



The Draft Genome of *Kochia scoparia* and the Mechanism of Glyphosate Resistance via Transposon-Mediated *EPSPS* Tandem Gene Duplication

Eric L. Patterson^{1,2}, Christopher A. Saski², Daniel B. Sloan ³, Patrick J. Tranel⁴, Philip Westra¹, and Todd A. Gaines ^{1,*}

¹Department of Bioagricultural Sciences and Pest Management, Colorado State University

²Department of Genetics and Biochemistry, Clemson University

³Department of Biology, Colorado State University

⁴Department of Crop Sciences, University of Illinois, Urbana

*Corresponding author: E-mail: todd.gaines@colostate.edu

Accepted: September 8, 2019

Data deposition: All raw sequence read files for the whole genome sequencing have been deposited in the Sequence Read Archive database at NCBI under BioProject ID PRJNA526487 (SRR8835960–SRR8835963). The genome assembly was submitted to the NCBI genomes database with the accession SNQN00000000.

Abstract

Increased copy number of the 5-enolpyruvylshikimate-3-phosphate synthase (*EPSPS*) gene confers resistance to glyphosate, the world's most-used herbicide. There are typically three to eight *EPSPS* copies arranged in tandem in glyphosate-resistant populations of the weed kochia (*Kochia scoparia*). Here, we report a draft genome assembly from a glyphosate-susceptible kochia individual. Additionally, we assembled the *EPSPS* locus from a glyphosate-resistant kochia plant by sequencing select bacterial artificial chromosomes from a kochia bacterial artificial chromosome library. Comparing the resistant and susceptible *EPSPS* locus allowed us to reconstruct the history of duplication in the structurally complex *EPSPS* locus and uncover the genes that are coduplicated with *EPSPS*, several of which have a corresponding change in transcription. The comparison between the susceptible and resistant assemblies revealed two dominant repeat types. Additionally, we discovered a mobile genetic element with a FHY3/FAR1-like gene predicted in its sequence that is associated with the duplicated *EPSPS* gene copies in the resistant line. We present a hypothetical model based on unequal crossing over that implicates this mobile element as responsible for the origin of the *EPSPS* gene duplication event and the evolution of herbicide resistance in this system. These findings add to our understanding of stress resistance evolution and provide an example of rapid resistance evolution to high levels of environmental stress.

Key words: genomics, weed biology, molecular evolution, herbicide resistance, mobile genetic element, gene duplication.

Introduction

Gene copy number variation is an important source of genetic variation that can be deleterious in some cases, such as causing cancer in humans, but can also increase genetic variation and lead to adaptations (Schimke et al. 1985; Lynch and Conery 2000; DeBolt 2010; Xi et al. 2011; Hull et al. 2017). This is especially true in plants where novel genetic variation is essential in the face of rapidly changing environments (DeBolt 2010). Increases in copy number of the 5-enolpyruvylshikimate-3-phosphate synthase (*EPSPS*) gene confer resistance

to glyphosate, the world's most-used herbicide, in several plant species (reviewed in Sammons and Gaines 2014). Increased *EPSPS* gene copy number results in the overproduction of the *EPSPS* protein, glyphosate's target (Gaines et al. 2010; Wiersma et al. 2015), making it necessary for the application of more glyphosate to have the same lethal effect (Vila-Aiub et al. 2014; Godar et al. 2015; Gaines et al. 2016; Koo et al. 2018). This phenomenon has been observed in eight weed species to date; however, the DNA sequence surrounding the *EPSPS* gene duplication has only been

resolved in one species, *Amaranthus palmeri* (Molin et al. 2017; Patterson et al. 2018), as most weed species do not have sequenced genomes. In the case of *A. palmeri*, *EPSPS* gene duplication is caused by a large, circular, extrachromosomal DNA element that replicates autonomously from the nuclear genome (Molin et al. 2017; Koo et al. 2018). This mechanism results in *A. palmeri* plants with up to hundreds of *EPSPS* copies (Gaines et al. 2010).

Recently, *EPSPS* gene duplication has been described in the weed species *Kochia scoparia* (*kochia*, syn. *Bassia scoparia*), one of the most important weeds in the Central Great Plains of United States and Canada (Beckie et al. 2013, 2015, 2018; Jugulam et al. 2014; Kumar et al. 2015; Wiersma et al. 2015; Gaines et al. 2016; Martin et al. 2017). In glyphosate-resistant *kochia*, *EPSPS* copy numbers typically range from three to eight with the highest reports at 11 copies (Gaines et al. 2016). In contrast to the extrachromosomal element observed in *A. palmeri*, fluorescence in situ hybridization (FISH) has shown that the *EPSPS* copies in *kochia* are arranged in tandem at a single chromosomal locus and are most likely generated by unequal crossing over (Jugulam et al. 2014). More detailed cytogenetics studies using Fiber-FISH estimated that most repeats of the *EPSPS* loci are either 45 or 66 kb in length. Inverted repeats and repeats of 70 kb in length were also observed (Jugulam et al. 2014). The initial event that started *EPSPS* gene duplication, the fine-scale sequence variation between the various types of repeats, and the other genes that may be coduplicated with *EPSPS* remain unresolved.

Understanding how gene copy number variants form and their potential phenotypic consequences is essential for determining how plants adapt to their environment and thrive in adverse conditions. In this article, we assembled a rough-draft genome of a glyphosate-susceptible *kochia* plant. We then identified the contig containing the *EPSPS* locus and investigated the genes that are coduplicated with *EPSPS*, their transcription in glyphosate-resistant and susceptible plants, and through whole-genome resequencing of a glyphosate-resistant plant, discovered the up- and downstream borders of the duplicated region. We also sequenced and assembled the *EPSPS* locus from a glyphosate-resistant *kochia* plant using bacterial artificial chromosomes (BACs) probed for 1) the *EPSPS* gene, 2) the downstream junction, and 3) the upstream junction. After assembling four BACs, we generated a model sequence of the *EPSPS*-duplicated locus containing six instances of the *EPSPS* gene. We discovered two dominant repeat types, an inversion, and rarer repeats of different sizes using a combination of qPCR markers, genomic resequencing, and RNA-Seq data. Through this analysis, we also discovered a 16-kb mobile genetic element (MGE) that is associated with the gene duplication event. This MGE contains four putative coding sequences. We hypothesize that the insertion of this MGE downstream of the *EPSPS* gene is responsible for a disruption of this region and the origin of the *EPSPS* gene duplication event.

Materials and Methods

Tissue Collection and Nucleic Acid Extraction

The herbicide-susceptible *K. scoparia* line “7710” (Preston et al. 2009; Pettinga et al. 2018) was used for genomic sequencing. All plants in this line were consistently controlled by glyphosate treatments at field rates of 860 g a.e. ha⁻¹. Plants were grown in a greenhouse at Colorado State University. After seeds germinated, they were transferred into 4-l pots filled with Fafard 4 P Mix supplemented with Osmocote fertilizer (Scotts Co. LLC), regularly watered, and grown under a 16-h photoperiod. Temperatures in the greenhouse cycled between 25 °C day and 20 °C nights. A single, healthy individual was selected for tissue collection.

A glyphosate-resistant line (M32) was obtained from a field population near Akron, Colorado (40.162382, -103.172849) in the autumn of 2012. After glyphosate failed to control these plants in the field, seed was collected from ten surviving individuals. Seeds were germinated and treated with 860 g a.e. ha⁻¹ of glyphosate and ammonium sulfate (2% w/v). Survivors were then collected, crossed, and seed was collected. This process was repeated for three generations until no susceptible individuals were observed in the progeny. All plants were confirmed to have elevated *EPSPS* copy number using genomic qPCR (Gaines et al. 2016).

For shotgun genome Illumina sequencing of the two lines, DNA was extracted from samples using a modified CTAB extraction protocol (see [Supplementary Material](#) online). For large-fragment, genomic PacBio sequencing of the glyphosate-susceptible line, the CTAB protocol was further modified to obtain more DNA of sufficiently large size (>10 kb) (see [Supplementary Material](#) online). For RNA-Seq, four susceptible and four resistant plants were grown in the greenhouse as described above, until they were ~10 cm tall and 100 mg of young expanding leaf tissue was taken from each plant. RNA was extracted from young leaf tissue from four plants from each of the glyphosate-susceptible and resistant lines using the Qiagen RNeasy Plus Mini Kit. Each replicate sample was normalized to a total mass of 200 ng total RNA.

Sequencing Libraries

Three genomic DNA libraries of glyphosate-susceptible *kochia* DNA were prepared for Illumina sequencing on a HiSeq 2500 at the University of Illinois, Roy J. Carver Biotechnology Center for genome assembly. First, DNA was size selected to 240 bp so that there was overlap between the read pairs in a high-coverage, short-insert library sequenced on one full flow cell (eight lanes) for use with ALLPATHS-LG. Second, two large insert, mate-pair libraries (5 and 10 kb) were each run on one lane at 2 × 150 bp.

