we cannot rule out the possibility of any of these genes, though based on the data available, *ambWUS* may be the strongest candidate for spurring divergence in male and female flower development.

The significant difference in gene expression of *ambMTN* is especially interesting given that it is the gene model with the highest Ks value that is located in the SDR inversion. *MTN1-2* genes encode 5'-methylthioadenosine (MTA) nucleosidase<sup>68</sup>, and double mutant *mtn1-1mtn2-1* flowers in *Arabidopsis* have indehiscent anthers and malformed pollen grains<sup>69</sup>. Double mutants also affected carpels and ovules, although the structures were aberrant but not necessarily non-functional, and 10% looked like wild type<sup>69</sup>. The observed anther phenotype in *Arabidopsis* is consistent with the staminode development in female flowers in *Amborella*, and together these lines of evidence suggest that *ambMTN* may be the male-sterility gene. Based on our analyses, we hypothesize that the W-linked *ambMTN* was the initial male-sterility mutation, creating the proto-W, followed by a loss-of-function mutation on the W-*ambWUS* and a Z-copy shift to dosage-dependant gynoeceium suppression. The genes we have identified here make ideal candidates for further functional genomic investigation and validation.

## DISCUSSION

Advances in sequencing technologies and assembly algorithms have enabled the construction of telomere-to-telomere genome assemblies for humans, including the X and Y sex chromosomes<sup>70,71</sup>. The sex chromosomes in humans and other mammals are often highly heteromorphic and can be the most challenging chromosomes to sequence and assemble<sup>72</sup>. Moreover, given their antiquity, it is not possible to reconstruct events dating back to the origin and early evolution of mammalian sex chromosomes. In some plants and animals, however, sex chromosomes have repeatedly evolved from different ancestral autosomes, with different sex-determining mutations<sup>2,3,73</sup> and with various mechanisms to impede recombination between the sex chromosome pair. Here we show that we can fully phase structurally similar sex chromosomes within a heterogametic individual. Our analyses highlight the utility of phased sex

chromosomes, and diversity sequencing, to develop models of sex chromosome evolution when experimental investigation of gene function is currently intractable. This research lays the foundation for examining sex chromosome evolution in all angiosperms.

## METHODS

### **DNA/RNA extraction, library prep, and sequencing.**

We sequenced *Amborella trichopoda* (var. Santa Cruz 75) using a whole genome shotgun sequencing strategy and standard sequencing protocols. High-molecular-weight DNA was extracted from young tissue using the protocol of Doyle and Doyle<sup>74</sup> with minor modifications. Flash-frozen young leaves were ground to a fine powder in a frozen mortar with liquid nitrogen followed by very gentle extraction in a 2% CTAB buffer (that included proteinase K, PVP-40, and beta-mercaptoethanol) for 30 minutes to 1 hour at 50 °C. After centrifugation, the supernatant was gently extracted twice with 24:1 Chloroform: Isoamyl alcohol. The upper phase was transferred to a new tube, and 1/10th volume of 3 M Sodium acetate was added, the solution gently mixed, and DNA precipitated with iso-propanol. The DNA precipitate was collected by centrifugation, washed with 70% ethanol, air dried for 5-10 minutes, and dissolved thoroughly in an elution buffer at room temperature followed by RNase treatment. DNA purity was measured with a Nanodrop, DNA concentration was measured with Qubit HS kit (Invitrogen, Waltham, MA), and DNA size was validated by CHEF-DR II system (Bio-Rad Laboratories, Hercules, CA). The A PacBio HiFi library was constructed using DNA that was sheared using a Diagenode Megaruptor 3 instrument. Libraries were constructed using a SMRTbell Template Prep Kit 2.0 and tightly sized on a SAGE ELF instrument (1-18kb) to a final library average insert size of 24-Kb. PacBio sequencing was completed using the SEQUEL II platform at the HudsonAlpha Institute for Biotechnology in Huntsville, Alabama, yielding 83.3 Gb of raw sequence, with a total coverage of 58.81x per haplotype (Supplementary Table 12).

Illumina Hi-C sequencing for Santa Cruz 75 was conducted at Phase Genomics with a single 2x80 Dovetail Hi-C library (42.31x; Supplementary Table 1). DNA for the Illumina PCR-free library was extracted using a Qiagen DNeasy kit (Qiagen, Hilden, Germany) and was sequenced at the HudsonAlpha Institute for Biotechnology in Huntsville, Alabama. Illumina reads were sequenced on the Illumina NovaSeq 6000 platform using a 400-bp-insert TruSeq PCR-free fragment library (49.62x). Prior to assembly, Illumina fragment reads were screened for phix contamination. Reads composed of >95% simple sequence and those <50 bp after trimming for adapter and quality (q<20) were removed. The final read set consists of 158,007,088 reads for a total of 49.62x of high-quality Illumina bases.

