## UC Irvine UC Irvine Previously Published Works

## Title

The genomic and epigenomic landscapes of hemizygous genes across crops with contrasting reproductive systems.

## Permalink

https://escholarship.org/uc/item/7hf835rk

### Journal

Proceedings of the National Academy of Sciences of the United States of America, 122(6)

## **Authors**

Peng, Yanling Wang, Yiwen Liu, Yuting <u>et al.</u>

Publication Date

2025-02-11

## DOI

10.1073/pnas.2422487122

Peer reviewed



# The genomic and epigenomic landscapes of hemizygous genes across crops with contrasting reproductive systems

Yanling Peng<sup>a</sup>, Yiwen Wang<sup>a</sup>, Yuting Liu<sup>a</sup>, Xinyue Fang<sup>a</sup>, Lin Cheng<sup>a</sup>, Qiming Long<sup>a</sup>, Dalu Su<sup>a</sup>, Tianhao Zhang<sup>a</sup>, Xiaoya Shi<sup>a</sup>, Xiaodong Xu<sup>a</sup>, Qi Xu<sup>a</sup>, Nan Wang<sup>a</sup>, Fan Zhang<sup>a</sup>, Zhongjie Liu<sup>a</sup>, Hua Xiao<sup>a</sup>, Jin Yao<sup>b</sup>, Ling Tian<sup>b</sup>, Wei Hu<sup>c</sup>, Songbi Chen<sup>c</sup>, Haibo Wang<sup>d</sup>, Sanwen Huang<sup>a,c</sup>, Randon S. Gaut<sup>e,1</sup>, and Yongfeng Zhou<sup>a,c,1</sup>

Affiliations are included on p. 11.

Edited by Rod Wing, The University of Arizona, Tucson, AZ; received November 8, 2024; accepted January 6, 2025 by Editorial Board Member Douglas E. Soltis

Hemizygous genes, which are present on only one of the two homologous chromosomes of diploid organisms, have been mainly studied in the context of sex chromosomes and sex-linked genes. However, these genes can also occur on the autosomes of diploid plants due to structural variants (SVs), such as a deletion/insertion of one allele, and this phenomenon largely unexplored in plants. Here, we investigated the genomic and epigenomic landscapes of hemizygous genes across 22 genomes with varying propagation histories: eleven clonal lineages, seven outcrossed samples, and four inbred and putatively homozygous genomes. We identified SVs leading to genic hemizygosity. As expected, very few genes (0.01 to 1.2%) were hemizygous in the homozygous genomes, representing negative controls. Hemizygosity was appreciable among outcrossed lineages, averaging 8.7% of genes, but consistently elevated for the clonal samples at 13.8% genes, likely reflecting heterozygous SV accumulation during clonal propagation. Compared to diploid genes, hemizygous genes were more often situated in centromeric than telomeric regions and experienced weaker purifying selection. They also had reduced levels of expression, averaging  $\sim 20\%$  of the expression levels of diploid genes, violating the evolutionary model of dosage compensation. We also detected higher DNA methylation levels in hemizygous genes and transposable elements, which may contribute to their reduced expression. Finally, expression profiles showed that hemizygous genes were more specifically expressed in contexts related to fruit development, organ differentiation, and stress responses. Overall, hemizygous genes accumulate in clonally propagated lineages and display distinct genetic and epigenetic features compared to diploid genes, shedding unique insights into genetic studies and breeding programs of clonal crops.

grapevine | clonal propagation | structural variation | integrative genomics | heterozygosity

Hemizygous genes are present on only one of the two homologous chromosomes of a diploid organism (1-3). The most prominent examples of hemizygous genes are on the sex chromosomes of male mammals (XY) or female birds (ZW) (1-3). Similar sets of hemizygous genes are present in the sex-linked regions of dioecious plants with X/Y sex determination, such as palms (4, 5), asparagus (6, 7), kiwifruit (8, 9). Numerous studies have focused on the evolution, gene expression, and epigenetic regulation of sex-linked hemizygous genes compared to diploid genes (2, 10-15). For example, genomic studies in mammalian males have consistently revealed lower mutation rates and more efficient selection in sex-linked hemizygous genes than diploid genes (10), due in part to the fact that hemizygosity uncovers recessive alleles and makes them visible to selection (16). Similar studies have generally shown that the ratio of sex-linked to diploid (X:AA) (where X represents sex-linked genes and AA represents autosome genes) gene expression is  $\sim 0.5$  in animals and plants (11, 17). This ratio is inconsistent with the hypothesis that dosage compensation re-equalizes male and female expression to restore XY male expression back to its ancestral level (17).

Interestingly, estimated expression levels of XY sex chromosome alleles in males show an overall trend of reduced expression of Y-linked alleles relative to X-linked alleles in both animals and plants (12–15). Some of these expression effects are mediated by epigenetic marks, including histone modifications and DNA methylation (18–21). For example, the male-specific region of the papaya Y chromosome is associated with knob-like heterochromatic structures that are heavily methylated, suggesting that DNA methylation has played a role in the evolution of this Y chromosome (18). These observations indicate that sex-linked hemizygous genes often have distinct epigenetic and regulatory features. In addition to sex-linked regions, the absence of one paired allele is frequently observed

#### Significance

This study provides a comprehensive investigation of hemizygous genes in diploid plants, highlighting their distinct characteristics compared to diploid genes. By examining 22 genomes from contrasting reproductive systems, we identify distinct evolutionary and functional features of these genes. Clonally propagated lineages are particularly replete with hemizygous genes, presumably due to the accumulation of heterozygous structural variants, but they are also common in outcrossing diploids. The hemizygous genes are expressed at lower levels than the expected 50% of diploid genes, perhaps due to enhanced DNA methylation of hemizygous genes and transposable elements. These insights enhance our understanding of plant genetics and offer valuable implications for breeding strategies of clonal crops.

Published February 7, 2025.

The authors declare no competing interest.

This article is a PNAS Direct Submission. R.A.W. is a guest editor invited by the Editorial Board.

Copyright © 2025 the Author(s). Published by PNAS. This open access article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

<sup>&</sup>lt;sup>1</sup>To whom correspondence may be addressed. Email: bgaut@uci.edu or zhouyongfeng@caas.cn.

This article contains supporting information online at https://www.pnas.org/lookup/suppl/doi:10.1073/pnas. 2422487122/-/DCSupplemental.

in non-sex-linked regions of homologous chromosomes in diploid plants, leading to a significant presence of hemizygous genes (22). However, the extent and function of these genes remain largely uncharacterized.

Here, we study hemizygous genes across a collection of plant genomes with contrasting propagation histories. The identification of these hemizygous (or haploid) genes in diploid plant genomes has become feasible with the emergence of long-read sequencing technologies and the advancements in assembly algorithms. Precise genome assemblies facilitate the identification of structural variants (SVs) in heterozygous diploid plant genomes, thereby permitting genome-wide identification of hemizygous genes caused by SVs (22–26). For example, by remapping long-reads to a reference genome assembly, it has been inferred that ~13.5% and ~15% of genes are hemizygous in two clonal grapevine (Vitis vinifera ssp. vinifera) cultivars (22). This high value may in part reflect unique features of long-term clonal lineages, because recessive deleterious mutations are expected to accumulate in these lineages (22, 27). Nonetheless, hemizygosity is not limited solely to clonal lineages, because ~8.89% and ~4% of genes are estimated to be hemizygous in an outcrossing wild rice species (Oryza longistaminata) and in avocado (Persea americana), respectively (28, 29). In contrast, as expected, only a few genes have been inferred to be hemizygous in inbred, self-fertilized accessions. For example, only 0.73% and 0.35% of genes were inferred as hemizygous in rice cultivars Nipponbare and 93-11, respectively (28).

The extent of hemizygosity in plant genomes has only begun to be appreciated, largely because genome projects have historically focused on self-fertilized or homozygous materials (22). As a result, there is currently little information about natural variation in the number of hemizygous genes, about potential correlations between hemizygosity and life history traits, such as reproductive systems and historical population sizes. There is also limited information about the evolutionary dynamics and putative functions of hemizygous genes. However, we are aware that hemizygosity can affect function. For example, white berry color in grapevines is related to a complex series of mutations, which includes hemizygosity of a large genomic region (30). In this case, the key feature of hemizygosity is that it uncovered a recessive, nonfunctional allele that interrupts anthocyanin biosynthesis. The *ClpC*-like gene in *Citrus clementina* is also hemizygous and exhibits lower expression compared to the wild-type plants, resulting in a reduced chlorophyll a/b ratio in green tissues (31). The MdACT7 gene, located within a 2.8-Mb heterozygous deletion, is expressed at lower levels in Malus domestica "AGala" compared to "Gala", causing delayed fruit maturation in the former (32). Functional effects are not limited to plants; hemizygous deletions in human autosomes are often associated with diseases and cancer. These deletions are often accompanied by decreased gene expression and increased DNA methylation, as illustrated by the examples of genes like RUNX3, KLF4, and TP53 (33-35).

