Research Article



Gene capture by transposable elements leads to epigenetic conflict in maize

Aline Muyle¹, Danelle Seymour^{1,2}, Nikos Darzentas³, Elias Primetis⁴, Brandon S. Gaut^{1,*} and Alexandros Bousios^{4,*}

https://doi.org/10.1016/j.molp.2020.11.003

ABSTRACT

Transposable elements (TEs) regularly capture fragments of genes. When the host silences these TEs, siRNAs homologous to the captured regions may also target the genes. This epigenetic crosstalk establishes an intragenomic conflict: silencing the TEs has the cost of silencing the genes. If genes are important, however, natural selection may maintain function by moderating the silencing response, which may also advantage the TEs. In this study, we examined this model by focusing on Helitrons, Pack-MULEs, and Sirevirus LTR retrotransposons in the maize genome. We documented 1263 TEs containing exon fragments from 1629 donor genes. Consistent with epigenetic conflict, donor genes mapped more siRNAs and were more methylated than genes with no evidence of capture. However, these patterns differed between syntelog versus translocated donor genes. Syntelogs appeared to maintain function, as measured by gene expression, consistent with moderation of silencing for functionally important genes. Epigenetic marks did not spread beyond their captured regions and 24nt crosstalk siRNAs were linked with CHH methylation. Translocated genes, in contrast, bore the signature of silencing. They were highly methylated and less expressed, but also overrepresented among donor genes and located away from chromosomal arms, which suggests a link between capture and gene movement. Splitting genes into potential functional categories based on evolutionary constraint supported the synteny-based findings. TE families captured genes in different ways, but the evidence for their advantage was generally less obvious; nevertheless, TEs with captured fragments were older, mapped fewer siRNAs, and were slightly less methylated than TEs without captured fragments. Collectively, our results argue that TE capture triggers an intragenomic conflict that may not affect the function of important genes but may lead to the pseudogenization of less-constrained

Key words: transposable elements, intragenomic conflict, gene capture, epigenetic silencing, methylation, synteny

Muyle A., Seymour D., Darzentas N., Primetis E., Gaut B.S., and Bousios A. (2021). Gene capture by transposable elements leads to epigenetic conflict in maize. Mol. Plant. 14, 237-252.

INTRODUCTION

Transposable elements (TEs) constitute the majority of plant genomes and are major drivers of both genomic and phenotypic evolution (Lisch, 2013). TEs are generally not active under normal conditions. Based mostly on studies in Arabidopsis thaliana, it is known that this inactivation is achieved by host epigenetic silencing mechanisms that suppress TE activity both before and after transcription (Matzke and Mosher, 2014; Cuerda-Gil and Slotkin, 2016). These mechanisms rely on small interfering RNAs (siRNAs) that guide RNAi and RNA-directed DNA methylation (RdDM) against homologous sequences at the RNA and DNA level, respectively. RdDM is a feedback loop that initiates and spreads cytosine methylation in the CG, CHG, and CHH contexts (H = A, C, or T) of TE sequences. Symmetric CG and CHG methylation can then be maintained through cell division independently of RdDM, but asymmetric CHH methylation

Published by the Molecular Plant Shanghai Editorial Office in association with Cell Press, an imprint of Elsevier Inc., on behalf of CSPB and CEMPS, CAS.

¹Department of Ecology and Evolutionary Biology, UC Irvine, Irvine, CA 92697, USA

²Department of Botany and Plant Sciences, UC Riverside, Riverside, CA 92521, USA

³Central European Institute of Technology, Masaryk University, Brno, Czech Republic

⁴School of Life Sciences, University of Sussex, Brighton, UK

^{*}Correspondence: Brandon S. Gaut (bgaut@uci.edu), Alexandros Bousios (alexandros.bousios@gmail.com)

requires continuous *de novo* siRNA targeting (Matzke and Mosher, 2014; Cuerda-Gil and Slotkin, 2016). As a result, silenced TEs are usually heavily methylated only in the CG and CHG contexts, with CHH methylation typically at much lower levels. Once methylated, silenced TEs are often associated with a closed heterochromatic state (Sigman and Slotkin, 2016), which can influence the function and expression of genes, especially when TEs and genes reside in close proximity. For example, methylated and siRNA-targeted TEs can be associated with altered expression of neighboring genes (Hollister et al., 2011; Maumus and Quesneville, 2014; Lee and Karpen, 2017) and, as a result, may be subject to stronger purifying selection compared with unsilenced TEs or TEs far from genes (Hollister and Gaut, 2009; Lee and Karpen, 2017).

In contrast to the epigenetic effects of TEs near genes, much less is known about epigenetic interactions between TEs and genes over long distances, particularly through the trans-activity of siRNAs (Cho, 2018). For siRNAs to mediate long-distance interactions, there must be sequence similarity between genes and TEs, so that siRNAs are homologous to both. The requirement of sequence similarity can be satisfied by varied evolutionary scenarios, such as the exaptation of portions of TEs into coding genes (Lockton and Gaut, 2009), but it is especially relevant in the phenomenon of gene capture by TEs. Gene capture has been investigated widely in both animals and plants (Thomas and Pritham, 2015; Zhao et al., 2018). Within plant genomes, capture has been best characterized for Helitrons and Pack-MULE DNA transposons, which together have captured thousands of gene fragments (Yang and Bennetzen, 2009; Zhao et al., 2018). Capture is common enough that a single TE often contains fragments of multiple host genes from unlinked genomic locations (Jiang et al., 2004; Thomas and Pritham, 2015). Although it is clear that gene capture is common, the mechanisms remain uncertain. However, several mechanisms have been proposed (Grabundzija et al., 2016; Catoni et al., 2019), and evidence suggests that capture can occur through both DNA- and RNAmediated processes (Jiang et al., 2004; Morgante et al., 2005).

The evolutionary consequences of gene capture are also not well characterized. One potential consequence is that the shuffling and rejoining of coding information within a TE leads to the emergence of a novel gene (Thomas and Pritham, 2015). Although this conjecture has been disputed (Juretic et al., 2005), a substantial proportion of TE-captured gene sequences are expressed (Wang et al., 2016; Zhao et al., 2018), a subset of those are translated (Hanada et al., 2009; Zhao et al., 2018), and a few exhibit signatures of selective constraint (Juretic et al., 2005; Hanada et al., 2009; Yang and Bennetzen, 2009). Another distinct possibility is that gene capture is a neutral mutational process that has few downstream evolutionary ramifications. Finally, gene capture may establish evolutionary conflicts between TEs and genes. Lisch (2009) has argued that gene capture is in a TE's evolutionary interest, because it blurs the line between host and TE "by combining both transposon and host sequences ... to increase the cost of efficiently silencing those transposons" (Lisch, 2009). This argument suggests a model of genomic conflict in which a TE captures a fragment from a gene, and the host mounts an siRNA-mediated response against the TE. Because the siRNAs from the captured fragment within

the TE can also target the captured region of the "donor" gene (i.e., the gene from which the fragment has been captured), the host response to the TE can simultaneously act in *trans* against the donor gene.

Under this scenario, transcriptional silencing of the TE may have collateral effects on the donor gene, including targeting by siRNAs that lead to DNA methylation and subsequent silencing (Figure 1A). If the donor gene has an important function, then natural selection is likely to either remove affected individuals from the population or limit potential silencing effects on the gene. The latter creates an intragenomic conflict, whereby the advantage of silencing the TE is balanced by potential damage to donor gene function. Conversely, selection to moderate the host response potentially advantages the TE with the captured gene fragment. Notably, this conflict model makes testable predictions that: (1) donor genes bear the signature of trans-epigenetic effects, including increased siRNA targeting and consequent methylation, (2) selection may limit these trans-epigenetic effects for important versus less functionally important genes, and (3) TEs benefit from capture via decreased host response.

The possibility of epigenetic links between TEs and donor genes has been discussed previously (Thomas and Pritham, 2015), but to our knowledge only one study has examined how often siRNAs map to both donor genes and to their captured fragments (Hanada et al., 2009). This study focused on Pack-MULEs in rice (Oryza sativa) and found siRNAs that map to both TEs and donor genes, thus providing the potential for siRNA "crosstalk" between donor genes and captured gene fragments. The study also found that genes with crosstalk are less expressed compared with genes without any mapped siRNAs. Two recent studies of rice Pack-MULEs extended this line of enquiry by investigating whether donor genes are methylated (Wang et al., 2016; Zhao et al., 2018), which could be indicative of epigenetic effects consistent with the conflict model. They found, however, that donor genes have low methylation levels that do not differ substantially from genes with no apparent history of capture by TEs (hereafter termed free genes). These studies provide some, but limited, evidence for epigenetic conflict.