To annotate gene models, we generated RNAseq and Iso-Seq data for several stages of leaf, flower, and fruit for Santa Cruz 75 and two male isolates, ABG 2006-2975 and ABG 2008-1967 (Supplementary Table 11). Total RNAs were extracted using a Qiagen RNeasy kit. The PacBio Iso-Seq libraries were constructed using a PacBio Iso-Seq Express 2.0 kit. Libraries were either sized (0.66x bead ratio) or unsized (1.2x bead ratio) to give final libraries with average transcript sizes of 2kb or 3kb, respectively. Libraries were sequenced using polymerase

V2.1 on a PacBio Sequel II Platform. The RNASeq libraries were constructed using an Illumina TruSeq Stranded mRNA Library Prep Kit using standard protocols and were sequenced using a NovaSeq 6000 Instrument PE150 to 40 million reads per library.

To identify the sex chromosomes, we additionally sequenced the whole genomes of 52 *Amborella* individuals sampled from natural populations (Supplementary Table 11). DNA extractions were performed using a standard CTAB protocol. Illumina sequencing was performed on NovaSeq and HiSeq platforms at RAPID Genomics in Gainesville, Florida, using a 2x150 paired-end library. The voucher specimens are deposited at the New Caledonia Herbarium in Nouméa (Herbarium code: NOU) and Indiana University (IND). Existing data used to support this manuscript are found in Supplementary Table 11.

### Genome assembly

The version 2.0 HAP1 and HAP2 assemblies were generated by assembling the 3,605,703 PacBio CCS reads (58.81x per haplotype) using the HiFiAsm+HIC assembler<sup>75</sup> and subsequently polished using RACON<sup>76</sup>. This approach produced initial assemblies of both haplotypes. The HAP1 assembly consisted of 1,522 scaffolds (1,522 contigs), with a contig N50 of 25.5 Mb, and a total genome size of 800.6 Mb (Supplementary Table 13). The HAP2 assembly consisted of 1,043 scaffolds (1,043 contigs), with a contig N50 of 43.0 Mb, and a total genome size of 773.5 Mb (Supplementary Table 13).

Hi-C Illumina reads from *Amborella trichopoda* isolate Santa Cruz 75 were separately aligned to the HAP1 and HAP2 contig sets with Juicer<sup>77</sup>, and chromosome-scale scaffolding was performed with 3D-DNA<sup>78</sup>. No misjoins were identified in either the HAP1 or HAP2 assemblies. The contigs were then oriented, ordered, and joined together into 13 chromosomes per haplotype using the Hi-C data. A total of 31 joins was applied to the HAP1 assembly, and 20 joins for the HAP2 assembly. Each chromosome join is padded with 10,000 Ns. Contigs terminating in significant telomeric sequence were identified using the (TTTAGGG)<sub>n</sub> repeat, and care was taken to make sure that the repeats were properly oriented in the production assembly. The remaining scaffolds were screened against bacterial proteins, organelle sequences, and GenBank non-redundant database, and any scaffold found to be a contaminant was removed. After the chromosomes were formed, it was observed that some small (<20Kb) redundant sequences were present on adjacent contig ends within chromosomes. To resolve this issue, adjacent contig ends were aligned to one another using BLAT<sup>79</sup>, and duplicate sequences were collapsed to close the gap between them. A total of 5 adjacent contig pairs were collapsed in the HAP1 assembly and 4 in the HAP2 assembly.

Finally, homozygous SNPs and INDELS were corrected in the HAP1 and HAP2 releases using ~49x of Illumina reads (2x150, 400-bp insert) by aligning the reads using BWA-MEM<sup>80</sup> and identifying homozygous SNPs and INDELS with GATK's UnifiedGenotyper tool<sup>81</sup>. A total of 465 homozygous SNPs and 15,763 homozygous INDELS were corrected in the HAP1 release, while a total of 473 homozygous SNPs and 17,208 homozygous INDELS were corrected in the HAP2 release. The final version 2.0 HAP1 release contained 707.9 Mb of sequence, consisting

of 59 contigs with a contig N50 of 36.3 Mb and a total of 99.69% of assembled bases in chromosomes. The final version 2.0 HAP2 release contained 700.3 Mb of sequence, consisting of 45 contigs with a contig N50 of 44.5 Mb and a total of 99.87% of assembled bases in chromosomes.

### **Genome annotation**

Transcript assemblies were made from ~757M pairs of 2x150 stranded paired-end Illumina RNAseq reads using PERTRAN, which conducts genome-guided transcriptome short-read assembly via GSNAP<sup>82</sup> and builds splice alignment graphs after alignment validation, realignment, and correction. To obtain 825K putative full-length transcripts, about 20M PacBio Iso-Seq CCSs were corrected and collapsed by a genome-guided correction pipeline, which aligns CCS reads to the genome with GMAP<sup>82</sup> with intron correction for small indels in splice junctions, if any, and clusters alignments when all introns are the same or have 95% overlap for a single exon. Subsequently 563,694 transcript assemblies were constructed using PASA<sup>83</sup> from ESTs and RNAseq transcript assemblies described above. Loci were determined by transcript assembly alignments and/or EXONERATE alignments of proteins from *Arabidopsis thaliana*, *Glycine max*, *Sorghum bicolor*, *Oryza sativa*, *Lactuca sativa*, *Helianthus annuus*, *Cynara cardunculus*, *Selaginella moellendorffii*, *Physcomitrella patens*, *Nymphaea colorata*, *Solanum lycopersicum*, and *Vitis vinifera*, and Swiss-Prot eukaryote proteomes to the repeat-soft-masked *Amborella trichopoda* HAP1 genome using RepeatMasker<sup>84</sup> with up to 2kb extension on both ends unless extending into another locus on the same strand. Gene models were predicted by homology-based predictors, FGENESH+<sup>85</sup>, FGENESH\_EST (similar to FGENESH+, but using EST to compute splice site and intron input instead of protein/translated ORF), EXONERATE<sup>86</sup>, PASA assembly ORFs (in-house homology-constrained ORF finder), and AUGUSTUS<sup>87</sup> trained by the high-confidence PASA assembly ORFs and with intron hints from short-read alignments. The best-scored predictions for each locus were selected using multiple positive factors, including EST and protein support, and one negative factor: overlap with repeats. The selected gene predictions were improved by PASA, and the optimal set was selected using several curated gene quality metrics<sup>88</sup>. We assessed the gene annotations using compleasm v0.2.6<sup>89</sup> using the Embryophyta database.