Thus far, the extent of hemizygosity has only been estimated in a handful of plant genomes, and there have been no accompanying genome-wide analyses of hemizygous gene function and epigenetic states. In this study, we build or gather haplotype-resolved genomes and primary assemblies for 22 genome samples, all based on PacBio HiFi data, and subsequently identify SVs that define hemizygous genes. The 22 samples represent a range of reproductive histories, including eleven clonally propagated plants representing grapevines (*V. vinifera*), apples (*M. domestica*), and cassava (*Manihot esculenta*). Since extensive hemizygosity may be elevated in clonal lineages, we have also included comparisons from seven outcrossing samples—including wild grapevines, wild apples, and one wild rice species (*Oryza rufipogon*)—and four inbred cultivars (*SI Appendix*, Table S1). To compare results across samples, we have assembled PacBio HiFi sequencing data for each genome and identified SVs that define hemizygous genes.

Given the identification of hemizygous genes, along with the availability of transcriptomic and epigenomic data for a subset of the samples, we ask the following questions: i) Are hemizygous genes widespread in diploid plant genomes, or are they particularly abundant in clonal lineages? ii) Do hemizygous genes have distinct sequences and evolutionary features compared to diploid genes? For example, are they enriched for specific biological processes? If they are, are they expressed at half the average expression levels of diploid genes? iii) Hemizygous regions can include genes as well as other sequence features, like transposable elements (TEs). How extensive are hemizygous TEs, and do they have detectable correlations with the expression of nearby diploid genes? Finally, iv) Do hemizygous genes exhibit distinct epigenetic patterns when compared to diploid genes? If they do, is expression related to these epigenetic effects? By addressing these questions, our goal is to further understand the evolutionary and functional consequences of genic hemizygosity. Ultimately this knowledge will be beneficial for understanding the genetics, breeding, and evolution of plants with heterozygous genomes.

#### Results

The Prevalence of Hemizygous Genes in Clonal Plant Genomes. To identify hemizygous genes, we either built or gathered haplotype-resolved genome assemblies for 11 clonal and seven outcrossing plants, as well as primary genome assemblies for four inbred samples (SI Appendix, Table S1). The four inbred samples, representing grapevine, tomato, and rice, were included as a control, because hemizygosity should be near zero in these lineages. All 22 surveyed genomes were assembled with PacBio HiFi data (SI Appendix, Table S2). Among the assemblies, three Vitis genomes were generated based on PacBio HiFi and Hi-C data, and haplotypes were resolved with ultralong Oxford Nanopore Technologiesdata (> 100 kb) for gap-closing (for Chardonnay) (SI Appendix, Fig. S1). The three haplotype-resolved genome assemblies were anchored to 38 chromosomes and highly contiguous, with scaffold N50 sizes ranging from 25.1 to 26.3 Mb and Benchmarking Universal Single-Copy Orthologs (BUSCO) completeness scores of 97.3 to 98.5% (SI Appendix, Table S3). The data and reference assemblies for the remaining 19 accessions were retrieved from public repositories (SI Appendix, Table S2). The 19 assemblies were anchored to chromosome level, had scaffold N50 values of 25.1 to 67.6 Mb and BUSCO completeness scores of 93.0 to 99.2% (SI Appendix, Table S3).

Given these high-quality references, we identified hemizygous regions by remapping PacBio HiFi reads longer than 10 kb to genome assemblies, focusing on SVs longer than 50 bp. The SVs were identified using the Sniffles pipeline, followed by several filtering steps, including thresholds for quality and coverage (Materials and Methods). Because the depth of coverage varied from ~40× to 80× or higher across the 22 samples, we downsampled the PacBio HiFi data (reads longer than 10 kb) to 40× coverage before mapping to facilitate fair comparisons (SI Appendix, Fig. S4). Focusing on grapevine, this approach yielded between 59,015 and 61,066 heterozygous SVs (hSVs) in the three clonal samples. Of these, more than one-third (i.e., 26,845 to 27,982, 45.1 to 45.8% of total hSVs) were heterozygous deletions (hDELs) relative to the reference (SI Appendix, Fig. S4A); among the remaining SV types, 13.8 to 16.8% were classified as heterozygous breakends (hBND), duplications (hDUP), or inversions (hINV), while heterozygous insertions (hINS) accounted for 38.1 to 40.4% of SVs. Similar patterns were observed in the remaining clonal samples (*SI Appendix*, Fig. S4A). We excluded the hINS class from further analyses because the reference lacks information about content within these SVs.

Even after excluding hINSs, our focus on hDELs, hBNDs, hDUPs, and hINVs revealed extensive hemizygosity within clonal plants compared to outcrossing and inbred lineages. In the three clonal genomes of grapevine, for example, hemizygous regions ranged from 76.0 to 124.6 Mb, corresponding to 15.4 to 25.6% of the total genome size (*SI Appendix*, Table S1). The genome-wide extent of hSVs in clonally propagated apples, potatoes, and cassava was even higher than that of the grapevine varietals, perhaps reflecting their larger genome sizes (*SI Appendix*, Table S1) and potentially reflecting features of their reproductive and life histories, such as the duration of the clonal lineage.

The hemizygous genome proportion was substantially higher, averaging at ~25%, for clonal samples compared to the outcrossing (~16% on average) and inbred (< 1.1%) samples (*SI Appendix*, Table S1). The low proportion for the inbred samples suggests that our methods have low false positive errors—i.e., representing < 1.1 % of the genome. The differences in proportions of hemizygous genes were highly significant among clonal and outcrossing propagation histories (e.g., P < 0.05 between clonal and outcrossing samples; Mann–Whitney U-test). Although one must exercise caution with this comparison, as our samples were not phylogenetically independent, the data consistently indicated that clonally propagated lineages have elevated hemizygosity.

We further characterized hemizygous regions by examining the presence of genes and TEs within them (*SI Appendix*, Fig. S3). Focusing again on grapevines as an example, we detected between 5,131 and 6,573 hemizygous genes (representing between 12.6 and 16.8% of total genes) in the three clonal samples, with Thompson Seedless exhibiting the highest proportion (Fig. 1 and *SI Appendix*, Table S1). These genomes also contained numerous hemizygous TEs, between 174,850 and 194,536, representing 16.2 to 18.1% of annotated TEs (*SI Appendix*, Fig. S5 and Table S6), with Thompson Seedless again exhibiting the highest proportion. In contrast to the clonal samples, the outcrossing *Vitis* 

genomes had fewer hemizygous genes and TEs, with between 3,745 and 4,356 hemizygous genes representing 8.5 to 10.5% of total genes (SI Appendix, Table S1) and between 60,949 and 137,246 hemizygous TEs (between 5.4 and 12.2% of total TEs) (SI Appendix, Fig. S5 and Table S6). As expected, the inbred, nearly homozygous PN40024 V. vinifera genome had far fewer hemizygous genes and TEs, with 166 hemizygous genes (0.4% of annotated genes) and 4,828 hemizygous TEs (0.4% of total TEs). The analysis of nine Malus, one Manihot, two Solanum, and three Oryza accessions showed similar patterns (Fig. 1 and SI Appendix, Fig. S5 and Tables S1 and S6)—i.e., elevated gene and TE hemizygosity in clonal accessions. Altogether, these observations generalize previous findings, based on only a few genomes, which estimated that i) ~13.5% and 15% of genes are hemizygous in clonal grapevine cultivars (23, 27); ii) a lower but notable percentage of hemizygous genes in outcrossing plants (i.e., 8.89% hemizygous genes in O. longistaminata (28), 4% in avocado (29), and now between 5.9 and 11.7% in wild outcrossing samples (SI Appendix, Table S1); and iii) consistently low rates of genic hemizygosity (<1%) in putatively homozygous materials (28).

#### Refinement of the Set of Hemizygous Genes in Vitis Assemblies.

We anticipated that our observations based on Sniffles represented bona fide SVs. However, for downstream characterization, we deemed it critical to define a subset of genes with additional evidence supporting hemizygosity. For these analyses, we focused on the *Vitis* accessions because the genome and annotation pipeline were consistent (thereby limiting technical variation among groups or sequencing sites). Additionally, expression and epigenetic data were available for most samples, and the selected samples represent the range of reproductive histories considered in this study (*SI Appendix*, Table S1).

We identified a refined set of hSVs for the three clonal and three outcrossing *Vitis* samples by aligning the haplotypes with each other, inferring hemizygous regions, and then taking the intersection with the Sniffles-based inferences (*Materials and Methods*). This additional filter reduced the hemizygous gene set by ~50% across the six samples to an average of 2,365 genes



Reproductive system 🔶 Clonal 🔶 Inbreeding 🔶 Outcrossing

Fig. 1. Proportion of hemizygous genes in crop genomes with contrasting reproductive systems.

(range: 1,409 to 3,258) (*SI Appendix*, Fig. S6A). We focused on these genic sets for downstream analyses and contrasted them with diploid genes—i.e., genes without any evidence of an overlapping hSV. We did not include the inbred PN40024 sample in subsequent analyses, however, because it contained so few hemizygous genes.

Given this filtered set of hemizygous genes, we examined statistics such as the proportion of genes with a single exon, the number of exons per gene, exon length, and overall gene length (SI Appendix, Fig. S6 B-E). The percentage of one-exon genes was significantly higher for hemizygous than diploid genes (Fisher exact test, *P* < 0.05; *SI Appendix*, Fig. S6*B*), but the average number of exons, exon length, and gene length of hemizygous genes were all significantly lower (P < 0.05, *SI Appendix*, Fig. S6 *C*–*E*). These results were not unexpected, because one might naively expect that shorter genes have a higher probability of being encompassed by an SV event. It is also possible that the SV and alignment algorithms were biased toward identifying shorter genes, but we attempted to obviate potential biases by using only >10 kb reads and adding an additional filter based on whole genome alignments. Altogether, these results suggest that hemizygous genes are shorter and structurally simpler than diploid genes.