Additionally, genomic DNA from the glyphosate-resistant line was prepared for Illumina sequencing using the Genomic DNA Sample Prep Kit from Illumina following the

manufacturer's protocols and sequenced on one entire lane of a HiSeq 2500 flow cell. Quality of the raw Illumina sequence reads was assessed using FASTQC v0.10.1. Adapters were removed using Trimmomatic version 0.60 with the parameters "ILLUMINACLIP: tranel_adaptors.fa: 2:30:10 TRAILING: 30 LEADING: 30 MINLEN: 45" using these adapters: "AGATCGGAAGAGCAC" and "AGATCGGAAGAGCGT."

A large insert DNA library for PacBio sequencing was generated at the UC Davis Genome Center using the PacBio SMRT Library Prep for RSII followed by BluePippin size selection for fragments >10 kb. The library was sequenced with 12 PacBio SMRT cells using the RSII chemistry after a titration cell to determine optimal loading. In total, 2,760,348 PacBio reads were generated with a read N50 of 6,576 bp with the largest read being 41,738 bp.

Strand-specific RNA-Seq libraries were prepared robotically on a Hamilton Star Microlab at the Clemson University Genomics and Computational Facility following in-house automation procedures generally based on the TruSeq Stranded mRNAseq preparation guide. The prepared libraries were pooled and 100-bp paired-end reads were generated using a NextSeq 500/550.

Susceptible Genome Assembly

Two different assemblies were generated that integrated the PacBio and Illumina data of the susceptible kochia line. These two assemblies were then compared and merged by consensus for a single final assembly referred to as KoSco-1.0. For the first assembly, raw PacBio reads were error corrected using the high-coverage, paired-end Illumina library with the error correcting software Proovread 2.13.11 (Hackl et al. 2014). Proovread was run with standard parameters, using the high-coverage 150 bp, paired-end Illumina library on each SMRT cell individually. Error corrected reads were then assembled using the Celera Assembler fork for long reads, Canu 1.0 (Koren et al. 2017). Canu was run with a predicted genome size of 1 Gb, and the PacBio-corrected settings. For the second assembly, an initial ALLPATHS-LG v r52488 assembly was made with all three Illumina libraries (Butler et al. 2008). ALLPATHS was run assuming a haploid genome of 1 Gb. The resulting contigs were then scaffolded using the uncorrected PacBio reads and the software PBJelly 15.8.24 (English et al. 2012). PBJelly was run with the following blasr settings: "-minMatch 8 -sdpTupleSize 8 -minPctIdentity 75 -bestn 1 -nCandidates 10 -maxScore -500 -nproc 19 -noSplitSubreads." The two assemblies were then merged with GARM Meta assembler 0.7.3 to get a final version of the genome assembly for our analysis (Mayela Soto-Jimenez et al. 2014). The assembly from ALLPATHS was set to assembly "A" and the assembly from Canu was set as genome "B." All other parameters were kept standard. We refer to the resulting meta-assembly as KoSco-1.0.

Genome Annotation

The merged assembly was annotated with the WQ-Maker 2.31.8 pipeline in conjunction with CyVerse (Cantarel et al. 2008; Thrasher et al. 2014). WQ-Maker was informed with kochia transcriptome from Wiersma et al. (2015), all expressed sequence tags from the Chenopodiaceae downloaded from NCBI, all protein sequence from the Chenopodiaceae family downloaded from NCBI, and Augustus using *Arabidopsis thaliana* gene models. The resulting predictions were then used to train SNAP (February 16, 2013) through two rounds for final gene model predictions. Gene space completeness was assessed using BUSCO v3 and the eudicotyledons *odb10* prerelease data set using standard parameters (Simão et al. 2015).

The predicted gene transcripts (mRNA) and predicted translated protein sequence were then annotated using Basic Local Alignment Search Tool Nucleotide (BlastN) and Protein (BlastP) 2.2.18+ for similarity to known transcripts and proteins, respectively. Alignments were made to the entire NCBI nucleotide and protein databases. For all BLAST homology searches, the e-value was set at 1e-25 and only the best match was considered. The predicted proteins were further annotated using InterProScan 5.28-67.0 for protein domain predictions (Mi et al. 2004; Camacho et al. 2009; Jones et al. 2014). InterProScan was run using standard settings. The complete assembly was analyzed using RepeatMasker 4.0.6 to search for small interspersed repeats, DNA transposon elements, and other known repetitive elements using the "Viridiplantae" repeat database and standard search parameters (Tarailo-Graovac and Chen 2009).

Genomic Resequencing of Glyphosate-Resistant Kochia and Differential Gene Expression

Genomic resequencing reads from the glyphosate-resistant plant were aligned to the KoSco-1.0 genome assembly using the BWA-backtrack alignment program with default parameters (Li and Durbin 2009). The boundaries of the *EPSPS* copy number variant were manually detected where coverage dramatically increased up- and downstream of the *EPSPS* gene.

RNA-Seq reads from susceptible and resistant plants were aligned to the gene models from the genome assembly using the mem algorithm from the BWA alignment program version 0.7.15 under standard parameters. Read counts for each gene were extracted from this alignment using the software featureCounts in the Subread 1.6.0 package and the gene annotation generated by WQ-Maker (Liao et al. 2014). Expression level and differential expression between the glyphosate susceptible and glyphosate-resistant plants for all genes were calculated with the EdgeR package using the quasi-likelihood approach in the generalized linear model (glm) framework as described in the user manual (Robinson et al. 2010).

Assembling the *EPSPS* Locus from a Glyphosate-Resistant Plant

A library of BACs was generated from a single glyphosate-resistant kochia plant selected from the glyphosate-resistant population following the protocol described in Luo and Wing (2003) with modifications as described in Molin et al. (2017). High-molecular weight (HMW) DNA was extracted from young leaf tissue from a single glyphosate-resistant plant using a modified CTAB DNA extraction protocol. This HMW DNA was ligated to a linearized vector and transformed into *Escherichia coli* using electroporation. Recombinant colonies were then grown on LB plates. Radiolabeled probes were designed for the *EPSPS* gene itself, a sequence upstream, and a sequence downstream of the *EPSPS* CNV. Predicted locations for the probes were determined by looking at the alignment of shotgun Illumina data from the glyphosate-resistant line against the contig containing *EPSPS* in the genome assembly. Several colonies containing the appropriate sequences were identified for each probe. These identified BACs were end sequenced to determine their approximate location and run on pulse-field gel electrophoresis to determine their approximate size. Colonies containing positive BACs of the correct position and size were isolated and cultured. HMW DNA was extracted from these colonies and prepared using a SMRTbell Template Prep Kit, 1.0 using the manufacturer-recommended protocols. Finally, the HMW DNA was sent for RSII PacBio sequencing on two SMRT cells performed at The University of Delaware, DNA Sequencing & Genotyping Center.

PacBio reads were assembled using the software Canu (Koren et al. 2017). The BAC vector sequence was then removed from the assembled contigs. Using the known size of the BACs, their end-sequences and the corresponding contig from the susceptible genome assembly, entire BAC sequences were reconstructed manually from the contigs produced by CANU. These “full-length” BACs were then aligned, and overlaps were used to generate the largest contiguous length possible. This BAC meta-assembly was aligned to the susceptible contig from the genome assembly containing the *EPSPS* gene using YASS. Additionally, the BAC insert sequences were run through the MAKER pipeline, informed with cDNA and protein annotations from the Chenopodiaceae and the gene models from the kochia genome (Cantarel et al. 2008) for gene annotation. This BAC assembly led to the discovery of two dominant repeat types (a full length 56.1-kb repeat and a smaller 32.9-kb repeat), the up- and downstream boundaries of the CNV, as well as a large MGE that was interspersed in the repeat structure.

Using the Illumina genomic resequencing data from the resistant line, we calculated the copy number of four regions from the CNV by read depth as follows: 1) the region directly upstream of the CNV; 2) the region directly downstream of the CNV; 3) the MGE; and 4) the full length, 56.1-kb repeat.

This 56.1-kb repeat was then subdivided into the region only present within the 56.1-kb repeat and the region that is shared between the 56.1-kb repeat and a smaller 32.9-kb repeat. Highly repetitive regions and those containing transposable elements were masked for the alignment of resequencing reads. Genomic resequencing reads from the glyphosate-resistant plant were aligned to these units using the BWA-backtrack alignment program using standard parameters. The number of reads mapping to each unit was calculated and divided by the length of that region to get the average number of reads per unmasked DNA length. The up- and downstream read depths were averaged and used to standardize the read depths of each of the four units. These standardized read depths correspond with the predicted copy number of each unit.

Constructing the MGE in the Susceptible Line

The raw PacBio reads from the genome assembly were aligned against the MGE from the resistant *EPSPS* loci using minimap2 with the “map-pb” preset parameters (Li 2018). Variants from this alignment were called using Samtools v 1.9 mpileup (Li et al. 2009). The resulting variant call file and the MGE sequence of resistant *EPSPS* loci were used to make a consensus sequence for the susceptible line using BCftools v 1.9 consensus command (Li et al. 2009). Commands were run using standard parameters unless otherwise noted. The MGE sequences from the susceptible and resistant lines were aligned using YASS (Noé and Kucherov 2005).