We further annotated repeats using EDTA v2.0.0<sup>90</sup> using the sensitive mode that runs RepeatModeler<sup>91</sup>. To identify tandem repeats, we used Tandem Repeats Finder v.4.09.1<sup>92</sup> (parameters 2 7 7 80 10 50 500 -f -d -m -h). We ran StainedGlass v0.5<sup>93</sup> to visualize the massive tandem repeat arrays for chromosomes in both haplotypes. To build the repeat landscapes for assessing recent expansion events, we followed the methods outlined in EDTA Github Issue #92: Draw Repeat Landscapes, utilizing a library generated from an independent annotation on the combined haplotypes with EDTA v2.0.1.

### **Comparisons between assembly haplotypes**

To plot comparisons between the two haplotypes, including genes and repeats, we used GENESPACE v.1.3.1<sup>94</sup>. To generate synteny between the two haplotypes, we first performed genome alignments. HAP1 and HAP2 were aligned using AnchorWave v1.0.1<sup>95</sup> using the 'genoAli' method and '-IV' parameter to allow for inversions. Alignment was performed using only "chromosome" sequence for each haplotype. The alignment was converted to SAM format using the 'maf-convert' tool provided in 'last' v460<sup>96</sup> and used for calling variants with SyRI v1.6.3<sup>97</sup>. The output from SyRI was used to make chromosome-level synteny and SV plots using plotsr v0.5.4<sup>98</sup>.

### **Identification of the sex chromosome non-recombining region**

We used whole-genome sequencing data to identify the sex-determining region (SDR) of the W. All paired-end Illumina data had adapters removed and were quality filtered using TRIMMOMATIC v0.39<sup>99</sup> with leading and trailing values of 3, sliding window of 30, jump of 10, and a minimum remaining read length of 40. We next found all canonical 21-mers in each isolate using Jellyfish v2.3.0<sup>100</sup> and used the bash *comm* command to find all *k*-mers shared in all female isolates and not found in any male isolate (W-mers). We mapped the W-mers to both haplotype assemblies using BWA-MEM v0.7.17<sup>80</sup>, with parameters '-k 21' '-T 21' '-a' '-c 10'. W-mer mapping was visualized by first calculating coverage in 100,000-bp sliding windows (10,000 bp jump) using BEDTools v2.28.0<sup>101</sup> and plotted using karyoploteR v1.26.0<sup>102</sup>.

### **Structural variation**

To identify structural variants between the haplotypes, we mapped PacBio reads using minimap2 v2.24<sup>103</sup> in HiFi mode, added the MD tag using samtools v1.10 *calmd*, and called structural variants using Sniffles v2.0.7<sup>104</sup>. We also performed whole-genome alignments using minimap2 v2.24<sup>103</sup> and visualized the dotplot using pafR v0.0.2<sup>105</sup>.

### **Gene homology and protein evolution**

To identify one-to-one orthologs on the ZW to examine protein evolution, we ran OrthoFinder v2.5.2<sup>106,107</sup> using only the *Amborella* haplotypes. We calculated synonymous (Ks) and nonsynonymous (Ka) changes in codons using Ka/Ks Calculator v2.0<sup>108</sup>.

To identify the boundaries of evolutionary strata, we used the R package mcp v0.3.4<sup>109</sup> on dXY and Ks. For Ks, we first ran a test for outliers using PMCMRplus v1.9.10<sup>110</sup> to run Rosner's generalized extreme studentized deviate many-outlier test<sup>111</sup>. For mcp, we used the model, 'y ~ 1, ~ 1' to identify the change point between two plateaus, and we used 100,000 iterations, 3 chains, and a burn-in of 100,000 (i.e., 'adapt').

### **Nucleotide differences between the sexes**

BWA v0.7.17<sup>80</sup> was used to map reads, and bcftools v1.9 *mpileup* and *call*<sup>112</sup> functions were used to call variants using the Island-wide sampling (nine male and six female plants; Supplementary Table 11). We filtered the vcf file using 'QUAL>20 & DP>5 & MQ>30', minor

allele frequency of 0.05, and dropped sites with > 25% missing data. To calculate Nei's nucleotide diversity between the sexes (dXY), we used pixy v1.2.7.beta1<sup>113</sup>. dXY was calculated using 100,000-bp windows with a 10,000-bp jump, and on the gene models only separately.