The study of Pack-MULEs in rice suffers from two potential shortcomings with respect to investigating epigenetic interactions. The first is Pack-MULEs themselves. They commonly capture genes and therefore provide a rich template for study, but often have lower methylation levels than other TE families (Zhao et al., 2018; Stitzer et al., 2019), possibly because they preferentially insert near the 5' termini of genes (Jiang et al., 2011). This tendency may lessen the potential for intragenomic conflict with their donor genes. The second shortcoming is the small genome size of rice. Large genomes differ from small genomes in their TE content and also their genic methylation patterns. For example, Takuno et al. (2016) showed that only 6% of genes in rice (490 Mb) and 2% of genes in A. thaliana (156 Mb) have high levels of methylation (≥90% of methylated cytosines) in the CG context compared with 24% of genes in the much larger (2646 Mb) genome of maize (Takuno et al., 2016). Genic methylation in the CG context does not suppress expression, but methylation in the CHG context likely does; 12%, 1%, and <1% of maize, rice and A. thaliana genes,

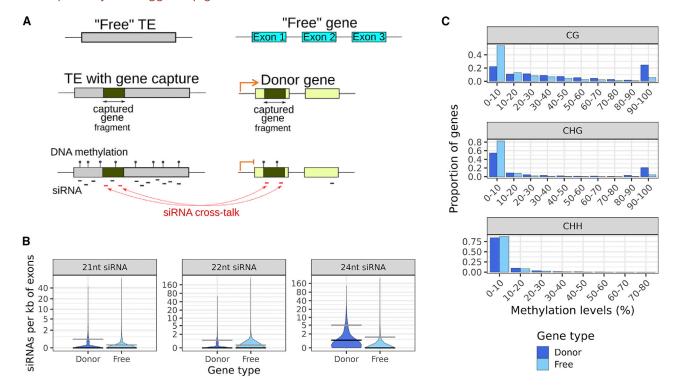


Figure 1. Epigenetic Effects of TE Capture on Donor Genes.

(A) Schematic of a capture event by a TE and ensuing epigenetic interactions. Definitions used in the text are shown, including donor and free genes, free TEs and TEs with captured fragments, and crosstalk siRNAs that may act in *trans*. The orange arrows indicate expression.

(B) Number of 21nt, 22nt, and 24nt distinct siRNA sequences per kb of exonic mapping to donor and free genes. The gray lines indicate the mean, and the black lines represent the median.

(C) Distribution of the proportion of CG, CHG, and CHH exonic methylation of donor and free genes. Data are from ear tissue.

respectively, have high CHG methylation levels (\geq 90%), reflecting the positive correlation between genic CHG methylation and genome size (Niederhuth et al., 2016; Takuno et al., 2016; Seymour and Gaut, 2020). As a result of these differences, large genomes like maize may represent better systems to study gene capture by TEs and the impact on donor genes.

Here, we hypothesize that gene capture may have epigenetic consequences for endogenous genes in maize. To test this hypothesis, we identify capture events representing all three major TE classes, i.e., Helitron rolling circle transposons, Pack-MULE class II DNA transposons, and a representative of class I retroelements, Sirevirus LTR retrotransposons (Bousios and Darzentas, 2013). Sireviruses are crucial because they comprise ~20% of the maize genome (Bousios et al., 2012a), are targeted by large numbers of siRNAs, and are highly methylated (Bousios et al., 2016). Given sets of TEs with gene capture events, we integrate evolutionary analyses with siRNA, methylation, and gene expression data to address two sets of predictions. The first set focuses on the genic viewpoint. If the conflict model holds, we predict that donor genes bear the signature of trans-epigenetic effects compared with free genes. We also predict that natural selection will moderate these epigenetic effects on functionally important genes relative to less-important genes. In the second set of predictions, we focus on TEs with captured gene fragments. Is there any evidence that they benefit from gene capture via decreased host response?

RESULTS

Identifying TE-Captured Gene Fragments and Their Donor Genes

To investigate the potential for intragenomic conflict, we first identified gene capture events. Identifying true events is a challenging task, because annotation errors can lead to false positives that mislead downstream analyses. Annotation errors can be particularly pernicious for TEs, because a proportion of putatively full-length elements may represent partial sequences or mosaics of different TEs. This ambiguity is evident in TE annotations of the recent B73 RefGen_v4 genome that predict, for example, different numbers of Helitron sequences, 49 235 (Jiao et al., 2017) versus 22 339 (Stitzer et al., 2019), even though both used HelitronScanner (Xiong et al., 2014). To address this concern, we favored specificity over sensitivity by using previously published and carefully curated smaller datasets of full-length elements for Helitrons (Xiong et al., 2014), Pack-MULEs (Jiang et al., 2011), and Sireviruses (Bousios et al., 2012b). These datasets were mostly based on RefGen_v2 and contained 1,351, 275, and 13 833 elements, respectively. For example, the 1351 Helitrons represented a high-quality subset of 31 233 full-length elements identified in the original Helitron-Scanner manuscript that were, however, additionally validated in the same study with in silico comparisons with the genome of the Mo17 inbred line (Xiong et al., 2014). We curated these datasets to further remove problematic elements and converted their chromosomal coordinates to RefGen v4 to ensure that

these TEs are physically present in the most recent genome version (see Materials and Methods). Overall, we generated a sample of 7473 TEs, which consisted of 1035 Helitrons, 238 Pack-MULEs, and 6200 Sireviruses. We implemented a similarly strict pipeline for the 39 423 genes of the filtered gene set (FGS) to remove low-quality candidates (e.g., possible misannotated TEs) and 5495 genes that were no longer annotated in RefGen_v4 (see Materials and Methods). The final dataset consisted of 27 056 genes.

We then performed strict BLASTN comparisons (E value cutoff of 1×10^{-40}) between the TEs and the exons of the genes to identify both captured gene fragments within TEs and their donor genes. We only kept hits that belonged to the longest alternative transcript of each gene, removed cases of physical overlaps between full-length TEs and complete genes, and used the BLASTN bit score to select the true donor gene when exons from multiple candidates generated overlapping hits within a TE (Hanada et al., 2009; Jiang et al., 2011). This approach derived a final set of 1629 donor genes (out of 4117 candidate genes with hits to TEs), with the remaining 22 939 genes characterized as free genes. Several features of the donor genes suggest that they are neither pseudogenes nor small gene fragments located within TEs. For example, similar proportions of donor and free genes were assigned a specific function in RefGen_v4 (82.2% versus 84.4%), and donor genes had more exons and total exonic length than free genes (Supplemental Figure 1A and 1B). The donor genes were captured by 1263 distinct TEs. Most Helitrons (873; 84%) and Pack-MULEs (186; 78%) contained gene fragments, in contrast to a much smaller proportion of Sireviruses (204; 3%). Like previous studies (Jiang et al., 2004; Thomas and Pritham, 2015), we found that individual elements often contained multiple independent capture events: 68% of Helitrons harbored ≥2 captured fragments, as did 50% of Pack-MULEs and 15% of Sireviruses. Finally, we confirmed that the elements of each family had sequence or structural features that correspond to full-length TEs, i.e., the conserved 5' and 3' termini for Helitrons (Xiong et al., 2014), terminal inverted repeats for Pack-MULEs, conserved LTR termini for Sireviruses, and short target site duplications for both Pack-MULEs and Sireviruses (Supplemental Figure 2).