Markers for Confirming the Structure of the *EPSPS* CNV

Primers were designed that were spaced at regular intervals (~5–15 kb) along the susceptible contig that spanned the putative CNV area for genomic qPCR analysis (supplementary table 1, Supplementary Material online). Additionally, qPCR primers were designed that spanned the junctions of the two dominant repeat types, the up- and downstream boundaries of the CNV, as well as for the MGE (supplementary table 1, Supplementary Material online). Primers were designed to closely mimic the primers already published for the *EPSPS* gene (Wiersma et al. 2015), including a melting temperature between 51 and 56 °C, a GC content between 40% and 50%, and a length of between 20 and 24 bp. Furthermore, the resulting amplicon had to be between 100 and 200 bp long. All genomic PCR was performed using the same protocol established for *EPSPS* copy number assay (Gaines et al. 2016).

For genomic PCR screening of kochia populations for these repeat features, both susceptible and resistant plants were grown in the greenhouse until they were ~10 cm tall and 100 mg of young expanding leaf tissue was taken from each plant. DNA was extracted from this tissue using the recommended protocol from the DNeasy Plant Mini Kit.

The DNA quality and concentration were checked using a NanoDrop 1000 and diluted to 5 ng/μl. For qPCR, two genes were used as single-copy controls: acetolactate synthase (*ALS*) and copalyl diphosphate synthetase 1 (*CPS*). Each qPCR reaction consisted of 12.5 μl PerfeCTa SYBR green Super Mix (Quanta Biosciences), 1 μl of the forward and reverse primers at 10 μM, 10 ng gDNA (2 μl), and 9.5 μl of sterile water for a total volume of 25 μl.

A BioRad CFX Connect Real-Time System was used for qPCR. The temperature cycle for all reactions was as follows: an initial 3 min at 95 °C followed by 35 rounds of 95 °C for 30 s and 53 °C for 30 s with a fluorescence reading at 497 nm after each round. A melt curve was performed from 65 to 95 °C in 0.5 °C increments for each reaction to verify the production of a single PCR product. Additionally, all products from a susceptible line were run on a 1.5% agarose gel to verify a single product with low to no primer dimerization. Relative quantification was calculated using the comparative C_t method: $2^{\Delta C_t}$ ($\Delta C_t = (C_t^{(ALS)} + C_t^{(CPS)})/2 - C_t^{(EPSPS)}$) (Schmittgen and Livak 2008).

Results

Genome Assembly and Annotation

The KoSco-1.0 assembly consisted of 19,671 scaffolds, spanning 711 Mb. The longest scaffold was 770 kb and the N50 was 62 kb for this assembly. Approximately 9.43% of the base pairs were unknown “N” bases that serve only as scaffolding and distance information (supplementary table 2, Supplementary Material online). After annotation with Maker, 47,414 genes were predicted in KoSco-1.0 with an average transcript length of 943 bp (supplementary table 3, Supplementary Material online), compared with the 27,429 genes in *Beta vulgaris* (Dohm et al. 2014). KoSco-1.0 was analyzed using BUSCO for completeness, which found 1,490 out of 2,121 (70.3%) ultraconserved genes from the eudicotyledons *odb10* data set (supplementary table 4, Supplementary Material online). Approximately 62% of predicted kochia genes found one or more matches in the NCBI database(s) using a BLAST e-value < 1e-25 and almost 82% of predicted proteins were assigned one or more functional InterPro domain(s) (supplementary table 3, Supplementary Material online). RepeatMasker uncovered 6.25% of the genome assembly consisting of interspersed repeats with the largest proportion consisting of LTR elements of either the Ty1/Copia or Gypsy/DIRS1 variety. Simple repeats made up ~2.5% of the assembly (supplementary table 5, Supplementary Material online). The genome assembly is of lower quality than expected (e.g., scaffold length, gene content coverage) based on the hybrid assembly approach incorporating both Illumina and PacBio reads. The fragmented assembly may be partly due to remaining heterozygosity across the genome in the sequenced line, and/or due to sequence complexity in the kochia genome that remains to be

resolved. However, the genome assembly did enable detailed analysis of the *EPSPS* gene duplication.

The *EPSPS* Locus and Differential Gene Expression

The contig containing the *EPSPS* locus from the susceptible genome assembly was 399,779 bp long. The *EPSPS* gene model was 5,551 bp long (UTRs, exons, and introns included) and located between base pairs 91,663–97,214 of the contig. When this contig was aligned to *Beta vulgaris* near perfect synteny was observed; however, when compared with the sequence responsible for duplicating *EPSPS* from *A. palmeri*, little similarity existed outside of the *EPSPS* gene itself (fig. 1).

When shotgun Illumina genomic reads from the glyphosate-resistant line were aligned to the contig, the read depth of *EPSPS* and its surrounding area was much greater (>7.26-fold) than the background read depth (supplementary fig. 1, Supplementary Material online). Using this alignment, it was possible to predict the exact boundaries of the *EPSPS* CNV starting at base pair 41,684 and continuing to base pair 101,128. This region contains seven coding genes of various functions including *EPSPS* itself (table 1). When differential expression of all genes in the genome was calculated using RNA-Seq data, five of the genes in this region showed over expression in the glyphosate-resistant line, one gene showed underexpression in the glyphosate-resistant line, and one showed no significant difference (FDR adjusted *P* value < 0.05) (table 1). The lack of differential expression for some duplicated genes in the repeat between the resistant and susceptible plants may be because these genes are developmentally regulated and expressed in other developmental stages than the leaf vegetative stage sampled in this experiment, or they may be regulated in response to specific environmental conditions. When the *EPSPS* contig was aligned to itself, there was no evidence for sequence complexity (simple sequence repeats, inverted repeats, and self-homology) at the predicted boundaries of the CNV (supplementary fig. 2, Supplementary Material online).

The *EPSPS* Locus from a Glyphosate-Resistant Plant

Using PacBio data of four BACs from a glyphosate-resistant plant, we assembled four contigs that were 129.0 kb for the BAC detected with the upstream probe, 134.2 kb for the BAC detected with the downstream probe, and 140.5 and 78.0 kb for two BACs detected with the *EPSPS* probe. The whole BAC assembly was 429,317 bp long and encompassed six repeats of the *EPSPS* gene and a significant portion of the up- and downstream sequence (sequence available as supplementary file “all_fasta.txt,” Supplementary Material online).

The largest and most complete repeat was 56.1 kb long and contained the entire region predicted from the alignment of resistant Illumina data against the susceptible *EPSPS* contig, including all seven of the predicted genes in this region. The