### **Presence-absence variation**

Presence-absence variation (PAV) was identified following the methods of Hu et al.<sup>114</sup> mapping reads from the Island-wide sampling (eight male and six female plants; the Atlanta Botanical Gardens isolate was removed due to low resequencing depth; Supplementary Table 11) to our new reference genome and annotation. Briefly, reads for the samples were aligned to each haplotype using BWA v0.7.17<sup>80</sup>. Sorted BAM files were converted to bedgraph format using bedtools v2.30.0<sup>101</sup>. Genes were called absent if the horizontal coverage of exons was <5% and the average depth was <2x. A test for equality in the proportion of PAV rate across chromosomes was performed in R using the 'prop.test()' function.

### **Gene expression analyses**

To examine gene expression and identify candidate sex-determining genes, we used existing RNAseq data from 10 females and 10 males<sup>11</sup>. From the reads, we first filtered using TRIMMOMATIC (same parameters as above). Filtered reads were mapped to the HAP1 genome assembly using STAR v2.7.9a<sup>115</sup> and expression estimated for the annotated gene models using StringTie v2.1.7 (-e, -G)<sup>116</sup>. We performed differential gene expression analyses using DESeq2 v1.32.0<sup>117</sup>, with the contrast being between the sexes.

## **DATA AVAILABILITY**

The genome assemblies and annotations (v.2.1) are available on Phytozome v13 (<https://phytozome-next.jgi.doe.gov/>) and have been deposited on NCBI under BioProjects PRJNA1100625 and PRJNA1167780. Sequencing libraries for the genome assembly and annotation are publicly available on NCBI under BioProject PRJNA1100625 and the whole-genome sequencing of additional isolates are under PRJNA1161132. Individual accession numbers are provided in Supplementary Tables S10-11.

## **ACKNOWLEDGEMENTS**

The work (proposal no. 10.46936/10.25585/60001405) conducted by the U.S. Department of Energy (DOE) Joint Genome Institute (<https://ror.org/04xm1d337>), a DOE Office of Science User Facility, is supported under contract no. DE-AC02-05CH11231. Additional support for analysis was provided by the United States Department of Agriculture National Institute of Food and Agriculture Postdoctoral Fellowship no. 2022-67012-38987 (S.B.C.), National Science Foundation (NSF) IOS-PGRP CAREER no. 2239530 (A.H.), and National Science Foundation GRFP (L.A.). We thank the Atlanta Botanical Garden for providing *Amborella* material used in this study and Adam Bewick for the images of *Amborella* flowers.

### **AUTHOR CONTRIBUTIONS**

Concept and research design: S.B.C., J.S., J.L.-M., A.H.

Sample collection, data collection, sequencing: A.S., T.J., K.B., S.R., J.T., P.P.L., J.M., E.B.K., D.E.S., P.S.S., J.G., J.L.-M.

Genome assembly and annotation: S.B.C., J.J., S.S.

Computational and statistical analyses: S.B.C., L.A., J.T.L., A.L.H., P.G., A.Y.

Wrote the paper (with contributions from all authors): S.B.C., L.A., J.T.L., J.J., A.L.H., C.S., D.E.S., P.S.S., J.L.-M., A.H.

### **COMPETING INTERESTS STATEMENT**

The authors declare no competing interests.



## FIGURE LEGENDS

**Fig. 1. *Amborella* and its genome structure.** A-B) Female and male *Amborella* flowers. The *Amborella* genome (C) and chromosome 9 (D) are typical of flowering plants: gene-rich chromosome arms and repeat-dense, large pericentromeric regions. Gene positions were extracted from the protein-coding gene annotations, repeats from EDTA, and exact matches of 536,985 female-specific *k*-mers (W-mers). Syntenic mapping was calculated by AnchorWave and processed by SyRI, only plotting inversions and insertions and deletions > 10 kb. Visualization of synteny was accomplished with GENESPACE and sliding windows with gscTools. Panel B highlights the sex determination region of chromosome 9 with W-mers. All chromosomes in haplotype 1 and all but four in haplotype 2 have both left and right telomeres in the assembly (flagged with red \*), defined as a region of  $\geq 150$  bp made up of  $\geq 90\%$  plant telomere *k*-mers (CCCGAAA, CCCTAAA, RC) separated by no more than 100 bp.

**Fig 2. Sex chromosome location in *Amborella*.** A) W-mer coverage in the sex-determining region (SDR) and (B) homologous region of the Z (HZR) using four different sampling strategies for isolates. SDR (C) and HZR (D) location and their proximity to the Chr09 centromere. Ty3 elements (dark blue) are often enriched in the pericentromeric regions of plants and correspond to the low-complexity block of tandem repeat arrays (gray) that also contain the high-complexity centromeric block, indicated by the satellite monomer density (light blue). Gene density (orange) also predictably decreases near the pericentromeric region. The SDR (red) is notably outside of the putative pericentromeric region and distant from the centromere.