Donor Genes Are Targets of siRNAs and Are Highly Methylated

Under our conflict model, the first prediction is that gene capture should lead to siRNA crosstalk between genes and TEs, potentially leading to increased methylation of donor genes. Accordingly, we contrasted the exons of donor and free genes for siRNA mapping and methylation characteristics. Throughout this study, we relied on published siRNA and bisulfite sequencing (BS-seg) datasets, focusing on libraries from unfertilized ears (Nobuta et al., 2008; Gent et al., 2013), leaves of maize seedlings (Diez et al., 2014; Li et al., 2014), and tassels (Zhang et al., 2009) (see Materials and Methods). We analyzed 21nt, 22nt, and 24nt siRNAs (both uniquely and multiply mapped in the genome), because these lengths are involved in TE silencing. Considering, however, that 21nt/22nt siRNAs mostly participate in RNAi/post-transcriptional silencing while 24nt siRNAs participate in RdDM/transcriptional silencing (Matzke and Mosher, 2014; Cuerda-Gil and Slotkin, 2016), we analyzed each

Gene capture by TEs triggers epigenetic conflict

length separately. For each gene, we calculated the number of distinct siRNA sequences per kb of all their exons combined, a metric that avoids the errors inherent to measuring siRNA expression (Bousios et al., 2017). Across all genes and libraries, exonic mapping was strongly correlated for 21nt versus 22nt siRNAs (mean Pearson coefficient r = 0.85, p = 0) but not as much for 21nt/22nt versus 24nt siRNAs (mean Pearson coefficient r = 0.58, p = 0), likely reflecting their roles in different epigenetic pathways. Results were generally consistent among tissues; hence, we report data from the ear in the main text, and provide results from the other two tissues mostly in Supplemental Information.

The comparison of the siRNA mapping profiles of the exons of donor and free genes revealed striking differences: the 1629 donor genes mapped more siRNAs per kb than the 22 939 free genes (Figure 1B and Supplemental Figure 3, one-sided Mann-Whitney U test $p < 2.2 \times 10^{-16}$ for all siRNA lengths and tissues). Across all tissues combined, donor genes mapped 3.0 times more 24nt siRNAs per kb on average than free genes, compared with 1.7 times for 21nt and 22nt siRNAs, respectively. Differences in siRNA mapping are expected to affect methylation patterns. We calculated the proportion of methylated cytosines in the CG, CHG, and CHH contexts of exons using only uniquely mapped BS-seq reads across the genome. Of the 1629 donor and 22 939 free genes, 1525 and 21 614 passed CG methylation filters (≥10 covered CG sites), representing ~94% of the genic dataset, with similar proportions retained for CHG and CHH methylation. We found that the distribution of CG and CHG methylation was notably bimodal, with most genes having either low (\leq 10%) or high (≥90%) methylation (Figure 1C). This pattern is consistent with previous work on several plant species (Niederhuth et al., 2016; Takuno et al., 2016). However, donor and free genes generated strikingly different distributions, showing a bias toward high and low methylation, respectively: in ear, 24.9% of donor genes had ≥90% of their cytosines methylated in the CG context, and 20.9% in the CHG context. In contrast, only 5.5% of free genes were highly methylated in the CG context, with 4.6% in the CHG context. In fact, most free genes had \leq 10% CG and CHG methylation, 54.1% and 82.6%, respectively (Figure 1C). As expected, methylation in the CHH context was much lower, because the majority of genes in both datasets had low (≤5%) levels of methylation. However, although not clearly evident in the histogram, donor genes were significantly more methylated than free genes, with a mean of 5% versus 3.7% (one-sided Mann-Whitney U test, $p < 2.2 \times 10^{-16}$), and a higher proportion with high (≥15%) CHH methylation—i.e., 10.6% versus 6.6% of free genes. Overall, the trends were clear and consistent across all tissues (Figure 1 and Supplemental Figures 3 and 4): donor genes mapped more siRNAs and were more highly methylated.

Dramatic Differences in the Epigenetic Profiles of Syntenic versus Non-syntenic Donor Genes

Our results support the predictions of the conflict model by showing that donor genes are heavily enriched for both siRNA mapping and methylation levels. However, the model specifically proposes that intragenomic conflict arises for functional genes, but many donor genes have high levels of methylation, especially in the CHG context (Figure 1C and Supplemental Figure 4), which

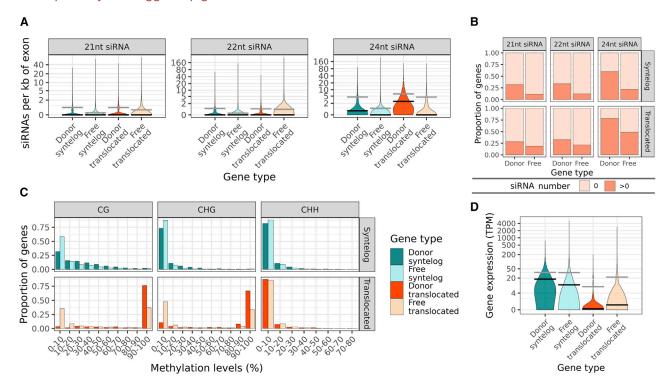


Figure 2. Epigenetic and Expression Profiles of Donor and Free Genes Split by Their Syntenic Status with Sorghum.

The four gene categories in all plots are donor syntelogs, free syntelogs, donor translocated genes, and free translocated genes.

- (A) Number of 21nt, 22nt, and 24nt distinct siRNA sequences per kb of exonic mapping.
- (B) Proportion of genes with no siRNA mapping.
- (C) Distribution of the proportion of CG, CHG, and CHH exonic methylation.
- (D) Gene expression measured in TPM. Data are from ear tissue. The gray lines in (A) and (D) indicate the mean, and the black lines represent the median.

is a potential signature of silencing. To better test the conflict model, we split the genic dataset according to syntenic relationships with *Sorghum bicolor* (Springer et al., 2018). Our reasoning was that syntenic orthologs are enriched for genes that are functional and associated with phenotypes (Schnable, 2015, 2019). We thus expect that these genes are more often subject to selective constraint and hence susceptible to epigenetic conflict. In contrast, non-syntenic genes are more likely to be dispensable or non-functional (Schnable, 2015, 2019) and thus less likely to be under strong selective constraint. We therefore predict that, as a general trend, the conflict model should be less obvious for non-syntenic genes.

We assigned the 1629 donor and 22 939 free genes into two categories: syntenic orthologs (hereafter syntelogs) and genes that have moved their location in maize relative to sorghum (hereafter translocated) (see Materials and Methods). The two categories yielded a striking observation: translocated genes had a higher probability than syntelogs to be captured by TEs. In total, 58.2% (948) of donor genes were syntelogs and 27.1% (442) were translocated, while their proportions among free genes were 78.7% (18 046) and 9.8% (2243), respectively (chisquare = 512.37, $p < 2.2 \times 10^{-16}$). We also note that a much higher proportion of donor and free syntelog genes (91.6% and 87.2%) was assigned a specific function compared with translocated genes (64% and 63.3%) based on the RefGen_v4 annotation, which supports our contention that syntelogs are

more likely to be functionally important. Most of the remaining donor (201; 12.3%) and free (2284; 10%) genes were either located in regions of maize chromosomes that were not identified in sorghum or completely lacked synteny information. We excluded these two categories from further analyses due to their ambiguous syntenic status.

We then contrasted the epigenetic profiles of syntelog and translocated genes, starting with siRNA mapping. Looking within each synteny-based category, the differences between donor and free genes remained consistent to our analysis based on all genes, i.e., donor genes were targeted by more siRNAs per kb than free genes across all tissues (Figure 2A and Supplemental Figure 5A, one-sided Mann–Whitney U test $p < 2.2 \times 10^{-16}$ for most combinations). This result was largely due to the higher fraction of free genes that did not map any siRNAs, with this difference being more prominent for 24nt siRNAs, where twice as many free genes had no mapping events compared with donor genes (Figure 2B). Removing genes with no siRNAs did not change the mapping differences within syntelogs but did so for translocated genes because donor and free genes were equally targeted by siRNAs (Supplemental Figure 5B). For syntelog genes and across all tissues combined, donor genes mapped 3.8-fold more 24nt siRNAs per kb on average than free genes, compared with 2.1-fold for both 21nt and 22nt siRNAs. These differences in average siRNA mapping were not as strong within translocated genes; for example, donor genes mapped 1.4-fold more 24nt siRNAs than free genes.

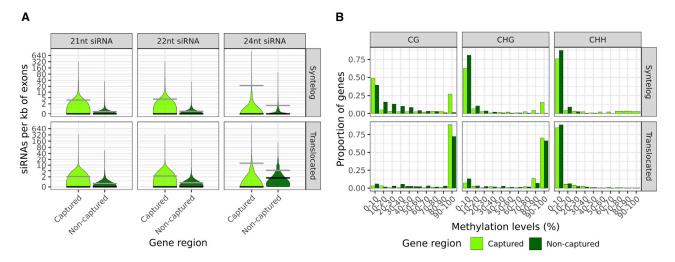


Figure 3. siRNA and Methylation Patterns of TE-Captured and Non-captured Regions of Syntelog and Translocated Donor Genes.
(A) Number of 21nt, 22nt, and 24nt distinct siRNA sequences per kb of captured and non-captured exonic regions. The gray lines indicate the mean, and the black lines represent the median.