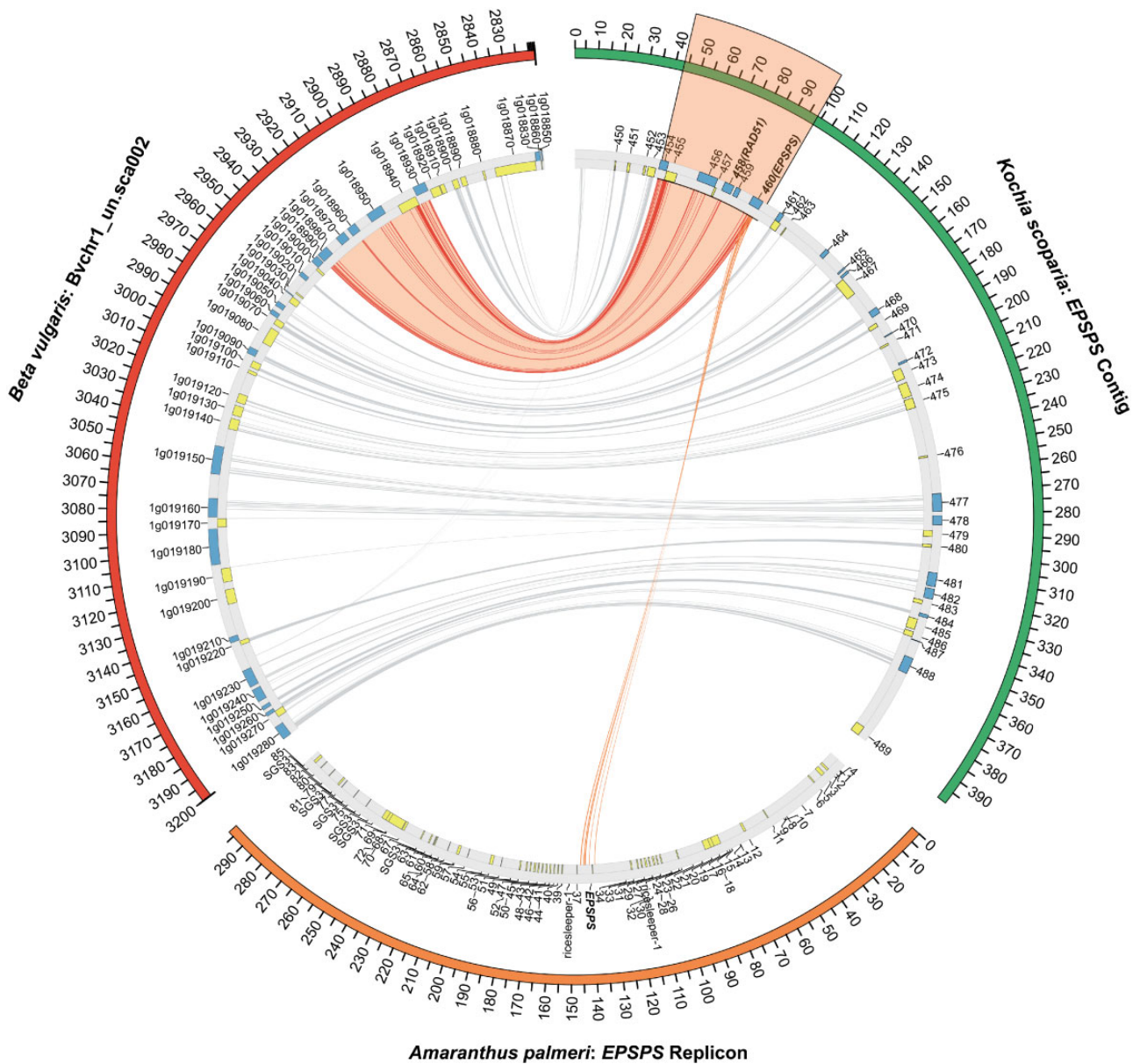


Fig. 1.—A comparison of the *EPSPS* contig from *Kochia* (Green), an genomic scaffold from chromosome 1 of the *Beta vulgaris* genome (Red) (GenBank ID: KQ090199.1) (Dohm et al. 2014), and the *EPSPS* replicon from *Amaranthus palmeri* (Orange) (Molin et al. 2017). Blue and yellow blocks indicate genes in the forward and reverse orientation, respectively. The *EPSPS* gene is highlighted in orange. Red, connecting lines, indicate areas of high similarity between *Beta vulgaris* and *Kochia*. Orange, connecting lines indicate areas of high similarity between *A. palmeri* and *Kochia*. Number of base pairs in the alignment is listed on the outside track. The links between *Beta vulgaris* and *Kochia* that fall within the *EPSPS*-duplicated region are highlighted in orange.

second type was 32.7 kb and contained only four of the seven coduplicated genes from the 56.1-kb repeat, including *EPSPS* and the three genes immediately upstream of it. The third repeat was a full-length inversion of the 56.1-kb repeat. The fourth type of repeat was an 18.2-kb inverted repeat that contained only *EPSPS* and a fraction of one upstream gene. The fifth and final repeat structure was identified as a forward repeat of 33.1 kb, containing *EPSPS* and the three genes immediately upstream of it (fig. 2). All repeats end at

the same downstream base pair, directly after *EPSPS*; however, the beginning upstream base pair of each repeat type is variable (figs. 2 and 3; supplementary table 6, Supplementary Material online).

Enough overlap existed among the BAC contigs to composite all BAC assemblies together to make a representative sequence (meta-assembly) that contained two full-length 56.1-kb repeats and one of each of the other repeat types. Additionally, the flanking single-copy up- and downstream

Table 1List of Genes Near *EPSPS* That Are in or Flanking the *EPSPS* CNV Event

Gene	Beginning	Ending	Length	Orientation	Description	Part of the CNV?	Read Depth	DE	P Value
KS_00451	27,406	28,674	1,268	Reverse	GRAVITROPIC IN THE LIGHT 1-like	No	0	-0.43	0.00
KS_00452	35,728	36,696	968	Reverse	IRK-interacting protein	No	0	-2.62	0.05
KS_00453	37,839	41,640	3,801	Reverse	Nitroreductase family	No	0	0.74	0.00
KS_00454	43,124	47,121	3,997	Forward	Arginase 1, mitochondrial	56.1 kb	2.86	2.23	0.00
KS_00455	47,240	52,651	5,411	Reverse	Protein NRT1/ PTR FAMILY 7.2-like	56.1 kb	2.86	0.72	0.58
KS_00456	63,014	72,467	9,453	Forward	tRNA N6-adenosine threonylcarbamoyltransferase	56.1 kb and 32.7 kb	3.49	3.03	0.00
KS_00457	72,617	73,531	914	Reverse	Golgin subfamily A member 6-like	56.1 kb and 32.7 kb	3.49	-3.18	0.00
KS_00458	76,342	81,181	4,839	Forward	DNA repair protein RAD51	56.1 kb and 32.7 kb	3.46	1.33	0.00
KS_00459	82,421	84,836	2,415	Forward	Transketolase, chloroplastic-like	56.1 kb and 32.7 kb	3.29	3.83	0.00
KS_00460	91,663	97,214	5,551	Forward	3-phosphoshikimate 1-carboxyvinyltransferase 2 (<i>EPSPS</i>)	56.1 kb and 32.7 kb	3.12	4.01	0.00
KS_00461	106,901	109,241	2,340	Forward	NAD-dependent epimerase	No	0	2.52	0.00
KS_00462	106,975	110,332	3,357	Reverse	Uncharacterized protein	No	0	2.54	0.06
KS_00463	113,504	114,006	502	Reverse	DUF861	No	0	0.05	0.85

NOTE.—Read depth is the \log_2 of the difference between the background read depth and the read depth of each gene from genomic Illumina sequencing of a glyphosate-resistant line. Base-pair coordinates are given relative to their position in the contig from the susceptible genome assembly. DE is the \log_2 differential expression between four resistant and four susceptible individuals from RNA-Seq. *P* value is the significance of DE and is adjusted for false discovery rate.

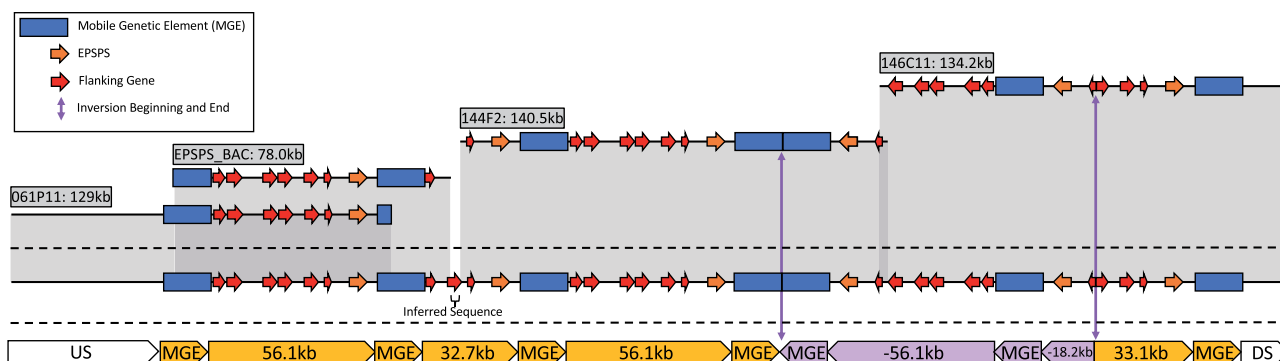


FIG. 2.—A diagram of the four assembled BACs and how they overlap to generate five different repeat types of the *EPSPS* CNV locus from glyphosate-resistant *Kochia*. The MGE is illustrated as a blue rectangle, the *EPSPS* gene is a green arrow, the coduplicated genes are orange arrows, and the beginning and end of the inverted repeat are vertical arrow lines.

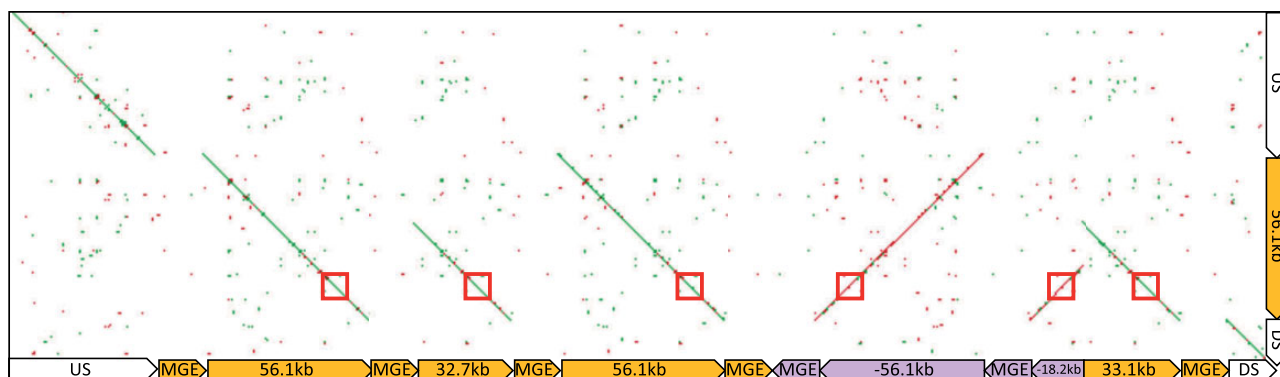


FIG. 3.—A dot-plot alignment of the assembled resistant *EPSPS* locus to the contig containing *EPSPS* from the susceptible genome assembly. The location of *EPSPS* is indicated by a red box. Large gaps in alignment are the insertion sites of the MGE.