#### **Evolutionary and Functional Properties of Hemizygous Genes.**

We then explored evolutionary and functional features of putatively hemizygous genes relative to diploid genes for the six Vitis samples, three clonal and three wild outcrossing samples (SI Appendix, Table S1). For instance, we analyzed the proportion of hemizygous genes in centromeric and telomeric regions, representing regions of lower and higher recombination, respectively. The proportion of hemizygous genes was higher in centromeric regions compared to diploid genes, whereas the opposite pattern was observed in telomeres (Fig. 2A). This suggests that hemizygous genes tend to be biased toward low recombination regions, where selection is less effective. We also measured nonsynonymous (Ka) and synonymous (Ks) substitution rates for each gene by comparing sequences to a paired outgroup (e.g., Muscadinia rotundifolia). In each of the samples, a lower percentage of hemizygous genes (2.5% across six genomes on average) were alignable to the outgroup than diploid genes (23.3% on average), suggesting that hemizygous genes either evolve more rapidly than diploid genes or are more dispensable (i.e., more often lost over evolutionary time). After aligning the available genes, the median Ks value in hemizygous genes was significantly lower than diploid genes in all six Vitis samples (Wilcoxon rank-sum test, P < 0.05; Fig. 2B), which could reflect that the alignable hemizygous genes are a conserved, biased subset. Nonetheless, these hemizygous genes have correspondingly higher median Ka/Ks values relative to diploid genes (Wilcoxon rank-sum test, P < 0.05; Fig. 2B), consistent with weaker purifying selection, lower mutation rates (as implied by lower Ks values), or some combination of these two processes.

We then examined the age of hemizygous genes by employing phylostratigraphic analyses (36). These analyses indicated that hemizygous genes originated more recently than diploid genes (Wilcoxon rank-sum test, P < 0.05; Fig. 2*C*). We also calculated the insertion times of hemizygous and diploid long terminal repeat (LTR) TEs, finding that hemizygous TEs are also evolutionary younger (Wilcoxon rank-sum test, P < 0.05; Fig. 2*D*).

We also investigated the proportion of single-copy and multicopy genes in both hemizygous genes and diploid genes. We hypothesized that hemizygous genes were more likely to belong to a multigene family, because gene family membership can provide functional redundancies that make hemizygosity potentially less detrimental. Our results supported our hypothesis, because Finally, we investigated the possible biological processes of hemizygous genes in the six *Vitis* accessions (Fig. 2F and *SI Appendix*, Fig. S8). In Pinot Noir and *V. retordii*, GO terms included pollen recognition, endosperm development, and defense against bacteria. In Chardonnay, enriched terms were related to pollen recognition, floral organ development, and defense response. In Thompson Seedless, enriched terms highlighted responses to abscisic acid, plant ovule development, and immune response. *V. piasezkii* and *V. davidii* showed enriched responses to salt stress and defense mechanisms. These results indicate that hemizygous genes can be involved in fundamental processes like reproduction and mitosis, but they are also consistently enriched for responses to biotic and abiotic stress.

**Unique Expression Patterns of Hemizygous Genes.** A simple null expectation for hemizygous genes is that they are expressed at 50% of the average level of diploid genes. To explore this hypothesis, we amassed RNA-seq datasets generated across accessions, developmental stages (e.g., fruit development), and experimental regimes, such as stress treatments (*SI Appendix,* Table S8). Our goals with these data were i) to investigate the level of expression in hemizygous genes relative to diploid genes and ii) to determine whether hemizygous genes had patterns of expression consistent with contributions to development and other processes.

We first assessed whether hemizygous genes were expressed. Across taxa and individual RNA-seq samples, a significantly higher proportion of hemizygous genes had no evidence of expression relative to diploid genes (Fig. 3A and SI Appendix, Fig. S9A). For example, across all samples in Pinot Noir, 49.2% (1,334 out of 2,711 genes) of hemizygous genes had evidence for expression, but that proportion was 74.2% (26,309 out of 35,459 genes) for diploid genes (chi-sq = 704.1, df = 1; P < 0.05). This trend held true in each tissue/treatment for all taxa—e.g., among the 21 RNA-seq tissues/treatments in Pinot Noir, hemizygous genes were expressed in lower proportions for all tissues/treatments. These results strongly suggest that hemizygous genes are enriched for pseudogenes or for a subset of genes that are expressed under fewer experimental and developmental conditions. However, not all hemizygous genes were pseudogenes; across all data samples, we detected expression for 49.2%, 66.4%, 44.5%, 24.7%, 52.7%, 29.4%, and 61.1% of hemizygous genes in Pinot Noir, Chardonnay, Thompson Seedless, V. piasezkii, V. davidii, and V. retordii, and Golden Delicious, respectively (Fig. 3A and SI Appendix, Fig. S9A).

We next investigated average levels of expression for the subset of genes with evidence for expression. Hemizygous genes were consistently expressed at significantly lower levels than diploid genes based on average expression across all tissues/treatments and within each tissue/treatment (Fig. 3B and SI Appendix, Fig. S9B). The hemizygous:diploid ratio of median expression was ~0.07 for Pinot Noir and V. davidii but higher for other Vitis samples (range 0.14 to 0.21) and highest for Golden Delicious apples (0.32) (SI Appendix, Fig. S9C). Nonetheless, all these values were significantly (P < 0.05, Wilcoxon signed rank test) lower than the 50% expected if hemizygous alleles were expressed at similar levels to diploid alleles. These results imply that there is a diminution of gene expression associated with hemizygosity, and this diminution typically results in average expression levels of hemizygous genes being less than 50% of those of diploid genes (Discussion).



**Fig. 2.** Characterization of hemizygous genes in clonal and outcrossing *Vitis* accessions. (*A*) Proportions of hemizygous and diploid genes located in centromeres (*Left*) and telomeres (*Right*). (*B*) *Ks* and *Ka/Ks* values comparing hemizygous and diploid genes; boxplots show medians, quartiles, and ranges. (*C*) Average phylostrata levels for both hemizygous and diploid genes, with larger values indicating younger genes. error bars indicate 95% bootstrap Cl. (*D*) Insertion times for hemizygous and diploid LTRs. (*E*) Proportion of single-copy genes among hemizygous and diploid genes. (*F*) Top 15 enriched biological processes for hemizygous genes in Pinot Noir. Fisher's exact test was used for (*A* and *E*); Wilcoxon rank-sum test for (*B*), (*C*), and (*D*). "ns" denotes nonsignificant (*P* > 0.05), significant differences (*P* < 0.05) are unmarked.

We then examined patterns of hemizygous gene expression across specific functional categories that were explored in the original RNA-seq experiments. These categories included stages of fruit development, organ differentiation, and responses to biotic and abiotic stresses. First, we estimated the proportion of hemizygous genes that were differentially expressed across comparisons. Hemizygous genes generally had a lower proportion of differentially expressed genes than diploid genes (Fig. 3C and *SI Appendix*, Fig. S9*D*). For example, in Pinot Noir, 28.3% (378) out of 1,334) of hemizygous genes and 77.5% (20,394 out of 26,309) of diploid genes were differentially expressed across all paired comparisons (chi-sq = 1712.0, df = 1; P < 0.05). The corresponding values for Chardonnay were 24.2% (454 out of 1,873) vs. 60.7% (16,987 out of 27,978) (chi-sq = 942.2, df = 1; P < 0.05, 24.8% (359 out of 1,450) vs. 68.6% (16,168 out of 23,578) (chi-sq = 1105.6, df = 1; P < 0.05) in Thompson Seedless, 25.5% (189 out of 742) vs. 66.3% (18,468 out of 27,838) (chi-sq = 515.9, df = 1; P < 0.05) in V. davidii, and 26.3% (366 out of 1,389) vs. 53.4% (17,261 out of 32,313) (chi-sq = 386.3, df = 1; P < 0.05) in Golden Delicious. Overall, fewer expressed hemizy-gous genes differed in expression during fruit development, organ differentiation, or during abiotic and biotic stress.