(B) Distribution of the proportion of CG, CHG, and CHH exonic methylation of captured and non-captured regions. Data are from ear tissue.

Noting that donor genes had higher levels of siRNA targeting than free genes within each synteny-based category, the differences between categories were intriguing. Across libraries, the most striking pattern was that the level of 24nt siRNA targeting of donor translocated genes was higher than any other category (Figure 2A and Supplemental Figure 5A; one-sided Mann–Whitney U test p < 2.2×10^{-16}). This was not the case for 21nt-22nt siRNAs that often did not map at levels statistically different between syntelog and translocated genes (Figure 2A and Supplemental Figure 5A). Taken together, these results lead to three main observations. First, gene capture by TEs is linked to increased levels of siRNA targeting in donor genes (e.g., donor syntelog versus free syntelog genes); second, these levels are lower in "important" compared with "less important" genes (e.g., donor syntelogs versus donor translocated genes); and third, 24nt siRNAs appear to be the crucial component of these differences.

The corresponding methylation patterns of the four gene categories supported the siRNA results. The distribution of CG and CHG methylation in donor and free syntelogs was no longer bimodal due to the absence of genes with high methylation (Figure 2C and Supplemental Figure 6). However, donor syntelogs still had higher methylation levels than free syntelogs (in ear mean CG 26.7% versus 15.5%, one-sided Mann-Whitney U test $p < 2.2 \times 10^{-16}$; CHG 9.6% versus 4.9%, p < 2.2×10^{-16}), including in the CHH context (mean 5.3% versus 3.6%, $p = 3.5 \times 10^{-11}$), where twice as many donor syntelogs had high methylation (12.6% versus 6.5%). The methylation profile of translocated genes differed from the above patterns (Figure 2C and Supplemental Figure 6). The majority of translocated donor genes had high CG (76.7%) and CHG (66.8%) methylation and virtually none had low methylation, while translocated free genes were clearly distinguished by having a bimodal distribution of methylation (Figure 2C and Supplemental Figure 6). Unlike syntelogs, CHH methylation was more similar between the two sets of translocated genes (in ear mean 4.9% versus 4.6%, p = 0.0052; 7.8% of donor versus 8.4% of free with high CHH methylation). Overall, the methylation patterns recapitulated the siRNA patterns by suggesting that donor genes tended to be more methylated than free genes, which in the case of donor translocated genes reached very high levels.

TE-Captured Regions of Syntelog Donor Genes Are Enriched for Repressive Epigenetic Marks

An additional prediction of the conflict model suggests that the epigenetic marks of siRNA mapping and methylation should be overrepresented in the regions that were captured by TEs, at least for functionally important genes (Figure 1A). To examine this prediction, we compared siRNA mapping and methylation levels between the captured versus non-captured exonic regions of donor genes. As predicted, significantly more siRNAs mapped to the captured than the non-captured regions of syntelogs, with 24nt siRNAs generating the strongest difference across libraries (one-sided Wilcoxon signed rank test $p < 2.2 \times 10^{-16}$; Figure 3A and Supplemental Figure 7A). The captured regions of syntelogs were also significantly more methylated in the CG (in ear one-sided Wilcoxon signed rank test $p = 9.45 \times 10^{-12}$), CHG ($p = 4.36 \times 10^{-12}$), and CHH (p = 1.49×10^{-5}) contexts (Figure 3B and Supplemental Figure 7B). Supporting the statistical tests, high levels of methylation were found only in captured regions of syntelogs; in ear, 27.3% of captured regions had ≥90% of their CG sites methylated versus only 1.4% for non-captured regions. These differences extended to high CHG (15.6% captured versus 0.2% non-captured regions) and high CHH (22% captured versus 6.5% non-captured regions) methylation. This result was robust when we increased the coverage filter in each locus from ≥ 10 to ≥ 40 covered sites, which tended to exclude captured regions with short lengths that could bias the results (Supplemental Table 1). In contrast to syntelogs, the captured and non-captured regions of translocated genes did not significantly differ in either siRNA targeting or methylation levels in most cases (Figure 3 and Supplemental Figure 7). Altogether, these results suggest that repressive epigenetic marks may

M	lecu	ılar	P	lar	٦ŧ

Region	Methylation context	Estimate	SE	t Value	p Value	Marginal R ²	
Captured	CG	2.82 ×10 ⁻³	5.41 ×10 ⁻⁴	5.21	2.21 ×10 ⁻⁷	0.0156	
	CHG	0.00577	0.000455	12.7	<2 ×10 ⁻¹⁶	0.107	
	СНН	0.00513	0.000280	18.3	<2 ×10 ⁻¹⁶	0.215	
Non-captured	CG	-0.001326	0.0003339	-3.97	0.0000746	0.0096	
	CHG	no significant effect (p = 0.1402)					
	СНН	no significant effect (p = 0.5697)					

Table 1. Correlation between the Number of 24nt Crosstalk siRNAs and Methylation Levels of TE-Captured and Non-captured Regions of Donor Syntelog Genes.

not spread beyond the captured regions of important syntelog genes but may do so in translocated genes.

24nt Crosstalk siRNAs Are Enriched in Donor Syntelog Genes and Affect CHH Methylation In Trans

Our conflict model is based on crosstalk siRNAs, which are siRNAs that map both to the captured fragment within the TE and the gene and may act in trans. To test whether crosstalk siRNAs represent an enriched fraction of the total number of siRNAs that mapped to donor genes, we employed a binomial test that compared the observed proportion of crosstalk siRNAs (crosstalk siRNAs/all siRNAs) to the proportion of captured gene length (captured exon length/total exon length) across all genes. In each tissue, the binomial test revealed a significant enrichment of crosstalk siRNAs in donor syntelogs (Supplemental Table 2), and this was especially strong for 24nt siRNAs (p ~ 0). Translocated genes had significantly fewer crosstalk siRNAs in leaf and tassel, and moderate statistical support for enriched crosstalk in ear (Supplemental Table 2). These contrasting patterns suggest that syntelogs are disproportionately targeted by crosstalk siRNAs, while translocated genes are targeted more generally.

A key prediction of the conflict model—i.e., that gene capture has the capacity to modify the epigenetic state of the donor gene presupposes that crosstalk siRNAs can act in trans. Hence, we used a linear model with mixed effects across all tissues (see Materials and Methods) to examine the relationship between crosstalk siRNAs and methylation of captured and noncaptured regions of donor genes. We found that the number of 24nt crosstalk siRNAs was positively correlated to the methylation levels of captured regions within donor syntelogs (Table 1). This was especially true for CHH methylation, where 21.5% of the variance across captured regions was explained by the abundance of 24nt crosstalk siRNAs, but 10.7% of the variance was also explained for CHG methylation and 1.6% for CG methylation. The fact that more variation was explained for CHH methylation makes biological sense, because methylation in this context is maintained de novo by 24nt siRNAs via RdDM (Matzke and Mosher, 2014). In contrast, the methylation levels of the non-captured regions were not positively associated with the abundance of 24nt crosstalk siRNAs (Table 1). In addition, these patterns did not hold as clearly for the other siRNA lengths for donor syntelogs or, generally, for translocated genes (Supplemental Table 3). For example, only 2.6% and 4.7% of the variance of CHH methylation within the captured regions of donor syntelogs was explained by 21nt and 22nt crosstalk siRNAs respectively. Taken together, these results establish an epigenetic link in donor syntelogs between gene capture, 24nt crosstalk siRNAs that act in trans, and methylation, particularly in the CHH context.