**Fig. 3. Molecular evolution of the *Amborella* sex chromosomes.** A) Evidence for two strata. For Ks, points above 0.06 were excluded. B-C) The repeat landscapes of the *Amborella* haplotypes indicate similar patterns of expansion and minimal evidence of recent TE proliferation. Relative time is determined by the Kimura substitution level with lower values closer to 0 representing more recent events and higher values approaching 40 representing older events.

## REFERENCES

1. Renner, S. S. The relative and absolute frequencies of angiosperm sexual systems: dioecy, monoecy, gynodioecy, and an updated online database. *Am. J. Bot.* **101**, 1588–1596 (2014).
2. Carey, S., Yu, Q. & Harkess, A. The Diversity of Plant Sex Chromosomes Highlighted through Advances in Genome Sequencing. *Genes* **12**, (2021).
3. Renner, S. S. & Müller, N. A. Plant sex chromosomes defy evolutionary models of expanding recombination suppression and genetic degeneration. *Nat Plants* **7**, 392–402 (2021).
4. Soltis, P. S., Soltis, D. E. & Chase, M. W. Angiosperm phylogeny inferred from multiple genes as a tool for comparative biology. *Nature* **402**, 402–404 (1999).
5. Moore, M. J., Bell, C. D., Soltis, P. S. & Soltis, D. E. Using plastid genome-scale data to resolve enigmatic relationships among basal angiosperms. *Proc. Natl. Acad. Sci. U. S. A.* **104**, 19363–19368 (2007).
6. Burleigh, J. G. *et al.* Genome-scale phylogenetics: inferring the plant tree of life from 18,896 gene trees. *Syst. Biol.* **60**, 117–125 (2011).
7. Soltis, D. E. *et al.* Angiosperm phylogeny: 17 genes, 640 taxa. *Am. J. Bot.* **98**, 704–730 (2011).
8. One Thousand Plant Transcriptomes Initiative. One thousand plant transcriptomes and the phylogenomics of green plants. *Nature* **574**, 679–685 (2019).
9. Sauquet, H. *et al.* The ancestral flower of angiosperms and its early diversification. *Nat. Commun.* **8**, 16047 (2017).
10. Anger, N., Fogliani, B., Scutt, C. P. & Gâteblé, G. Dioecy in *Amborella trichopoda*: evidence for genetically based sex determination and its consequences for inferences of the

- breeding system in early angiosperms. *Ann. Bot.* **119**, 591–597 (2017).
11. Käfer, J. *et al.* A derived ZW chromosome system in *Amborella trichopoda*, representing the sister lineage to all other extant flowering plants. *New Phytol.* **233**, 1636–1642 (2022).
  12. Akagi, T., Henry, I. M., Tao, R. & Comai, L. A Y-chromosome–encoded small RNA acts as a sex determinant in persimmons. *Science* (2014).
  13. Torres, M. F. *et al.* Genus-wide sequencing supports a two-locus model for sex-determination in Phoenix. *Nat. Commun.* **9**, 3969 (2018).
  14. Akagi, T. *et al.* Two Y-chromosome-encoded genes determine sex in kiwifruit. *Nat Plants* **5**, 801–809 (2019).
  15. Harkess, A. *et al.* Sex Determination by Two Y-Linked Genes in Garden Asparagus. *Plant Cell* **32**, 1790–1796 (2020).
  16. Kazama, Y. *et al.* A CLAVATA3-like Gene Acts as a Gynoecium Suppression Function in White Campion. *Mol. Biol. Evol.* **39**, (2022).
  17. Müller, N. A. *et al.* A single gene underlies the dynamic evolution of poplar sex determination. *Nat Plants* **6**, 630–637 (2020).
  18. Amborella Genome Project. The Amborella genome and the evolution of flowering plants. *Science* **342**, 1241089 (2013).
  19. Oginuma, K., Jaffré, T. & Tobe, H. The Karyotype Analysis of Somatic Chromosomes in *Amborella trichopoda* (Amborellaceae). *J. Plant Res.* **113**, 281–283 (2000).
  20. Rhie, A., Walenz, B. P., Koren, S. & Phillippy, A. M. Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies. *Genome Biol.* **21**, 245 (2020).
  21. Magallón, S., Gómez-Acevedo, S., Sánchez-Reyes, L. L. & Hernández-Hernández, T. A