We also explored expression specificity by detecting genes that were expressed in specific tissues or treatments. That is, we counted the number of genes that had significant evidence for expression in only one tissue/treatment of a paired comparison. Altogether, hemizygous genes had a higher proportion of tissue/treatment-specific genes than the diploid genes (*SI Appendix*, Fig. S10). Across 10 tissue/treatment comparisons in Pinot Noir, 31.0 to 50.0% of hemizygous genes were expressed in only one of the paired tissues or treatments, but these values were substantially lower for diploid genes (*SI Appendix*, Fig. S10). Similar patterns—i.e., more tissue/treatment-specific expression for hemizygous genes—were also found in Chardonnay, Thompson Seedless, *V. davidii*, and Golden Delicious (*SI Appendix*, Fig. S10). Given these results, we anticipated that a lower proportion of



**Fig. 3.** Transcriptomic landscapes of hemizygous and diploid genes in clonal and outcrossing *Vitis* accessions. (*A*) Proportion of expressed hemizygous and diploid genes in each tissue/treatment and overall across six samples. (*B*) Expression levels (FPKM) of hemizygous and diploid genes in each tissue/treatment from the six *Vitis* samples. Error bars show 95% bootstrap-based Cl. (*C*) Proportion of differentially expressed genes for hemizygous and diploid genes for four *Vitis* accessions with control-treatment contrasts. (*D*) Proportion of common and unique DEGs of hemizygous and diploid genes among three processes, including fruit development (Fr), organ differentiation (Or), and abiotic and biotic stress stimulus processes (Stress). Fisher's test was used for (*A*) and (*C*; Wilcoxon rank-sum store (*B*). Significant differences (*P* < 0.05) were observed throughout. Circles in (*A*) and (*C*) represent increments of 25%, from 25% (innermost) to 100% (outermost).

hemizygous genes would be expressed across experiments investigating fruit development, organ differentiation, abiotic and stress; indeed, this was true for all relevant samples (Fig. 3*D* and *SI Appendix*, Fig. S9*E*).

The Cis-Regulatory Effects of Hemizygous and Diploid TEs on Gene Expression. We then explored the cis-regulatory effects of TEs on gene expression. To do so, we used both the RepeatModeler and EDTA pipelines to identify TEs for the three clonal and three outcrossing *Vitis* samples, detecting from 337,284 to 502,618 TEs across the six accessions based on both pipelines (*SI Appendix*, Tables S5 and S6). We then classified genes into four categories based on their proximity to annotated TEs. The four categories were i) hemizygous genes with nearby TEs (i.e., within 2 kb of either the 5' or 3' ends of genes), ii) hemizygous genes without nearby TEs, iv) diploid genes without nearby TEs.

Focusing on diploid genes, the pattern was consistent and clear: Among the group of genes without TEs, a higher percentage were expressed (Fig. 4A) and expressed at higher levels (Fig. 4B) than genes with nearby TEs. This observation held across the six taxa and across individual RNA-seq samples (*SI Appendix*, Figs. S11 and S12). The difference could be striking; for example, in one leaf sample of Pinot Noir, 93.3% of diploid genes without a nearby TE were expressed, while only 71.4% of diploid genes with a nearby TE were expressed. The pattern in diploid genes was consistent with the findings that host silencing of TEs near genes often negatively affects expression of a neighboring gene (37)—e.g., siRNA-targeted TEs are associated with reduced gene expression (38)—and TEs close to genes may disrupt cis-regulatory elements such as enhancers and silencers that affect gene expression (39, 40).

However, these patterns were less pronounced for hemizygous genes (Fig. 4*A* and *SI Appendix*, Figs. S11 and S12). For example, hemizygous genes near TEs tended to express more often in two taxa, *V. davidii* and *V. retordii* (Fig. 4*A*, Wilcoxon rank-sum test, P < 0.05), while the other four taxa showed no significant differences (Fig. 4*A*, Wilcoxon rank-sum test, P > 0.05). Meanwhile, there were no significant differences in the expression levels of hemizygous genes near TEs compared to those without nearby TEs in the six taxa (Fig. 4*B*, Wilcoxon rank-sum test, P > 0.05).

Thus, the relationship between hemizygous genes and the presence of nearby TEs is less evident compared to that of diploid genes

Like genes, TEs can be diploid or hemizygous, so we also explored this effect on gene expression. Across six taxa, the percentage of expressed diploid genes with nearby diploid TEs was generally higher compared to diploid genes with nearby hemizygous TEs (Fig. 4*C* and *SI Appendix*, Fig. S13) and at higher levels (Fig. 4*D* and *SI Appendix*, Fig. S14). Although the pattern for hemizygous genes was less obvious (Fig. 4*D* and *SI Appendix*, Fig. S14), the results generally suggest that SVs near genes (i.e., those resulting in hemizygous TEs) tend to reduce the expression of diploid genes more than nearby diploid TEs.

Higher DNA Methylation in Hemizygous Genes and TEs. One potential explanation for the effect of hemizygous vs. diploid TEs has to do with epigenetic patterns—i.e., if hemizygous TEs are more highly methylated, they may more effectively dampen gene expression. We thus investigated DNA methylation patterns. We began by surveying hemizygous vs. diploid genes from leaves of four *Vitis* samples (Pinot Noir, Chardonnay, *V. piasezkii*, and *V. retordii*; *SI Appendix*, Table S9). For hemizygous genes in Pinot Noir leaves, we detected average weighted genomic DNA methylation levels of 45.8%, 23.0%, and 2.3% in the CG, CHG, and CHH sequence contexts, respectively (Fig. 5 *A* and *B* and *SI Appendix*, Fig. S15); Similar to previous reports (41, 42), genic methylation levels were lower than the genome-wide methylation levels. For hemizygous vs.

diploid genes, the average DNA methylation level was 52.2% vs. 40.8%, 24.6% vs. 14.8%, and 2.3% vs. 1.7% in the CG, CHG, and CHH contexts, respectively. These patterns were largely consistent across taxa, and they generally reflect higher methylation levels for hemizygous genes compared to diploid genes.

As expected, TEs were methylated at higher levels than genome-wide averages. However, it is interesting to note that hemizygous TEs close to diploid or hemizygous genes tended to be methylated at higher levels than diploid TEs close to diploid or hemizygous genes. For example, in the Pinot Noir sample, hemizygous and diploid TEs close to hemizygous genes have methylation levels of 81.6 and 74.2% in the CG context, 73.2 and 61.7% in the CHG context and 5.3 and 4.5% in the CHH context. Similar patterns were found in Chardonnay, *V. piasezkii*, and *V. retordii* (Fig. 5*A*). Hence, hemizygous TEs generally have higher methylation levels than diploid TEs, which may explain their stronger effect on nearby gene expression.

Hemizygous Gene Expression Levels Correlated with Gene Body Methylation. We then turned to the methylation status of individual genes across different taxa. Hemizygous genes had a lower proportion of gbM genes than diploid genes in Chardonnay and *V. retordii*. For instance, in Chardonnay, 18.1% of hemizygous genes (511 out of 2,821) were classified as gbM, in contrast to 28.8% of diploid genes (9,797 out of 33,980) (Fig. 6A). A similar pattern was found in *V. retordii*. The difference between gbM



**Fig. 4.** Cis-regulation of TE effects on hemizygous and diploid genes in *Vitis* accessions. (*A*) Proportion of expressed hemizygous and diploid genes with/without nearby TEs. (*B*) Expression levels  $[log_2(FPKM+1)]$  of expressed hemizygous and diploid genes with/without nearby TEs. (*C*) Proportion of expressed genes with nearby hemizygous or diploid TEs. (*D*) Expression levels  $[log_2(FPKM+1)]$  of genes with nearby hemizygous or diploid TEs. (*D*) Expression levels  $[log_2(FPKM+1)]$  of genes with nearby hemizygous or diploid TEs. Fisher's exact test was used for (*A*) and (*C*), Wilcoxon rank-sum test for (*B*) and (*D*). Is denotes nonsignificant (*P* > 0.05); significant differences (*P* < 0.05) are unmarked. Boxplots in (*B*) and (*D*) show median, quartiles, and range.

proportions in hemizygous and diploid genes (i.e., hemizygous < diploid) and average CG genic methylation ratio pattern (hemizygous > diploid) in Chardonnay and *V. retordii* can potentially be attributed to the higher proportion of mCHG genes in hemizygous regions, which may influence their methylation status. For example, in Chardonnay, 54.4% of hemizygous genes (1,534 out of 2,821) and 23.5% of diploid genes (7,985 out of 33,980) were classified as mCHG genes.

After identifying gbM, mCHG, and UM genes, we investigated their expression patterns and observed several distinct trends. First, a smaller proportion of mCHG genes were expressed compared to gbM and UM genes (Fig. 6C), regardless of whether they were hemizygous or diploid. This finding aligns with previous studies indicating that mCHG methylation suppresses gene expression (41). The high proportion of hemizygous mCHG genes contributed to the overall lower expression levels of hemizygous vs. diploid genes (Fig. 3*B*). Second, a higher proportion of gbM and UM genes were expressed in diploid genes compared to hemizygous genes (Fig. 6B). Third, the patterns based on the proportion of expressed genes were largely reflected in expression levels. That is, mCHG genes were relatively lowly expressed, no matter if they were hemizygous or diploid (Fig. 6C); gbM and UM genes were expressed at higher levels than mCHG genes (Fig. 6C); and hemizygous gbM and UM genes were consistently expressed at lower levels than diploid genes (Fig. 6C).

#### Discussion

Hemizygous genes have been studied extensively in sex-linked regions, but they can also occur beyond sex-linked regions of homologous chromosomes due to hSVs. Some SVs lead to the presence of a single allele on one homologous chromosome of an otherwise diploid organism. Here, we have integrated genomic, transcriptomic, and epigenomic analyses to estimate the frequency of these hemizygous genes, to explore their relationship to propagation history, and to characterize their features, expression, and epigenetic regulation.