The fact that siRNA crosstalk is significant for syntelog genes raises an interesting question: what is the relationship between siRNA crosstalk and the time since the capturing event took place? This is probably a complex relationship, for two reasons. First, the initiation of the host epigenetic response against a new capture event may not be immediate, so that very recent capture events may not generate enough siRNAs to detect crosstalk. Second, the opportunities for crosstalk are finite, because the sequences of the donor gene and the captured fragment within the TE diverge over time. As they diverge, crosstalk can no longer occur as efficiently because siRNAs no longer match both entities. We used synonymous divergence (dS) between the donor syntelog and the TE-captured exon as a proxy of the age of capture. We then examined the relationship between the abundance of siRNA crosstalk and time since gene capture. The tests were significant and positively correlated only when we combined all siRNA lengths. This correlation suggests that donor syntelogs with older capture events had more crosstalk siRNAs over time, despite the increased divergence of their captured sequences (generalized linear model with mixed effects across all tissues z value = 2.04, p = 0.0413, marginal $R^2 = 0.006$, Supplemental Figure 8, see Materials and Methods). The results held after taking into account the variability in dS across genes by using the dS of maize-sorghum syntelogs as an offset in the generalized linear model (z value = 1.901, p = 0.0573, marginal $R^2 = 0.00589$). Overall, we interpret these results to imply that it takes time for crosstalk to evolve after the capture event.

TE Capture Does Not Affect the Expression of Donor **Syntelog Genes**

Our analyses are consistent with the interpretation that crosstalk siRNAs drive, to some extent, methylation of donor genes. The conflict model predicts, however, that these epigenetic modifications will have minimal effects on important genes, because natural selection acts against changes that affect function. To test this conjecture, we contrasted expression patterns of donor and free genes using data from the ATLAS Expression database (see Materials and Methods). Indeed, we did not find significantly lower levels of expression in donor syntelog compared with free syntelog genes in ear, leaf, and 10 different cell types of the maize kernel (Figure 2D and Supplemental Figure 9). In fact, donor syntelogs were expressed at significantly higher levels

(log transformed average of transcripts per million [TPM] 1.88 versus 1.23 across tissues, one-sided Mann-Whitney U test p between 2.7×10^{-8} and 2.2×10^{-16}) and had a lower proportion of genes with no expression (average of 11.5% versus 24% across tissues) than free syntelogs. In contrast, both categories of translocated genes had lower levels of expression compared with syntelog genes (Figure 2D and Supplemental Figure 9), which is in agreement with previous studies that showed translocated genes to have pseudogene-like characteristics (Schnable, 2015, 2019; El Baidouri et al., 2018). Donor translocated genes appeared to be driving this difference, because they had significantly lower expression than free translocated genes (log transformed average TPM across tissues -0.52 versus 0.32, one-sided Mann-Whitney U test p between 0.0009 and 8.23×10^{-12}). Therefore, donor translocated genes exhibit a signal consistent with runaway epigenetic interactions with TEs that are not moderated by functional constraints, hence dramatically reducing expression.

Finally, we examined if gene expression is affected by the position of the captured fragment within the gene. For example, it is possible that capture and subsequent epigenetic interactions at the 5' or 3' untranslated regions (UTRs) may affect expression levels, because these are regions of major importance for gene regulation (Barrett et al., 2012; Dvir et al., 2013). To investigate this, we classified each gene based on whether the captured fragment(s) were part of the 5' or 3' exons to approximate the location of UTRs, or any internal exon. By analyzing genes whose captured fragment(s) were derived from a single position only, we found that capture of the 5' or 3' exons of syntelogs significantly reduced expression across all tissues compared with capture of internal exons (one-sided Mann-Whitney U test p between 1.9 \times 10⁻⁸ and 2.2 \times 10⁻¹⁶), while this was not the case for translocated genes (Supplemental Figure 10). We note, nevertheless, that the expression levels of these donor syntelogs were still higher than free syntelogs (one-sided Mann-Whitney U test p between 0.02951 and 2.1×10^{-6} across tissues).

Genes under Weak Selective Constraint Suffer from a Broader Epigenetic Impact of TE Capture

One of the reasons for separating genes into syntelogs and translocated genes was to compare the consequences of gene capture between more- and less-constrained genes. It has been shown that translocated genes tend to be under weaker selective constraint (Schnable, 2015, 2019; El Baidouri et al., 2018), but we sought to directly test this by calculating the dN (nonsynonymous divergence)/dS ratio of maize-sorghum orthologs (see Materials and Methods). Our findings corroborate the previous studies: translocated genes had a significantly higher dN/dS ratio compared with syntelogs (Supplemental Figure 11). We also used the dN/dS values to independently categorize donor and free genes into constrained (dN/dS < 0.4) or less constrained (dN/dS > 0.4) (see Materials and Methods). This represents an alternative method to contrast genes with different functional constraints, irrespective of their status as syntelogs or translocated genes. We repeated all previous analyses and, overall, confirmed the results of the synteny-based approach: first, donor genes mapped significantly more siRNAs than free genes within each dN/dS-based category (Supplemental

Gene capture by TEs triggers epigenetic conflict

Figure 12A); second, constrained donor genes accumulated significantly more methylation in all three contexts compared with constrained free genes; third, less-constrained donor genes had the highest CG and CHG methylation levels among gene categories (Supplemental Figure 12B); fourth, constrained donor genes mapped significantly more siRNAs and methylation inside the captured region only, but lessconstrained donor genes do so across their entire sequence (Supplemental Figure 13); finally, the categories of constraint recapitulated the synteny-based analyses with respect to gene expression, i.e., constrained donor genes were expressed at significantly higher levels than constrained free genes (Supplemental Figure 14A). Overall, these results support the existence of intragenomic conflict with outcomes that vary according to the putative functional importance of genes, but without relying on synteny-based definitions.

Potential Advantages for TEs to Capture Gene Fragments

As well as the impact on genes, the conflict model also predicts that TEs with captured gene fragments gain an advantage due to a moderation of the host response. To explore this possibility, we focused on 852 TEs that captured at least one fragment from a syntelog and contrasted them to 5931 free TEs that had no BLASTN hit to the gene dataset. Given these two groups, we considered three potential measures of advantage for TEs with syntelog capture: (1) they may be retained within the genome for longer lengths of time, (2) they may be targeted by fewer siRNAs, and (3) they may have lower levels of methylation.

To test the first idea, we used age estimates from terminal branch lengths of TE phylogenetic trees generated by Stitzer et al. (2019). We found that TEs with syntelog capture are older than free TEs (mean of 0.135 versus 0.066 million years, one-sided Mann-Whitney U test $p < 2.2 \times 10^{-16}$; Figure 4A), suggesting that they have remained intact within the genome for longer periods. They also mapped significantly less siRNAs of all lengths based on a linear model across all tissues and after removing the captured regions from TEs with captured fragments (for example, contrast for 24nt siRNAs z ratio = -57.59, p < 0.0001, marginal R² = 15.55%, see Materials and Methods) (Figure 4B and Supplemental Figure 15A). This result remained significant after including TE age in the model (Supplemental Table 4). We also found that TEs with syntelog capture were less methylated than free TEs in both the CG (ear mean 95.5% versus 98.5%) and CHG (89.4% versus 91.4%) contexts (Figure 4C and Supplemental Figure 15B). These differences were small but significant for CG methylation in a linear model across all tissues and held after controlling for TE age (CG: contrast z value = 4.71, p = 2.44×10^{-6} ; CHG contrast p = 0.695; Supplemental Table 5). However, TEs with syntelog capture had significantly more CHH methylation compared with free TEs (12.0% versus 3.1%, contrast t value = -25.86, $p < 2 \times 10^{-6}$, Figure 4C and Supplemental Figure 15B, Supplemental Table 5).

Patterns of Conflict Are Reproducible at the TE Family Level

Thus far we have primarily reported analyses based on all three TE families together. We combined families to provide a large

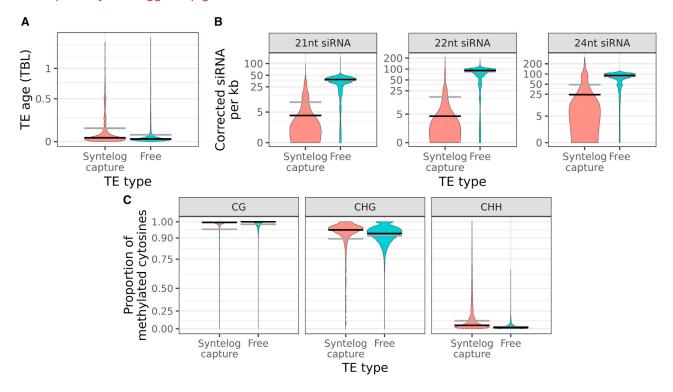


Figure 4. Characteristics of TEs with Syntelog Capture versus Free TEs.