- metacalibrated time-tree documents the early rise of flowering plant phylogenetic diversity. *New Phytol.* **207**, 437–453 (2015).
22. Marchant, D. B. *et al.* Dynamic genome evolution in a model fern. *Nat Plants* **8**, 1038–1051 (2022).
  23. Niu, S. *et al.* The Chinese pine genome and methylome unveil key features of conifer evolution. *Cell* **185**, 204–217.e14 (2022).
  24. Healey, A. L. *et al.* Newly identified sex chromosomes in the Sphagnum (peat moss) genome alter carbon sequestration and ecosystem dynamics. *Nat Plants* **9**, 238–254 (2023).
  25. Neumann, P. *et al.* Plant centromeric retrotransposons: a structural and cytogenetic perspective. *Mob. DNA* **2**, 4 (2011).
  26. Sigman, M. J. & Slotkin, R. K. The first rule of plant transposable element silencing: location, location, location. *Plant Cell* (2016).
  27. Lappin, F. M. *et al.* A polymorphic pseudoautosomal boundary in the *Carica papaya* sex chromosomes. *Mol. Genet. Genomics* **290**, 1511–1522 (2015).
  28. Cotter, D. J., Brotman, S. M. & Wilson Sayres, M. A. Genetic Diversity on the Human X Chromosome Does Not Support a Strict Pseudoautosomal Boundary. *Genetics* **203**, 485–492 (2016).
  29. Palmer, D. H., Rogers, T. F., Dean, R. & Wright, A. E. How to identify sex chromosomes and their turnover. *Mol. Ecol.* **28**, 4709–4724 (2019).
  30. Tennessen, J. A. *et al.* Repeated translocation of a gene cassette drives sex-chromosome turnover in strawberries. *PLoS Biol.* **16**, e2006062 (2018).
  31. Yu, Q. *et al.* A physical map of the papaya genome with integrated genetic map and genome sequence. *BMC Genomics* **10**, 371 (2009).

32. Lahn, B. T. & Page, D. C. Four evolutionary strata on the human X chromosome. *Science* **286**, 964–967 (1999).
33. Rice, W. R. THE ACCUMULATION OF SEXUALLY ANTAGONISTIC GENES AS A SELECTIVE AGENT PROMOTING THE EVOLUTION OF REDUCED RECOMBINATION BETWEEN PRIMITIVE SEX CHROMOSOMES. *Evolution* **41**, 911–914 (1987).
34. Charlesworth, D., Charlesworth, B. & Marais, G. Steps in the evolution of heteromorphic sex chromosomes. *Heredity* **95**, 118–128 (2005).
35. Papadopulos, A. S. T., Chester, M., Ridout, K. & Filatov, D. A. Rapid Y degeneration and dosage compensation in plant sex chromosomes. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 13021–13026 (2015).
36. Wu, M. & Moore, R. C. The Evolutionary Tempo of Sex Chromosome Degradation in *Carica papaya*. *J. Mol. Evol.* **80**, 265–277 (2015).
37. Hobza, R. *et al.* Impact of Repetitive Elements on the Y Chromosome Formation in Plants. *Genes* **8**, (2017).
38. Sacchi, B. *et al.* Phased assembly of neo-sex chromosomes reveals extensive Y degeneration and rapid genome evolution in *Rumex hastatulus*. *bioRxiv* 2023.09.26.559509 (2023) doi:10.1101/2023.09.26.559509.
39. Jedlicka, P., Lexa, M. & Kejnovsky, E. What Can Long Terminal Repeats Tell Us About the Age of LTR Retrotransposons, Gene Conversion and Ectopic Recombination? *Front. Plant Sci.* **11**, 644 (2020).
40. Cornet, C. *et al.* Holocentric repeat landscapes: From micro-evolutionary patterns to macro-evolutionary associations with karyotype evolution. *Mol. Ecol.* (2023)

doi:10.1111/mec.17100.

41. Bachtrog, D. Y-chromosome evolution: emerging insights into processes of Y-chromosome degeneration. *Nat. Rev. Genet.* **14**, 113–124 (2013).
42. Charlesworth, D. The timing of genetic degeneration of sex chromosomes. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **376**, 20200093 (2021).
43. Hibbins, M. S. *et al.* Phylogenomics resolves key relationships in *Rumex* and uncovers a dynamic history of independently evolving sex chromosomes. *bioRxiv* 2023.12.13.571571 (2023) doi:10.1101/2023.12.13.571571.
44. Krasovec, M., Chester, M., Ridout, K. & Filatov, D. A. The Mutation Rate and the Age of the Sex Chromosomes in *Silene latifolia*. *Curr. Biol.* **28**, 1832–1838.e4 (2018).
45. Akagi, T. *et al.* Rapid and dynamic evolution of a giant Y chromosome in *Silene latifolia*. *bioRxiv* 2023.09.21.558759 (2023) doi:10.1101/2023.09.21.558759.
46. Ma, X. *et al.* The spinach YY genome reveals sex chromosome evolution, domestication, and introgression history of the species. *Genome Biol.* **23**, 75 (2022).
47. She, H. *et al.* Evolution of the spinach sex-linked region within a rarely recombining pericentromeric region. *Plant Physiol.* **193**, 1263–1280 (2023).
48. Bachtrog, D. *et al.* Are all sex chromosomes created equal? *Trends Genet.* **27**, 350–357 (2011).
49. Charlesworth, B. & Charlesworth, D. A model for the evolution of dioecy and gynodioecy. *Am. Nat.* (1978).
50. Jay, P., Tezenas, E., Véber, A. & Giraud, T. Sheltering of deleterious mutations explains the stepwise extension of recombination suppression on sex chromosomes and other supergenes. *PLoS Biol.* **20**, e3001698 (2022).