Hemizygous Genes Are Most Common in Clonal Lineages. Consistent with previous work, we have found that hemizygous genes are more common in clonal, as opposed to outcrossing lineages. Although hemizygosity has already been measured in a handful of plant taxa—i.e., primarily grape varieties and rice species-we have expanded observations to encompass nine clonally propagated cultivars from grapevine, apple, cassava, and potato, seven outcrossed samples from wild grapevine, wild apple, and wild rice, and four genomes expected to be fully homozygous (SI Appendix, Table S1). By focusing on hDELs, hBNDs, hDUPs, and hINVs relative to the reference assembly, we have documented, as expected, little evidence for hemizygosity in the homozygous samples, with estimates representing <1.2%of the genome (SI Appendix, Table S1). These results are not particularly surprising, but they show that we do not estimate high hemizygosity where there should be none.

In contrast to the homozygous samples, our work substantiates a growing consensus that outcrossing species can harbor a substantive portion of their genome as hemizygous. Among the seven outcrossed samples, 5.9 to 11.7% of their genes are captured within hSVs, mimicking levels found in outcrossing rice and avocado. (Avocado is clonally propagated in cultivation, but the investigated tree had been produced by a recent outcrossing event.) In contrast, long-term clonal lineages consistently have an even more substantial fraction of their genomes and genes captured in a hemizygous state. Most of the observations to date have been based on grapevine clones, some of which have been propagated for 1000 or more years (43). However, by including cultivated apple, cassava, and potato, we have shown that this clonal phenomenon is not limited to grapevines (*SI Appendix*, Table S1). Moreover, the results accentuate how a traditional focus on inbred plants like *Arabidopsis thaliana*, rice, and tomato has biased our understanding of genetic variation. The inbred plants are typically highly homozygous with few sequence variants, but the genomes of clonal plants are highly heterozygous with genetic diversity that includes SVs and hemizygous genes (44).

High genetic variation in clonal lineages is not particularly unexpected, for two reasons. First, previous work on SVs has inferred, based on population samples, that they tend to be deleterious (22, 28). Second, forward simulations have consistently revealed that heterozygous, deleterious variants are expected to accumulate over time in clonal lineages, a phenomenon not observed in outcrossing plants (22, 45, 46). This accumulation reflects the fact that recessive deleterious alleles can hide as heterozygotes within a clonal lineage; whereas in outcrossing systems, they are expected to occasionally become homozygous and thus subject to selection. This accumulation also reflects that recombination is limited (i.e., effectively zero) in strictly clonal lineages, meaning that deleterious mutations cannot recombine onto different genetic backgrounds. Consistent with this argument, we find that hemizygous genes tend to be biased toward low-recombination, centromeric regions where selection is likely to be less efficient (Fig. 2A). Finally, it is interesting to note that previous work showed that domesticated, clonally propagated cassava has a marked 26% higher genomic burden of putatively deleterious nucleotides compared with its wild congener (47), consistent with its substantial burden of hemizygous genes (SI Appendix, Table S1).

Despite previous studies about the accumulation of deleterious variants in clonal lineages, the large number of hemizygous genes in clonal lineages is still somewhat surprising, because functionally hemizygous genes cannot (by definition) be recessive. Hence, the dynamics of the accumulation of hemizygous genes are likely to differ somewhat from the deleterious recessive case studied by forward simulation. Assuming that many (but not all; see below) of the SV events are slightly deleterious, several functional and evolutionary processes likely contribute to the accumulation of hemizygous genes in clonal lineages. One is a ratchet mechanism—i.e., once an SV occurs in a clonal lineage, it has only one possible fate, so long as it is not lethal, which is to remain in the clonal lineage. By this process, clonal lineages are expected to accumulate SVs. In theory, this accumulation is more likely when the SV events do not severely affect fitness; for that reason, we expect deleterious SVs to often have moderate functional effects.

#### Hemizygous Gene Expression Is Moderated by Epigenetic Effects.

Indeed, we have accrued evidence that hemizygous genes have moderate functional effects, based on three pieces of evidence. First, hemizygous genes are more likely to be nonexpressed than diploid genes in our samples (Fig. 3A). That is, a higher proportion of hemizygous genes appear to be pseudogenes. Second, hemizygous genes are more likely to be members of gene families (Fig. 2*E*), implying that they are more likely to be functionally redundant. Thus, the loss of one copy of a multicopy gene is likely to carry fewer fitness consequences than the loss of one allele of a critical single-copy gene. Finally, and somewhat surprisingly, as a group, hemizygous genes tend to be expressed at less than half the level of average diploid genes, at about 20% (SI Appendix, Fig. S9C). This value is substantially less than the 50% expected of a single allele. It is hard to know the cause of this low expression pattern. It is possible, for example, that hemizygous genes are a biased sample that were lowly expressed in their diploid state before the SV event. Another possibility is



**Fig. 5.** Epigenomic landscapes of hemizygous genes and their nearby diploid and hemizygous TEs. (A) Average DNA methylation levels in CG, CHG, or CHH contexts for different sequence types in four *Vitis* accessions. (B) DNA methylation distribution for each sequence type for Pinot Noir, with "TSS" and "TES" marking transcription start/end sites of genes or TE boundaries, including 1-kb flanking regions.

that epigenetic effects act especially strongly on hemizygous genes to moderate their expression (see below).

In this context, it is worth accentuating that another set of hemizygous genes that have been studied intensively—i.e., sex-linked genes. Sex-linked genes tend to have an X:AA gene expression ratio of ~0.5 in the human, mouse, and nematode (11). Another possibility for sex-linked genes is dosage compensation, which predicts that hemizygous X-linked genes are expressed at twice the level of diploid genes per active allele to balance the gene dosage between the X chromosome and autosomes (12). The upregulation of the hemizygous copy may be sufficient to mitigate negative fitness effects, even if expression still falls significantly short of ancestral expression levels, and may also mitigate the effects of an euploidy (48–50). In contrast, we do not see any overarching evidence of complete or even partial dosage compensation of hemizygous genes. Instead, the opposite is true: The expression of hemizygous alleles is substantially less expressed than the average diploid allele.

We suspect that lower expression is at least partially due to epigenetic phenomena, for three reasons. First, in all four samples investigated, for both diploid and hemizygous genes, nearby hemizygous TEs have elevated levels of DNA methylation relative to their nearby diploid TEs (Fig. 5*A*). Several phenomena may contribute to this observation, including that hemizygous TEs may be more recent insertions (and therefore more actively targeted by host epigenetic responses) (Fig. 2*D*). Whatever the cause, the data hint that hemizygous TEs differ quantitatively in their methylation effects. Second, hemizygous genes also exhibit higher levels of methylation than diploid genes, specifically a higher proportion of mCHG alleles (Fig. 6*A*), which are usually a mark of low expression (Fig. 6*C*). Finally, we have shown that genes near TEs are consistently more lowly expressed than genes far from TEs (Fig. 4*B*), but this effect is more prominent for genes near hemizygous TEs (Fig. 4*D*). This may be a partial explanation as to why genes close to SVs are associated with reduced gene expression levels in other species, like tomato (51).

These observations have interesting parallels to previous studies that have suggested that DNA methylation is correlated with reduced gene expression levels for sex-limited genes on the Y or W chromosome (52). High levels of DNA methylation have also been associated with sex chromosomes in sticklebacks and papaya (18, 19). In addition, DNA methylation is a key feature in X-chromosome inactivation (20). These results suggest some similar features of DNA methylation patterns between sex-linked and non-sex-linked hemizygous genes. Hemizygosity in human autosomes is linked to decreased gene expression and increased methylation, and these phenomena contribute to cancer and disease. For example, 45 to 60% of human gastric cancer cells show reduced RUNX3 expression due to a hemizygous deletion and promoter hypermethylation (33). Similarly, promoter hypermethylation and a hemizygous deletion leads to KLF4 down-regulation and apoptosis induction, enhancing its antitumor activity (34). Finally, a hemizygous deletion containing the TP53 gene is found in over 10% of newly diagnosed multiple myeloma patients, and it is associated with decreased expression, impaired *p53* response, and resistance to apoptosis due to promoter hypermethylation (35). Clearly, we cannot be certain what, if any, epigenetic mechanisms might be shared between sex-linked and human disease hemizygosity and that which we have studied here, but it is an interesting question for further research.



**Fig. 6.** Epigenetic effects on hemizygous gene expression in four *Vitis* accessions. (*A*) Proportion of gbM, mCHG, and UM genes for hemizygous and diploid genes across four taxa. (*B*) Proportion of expressed gbM, mCHG, and UM genes. (*C*) Expression levels  $[log_2(FPKM+1)]$  of expressed gbM, mCHG, and UM genes. Fisher's exact test was used for (*A*) and (*B*); Wilcoxon rank-sum test for (*C*). ns denotes nonsignificant (*P* > 0.05), while significant differences (*P* < 0.05) are unmarked. Boxplots in (*C*) show median, quartiles, and range.

Are Hemizygous Genes Merely Functional Remnants? Given the evidence that hemizygous genes tend to be shorter than diploid genes (*SI Appendix*, Fig. S6 *D* and *E*), expressed at lower levels (Fig. 3*B*), potentially subjected to lower levels of purifying selection (as measured by *KalKs*; Fig. 2*B*), preferably distributed in centromeric regions representing low recombination regions, and more heavily methylated (Fig. 4 *A* and 4 *B*), it is tempting to conclude that hemizygous genes are typically pseudogenes. Are they merely functional remnants of previously functional genes? Might they also simply represent the dispensable component of the genome, which tend to be less expressed and more methylated (53)? While the answer to this question is likely "yes" for most hemizygous genes, there is some tantalizing evidence suggesting that the answer may often be "no".