(A) TE transposition age in terminal branch lengths (TBL).

(B) Number of 21nt, 22nt, and 24nt distinct siRNA sequences per kb mapping to TEs. This was computed after removing captured regions from TEs, but results were qualitatively identical when they were included.

(C) Proportion of methylated cytosines in CG, CHG, and CHH contexts of TEs. The gray lines indicate the mean, and the black lines represent the median.

number of observations for both syntelog and translocated genes, but also to provide a global view of the epigenetic effect of TE capture on host genes. We did, however, examine each family separately and found that the main results were reproducible at the family level (summarized in Supplemental Tables 6 and 7): (1) donor genes mapped more siRNAs and were more methylated than free genes in both syntenic categories; (2) the hotspot for these epigenetic patterns was the captured region for donor syntelogs: (3) the expression level was higher for donor syntelogs compared with free syntelogs and, conversely, it was lower for donor translocated versus free translocated genes. We note, however, that the expression of genes captured by Pack-MULEs was not statistically different to free genes either for syntelogs or translocated genes. Finally, we repeated the analysis for TE advantage. Only Sireviruses generated significant trends and only for age (mean of 0.0972 versus 0.0652 million years, one-sided Mann-Whitney U test p = 6.14 \times 10⁻⁶) and siRNA mapping (Supplemental Tables 4 and 5). The lack of significance for Helitrons and Pack-MULEs may reflect the fact that few of these elements lacked captured gene fragments.

Overall, these findings suggest that gene capture by any TE family is likely to trigger similar downstream epigenetic effects. However, we also generated evidence that TE families may be capturing genes in distinct ways. For example, Helitrons and, to a lesser extent, Sireviruses exhibited a preference for translocated genes: compared with their proportion of 9.8% among free genes, translocated genes represented 30.4% (412/1356)

of the genes captured by Helitrons (chi-square = 556.34, p < 2.2×10^{-16}) and 20.3% (30/148) of those captured by Sireviruses (chi-square = 17.075, $p < 1.8 \times 10^{-5}$) (Figure 4B and Supplemental Figure 16A). In contrast, only 7.5% (15/200) of the Pack-MULE genes were translocated. Furthermore, the three families tend to capture different parts of genes. Helitrons exhibited a preference for internal exons (45.3%), but Pack-MULEs most often captured 5' exons (42.1%) and Sireviruses most often captured 3' exons (64.2%) (Figure 4B and Supplemental Figure 16B). If the capture of internal exons has smaller effects on donor gene expression (Supplemental Figure 10), then Helitrons may cause less functionally impactful epigenetic conflicts than the other families. The orientation of the captured fragments also differed among families. While there was no orientation bias for Pack-MULEs and Sireviruses, 78.6% of fragments captured by Helitrons were in the sense orientation (Supplemental Figure 16C). If antisense capture events can trigger RNAi more readily during TE expression, then their deficit in Helitrons may again limit some aspects of their epigenetic effects on donor genes.

DISCUSSION

TEs are often in conflict with their plant hosts, because their proliferation tends to have a deleterious effect on host fitness. While this aspect of the TE-host conflict is well established, here we have studied a unique aspect of their conflict, which is driven by gene capture and ensuing epigenetic interactions between TEs and genes. To study this conflict, we have formalized a model

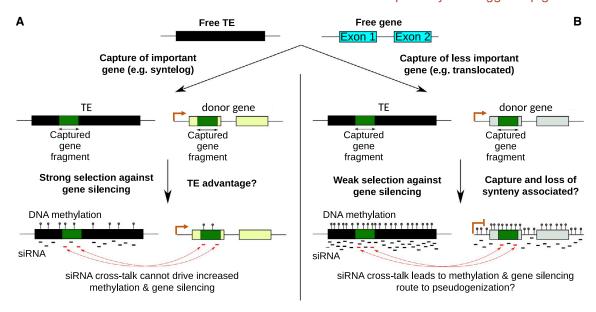


Figure 5. The Epigenetic Conflict Model of Gene Capture.

When TEs capture fragments of genes, siRNAs derived by the TEs may act in trans to accidentally mediate an epigenetic response against the gene, leading to increased methylation and reduced expression. The conflict comes from evolutionary pressure to silence TEs without simultaneously silencing functionally important genes, syntelogs in our example. As a result, (A) epigenetic effects on these genes are moderated by natural selection and expression is not affected. TEs may benefit from this moderation, although this remains unclear. In contrast, (B) for genes that are not under strong selective constraint, methylation can increase in the absence of conflict, leading to loss of expression and potential pseudogenization. This profile is characteristic of genes that have moved from their syntenic loci, which are overrepresented among donor genes, suggesting that capture may trigger movement.

suggested by Lisch (2009). He argued that gene capture can have a beneficial effect on TEs because they become "camouflaged" and, hence, are less apt to be silenced by the host epigenetic machinery. We have extended the model to also consider the effect of capture on donor genes, predicting that they should have higher siRNA mapping relative to genes with no history of capture. Moreover, if the TE is subjected to silencing, we predict that siRNA crosstalk between the TE and the gene drives epigenetic alterations to the gene itself. The epigenetic modification of the donor gene may eventually reach a threshold that affects gene function, ultimately driving intragenomic conflict, especially if the gene is functionally important. That is, when the silencing response against the TE becomes deleterious to the donor gene, then natural selection may favor a moderation of the silencing response (Figure 5A).

The Case for Conflict: Syntelog Genes

What is the evidence to support this model? Based on our dataset of syntelogs—to which the conflict model should apply because they are enriched for functionally important genes according to both previous studies (Schnable, 2015, 2019) and our own dN/dS analyses—we find that donor syntelogs map more siRNAs and are more highly methylated than free syntelogs (Figure 2A and 2C). These epigenetic markers are enriched in the captured regions of donor syntelogs (Figure 3), where a large fraction of siRNAs also map (crosstalk) to the captured fragment within the TE (Supplemental Table 2). In addition, there is a clear relationship between 24nt crosstalk siRNAs and methylation levels in the captured region (Table 1 and Supplemental Table 3). This relationship is stronger for CHH methylation, which is more reliant on RdDM than CG and CHG methylation

(Matzke and Mosher, 2014). The magnitude of the effect is not inconsequential, because the number of 24nt crosstalk siRNAs explains ~21% of CHH methylation variation across captured fragments. These results suggest that (1) substantial RdDM activity occurs predominantly in the captured fragments of donor syntelogs, and that (2) this activity is likely driven by epigenetic crosstalk with the TEs that contain the captured sequences.

The conflict model further predicts that the epigenetic interactions should not proceed to the extent that gene expression is altered, because natural selection will conserve the function of important genes. We assessed function by comparing gene expression between donor and free syntelogs, and indeed found no evidence of reduction in expression (Figure 2D and Supplemental Figure 9); in fact, donor syntelogs were more highly expressed than free syntelogs. We propose that this difference likely reflects biases in capture events. This hypothesis presupposes that TEs are better able to capture highly expressed genes in open chromatin, and it conforms to the integration preferences of several TE families across plants and animals for genic regions (Bousios et al., 2012a; Sultana et al., 2017; Zhang et al., 2020). Another interesting feature is the genic region that has been captured. Evidence suggests that methylation of the 5' and 3' UTRs-two regions with important regulatory roles for gene function (Barrett et al., 2012; Dvir et al., 2013)-significantly affects expression levels in humans (Luo et al., 2018). Our analysis supports this claim by showing that capture of the 5' or 3' side of syntelogs is associated with lower expression than capture of internal exons (Supplemental Figure 10). Based on this negative effect, one expects 5' and 3' capture events to be rare for important genes

or that natural selection quickly removes them from the population. However, we did not observe a scarcity of 5' or 3' captured fragments in donor syntelogs compared with donor translocated genes (17.3% versus 15.6% for 5' exons; 20.8% versus 21.5% for 3' exons); we hypothesize that their abundance might be linked to the time that it takes for siRNA crosstalk to establish.