51. Lenormand, T. & Roze, D. Y recombination arrest and degeneration in the absence of sexual dimorphism. *Science* **375**, 663–666 (2022).
52. Buzgo, M., Soltis, P. S. & Soltis, D. E. Floral Developmental Morphology of *Amborella trichopoda* (Amborellaceae). *Int. J. Plant Sci.* **165**, 925–947 (2004).
53. Flores-Tornero, M. *et al.* Transcriptomic and Proteomic Insights into *Amborella trichopoda* Male Gametophyte Functions. *Plant Physiol.* **184**, 1640–1657 (2020).
54. Julca, I. *et al.* Comparative transcriptomic analysis reveals conserved programmes underpinning organogenesis and reproduction in land plants. *Nat Plants* **7**, 1143–1159 (2021).
55. Deyhle, F., Sarkar, A. K., Tucker, E. J. & Laux, T. WUSCHEL regulates cell differentiation during anther development. *Dev. Biol.* **302**, 154–159 (2007).
56. Zúñiga-Mayo, V. M., Gómez-Felipe, A., Herrera-Ubaldo, H. & de Folter, S. Gynoecium development: networks in *Arabidopsis* and beyond. *J. Exp. Bot.* **70**, 1447–1460 (2019).
57. Schoof, H. *et al.* The stem cell population of *Arabidopsis* shoot meristems is maintained by a regulatory loop between the *CLAVATA* and *WUSCHEL* genes. *Cell* **100**, 635–644 (2000).
58. Parvathy, S. T., Prabakaran, A. J. & Jayakrishna, T. Author Correction: Probing the floral developmental stages, bisexuality and sex reversions in castor (*Ricinus communis* L.). *Sci. Rep.* **11**, 10504 (2021).
59. Zhang, S. *et al.* The control of carpel determinacy pathway leads to sex determination in cucurbits. *Science* **378**, 543–549 (2022).
60. Schlegel, J. *et al.* Control of *Arabidopsis* shoot stem cell homeostasis by two antagonistic CLE peptide signalling pathways. (2021) doi:10.7554/eLife.70934.
61. Kurakawa, T. *et al.* Direct control of shoot meristem activity by a cytokinin-activating

- enzyme. *Nature* **445**, 652–655 (2007).
62. Hardtke, C. S. & Berleth, T. The Arabidopsis gene MONOPTEROS encodes a transcription factor mediating embryo axis formation and vascular development. *EMBO J.* **17**, 1405–1411 (1998).
  63. Aida, M., Vernoux, T., Furutani, M., Traas, J. & Tasaka, M. Roles of PIN-FORMED1 and MONOPTEROS in pattern formation of the apical region of the Arabidopsis embryo. *Development* **129**, 3965–3974 (2002).
  64. Stortenbeker, N. & Bemer, M. The SAUR gene family: the plant’s toolbox for adaptation of growth and development. *J. Exp. Bot.* **70**, 17–27 (2019).
  65. He, S.-L., Hsieh, H.-L. & Jauh, G.-Y. SMALL AUXIN UP RNA62/75 Are Required for the Translation of Transcripts Essential for Pollen Tube Growth. *Plant Physiol.* **178**, 626–640 (2018).
  66. Chae, K. *et al.* Arabidopsis SMALL AUXIN UP RNA63 promotes hypocotyl and stamen filament elongation. *Plant J.* **71**, 684–697 (2012).
  67. van Mourik, H., van Dijk, A. D. J., Stortenbeker, N., Angenent, G. C. & Bemer, M. Divergent regulation of Arabidopsis SAUR genes: a focus on the SAUR10-clade. *BMC Plant Biol.* **17**, 245 (2017).
  68. Bürstenbinder, K. *et al.* Inhibition of 5’-methylthioadenosine metabolism in the Yang cycle alters polyamine levels, and impairs seedling growth and reproduction in Arabidopsis. *Plant J.* **62**, 977–988 (2010).
  69. Waduwara-Jayabahu, I. *et al.* Recycling of methylthioadenosine is essential for normal vascular development and reproduction in Arabidopsis. *Plant Physiol.* **158**, 1728–1744 (2012).



70. Nurk, S. *et al.* The complete sequence of a human genome. *Science* **376**, 44–53 (2022).
71. Rhie, A. *et al.* The complete sequence of a human Y chromosome. *bioRxiv* 2022.12.01.518724 (2022) doi:10.1101/2022.12.01.518724.
72. Rhie, A. *et al.* Towards complete and error-free genome assemblies of all vertebrate species. *Nature* **592**, 737–746 (2021).
73. Zhu, Z., Younas, L. & Zhou, Q. Evolution and regulation of animal sex chromosomes. *Nat. Rev. Genet.* 1–16 (2024) doi:10.1038/s41576-024-00757-3.
74. Doyle, J. J. & Doyle, J. L. A rapid DNA isolation procedure for small quantities of fresh leaf tissue. *Phytochemical bulletin* (1987).
75. Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat. Methods* **18**, 170–175 (2021).
76. Vaser, R., Sović, I., Nagarajan, N. & Šikić, M. Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Res.* **27**, 737–746 (2017).
77. Durand, N. C. *et al.* Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell Syst* **3**, 95–98 (2016).
78. Dudchenko, O. *et al.* De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds. *Science* **356**, 92–95 (2017).
79. Kent, W. J. BLAT—The BLAST-Like Alignment Tool. *Genome Res.* **12**, 656–664 (2002).
80. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv [q-bio.GN]* (2013).
81. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
82. Wu, T. D. & Nacu, S. Fast and SNP-tolerant detection of complex variants and splicing in