Table 2

Copy Number Data from All qPCR Markers on Three Glyphosate-Susceptible (7710) and Five Glyphosate-Resistant (M32) Individuals

Line	Biological Replicate	1	2	3	4	5	6	7	8	9	10	11
7710	1	0.9	0.7	N/A	1.1	1.6	1.1	1.3	1.2	0.7	1.9	0.8
	2	0.7	0.7	N/A	1.0	1.5	1.2	1.4	1.4	0.9	1.7	1.2
	3	0.7	0.6	N/A	0.9	1.0	1.2	0.7	1.3	1.0	1.6	1.1
M32	1	0.9	0.7	9.5	6.1	11.3	11.2	11.3	11.5	1.0	N/A	1.0
	2	0.8	0.7	9.5	6.0	12.6	12.1	12.4	13.3	1.0	N/A	1.1
	3	0.7	0.6	7.6	3.2	10.9	11.1	11.0	11.7	1.0	N/A	1.0
	4	0.7	0.7	8.1	5.1	10.8	9.9	10.4	9.9	0.9	N/A	0.9
	5	1.2	1.0	14.2	10.0	20.3	19.0	19.6	20.0	1.3	N/A	1.4

NOTE.—Copy number is calculated as $\Delta C_t = (C_t^{(ALS)} + C_t^{(CPS)})/2 - C_t^{\text{Marker}}$. N/A, no amplification.**Table 3**

Copy Number Data for the Number of 56.1-kb Repeats, 32.7-kb Repeats, and the MGE on Three Glyphosate-Susceptible (7710) and Five Glyphosate-Resistant (M32) Individuals

Line	Replicate	56.1 kb	32.7 kb	MGE
7710	1	N/A	N/A	3.9
	2	N/A	N/A	5.5
	3	N/A	N/A	4.7
M32	1	5.4	1.8	16.2
	2	5.1	1.9	17.4
	3	5.1	1.7	18.2
	4	5.3	1.7	14.1
	5	6.9	2.1	17.7

NOTE.—Copy number is calculated as $\Delta C_t = (C_t^{(ALS)} + C_t^{(CPS)})/2 - C_t^{\text{Marker}}$. N/A, no amplification.

sequences were included. When this BAC meta-assembly from glyphosate-resistant kochia was aligned to the susceptible contig from the genome assembly, we observed perfect agreement between the resistant and susceptible loci; however, a large disparity was evident at each repeat junction and on either end of the resistant repeat structure (fig. 3). A 16,037-bp sequence was inserted just down- and upstream of all repeats in the glyphosate-resistant BAC assemblies. This insert shows no homology with any part of the susceptible contig; furthermore, when this insertion was aligned against the entire susceptible genome assembly, this region was not found in its entirety.

Maker was run on this insertion to predict gene models and identified four loci with putative coding genes. The first predicted gene belonged to the family of genes known as FHY3/FAR1 (IPRO31052) and contained the domains: “AR1 DNA binding” and “zinc finger, SWIM-type” (IPRO04330F, IPRO07527, respectively). The second gene’s function was less clear but was identified to be part of the Ubiquitin-like domain superfamily (IPRO29071). The third gene’s function was also unclear and was generally identified as belonging to the Endonuclease/exonuclease/phosphatase superfamily (IPRO36691). The fourth and final gene had no identifiable InterPro domains, and had BLAST hits to uncharacterized

proteins in NCBI. We refer to this insertion as the MGE in all figures and discussion as it seems to have inserted only in resistant lines from an unknown *trans* location in the genome. Of the four predicted genes, none had any expression in resistant or susceptible plants in the RNA-Seq data set. The MGE was assembled by aligning PacBio reads from the susceptible genome sequencing to the resistant MGE and generating a consensus (supplementary file “all_fasta.txt,” Supplementary Material online), confirming that it does exist in the kochia genome, but it was not assembled during whole genome assembly. The MGE sequences from resistance and susceptible were aligned as a dotplot using YASS and showed very high homology (supplementary fig. 3, Supplementary Material online).

Markers for Confirming the Structure of the *EPSPS* CNV

Quantitative PCR markers were developed dispersed across the entire CNV, including markers on both sides in regions that show no evidence of CNV (supplementary table 1, Supplementary Material online). These markers performed, for the most part, as predicted based on the resequencing of the glyphosate-resistant plants and the BAC sequencing. All markers up- and downstream of the CNV are approximately single copy. Markers 3 and 4, predicted to be only in the longer, 56.1-kb repeat, both show increased copy number in resistant individuals. Markers 5, 6, 7, and 8, are in both 56.1- and 32.7-kb repeats. These four markers were tightly associated, covaried for each individual, and showed higher copy number than markers 3 and 4 (table 2).

Additional qPCR markers were developed that only amplified when the MGE was flanked by either the two dominant repeat types of 56.1 or 32.7 kb. Using these markers, we quantified the number of 56.1- or 32.7-kb repeats in several individuals. In our line, 32.7-kb repeats were less frequent than 56.1-kb repeats. The tested individuals each had approximately two 32.7-kb repeats and between five and seven 56.1-kb repeats (table 3). These markers did not amplify in any susceptible plants, which support the discovery that the MGE is not present at the beginning of the susceptible *EPSPS* locus.

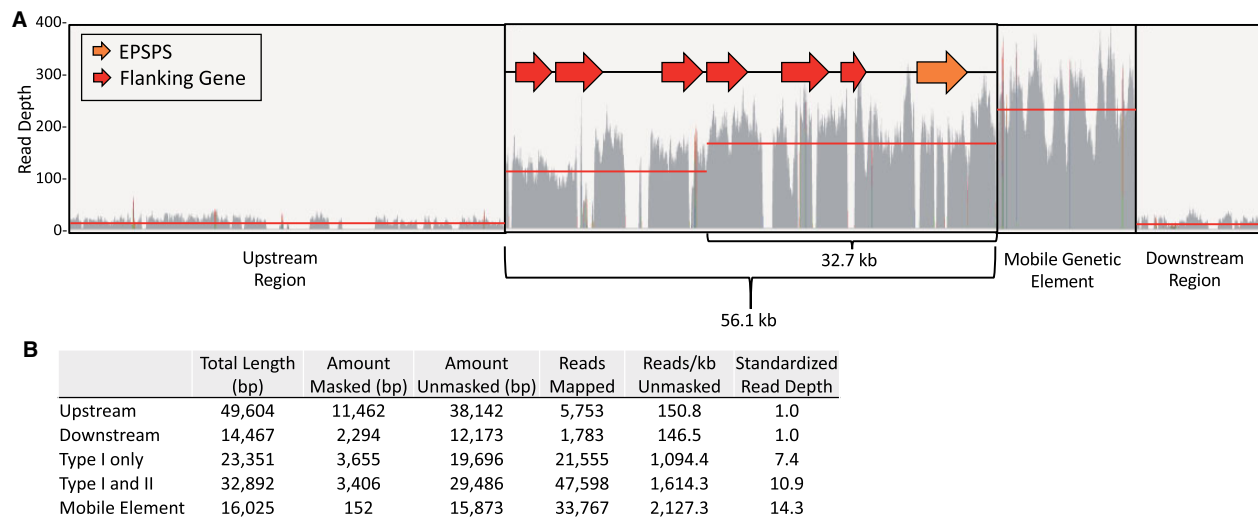


FIG. 4.—(A) Illumina shotgun genome resequencing data from a resistant kochia plant aligned to four distinct units from the BAC assembly: 1) The region directly upstream of the *EPSPS* tandem duplication, 2) the tandemly duplicated region of the genome containing *EPSPS*, 3) the MGE, and 4) the region directly downstream of the *EPSPS* tandem duplication. Red lines indicate the average read depth for that unit. Two averages are indicated for the tandemly duplicated region of the genome containing *EPSPS* due to two major repeat sites existing in the *EPSPS* CNV structure: the 56.1- and 32.7-kb repeat types. (B) A table outlining the calculation for copy number estimates for the four units. The total length of the region, the amount of repetitive DNA that was masked, the amount of DNA remaining unmasking, the number of reads mapped to the unmasked regions, the average reads per kilobase of unmasked DNA, and the read depth divided by the reads/kb unmasked of the nonduplicated region.

Additionally, we developed a marker internal to the MGE. All susceptible individuals had ~4–5 copies of this marker. The MGE assembly from the susceptible genome showed that the MGE primer sites were present in the susceptible sequence and identical to the MGE sequence at those positions in the resistant MGE. In resistant individuals, we detected 14–18 copies of the MGE. If we account for the 4–5 copies that are in the susceptible individuals and if we consider that an MGE exists at both the up- and downstream boundary, then we would predict 9–13 copies, which almost perfectly correlates with the copy number observed for qPCR markers 5, 6, 7, and 8. This would indicate that one copy of the MGE is associated with each repeat (table 3).

Illumina shotgun genome resequencing data from a resistant kochia plant aligned to four distinct units from the BAC assembly was used to calculate the copy number of each unit of the repeat structure and to confirm our qPCR results. After standardizing the read depth of each unit by the background read depth, we calculated 7.4 copies of the 56.1-kb repeat, 10.9 copies of the 32.7-kb repeat type, and 14.3 copies of the MGE (fig. 4A and B). It should be noted that the unit of the 32.7-kb repeat type includes reads from all repeats due to the sequence of this region being shared in all repeat types. With this information in conjunction with previously published cytogenetic work (Jugulam et al. 2014; Jugulam and Gill 2018), we propose a model for the structure of the *EPSPS* CNV from resistant kochia individuals (fig. 5).