We also need to address the case of directionality. Could it be that TEs simply tend to capture highly methylated genes? This notion can be rejected based on at least four pieces of information. First, most of the 1629 donor genes are syntelogs (948; 58.2%) whose methylation levels are much lower compared with donor translocated genes. Second, this argument does not easily explain why the epigenetic effects are found only within the captured region of donor syntelogs. This difference is not likely to be a simple function of statistical power, because noncaptured regions were consistently longer than captured regions. Third, the linear model (Table 1) establishes a positive relationship among capture, crosstalk siRNAs, and methylation. Since siRNAs facilitate methylation via RdDM and not vice versa, this result implies a directionality that contradicts the simple explanation of a high methylation capture bias. Finally, we can use orthology relationships with sorghum to assess whether capture is biased toward genes that are highly methylated (Supplemental Figure 17). We examined patterns of CG methylation between maize and sorghum syntelogs, separated between the free and donor categories based on our analysis in maize. Both free and donor genes show a positive correlation between species (Spearman coefficient r = 0.62 for free and r = 0.5 for donor genes, $p < 2.2 \times 10^{-16}$ for both), which reflects the well-established fact that genic methylation tends to be conserved over evolutionary time (Takuno and Gaut, 2013; Niederhuth et al., 2016; Takuno et al., 2016; Seymour and Gaut, 2020). More importantly, however, they imply that preexisting levels of methylation do not seem to trigger capture events, because genes located throughout the methylation spectrum have been captured in maize.

The Case for Conflict: Translocated Genes

Our results clearly illustrate the epigenetic effects of TE capture on donor syntelogs, but an additional feature that merits discussion is the curious case of translocated genes (Figure 5B). Translocated genes are overrepresented among donor genes, because 16.5% of all translocated genes were found to be donors compared with only 5% of all syntelog genes. This pattern was especially evident for Helitrons (Supplemental Figure 16A). Previous work has shown that TEs contribute to modifications of synteny (Wicker et al., 2010), suggesting that TE capture can trigger gene movement. It is therefore possible that the categories of donor and "translocated" are linked mechanistically, i.e., that gene capture and gene translocation happen concomitantly.

Donor translocated genes are also highly methylated as most have >90% CG and CHG methylation (Figure 2C). This pattern is consistent with gene silencing, which is supported by siRNA mapping that is not specific to the captured region (Figure 3 and Supplemental Figure 7), very low expression levels (Figure 2D), and low percentage with functional annotation

(64%). We propose that donor translocated genes are the exceptions that prove the rule-i.e., they illustrate the runaway effects of epigenetic interactions with TEs in the absence of selection for function. That said, it is worth emphasizing that the epigenetic patterns of donor translocated genes are not a feature of translocated genes in general, because only a subset of free translocated genes has the combination of high methylation and low expression levels (Figure 2C and 2D). In addition, donor translocated genes have a significantly higher dN/dS ratio than free translocated genes (Supplemental Figure 11), suggesting that these genes may be en route to pseudogenization.

We have suggested that TEs may tend to capture highly expressed genes, which seems to superficially contradict the observation that donor translocated genes are lowly expressed. These observations need not be at odds, however. If capture often triggers translocation, it can then lead to loss of expression as a downstream consequence for less-important genes. In contrast, critical genes may be constrained in location; that is, natural selection may filter translocation of important genes so that they never, in effect, become translocated. Our dN/dS analysis is consistent with this reasoning, because donor translocated genes have the highest proportion of genes in the less-constrained category (46% of genes with dN/dS > 0/4, Supplemental Table 8). This subset of genes has the lowest expression level among all gene categories (Supplemental Figure 14B), while nearly all of them (~90%) have high or intermediate CHG methylation levels (Supplemental Table 8).

Also, the chromosomal location of the translocation event is probably important. If the translocated gene "lands" in a region permissive to heterochromatin formation, then the gene is likely to reach higher levels of non-CG methylation that reduce expression. Intriguingly, we find that donor translocated genes are located in regions near the pericentromeres when they have high or intermediate levels of CHG methylation, while those with low CHG methylation are found in the chromosomal arms (Supplemental Figure 18). In contrast, all donor syntelogs, regardless of their CHG methylation levels, are found in the chromosome arms, which is the typical distribution of maize genes (Schnable et al., 2009). This difference suggests that translocation could lead to the acquisition of CHG methylation when it is directed toward heterochromatin-prone areas. This finding follows the results of a recent study in soybean that showed that non-syntenic paralogs are enriched for non-CG methylation as a result of their movement to pericentromeres (El Baidouri et al., 2018), and also provides evidence that this movement may be mediated by TE capture in some cases.

The Case for Conflict: TEs

The conflict model also predicts that TEs with captured fragments of important genes gain an advantage. It is an open question as to how to measure such an advantage, and so we investigated several potential options. We asked, for example, whether TEs with fragments of syntelog genes have a tendency for camouflage, as measured by siRNA mapping or methylation levels. Consistent with the conflict model, these TEs map fewer siRNAs than free TEs, even when the captured region was masked (Figure 4B) and when TE age was taken into account.

One caveat to this result is that we likely underestimated the size of the captured region; this could bias analyses if captured regions tend to map fewer siRNAs than TE-specific regions. TEs with syntelog capture also tend to have lower CG and CHG methylation than free TEs (Figure 4C). However, this is a nuanced result, for two reasons. First, we find that differences are small in magnitude and both TE types have >90% methylation on average. At these levels, any TE is probably effectively silenced. Second, TEs with syntelog capture events have ~3-fold higher levels of CHH methylation (Figure 4C), which is hard to reconcile with the lower number of matching siRNAs. The cause of this CHH difference remains elusive, but it contributes to the overall impression that these TEs have ongoing epigenetic interactions defined in large part by increased CHH methylation levels for both the TE and the donor genes. Finally, if gene fragments provide camouflage for TEs, one reasonable prediction is that they will exist within the genome for longer periods of time than free TEs. We found that this is indeed the case (Figure 4A), but this result was principally caused by Sireviruses, perhaps in part reflecting their higher proportion of free TEs. Notably, the slightly lower levels of CG/CHG methylation and longer periods of retention are similar to the findings of a recent study that examined gene capture by the GingerRoot DNA transposon in the clubmoss Selaginella lepidophylla (Cerbin et al., 2019). Altogether, we consider the case for TE advantage to be tantalizing and perhaps correct, but not yet fully convincing.

It is worth mentioning, however, two additional points. First, it is possible that the epigenetic response against a TE continues unabated after the capture event, so that it gains no advantage, but the epigenetic effects on donor genes are moderated by natural selection using other mechanisms, such as active CHG demethylation (Wendte et al., 2019). In such a scenario, further research could elucidate how demethylation may only occur in the subset of genes that are important for host function. Second, it is likely that there may be little advantage for a nonautonomous TE to capture genic sequences to avoid silencing. Most Helitrons and all Pack-MULEs are non-autonomous (Thomas and Pritham, 2015; Zhao et al., 2018), and for such elements the activity of the autonomous TE may be more important for proliferation than their own methylation levels. This relationship has been shown for the Ping/mPing family in rice (Lu et al., 2017). That said, capture by a non-autonomous TE can still establish epigenetic crosstalk with the donor gene and, hence, may offer some protection to the autonomous TE by increasing the cost of targeting a family as a whole.

Limitations of Our Analyses

We recognize that our set of donor genes does not represent all capture events throughout the history of the maize genome. This is because we did not examine all known TE families in maize and because we used criteria to identify capture events that were stricter than previous studies (e.g., E value cutoffs of 1 × 10⁻⁴⁰ versus 1 × 10⁻⁵) (Jiang et al., 2004, 2011; Du et al., 2009; Hanada et al., 2009; Yang and Bennetzen, 2009), a conservative approach that favors specificity over sensitivity. As a result of these methodological decisions, our set of free genes must contain false negatives, i.e., undetected capture events. Similarly, our set of donor genes may also contain false

Gene capture by TEs triggers epigenetic conflict

positives. The crucial point about both false negatives and false positives is that they should reduce—and not enhance—epigenetic differences between donor and free genes. Hence, we suspect that, if anything, we have systematically underestimated the magnitude of epigenetic effects of TE capture on donor genes, at least for maize.