Evidence supporting functionality of some hemizygous genes comes in a few forms. For example, a reasonable proportion of hemizygous genes have gbM patterns of methylation (Fig. 6*A*). In both hemizygous and diploid genes, gbM genes exhibit significantly higher expression levels compared to mCHG genes (Fig. 6*C*). Moreover, several studies have detected a correlation between the presence of gbM and the enhancement of gene or allelic expression (41), while others have found evidence that it is subject to natural selection based on population genetic arguments. In short, although the functional role of gbM (if any) is debated (54), it typically is a mark deposited and maintained on active genes (41). The fact that some hemizygous genes bear this epigenetic mark superficially suggests that they cannot be easily dismissed as nonfunctional.

In addition, hemizygous genes (as a group) demonstrate patterns of tissue/treatment-specific expression that are similar to those of diploid genes. This pattern does not hold at the single gene level, but nonetheless up to 50% of hemizygous genes exhibit tissue/treatment-specific expression in Pinot Noir (SI Appendix, Fig. S10). Of course, tissue/treatment-specific expression patterns are not proof of function, but it does indicate that some hemizygous genes are induced under different environmental and developmental conditions. Finally, there are some consistent patterns of GO enrichment, particularly for responses to biotic and abiotic stresses (Fig. 2F and SI Appendix, Fig. S8). Again, GO enrichment is not proof of function, but this evidence combines to make it reasonable to hypothesize that not all hemizygous genes are functional "junk". Of course, the mere act of uncovering a recessive allele can have important functional consequences; we invoke again the compelling case of hemizygosity and the white berry phenotype of grapes (22, 30).

#### **Materials and Methods**

Sample Selection and Genome Assembly and Annotation. We utilized PacBio-based genome assemblies for 22 diploid plant samples, including three haplotype-resolved assemblies for Pinot Noir (PN\_AGIS2\_hap1 and PN\_AGIS2\_hap2), Chardonnay (CHT2T\_AGIS1\_hap1 and CHT2T\_AGIS1\_hap2), and *V. piasezkii* (PIA\_AGIS1\_hap1 and PIA\_AGIS1\_hap2), which were assembled as part of this study, with Chardonnay genome achieving a haplotype-resolved telomere-to-telomere level (*SI Appendix*, Tables S1–S3). Plant materials were cultivated at the Agriculture Genomics Institute at Shenzhen (AGIS), CAAS. DNA extraction, SMRTbell

library construction, and sequencing on PacBio Sequel II (CCS mode), as well as ultralong ONT library preparation and Hi–C sequencing, followed established protocols in refs. 24, 25, and 55.

HiFi and Hi-C reads from three *Vitis* accessions (Pinot Noir, PN\_AGIS\_02; Chardonnay, CH\_AGIS\_01; *V. piasezkii*) were assembled into haplotigs using Hifiasm (v0.19.8). The assemblies were then anchored to 38 chromosomes based on PN\_T2T genome similarity (56) using RagTag (v2.1.0) (57), and scaffold with Juicer (v1.6) (58) and 3D-DNA (v190716) (59). Hi-C maps were visualized and manually adjusted in Juicebox (v2.17.00) (60). We aligned HiFi and ONT reads (as Chardonnay has ONT reads, while the others only have HiFi) to the genome using Minimap2 (v2.24-r1122) (61), and gaps were manually filled with Integrative Genomics Viewer (IGV) (v2.13.1) tool (62). The complete pipeline for genome assembly and gap filling is available on our lab GitHub@zhouyflab (*Data, Materials, and Software Availability*). The remaining 19 genome assemblies were retrieved from public resources (*SI Appendix*, Table S2). Detailed methods and sample information are provided in *SI Appendix*.

Gene annotations were generated using the MAKER (63) and Liftoff pipelines (64) (*SI Appendix*, Fig. S2 and Table S4) for all seven *Vitis* accessions, which included three accessions assembled in this study and four that were published previously (56, 65–67). TE annotation utilized RepeatModeler/RepeatMasker (RM) (68) and EDTA pipelines (69).

Identification and Characterization of Hemizygous Genes. To identify hemizygous genes, PacBio HiFi reads from the 22 genome assemblies were remapped to their respective genomes, and SVs were called using Sniffles (v2.0.6) (70). Reads longer than 10 kb were aligned with Minimap2 (v2.24) (61). SVs were filtered based on quality, size (>50 bp), and supporting reads (≥4). Hemizygous regions were defined as SV regions with 0/1 flags, and genes overlapping these regions (≥80% of their lengths) were classified as hemizygous genes. Further filtering was conducted by aligning haplotypes within six *Vitis* accessions using Mummer and Nucmer (71) to identify reliable hemizygous genes.

Sequence and evolutionary characteristics of hemizygous and diploid genes were compared, including exon numbers, gene lengths, synonymous mutation rate (*Ks*), and nonsynonymous/synonymous mutation ratio (*Ka/Ks*), gene ages, LTR insertion time, proportion of single-copy genes. Centromeric and telomeric repeats were annotated using Tandem Repeats Finder (TRF, version 4.09) (72), with regions extended 1 Mb for overlap analysis. *Ka* and *Ks* values were estimated using MCScanX (73) based on grapevine-*M. rotundifolia* genome sequence comparisons. Gene ages were determined using phylostratigraphy, mapping coding genes against proteomes of 12 representative species using BLASTP (36, 74) (*SI Appendix*, Fig. S7 and Table S7). LTR insertion times were calculated using EDTA (69) and MEGA-CC (75), and TEs were classified as hemizygous or diploid based on SV inferences. Single-copy orthologs were identified using OrthoFinder (76), with *M. rotundifolia* as outgroup. GO enrichment was analyzed using DAVID (77), with significant terms set at *P* value < 0.05. Additional details are provided in *SI Appendix*.

**Dissection of Hemizygous Gene Expression Patterns.** To explore how hemizygous genes respond to fruit development, organ differentiation, and stress stimuli, we analyzed 168 RNA-seq samples from six grapevine accessions and one apple, totaling 691 Gb (*SI Appendix*, Table S8). This included datasets from *V. piasezkii* and *V. retordii* leave generated in this study, along with 162 publicly available samples. The public available samples were categorized by experimental conditions, including fruit development, stress response, and organ differentiation.

Raw RNA-seq reads were processed using Trimmomatic (v0.39) (78) for quality trimming and subsequently mapped to their respective genomes with HISAT2 (v.2.2.1) (79). Gene counts were extracted with FeatureCounts (v2.0.1) (80), and expression levels were normalized to FPKM values using custom R scripts. Genes with FPKM > 0 were considered expressed. Detailed methods and sample information are provided in *SI Appendix*.

**Exploration of Cis-Regulatory Effects of TEs on Gene Expression.** Based on the identification of repeat sequences, we explored the cis-regulatory effects of TEs on gene expression. For this purpose, we first assigned each TE to its closest gene when

- 1. D. Charlesworth, Why and how do Y chromosome stop recombining? J. Evol. Biol. 36, 632-636 (2023).
- B. Charlesworth, D. Charlesworth, The degeneration of Y chromosomes. *Philos. Trans. R Soc. Lond B, Biol. Sci.* 355, 1563–1572 (2000).
- R. Bergero, D. Charlesworth, Preservation of the Y transcriptome in a 10-million-year-old plant sex chromosome system. *Curr. Biol.* 21, 1470–1474 (2011).

it was within 2 kb (the distance to either 5' or 3' end of gene with  $\ge 0$  kb and <2 kb) using command "bedtools closest -wo -a gene.bed -b TE.bed", and thus genes were separated in four classes: hemizygous genes with nearby TEs, hemizygous genes without nearby TEs, diploid genes with nearby TEs, diploid genes without nearby TEs. We also divided genes near TEs into four categories: hemizygous genes with either hemizygous or diploid TEs, and diploid genes with either hemizygous or diploid TEs.

**Unveiling DNA Methylation Patterns of Hemizygous Genes.** Bisulfite sequencing (BS-seq) was conducted on four samples, which were either generated in this study or obtained from public datasets (*SI Appendix*, Table S9). DNA was extracted using the Qiagen DNeasy Plant Mini kit, and bisulfite libraries were prepared as described previously (42). Libraries were sequenced on the Illumina HiSeq2500 platform, with lambda-DNA spiked in as a control for bisulfite conversion.

Trimmed reads were aligned to reference genomes using Bismark (v0.23.1) (81) with bowtie2 (v2.1.0) (82), and methylation status was determined using the bismark\_methylation\_extractor (minimum coverage = 2). Methylation levels were identified using a binomial test (FDR-corrected, P < 0.01) (83), and false methylation rates were computed using lambda-DNA or chloroplast DNA using MethylExtract (84). DNA methylation patterns across contexts (CG, CHG, CHH) were visualized with deepTools (85).