Discussion

Structure and Genetic Content of the *EPSPS* Tandem Duplication Region

Glyphosate resistance due to *EPSPS* gene duplication has independently evolved in multiple species within the Caryophyllales through very different genomic mechanisms, specifically tandem duplication in kochia and the proliferation of an extrachromosomal circular DNA containing *EPSPS* in *A. palmeri* (Jugulam et al. 2014; Molin et al. 2017; Koo et al. 2018; Patterson et al. 2018). We have determined the length and content of the repeat units in the tandem duplication found in one line of resistant kochia. Additionally, we discovered an MGE inserted between each repeat. Quantitative PCR shows that the most common repeats are 72.6 or 49.2 kb in length. These estimates are similar to, but slightly larger than, the previously described Fiber-FISH estimated sizes of 66 and 45 kb in another resistant kochia line (Jugulam et al. 2014). What accounts for the differences between our assemblies and the previously reported Fiber-FISH studies remains unclear, as Fiber-FISH can have a resolution of ~1 kb (Ersfeld 1994). It may be that different populations of kochia have different repeat sizes. Further testing and validation on the type and size of the *EPSPS* repeats in various, divergent populations is needed to confirm this. We did detect an inverted repeat near the downstream end of the CNV as shown by Jugulam et al. (2014).

RNA-Seq expression data show that four of the six genes within the conserved region of the tandem-repeat are

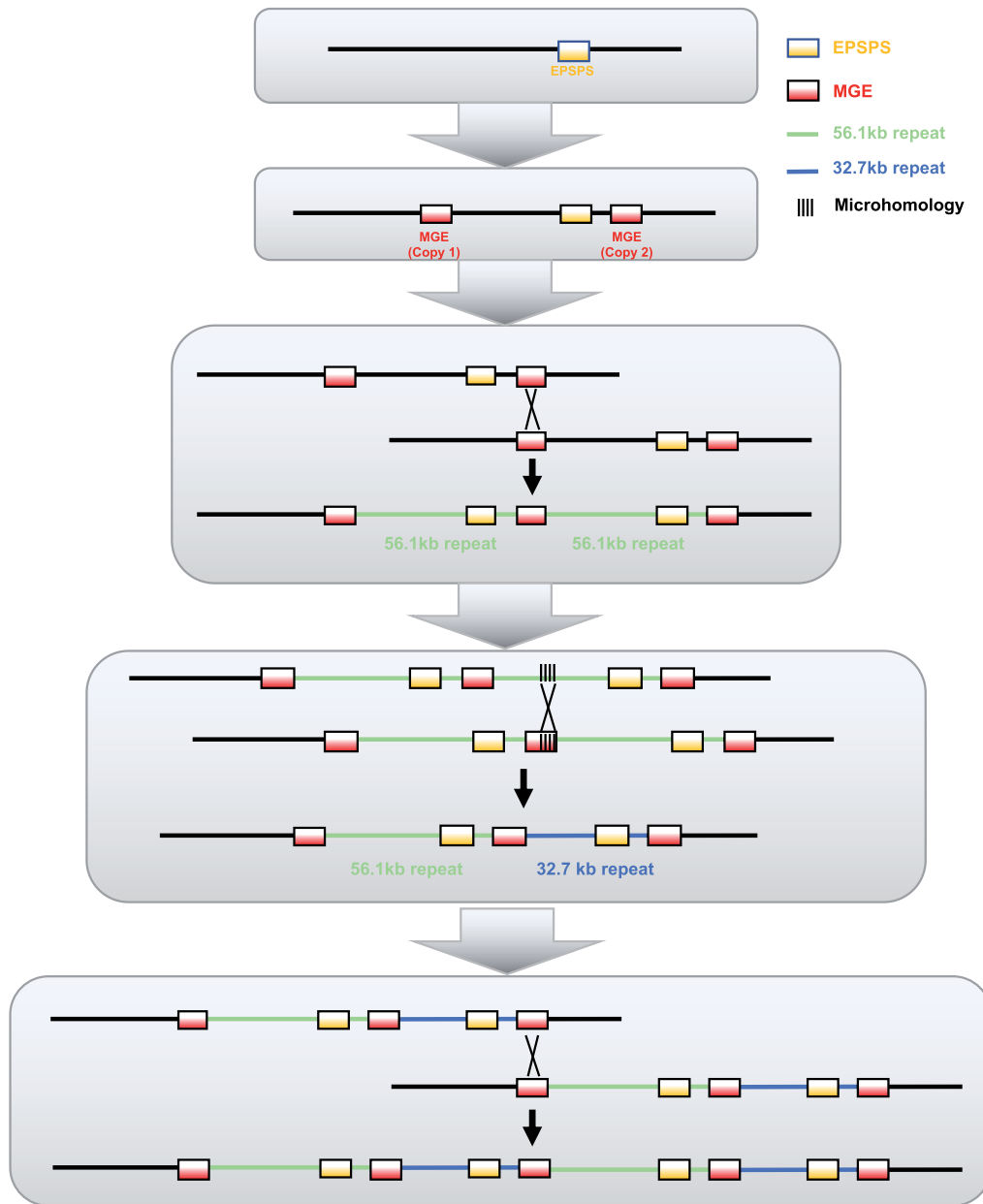


FIG. 5.—A model for the generation and continued increase of *EPSPS* copy number. The initial event that led to *EPSPS* gene duplication was the insertion of two mobile elements both up- and downstream of the *EPSPS* gene (MGE). After unequal crossing over, gametes were produced with >1 *EPSPS* gene copy. Subsequently, a double stranded break occurred at the MGE boundary that was incorrectly repaired using a microhomology-mediated mechanism within the middle of the repeat region, generating a shorter copy of this repeat region (32.7-kb repeat).

overexpressed at a rate commensurate with genomic resequencing read depth: *RAD51*, *transketolase*, *tRNA N6-adenosine threonylcarbamoyltransferase*, and *EPSPS* (FDR adjusted P value <0.05). The obvious benefit of *EPSPS* overexpression is glyphosate resistance, but the phenotypic effects due to increased expression of other genes in this CNV remain unclear.

The expression of the *RAD51* homolog is especially interesting due to its importance in regulating crossing over.

Misexpression, up or down, of *RAD51* has been shown to cause cancer in animal tissues as *RAD51* is involved in regulating homologous recombination of DNA during double stranded break repair (Maacke et al. 2000) (table 1). Additionally, *RAD51*, along with the recombinase DMC1, facilitate recombination of homologous chromosomes during meiosis in plants and animals (Crickard et al. 2018). In humans, *RAD51* expression is modulated by miRNAs and misregulation of these miRNAs are often associated with various

forms of cancer (Choi et al. 2014; Gasparini et al. 2014; Cortez et al. 2015; Liu, Xue, et al. 2015; Liu, Yang, et al. 2015). Therefore, we would predict that overexpression of *RAD51* in the resistant line would have a large impact phenotypic consequence and could change the recombination rates and double strand break repair.

The two other genes that are coduplicated and coexpressed with *EPSPS* and *RAD51* are annotated as *transketolase* and *tRNA N6-adenosine threonylcarbamoyltransferase*. Transketolase is an enzyme found in all organisms. Transketolase has two functions in plants, one in the pentose phosphate pathway (Horecker 2002) and a second in the Calvin cycle of photosynthesis (Flechner et al. 1996). Slight reductions in transketolase expression can have significant effects on the photosynthetic ability of tobacco plants (Henkes et al. 2001); however, overexpression of transketolase has been shown to have a mixture of phenotypes depending on species. Overexpression in rice does not result in an increase in an increase in CO₂ capture (Suzuki et al. 2017); however, in cucumbers, overexpression results in increased photosynthetic rate in transgenic cucumber leaves (Bi et al. 2013). It is possible that the duplication and overexpression of transketolase in kochia may have a significant phenotypic effect.

The tRNA N6-adenosine threonylcarbamoyltransferase protein is critical in the formation of the threonylcarbamoyl group on the adenosine at position 37 of tRNAs that read AXX codons (Jackman and Alfonzo 2013). The impact that the overexpression of this gene will have on phenotype is unclear; however, its function is universally conserved in all three kingdoms of life, with the specific gene *tsaD* being necessary in prokaryotes (Thiaville et al. 2015) but not yeast. The phenotypic impacts of *transketolase* and *tRNA N6-adenosine threonylcarbamoyltransferase* overexpression need to be investigated separately to be fully understood.

The expression of the two other genes in the *EPSPS* CNV (*golgin subfamily A member 6-like protein 6* and *NRT1/PTR Family 7.2-like*) is reduced in the resistant line. This reduction may be due to gene silencing, similar to what happens when multiple copies of transgenes are inserted in the same plant (Finnegan and McElroy 1994; Tang et al. 2006) (table 1).

We used qPCR genomic copy number primers to validate much of our BAC assembly. The results from a pair of primers that detected the presence and number of the MGE were surprising. In the susceptible plant, approximately 4–6 MGE copies were observed and the MGE was assembled separately from susceptible PacBio reads; therefore, this MGE is present in the susceptible plant but it was not assembled in the whole genome assembly. It may be that these background copies lie in repetitive or difficult to assemble regions. In the resistant plants, the number of MGE copies was always approximately equal to the *EPSPS* copy number plus 4–6 copies, indicating that the original copies found elsewhere in the genome are still present and the insert is being coduplicated with every

repeat of the *EPSPS* CNV. The fact that the MGE is present in the susceptible lineage implies that the insertion in the *EPSPS* region originated by transposition within the genome.