- short reads. *Bioinformatics* **26**, 873–881 (2010).
83. Haas, B. J. *et al.* Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. *Nucleic Acids Res.* **31**, 5654–5666 (2003).
  84. Smit, A. F. A., Hubley, R. & Green, P. RepeatModeler Open-1.0. 2008--2015. *Seattle, USA: Institute for Systems Biology. Available from: <httpwww.repeatmasker.org>, Last Accessed May 1, 2018* (2015).
  85. Salamov, A. A. & Solovyev, V. V. Ab initio gene finding in Drosophila genomic DNA. *Genome Res.* **10**, 516–522 (2000).
  86. Slater, G. S. C. & Birney, E. Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics* **6**, 31 (2005).
  87. Stanke, M., Schöffmann, O., Morgenstern, B. & Waack, S. Gene prediction in eukaryotes with a generalized hidden Markov model that uses hints from external sources. *BMC Bioinformatics* **7**, 62 (2006).
  88. Lovell, J. T. *et al.* The genomic landscape of molecular responses to natural drought stress in *Panicum hallii*. *Nat. Commun.* **9**, 5213 (2018).
  89. Huang, N. & Li, H. compleasm: a faster and more accurate reimplement of BUSCO. *Bioinformatics* **39**, (2023).
  90. Ou, S. *et al.* Benchmarking transposable element annotation methods for creation of a streamlined, comprehensive pipeline. *Genome Biol.* **20**, 275 (2019).
  91. Flynn, J. M. *et al.* RepeatModeler2 for automated genomic discovery of transposable element families. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 9451–9457 (2020).
  92. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.* **27**, 573–580 (1999).

93. Vollger, M. R., Kerpedjiev, P., Phillippy, A. M. & Eichler, E. E. StainedGlass: interactive visualization of massive tandem repeat structures with identity heatmaps. *Bioinformatics* **38**, 2049–2051 (2022).
94. Lovell, J. T. *et al.* GENESPACE tracks regions of interest and gene copy number variation across multiple genomes. *Elife* **11**, (2022).
95. Song, B. *et al.* AnchorWave: Sensitive alignment of genomes with high sequence diversity, extensive structural polymorphism, and whole-genome duplication. *Proc. Natl. Acad. Sci. U. S. A.* **119**, (2022).
96. Kielbasa, S. M., Wan, R., Sato, K., Horton, P. & Frith, M. C. Adaptive seeds tame genomic sequence comparison. *Genome Res.* **21**, 487–493 (2011).
97. Goel, M., Sun, H., Jiao, W.-B. & Schneeberger, K. SyRI: finding genomic rearrangements and local sequence differences from whole-genome assemblies. *Genome Biol.* **20**, 277 (2019).
98. Goel, M. & Schneeberger, K. plotsr: visualizing structural similarities and rearrangements between multiple genomes. *Bioinformatics* **38**, 2922–2926 (2022).
99. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
100. Marcais, G. & Kingsford, C. Jellyfish: A fast k-mer counter. *Tutorialis e Manuais* 1–8 (2012).
101. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
102. Gel, B. & Serra, E. karyoploteR: an R/Bioconductor package to plot customizable genomes displaying arbitrary data. *Bioinformatics* **33**, 3088–3090 (2017).

103. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
104. Sedlazeck, F. J. *et al.* Accurate detection of complex structural variations using single-molecule sequencing. *Nat. Methods* **15**, 461–468 (2018).
105. Winter, D., Lee, K. & Cox, M. pafr: read, manipulate and visualize ‘Pairwise mApping Format’ data. *The Comprehensive R Archive Network* (2020).
106. Emms, D. M. & Kelly, S. OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome Biol.* **16**, 157 (2015).
107. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biol.* **20**, 238 (2019).
108. Zhang, Z. *et al.* KaKs\_Calculator: calculating Ka and Ks through model selection and model averaging. *Genomics Proteomics Bioinformatics* **4**, 259–263 (2006).
109. Lindeløv, J. K. mcp: An R Package for Regression With Multiple Change Points. (2020) doi:10.31219/osf.io/fzqxv.
110. Pohlert, T. & Pohlert, M. T. Package ‘pmcmr’. *R package version 1*, (2018).
111. Rosner, B. Percentage Points for a Generalized ESD Many-Outlier Procedure. *Technometrics* **25**, 165–172 (1983).
112. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).
113. Korunes, K. L. & Samuk, K. pixy: Unbiased estimation of nucleotide diversity and divergence in the presence of missing data. *Mol. Ecol. Resour.* **21**, 1359–1368 (2021).

114. Hu, H. *et al.* Amborella gene presence/absence variation is associated with abiotic stress responses that may contribute to environmental adaptation. *New Phytol.* **233**, 1548–1555 (2022).
115. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
116. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat. Biotechnol.* **33**, 290–295 (2015).
117. Love, M., Anders, S. & Huber, W. Differential analysis of count data--the DESeq2 package. *Genome Biol.* **15**, 10–1186 (2014).