We defined body-methylated genes following the strategy of refs. 54 and 86. Briefly, the DNA methylation level of each protein-coding gene was quantified for all three contexts (CG, CHG, and CHH). *P* values were used to denote the deviation of methylation levels from the genomic averages for each context. We defined gene as BM ( $P_{CG} <= 0.05$ ,  $P_{CHG} > 0.05$  and  $P_{CHH} > 0.05$ ), mCHG ( $P_{CHG} <= 0.05$ , and  $P_{CHH} > 0.05$ ), mCHH ( $P_{CHH} <= 0.05$ ), and UM (( $P_{CG} > 0.05, P_{CHG} > 0.05$ , and  $P_{CHH} > 0.05$ ). In any other case, the methylation state was not inferred. Detailed methods and sample information are provided in *SI Appendix*.

Data, Materials, and Software Availability. The PacBio CCS, ONT, Hi-C, RNAseq, BS-seq data have been deposited to the NCBI short reads achieved under the project number: PRJNA1178252 (87). The genome assembly and annotation have been deposited to Zenodo: https://zenodo.org/records/14015567 (88). Code availability: All scripts and codes performed in this study are available on GitHub: https://github.com/zhouyflab/Genomic\_Epigenomic\_ Hemizygous\_Crops (89). All other data are included in the article and/or supporting information.

ACKNOWLEDGMENTS. We thank R. Gaut for generating the BS-seq data and all members in the Zhou lab for the useful comments and discussions. This work was supported by the National Natural Science Foundation of China (No. 32372662), the Science Fund Program for Distinguished Young Scholars of the National Natural Science Foundation of China (Overseas) to Yongfeng Zhou, the National Key Research and Development Program of China (2023YFD2200700), Hainan Province Key Research and Development Project (ZDYF2024XDNY156), Shenzhen Polytechnic University (SZPU) Research Project (6024330001K) to Lin Tian, and the NSF grant (No.1741627) to Brandon S. Gaut.

Author affiliations: <sup>a</sup>National Key Laboratory of Tropical Crop Breeding, Shenzhen Branch, Guangdong Laboratory of Lingnan Modern Agriculture, Key Laboratory of Synthetic Biology, Ministry of Agricultura and Rural Affairs, Agricultural Genomics Institute at Shenzhen, Chinese Academy of Agricultural Sciences, Shenzhen 518120, China; <sup>b</sup>School of Management, Shenzhen Polytechnic University, Shenzhen 518055, China; <sup>c</sup>National Key Laboratory of Tropical Crop Breeding, Tropical Crops Genetic Resources Institute, Chinese Academy of Tropical Agricultural Sciences, Haikou 571101, China; <sup>d</sup>Key Laboratory of Biology and Genetic Improvement of Horticultural Crops (Germplasm Resources Utilization) Ministry of Agriculture, Fruit Research Institute, Chinese Academy of Agricultural Sciences, Xingcheng 125100, Liaoning, China; and <sup>e</sup>Department of Ecology and Evolutionary Biology, University of California, Irvine, CA 92697

Author contributions: S.H., B.S.G., and Y.Z. designed research; Y.P., Y.W., Y.L., X.F., L.C., Q.L., Q.X., N.W., F.Z., Z.L., H.X., J.Y., L.T., W.H., S.C., H.W., and Y.Z. performed research; Y.P., D.S., T.Z., X.S., and X.X. analyzed data; Y.P. and Y.W. revised the manuscript; Y.P. and Y.L. and X.F. collect the data; L.C. revised the manuscript; Q.L. and X.X. collected the data; Q.X., N.W., F.Z., Z.L., H.X., J.Y., L.T., W.H., S.C., W.H., S.H., B.S.G., and Y.Z. revised the paper; and Y.P., B.S.G., and Y.Z. wrote the paper.

 E. Cherif et al., Male-specific DNA markers provide genetic evidence of an XY chromosome system, a recombination arrest and allow the tracing of paternal lineages in date palm. New Phytol. 197, 409-415 (2013).

 M. F. Torres et al., Genus-wide sequencing supports a two-locus model for sex-determination in Phoenix. Nat. Commun. 9, 3969 (2018).

- A. Harkess et al., The asparagus genome sheds light on the origin and evolution of a young Y chromosome. Nat. Commun. 8, 1279 (2017).
- K. Murase et al., MYB transcription factor gene involved in sex determination in Asparagus officinalis. Genes Cells 22, 115–123 (2017).
- 8. T. Akagi et al., Two Y-chromosome-encoded genes determine sex in kiwifruit. Nat. Plants 5, 801-809 (2019).
- J.Y.Yue et al., Origin and evolution of the kiwifruit Y chromosome. Plant Biotechnol. J. 22, 287-289 (2024).
  B. Vicoso, B. Charlesworth, Evolution on the X chromosome: Unusual patterns and processes. Nat.
- Rev. Genet. 7, 645–653 (2006).
  Y. Xiong et al., RNA sequencing shows no dosage compensation of the active X-chromosome. Nat.
- Y. Xiong et al., KNA sequencing shows no dosage compensation of the active X-chromosome. Nat. Genet. 42, 1043–1047 (2010).
- B. Charlesworth, The evolution of chromosomal sex determination and dosage compensation. *Curr. Biol.* 6, 149–162 (1996).
- R. Bergero, S. Qiu, D. Charlesworth, Gene loss from a plant sex chromosome system. *Curr. Biol.* 25, 1234–1240 (2015).
- Q. Zhou, D. Bachtrog, Chromosome-wide gene silencing initiates Y degeneration in *Drosophila*. *Curr. Biol.* 22, 522–525 (2012).
- D. Prentout et al., Plant genera Cannabis and Humulus share the same pair of well-differentiated sex chromosomes. New Phytol. 231, 1599–1611 (2021).
- P. A. Eyer, A. J. Blumenfeld, E. L. Vargo, Sexually antagonistic selection promotes genetic divergence between males and females in an ant. Proc. Natl. Acad. Sci. U.S.A. 116, 24157–24163 (2019).
- A. Muyle, G. A. B. Marais, V. Bacovsky, R. Hobza, T. Lenormand, Dosage compensation evolution in plants: Theories, controversies and mechanisms. *Philos Trans. R Soc. Lond B Biol. Sci.* 377, 20210222 (2022).
- W. Zhang, X. Wang, Q. Yu, R. Ming, J. Jiang, DNA methylation and heterochromatinization in the malespecific region of the primitive Y chromosome of papaya. *Genome Res.* 18, 1938–1943 (2008).
- D. C. H. Metzger, P. M. Schulte, The DNA methylation landscape of stickleback reveals patterns of sex chromosome evolution and effects of environmental salinity. *Genome Biol. Evol.* 10, 775–785 (2018).
- W. Reik, A. Lewis, Co-evolution of X-chromosome inactivation and imprinting in mammals. *Nat. Rev. Genet.* 6, 403–410 (2005).
- A. Muyle, D. Bachtrog, G. A. B. Marais, J. M. A. Turner, Epigenetics drive the evolution of sex
- chromosomes in animals and plants. *Philos Trans. R Soc. Lond B Biol. Sci.* 376, 20200124 (2021).
  Y. F. Zhou *et al.*, The population genetics of structural variants in grapevine domestication. *Nat. Plants* 5, 965–979 (2019).
- D. Tang *et al.*, Genome evolution and diversity of wild and cultivated potatoes. *Nature* **606**, 535–541 (2022).
- Y. Zhou et al., Graph pangenome explusions missing heritability and empowers tomato breeding. Nature 606, 527–534 (2022).
- Z. J. Liu *et al.*, Grapevine pangenome facilitates trait genetics and genomic breeding. *Nat. Genet.* 56, 2804–2814 (2024). 10.1038/s41588-024-01967-5.
- Q. M. Long et al., Population comparative genomics discovers gene gain and loss during grapevine domestication. Plant Physiol. 195, 1401–1413 (2024).
- A. M. Vondras *et al.*, The genomic diversification of grapevine clones. *BMC Genomics* 20, 972 (2019).
- Y. Kou et al., Evolutionary genomics of structural variation in asian rice (*Oryza sativa*) domestication. *Mol. Biol. Evol.* 37, 3507–3524 (2020).
- E. Solares et al., Insights into the domestication of avocado and potential genetic contributors to heterodichogamy. Genes Genom. Genet. 13, jkac323 (2023), 10.1093/g3journal/jkac323.
- P. Carbonell-Bejerano *et al.*, Catastrophic unbalanced genome rearrangements cause somatic loss of berry color in grapevine. *Plant Physiol.* **175**, 786–801 (2017).
- G. Ríos et al., Characterization of hemizygous deletions in using array-comparative genomic hybridization and microsynteny comparisons with the poplar genome. *BMC Genomics* 9, 1869143 (2008).
- S. Ban, I. El-Sharkawy, J. T. Zhao, Z. J. Fei, K. N. Xu, An apple somatic mutation of delayed fruit maturation date is primarily caused by a retrotransposon insertion-associated large deletion. *Plant* J. **111**, 1609–1625 (2022).
- Q. L. Li et al., Causal relationship between the loss of RUNX3 expression and gastric cancer. Cell 109, 113–124 (2002).
- D. Wei et al., Drastic down-regulation of Kruppel-like factor 4 expression is critical in human gastric cancer development and progression. *Cancer Res.* 65, 2746-2754 (2005).
- P. J. Teoh et al., p53 haploinsufficiency and functional abnormalities in multiple myeloma. Leukemia 28, 2066–2074 (2014).
- T. Domazet-Loso, J. Brajkovic, D. Tautz, A phylostratigraphy approach to uncover the genomic history of major adaptations in metazoan lineages. *Trends Genet.* 23, 533–539 (2007).
- J. D. Hollister, B. S. Gaut, Epigenetic silencing of transposable elements: A trade-off between reduced transposition and deleterious effects on neighboring gene expression. *Genome Res.* 19, 1419-1428 (2009).
- J. D. Hollister *et al.*, Transposable elements and small RNAs contribute to gene expression divergence between *Arabidopsis thaliana* and *Arabidopsis lyrata*. *Proc. Natl. Acad. Sci. U.S.A.* **108**, 2322–2327 (2011).
- C. D. Hirsch, N. M. Springer, Transposable element influences on gene expression in plants. *Biochim Biophys. Acta Gene Regul Mech* 1860, 157–165 (2017).
- J. M. Noshay *et al.*, Assessing the regulatory potential of transposable elements using chromatin accessibility profiles of maize transposons. *Genetics* 217, 1–13 (2021).
- A. M. Muyle, D. K. Seymour, Y. Lv, B. Huettel, B. S. Gaut, Gene body methylation in plants: Mechanisms, functions, and important implications for understanding evolutionary processes. *Genome Biol. Evol.* 14, evac038 (2022).
- D. K. Seymour, B. S. Gaut, Phylogenetic shifts in gene body methylation correlate with gene expression and reflect trait conservation. *Mol. Biol. Evol.* 37, 31–43 (2020).
- J. Ramos-Madrigal et al., Palaeogenomic insights into the origins of French grapevine diversity. Nat. Plants 5, 595–603 (2019).
- B. S. Gaut, D. K. Seymour, Q. Liu, Y. Zhou, Demography and its effects on genomic variation in crop domestication. *Nat. Plants* 4, 512–520 (2018).
- H. Xiao et al., Adaptive and maladaptive introgression in grapevine domestication. Proc. Natl. Acad. Sci. U.S.A. 120, e2222041120 (2023).
- Y. Zhou, M. Massonnet, J. S. Sanjak, D. Cantu, B. S. Gaut, Evolutionary genomics of grape (Vitis vinifera ssp. vinifera) domestication. Proc. Natl. Acad. Sci. U.S.A. 114, 11715–11720 (2017).
- P. Ramu et al., Cassava haplotype map highlights fixation of deleterious mutations during clonal propagation. Nat. Genet. 49, 959-963 (2017).