The Role of an MGE in *EPSPS* Gene Duplication

When the *EPSPS* contig from the susceptible genome assembly is aligned to itself, no complexities, such as SSRs or large homodimers of nucleotides, exist at the beginnings of any of the repeat types (supplementary fig. 1, Supplementary Material online). This would indicate that the sequence in the susceptible locus alone is insufficient for explaining why this region has become a site for copy number variation, which is inconsistent with earlier predictions that homology exists at the up- and downstream boundaries where an initial misalignment occurred (Jugulam et al. 2014); however MGEs, such as transposons, have been proposed to cause tandem repeats of sequences near their insertion point (Tsubota et al. 1989; Reams and Roth 2015).

We propose that the insertion of an MGE near the *EPSPS* locus in the resistant kochia line facilitated the subsequent history of tandem duplication in this region. The MGE contains a member of the *Fhy3/FAR1* gene family. Genes in this family are thought to be derived from MULE transposons and have been “domesticated” to have a role in the regulation of genes involved in circadian rhythm and light sensing in a wide phylogenetic distribution of angiosperms (Wang and Deng 2002; Hudson et al. 2003; Cowan et al. 2005; Tang et al. 2012). We hypothesize the insertion of the MGE near the *EPSPS* locus in resistant kochia line is evidence that *Fhy3/FAR1* elements may still be mobile and that they are not fully “domesticated.”

The MGE appears at both the up- and downstream borders of the CNV, therefore this MGE inserted at two locations flanking the *EPSPS* region. The insertion of two identical MGEs in close proximity could then facilitate misalignment and subsequent crossing-over events that would generate two alleles—one with two of the more common 56.1-kb repeats, and the other with no *EPSPS* gene, the latter of which would be lethal in the homozygous state. Such unequal crossing over could then facilitate further expansions of this region.

Interestingly, the beginning of the MGE shares a 7-bp stretch of perfect sequence identity with the exact beginning of the shorter, less common 32.7-kb repeat. We propose that a second recombination event took place between the MGE downstream boundary and the start site of the smaller 32.7-kb repeat, perhaps mediated by double-stranded break repair at the end of the MGE (fig. 5) (Ottaviani et al. 2014; Sfeir and Symington 2015). Short microhomology-mediated illegitimate recombination has been well studied in bacteria (Petes and Hill 1988; Nash 1996; Romero and Palacios 1997; de Vries and Wackernagel 2002; Reams and Neidle 2004).

The presence of an MGE end at the breakpoint of the large inversion in the tandem array (fig. 2) further implicates

double-stranded breaks at the MGE boundaries in this region. Homologous recombination and double strand break repair depend heavily on the enzyme RecA in bacteria and its homologue RAD51 in eukaryotes. These enzymes bind single-stranded DNA and promote strand invasion and therefore the exchange between homologous DNA molecules (Baumann and West 1998; Lin et al. 2006; Hastings et al. 2009). In *Kochia*, it remains unclear if the presence of *RAD51* in the duplicated region is coincidental or has affected the evolution of this tandem duplication event.

Conclusion

Widespread and repeated use of the herbicide glyphosate represents an intense abiotic selective pressure across large areas. Several weed species have evolved resistance to this pressure by means of increased copies of the target-site gene *EPSPS*. We identified an MGE at the duplicated *EPSPS* locus and hypothesize that the insertion of two or more of these MGEs initiated a tandem duplication event. Once the initial gene duplication occurred, the locus underwent several rounds of unequal recombination producing gametes with increased and decreased copy numbers. This interplay between transposable elements and target site copy number variation provides valuable insight into how genomic plasticity may contribute to rapid evolution of abiotic stress tolerance. Continuing to investigate the roles transposable elements and gene duplication play in shaping plant resilience is essential for understanding evolution and how plant genomes are changing in response to human activities.

Supplementary Material

Supplementary data are available at *Genome Biology and Evolution* online.

Acknowledgments

This work was partially supported by the [Colorado Wheat Administrative Committee](#), [Dow AgroSciences](#), and by the [USDA National Institute of Food and Agriculture, Hatch Project COL00783](#), accession number 1016207, to the Colorado State University Agricultural Experiment Station.

Literature Cited

- Baumann P, West SC. 1998. Role of the human RAD51 protein in homologous recombination and double-stranded-break repair. *Trends Biochem Sci.* 23(7):247–251.
- Beckie HJ, Blackshaw RE, Low R, Hall LM, Sauder CA. 2013. Glyphosate- and acetolactate synthase inhibitor-resistant *Kochia* (*Kochia scoparia*) in Western Canada. *Weed Sci.* 61(2):310–318.
- Beckie HJ, et al. 2015. Glyphosate-resistant *Kochia* (*Kochia scoparia* L. Schrad.) in Saskatchewan and Manitoba. *Can J Plant Sci.* 95(2):345–349.
- Beckie HJ, et al. 2018. Seed bank persistence, germination and early growth of glyphosate-resistant *Kochia scoparia*. *Weed Res.* 58(3):177–187.
- Bi H, Wang M, Dong X, Ai X. 2013. Cloning and expression analysis of transketolase gene in *Cucumis sativus* L. *Plant Physiol Biochem.* 70:512–521.
- Butler J, et al. 2008. ALLPATHS: de novo assembly of whole-genome shotgun microreads. *Genome Res.* 18(5):810–820.
- Camacho C, et al. 2009. BLAST+: architecture and applications. *BMC Bioinformatics* 10:421.
- Cantarel BL, et al. 2008. MAKER: an easy-to-use annotation pipeline designed for emerging model organism genomes. *Genome Res.* 18(1):188–196.
- Choi YE, et al. 2014. MicroRNAs down-regulate homologous recombination in the G1 phase of cycling cells to maintain genomic stability. *Elife* 3:e02445.
- Cortez MA, et al. 2015. In vivo delivery of miR-34a sensitizes lung tumors to radiation through RAD51 regulation. *Mol Ther Nucleic Acids.* 4:e270.
- Cowan RK, Hoen DR, Schoen DJ, Bureau TE. 2005. MUSTANG is a novel family of domesticated transposase genes found in diverse angiosperms. *Mol Biol Evol.* 22(10):2084–2089.
- Crickard JB, Kaniecki K, Kwon Y, Sung P, Greene EC. 2018. Spontaneous self-segregation of Rad51 and Dmc1 DNA recombinases within mixed recombinase filaments. *J Biol Chem.* 293(11):4191–4200.
- de Vries J, Wackernagel W. 2002. Integration of foreign DNA during natural transformation of *Acinetobacter* sp. by homology-facilitated illegitimate recombination. *Proc Natl Acad Sci U S A.* 99(4):2094–2099.
- DeBolt S. 2010. Copy number variation shapes genome diversity in *Arabidopsis* over immediate family generational scales. *Genome Biol Evol.* 2:441–453.
- Dohm JC, et al. 2014. The genome of the recently domesticated crop plant sugar beet (*Beta vulgaris*). *Nature* 505(7484):546–549.
- English AC, et al. 2012. Mind the gap: upgrading genomes with Pacific Biosciences RS long-read sequencing technology. *PLoS One* 7(11):e47768.
- Ersfeld K. 1994. Fiber-FISH: fluorescence in situ hybridization on stretched DNA. In: Melville SE, editor. *Parasite genomics protocols*. Totowa, New Jersey: Humana Press. p. 395–402.
- Finnegan J, McElroy D. 1994. Transgene inactivation: plants fight back! *Nat Biotechnol.* 12(9):883.
- Flechner A, et al. 1996. Molecular characterization of transketolase (EC 2.2. 1.1) active in the Calvin cycle of spinach chloroplasts. *Plant Mol Biol.* 32(3):475–484.
- Gaines TA, et al. 2010. Gene amplification confers glyphosate resistance in *Amaranthus palmeri*. *Proc Natl Acad Sci U S A.* 107(3):1029–1034.
- Gaines TA, et al. 2016. *EPSPS* gene copy number and whole-plant glyphosate resistance level in *Kochia scoparia*. *PLoS One* 11(12):e0168295.
- Gasparini P, et al. 2014. Protective role of miR-155 in breast cancer through RAD51 targeting impairs homologous recombination after irradiation. *Proc Natl Acad Sci U S A.* 111(12):4536–4541.
- Godar AS, Stahlman PW, Jugulam M, Dille JA. 2015. Glyphosate-resistant *Kochia* (*Kochia scoparia*) in Kansas: ePSPS gene copy number in relation to resistance levels. *Weed Sci.* 63(3):587–595.
- Hackl T, Hedrich R, Schultz J, Förster F. 2014. proofread: large-scale high-accuracy PacBio correction through iterative short read consensus. *Bioinformatics* 30(21):3004–3011.
- Hastings PJ, Lupski JR, Rosenberg SM, Ira G. 2009. Mechanisms of change in gene copy number. *Nat Rev Genet.* 10(8):551.
- Henkes S, Sonnwald U, Badur R, Flachmann R, Stitt M. 2001. A small decrease of plastid transketolase activity in antisense tobacco transformants has dramatic effects on photosynthesis and phenylpropanoid metabolism. *Plant Cell* 13(3):535–551.
- Horecker BL. 2002. The pentose phosphate pathway. *J Biol Chem.* 277(50):47965–47971.