- P. Stenberg, J. Larsson, Buffering and the evolution of chromosome-wide gene regulation. Chromosoma 120, 213-225 (2011).
- J. H. Malone et al., Mediation of Drosophila autosomal dosage effects and compensation by network interactions. Genome Biol. 13, r28 (2012).
- J. E. Mank, Sex chromosome dosage compensation: Definitely not for everyone. *Trends Genet.* 29, 677–683 (2013).
- M. Alonge et al., Major impacts of widespread structural variation on gene expression and crop improvement in Tomato. Cell 182, 145–161.e123 (2020).
- D. E. Shaw, M. A. White, The evolution of gene regulation on sex chromosomes. *Trends Genet.* 38, 844–855 (2022).
- T. Shi *et al.*, The super-pangenome of *Populus* unveils genomic facets for its adaptation and diversification in widespread forest trees. *Mol. Plant* 17, 725-746 (2024).
- A. Muyle, J. Ross-Ibarra, D. K. Seymour, B. S. Gaut, Gene body methylation is under selection in Arabidopsis thaliana. Genetics 218, iyab061 (2021).
- J. M. Belton et al., Hi-C: A comprehensive technique to capture the conformation of genomes. Methods 58, 268–276 (2012).
- X. Y. Shi et al., The complete reference genome for grapevine (Vitis vinifera L.) genetics and breeding. Hortic Res-England 10, uhad061 (2023).
- M. Alonge et al., RaGOO: Fast and accurate reference-guided scaffolding of draft genomes. Genome Biol. 20, 224 (2019).
- N. C. Durand *et al.*, Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell Syst.* 3, 95–98 (2016).
- O. Dudchenko et al., De novo assembly of the Aedes aegypti genome using Hi-C yields chromosome-length scaffolds. Science 356, 92–95 (2017).
- J. T. Robinson *et al.*, Juicebox.js provides a cloud-based visualization system for Hi-C data. *Cell Syst.* 6, 256–258.e251 (2018).
- H. Li, Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics* 34, 3094–3100 (2018).
- J. T. Robinson, H. Thorvaldsdottir, D. Turner, J. P., Mesirov, igv.js: An embeddable javascript implementation of the integrative genomics viewer (IGV). *Bioinformatics* 39, btac830 (2023).
- C. Holt, M. Yandell, MAKER2: An annotation pipeline and genome-database management tool for second-generation genome projects. *BMC Bioinformatics* 12, 491 (2011).
- A. Shumate, S. L. Salzberg, Liftoff: Accurate mapping of gene annotations. Bioinformatics 37, 1639–1643 (2021).
- X. Wang et al., Integrative genomics reveals the polygenic basis of seedlessness in grapevine. Curr. Biol. 34, 3763–3777.e3765 (2024).
- T. Zhang et al., Population genomics highlights structural variations in local adaptation to saline coastal environments in woolly grape. J. Integr Plant Biol. 66, 1408–1426 (2024).
- Y. Luo et al., Phased T2T genome assemblies facilitate the mining of disease-resistance genes in Vitis davidii. Hortic. Res. uhae306 (2024).
- J. M. Flynn et al., RepeatModeler2 for automated genomic discovery of transposable element families. Proc. Natl. Acad. Sci. U.S.A. 117, 9451–9457 (2020).
- S. Ou *et al.*, Benchmarking transposable element annotation methods for creation of a streamlined, comprehensive pipeline. *Genome Biol.* 20, 275 (2019).
- F. J. Sedlazeck *et al.*, Accurate detection of complex structural variations using single-molecule sequencing. *Nat. Methods* 15, 461–468 (2018).
- G. Marçais et al., MUMmer4: A fast and versatile genome alignment system. Plos Comput. Biol. 14, e1005944 (2018).
- G. Benson, Tandem repeats finder: A program to analyze DNA sequences. Nucleic Acids Res. 27, 573–580 (1999).
- Y. Wang et al., MCScanX: A toolkit for detection and evolutionary analysis of gene syntemy and collinearity. Nucleic Acids Res. 40, e49 (2012).
- 74. C. Camacho et al., BLAST+: Architecture and applications. BMC Bioinformatics 10, 421 (2009).
- S. Kumar, G. Stecher, D. Peterson, K. Tamura, MEGA-CC: Computing core of molecular evolutionary genetics analysis program for automated and iterative data analysis. *Bioinformatics* 28, 2685–2686 (2012).
- D. M. Emms, S. Kelly, OrthoFinder: Phylogenetic orthology inference for comparative genomics. Genome Biol. 20, 238 (2019).
- B. T. Sherman *et al.*, DAVID: A web server for functional enrichment analysis and functional annotation of gene lists (2021 update). *Nucleic Acids Res.* 50, W216–W221 (2022).
- A. M. Bolger, M. Lohse, B. Usadel, Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics* 30, 2114–2120 (2014).
- D. Kim, J. M. Paggi, C. Park, C. Bennett, S. L. Salzberg, Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* 37, 907–915 (2019).
- Y. Liao, G. K. Smyth, W. Shi, featureCounts: An efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* **30**, 923–930 (2014).
- F. Krueger, S. R. Andrews, Bismark: A flexible aligner and methylation caller for Bisulfite-Seq applications. *Bioinformatics* 27, 1571–1572 (2011).
- B. Langmead, S. L. Salzberg, Fast gapped-read alignment with Bowtie 2. Nat. Methods 9, 357–359 (2012).
- R. Lister et al., Highly integrated single-base resolution maps of the epigenome in Arabidopsis. Cell 133, 523-536 (2008).
- G. Barturen, A. Rueda, J. L. Oliver, M. Hackenberg, MethylExtract: High-Quality methylation maps and SNV calling from whole genome bisulfite sequencing data. *F1000Res.* 2, 217 (2013).
- F. Ramirez, F. Dundar, S. Diehl, B. A. Gruning, T. Manke, deepTools: A flexible platform for exploring deep-sequencing data. *Nucleic Acids Res.* 42, W187-W191 (2014).
- S. Takuno, B. S. Gaut, Body-methylated genes in Arabidopsis thaliana are functionally important and evolve slowly. Mol. Biol. Evol. 29, 219-227 (2012).
- Y. Peng et al., The genomic and epigenomic landscapes of hemizygous genes across crops with contrasting reproductive systems. NCBI BioProject. https://dataview.ncbi.nlm.nih.gov/object/ PRJNA1178252. Deposited 27 October 2024.
- Y. Peng, The genomic and epigenomic landscapes of hemizygous genes across crops with contrasting mating systems [Data set]. Zenodo. https://doi.org/10.5281/zenodo.14015567. Deposited 20 October 2024.
- Y. Peng et al., The genomic and epigenomic landscapes of hemizygous genes. GitHub. https://github. com/zhouyflab/Genomic\_Epigenomic\_Hemizygous\_Crops. Deposited 20 October 2024.