- Hudson ME, Lisch DR, Quail PH. 2003. The FHY3 and FAR1 genes encode transposase-related proteins involved in regulation of gene expression by the phytochrome A-signaling pathway. *Plant J*. 34(4):453–471.
- Hull RM, Cruz C, Jack CV, Houseley J. 2017. Environmental change drives accelerated adaptation through stimulated copy number variation. *PLoS Biol*. 15(6):e2001333.
- Jackman JE, Alfonzo JD. 2013. Transfer RNA modifications: nature's combinatorial chemistry playground. *Wires RNA*. 4(1):35–48.
- Jones P, et al. 2014. InterProScan 5: genome-scale protein function classification. *Bioinformatics* 30(9):1236–1240.
- Jugulam M, et al. 2014. Tandem amplification of a chromosomal segment harboring EPSPS locus confers glyphosate resistance in *Kochia scoparia*. *Plant Physiol*. 166(3):1200–1207.
- Jugulam M, Gill BS. 2018. Molecular cytogenetics to characterize mechanisms of gene duplication in pesticide resistance. *Pest Manag Sci*. 74(1):22–29.
- Koo D-H, et al. 2018. Extrachromosomal circular DNA-based amplification and transmission of herbicide resistance in crop weed *Amaranthus palmeri*. *Proc Natl Acad Sci U S A*. 115(13):3332–3337.
- Koren S, et al. 2017. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res*. 27(5):722–736.
- Kumar V, Jha P, Giacomini D, Westra EP, Westra P. 2015. Molecular basis of evolved resistance to glyphosate and acetolactate synthase-inhibitor herbicides in kochia (*Kochia scoparia*) accessions from Montana. *Weed Sci*. 63(4):758–769.
- Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* 25(14):1754–1760.
- Li H, et al. 2009. The sequence alignment/map format and SAMtools. *Bioinformatics* 25(16):2078–2079.
- Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34(18):3094–3100.
- Liao Y, Smyth GK, Shi W. 2014. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* 30(7):923–930.
- Lin Z, Kong H, Nei M, Ma H. 2006. Origins and evolution of the recA/RAD51 gene family: evidence for ancient gene duplication and endosymbiotic gene transfer. *Proc Natl Acad Sci U S A*. 103(27):10328–10333.
- Liu G, Xue F, Zhang W. 2015. miR-506: a regulator of chemo-sensitivity through suppression of the RAD51-homologous recombination axis. *Chin J Cancer*. 34(3):44.
- Liu G, Yang D, et al. 2015. Augmentation of response to chemotherapy by microRNA-506 through regulation of RAD51 in serous ovarian cancers. *J Natl Cancer Inst*. 107:djv108.
- Luo M, Wing RA. 2003. An improved method for plant BAC library construction. In: Grotewald E, editor. *Plant functional genomics, methods in molecular biology*. Vol. 236. Totowa (NJ): Humana Press. p. 3–19.
- Lynch M, Conery JS. 2000. The evolutionary fate and consequences of duplicate genes. *Science* 290(5494):1151–1155.
- Maacke H, et al. 2000. Over-expression of wild-type Rad51 correlates with histological grading of invasive ductal breast cancer. *Int J Cancer*. 88(6):907–913.
- Martin SL, et al. 2017. Glyphosate resistance reduces kochia fitness: comparison of segregating resistant and susceptible F2 populations. *Plant Sci*. 261:69–79.
- Mi H, Lazareva-Ulitsky B, Loo R, Kejariwal A, Vandergriff J. 2004. The PANTHER database of protein families, subfamilies, functions and pathways. *Nucleic Acids Res*. 33(Database issue):D284–D288.
- Molin WT, Wright AA, Lawton-Rauh A, Saski CA. 2017. The unique genomic landscape surrounding the EPSPS gene in glyphosate resistant *Amaranthus palmeri*: a repetitive path to resistance. *BMC Genet*. 18:91.
- Nash HA. 1996. Site-specific recombination: integration, excision, resolution, and inversion of defined DNA segments. In: Neidhardt FC, editor. *Escherichia coli and Salmonella: cellular and molecular biology*. Washington (DC): ASM Press. p. 2363–2376.
- Noé L, Kucherov G. 2005. YASS: enhancing the sensitivity of DNA similarity search. *Nucleic Acids Res*. 33(Web Server issue):W540–W543.
- Ottaviani D, LeCain M, Sheer D. 2014. The role of microhomology in genomic structural variation. *Trends Genet*. 30(3):85–94.
- Patterson EL, Pettinga DJ, Ravet K, Neve P, Gaines TA. 2018. Glyphosate resistance and EPSPS gene duplication: convergent evolution in multiple plant species. *J Hered*. 109(2):117–125.
- Petes TD, Hill CW. 1988. Recombination between repeated genes in microorganisms. *Annu Rev Genet*. 22:147–168.
- Pettinga DJ, et al. 2018. Increased Chalcone Synthase (CHS) expression is associated with dicamba resistance in *Kochia scoparia*. *Pest Manag Sci*. 74(10):2306–2315.
- Preston C, Belles DS, Westra PH, Nissen SJ, Ward SM. 2009. Inheritance of resistance to the auxinic herbicide dicamba in kochia (*Kochia scoparia*). *Weed Sci*. 57(1):43–47.
- Reams AB, Neidle EL. 2004. Gene amplification involves site-specific short homology-independent illegitimate recombination in *Acinetobacter* sp. strain ADP1. *J Mol Biol*. 338(4):643–656.
- Reams AB, Roth JR. 2015. Mechanisms of gene duplication and amplification. *Cold Spring Harb Perspect Biol*. 7(2):a016592.
- Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 26(1):139–140.
- Romero D, Palacios R. 1997. Gene amplification and genomic plasticity in prokaryotes. *Annu Rev Genet*. 31:91–111.
- Sammons DR, Gaines TA. 2014. Glyphosate resistance: state of knowledge. *Pest Manag Sci*. 70(9):1367–1377.
- Schimke R, Hill A, Johnston R. 1985. Methotrexate resistance and gene amplification: an experimental model for the generation of cellular. *Br J Cancer*. 51(4):459–465.
- Schmittgen TD, Livak KJ. 2008. Analyzing real-time PCR data by the comparative C_T method. *Nat Protoc*. 3(6):1101–1108.
- Sfeir A, Symington LS. 2015. Microhomology-mediated end joining: a back-up survival mechanism or dedicated pathway? *Trends Biochem Sci*. 40(11):701–714.
- Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31(19):3210–3212.
- Soto-Jimenez L, Estrada K, Sanchez-Flores A. 2014. GARM: genome assembly, reconciliation and merging pipeline. *Curr Top Med Chem*. 14(3):418–424.
- Suzuki Y, Kondo E, Makino A. 2017. Effects of co-overexpression of the genes of Rubisco and transketolase on photosynthesis in rice. *Photosyn Res*. 131(3):281–289.
- Tang W, et al. 2012. Transposase-derived proteins FHY3/FAR1 interact with PHYTOCHROME-INTERACTING FACTOR1 to regulate chlorophyll biosynthesis by modulating HEMB1 during deetiolation in *Arabidopsis*. *Plant Cell* 24(5):1984–2000.
- Tang W, Newton RJ, Weidner DA. 2006. Genetic transformation and gene silencing mediated by multiple copies of a transgene in eastern white pine. *J Exp Bot*. 58(3):545–554.
- Tarailo-Graovac M, Chen N. 2009. Using RepeatMasker to identify repetitive elements in genomic sequences. *Curr Prot Bioinformatics*. 25:4.10.1–4.10.14.
- Thiaville PC, et al. 2015. Essentiality of threonylcarbamoyladenosine (t6 A), a universal t RNA modification, in bacteria. *Mol Microbiol*. 98(6):1199–1221.
- Thrasher A, Musgrave Z, Kachmarck B, Thain D, Emrich S. 2014. Scaling up genome annotation using MAKER and work queue. *Int J Bioinform Res Appl*. 10(4/5):447–460.

- Tsubota SI, Rosenberg D, Szostak H, Rubin D, Schedl P. 1989. The cloning of the Bar region and the B breakpoint in *Drosophila melanogaster*: evidence for a transposon-induced rearrangement. *Genetics* 122(4):881–890.
- Vila-Aiub MM, et al. 2014. No fitness cost of glyphosate resistance endowed by massive *EPSPS* gene amplification in *Amaranthus palmeri*. *Planta* 239(4):793–801.
- Wang H, Deng XW. 2002. Arabidopsis FHY3 defines a key phytochrome A signaling component directly interacting with its homologous partner FAR1. *EMBO J.* 21(6):1339–1349.
- Wiersma AT, et al. 2015. Gene amplification of 5-enol-pyruvylshikimate-3-phosphate synthase in glyphosate-resistant *Kochia scoparia*. *Planta* 241(2):463–474.
- Xi R, et al. 2011. Copy number variation detection in whole-genome sequencing data using the Bayesian information criterion. *Proc Natl Acad Sci U S A.* 108(46):E1128–E1136.

Associate editor: Brandon Gaut