Another issue specifically concerns translocated genes. It is possible that some of these genes are misannotated TEs, which could explain their overrepresentation among donor genes (as a result of TE-TE hits during the BLASTN run). Although we cannot exclude this possibility, we believe that it is not a widespread phenomenon. This is based on our strict methodology, but it is also supported by recent data using long-read cDNA sequencing and gene-quality annotation of TEs, which showed that TEs have fewer and longer exons than genes (Panda and Slotkin, 2020). If donor translocated genes are TEs, they should have fewer and longer exons, but they do not (Supplemental Figure 1C and 1D). Furthermore, near identical proportions of donor translocated and free translocated genes had a sorghum ortholog identified and passed the dN/dS inference filters (Supplemental Table 8), suggesting that donor translocated genes are not preferentially enriched with misannotated TEs. Overall, our data suggest that misannotation bias is unlikely to drive the observed epigenetic and expression differences of donor translocated genes in relation to the other gene categories, especially the free translocated genes. Finally, we also acknowledge that some of the BLASTN hits may be the result of TE exaptation (i.e., a TE fragment becoming part of a gene) rather than a gene capture event. Crucially, TE exaptation should also trigger conflict as long as there is sufficient sequence similarity between the gene and another element of the same TE family. Examining the intensity and evolution of conflict for both these events within the same system is an interesting direction.

Concluding Remarks

The intragenomic conflict between TEs and host genomes described here raises several questions for future investigation. For example, it is likely that our model applies generally to plant genomes because gene capture by TEs is a common occurrence (Jiang et al., 2004; Morgante et al., 2005; Holligan et al., 2006; Thomas and Pritham, 2015; Zhao et al., 2018; Catoni et al., 2019), but it remains to be seen if the conflict is more pervasive in species with higher methylation levels and TE load, as is often the case for large genomes (Niederhuth et al., 2016; Takuno et al., 2016), or if it varies across TE types depending on their intrinsic transposition and capturing mechanisms. In addition, genes that have translocated from their syntenic loci account for a substantial proportion of the gene content of plants; for example, thousands of genes have lost synteny between maize inbred lines, such as B73, Mo17, and W22 (Springer et al., 2018; Sun et al., 2018). Further research is needed to show if capture by TEs is mechanistically linked with this gene movement; if true, then it may represent a main route toward pseudogenization.

MATERIALS AND METHODS

TE and Gene Datasets

For TEs, we utilized three published datasets that were carefully curated, representing full-length Helitrons, Pack-MULEs, and Sireviruses. For

Helitrons, we downloaded the coordinates on the B73 RefGen_v2 genome of 1351 high-quality elements that were in silico validated with the Mo17 inbred line in Xiong et al. (2014). Reflecting their sequence quality, most of these elements have a local combinational value score of >50. For Pack-MULEs, the coordinates of 275 full-length elements from Jiang et al. (2011) were based on the RefGen_v1 genome; hence, we aligned their sequences (BLASTN, E value 1 \times 10⁻¹⁸⁰) on the RefGen_v2 genome requiring 100% identity on the complete length of each element. This approach yielded 251 Pack-MULEs. We note that the official RefGen_v4 TE annotation (B73v4.TE.filtered.gff3) contains 1246 elements of the DNA transposon mutator (DTM) superfamily, some of which are presumed to be Pack-MULEs. A similar number of DTM elements (1300) were identified by TIR_Learner, a new tool for the identification of DNA transposons (Su et al., 2019). This low number suggests that DTM and Pack-MULE elements are not abundant in the maize genome and that our dataset captures a substantial proportion of them. Finally, we downloaded from MASiVEdb (Bousios et al., 2012b) the sequences of 13 833 Sireviruses identified in RefGen_v2 using the MASiVE algorithm (Darzentas et al., 2010). MASiVE is specifically built to detect full-length Sireviruses with single-nucleotide resolution by using highly conserved motifs located in the junctions of the LTRs with the internal domain. We then filtered out elements from all families that overlapped with each other and those with >5 consecutive "N" nucleotides, based on evidence that BLASTN hits between genes and TEs often mapped precisely at the border of these stretches, indicating potential errors during scaffold assembly. To ensure that TEs are physically present in RefGen_v4, we converted their chromosomal coordinates from RefGen_v2 to RefGen_v4 using the Assembly Converter tool (http://www.gramene.org/) and only kept TEs with ≥90% of length converted on the same chromosome as RefGen_v2. We note that ~97% of the TEs that passed this filter had \geq 99% length conversion. Our final TE population consisted of 1035 Helitrons, 238 Pack-MULEs, and 6200 Sireviruses (Supplemental Table 9).

Our input for genes was the RefGen v2 FGS (http://ftp.gramene.org/ maizesequence.org/). We only included evidence-based genes and filtered for TE-related keywords using the annotation files ZmB73_5b_FGS_ info.txt, ZmB73_5b_FGS.gff, ZmB73_5b_WGS_to_FGS.txt, ZmB73_5a_ gene_descriptors.txt, and ZmB73_5a_xref.txt. We also filtered for similarity (BLASTN, E value 1 \times 10⁻²⁰) of the exons to the conserved domains of the reverse transcriptase and integrase genes of LTR retrotransposons using hidden Markov models (PF07727 and PF00665) from Pfam (https://pfam. xfam.org/). To remove genes that were no longer annotated in RefGen_v4, we linked the RefGen_v2 and RefGen_v4 gene IDs using files "updated_ models" in https://download.maizegdb.org/B73_RefGen_v3/ and "maize.v3TOv4.genelDhistory.txt" in http://ftp.gramene.org, and accessed information on function with "Zea_mays.B73_RefGen_v4.43.chr.gff3" in http://ftp.gramene.org. These steps produced a final dataset of 27 056 genes. We finally assessed syntenic relationships with sorghum using data kindly provided to us by Dr. Margaret Woodhouse of MaizeGDB and produced for Springer et al. (2018). Supplemental Table 10 includes the list of donor and free genes together with their RefGen v2 and RefGen_v4 IDs and syntenic relationships with sorghum.

Identification of Capture Events

We first removed all cases of physical overlaps between genes and TEs and then ran a BLASTN search between the exons of the longest transcript of each gene and our TE dataset. We opted for a strict E value cutoff of 1 \times 10^{-40} because we intended to minimize false-positive events. The average capture length was 280nt, with a minimum of 90nt and a maximum of 1932nt. When exons from multiple genes overlapped partially or fully with a TE, we selected the highest BLASTN bit score to define the true donor gene (Hanada et al., 2009; Jiang et al., 2011). If exons from multiple genes had the same bit score, they were all regarded as true donors and kept for downstream analyses. Often, a TE contained multiple independent capture events, defined as non-overlapping areas

within the TE. In total, we identified 6838 such areas across all our TEs. We tested how this number changed after merging areas located in close proximity to each other, with the assumption that they may in reality represent a single capture event that BLASTN failed to identify in its entirety. By allowing a window of 10nt or 50nt, the number only slightly reduced to 6724 and 6379, respectively, suggesting that the majority represent truly independent capture events.

siRNA, Methylation, and Expression Data

For siRNA mapping, we retrieved short read libraries for ear (GSM306487), leaf (GSM1342517), and tassel (GSM448857). We used Trimmomatic (Bolger et al., 2014) to trim adaptor sequences, and FASTX toolkit (http:// hannonlab.cshl.edu/fastx_toolkit/) to remove low-quality nucleotides until reads had ≥3 consecutive nucleotides with a phred Q score >20 at the 3' end. Reads of 21nt, 22nt, and 24nt in length were kept and filtered for tRNAs (http://gtrnadb.ucsc.edu/), miRNAs (http://www.mirbase.org/), and rRNAs and snoRNAs (http://rfam.xfam.org/), and then mapped to the RefGen_v2 genome using BWA with default settings and no mismatches (Li and Durbin, 2010). Both uniquely and multiply mapping siRNAs were considered, and all loci of multiply mapping siRNAs were counted. Using a custom Perl script, we retrieved the number and IDs of all distinct siRNA sequences that mapped to a locus (e.g., captured region within the TE or an exon) to calculate mapping of distinct siRNA sequences per kb as suggested previously (Bousios et al., 2017). This metric collapses the number of reads in a library for a distinct sRNA sequence and it therefore permits the efficient calculation of the diversity and density of different siRNA sequences that map to a locus. Finally, using the siRNA IDs, we were able to identify siRNA crosstalk events.

For DNA methylation analysis, we used previously published BS-seq data from ear (SRA050144) and leaf (SRR850328). Reads were trimmed for quality and adapter sequences with Trimmomatic using default parameters and a minimum read length of 30nt (Bolger et al., 2014). Trimmed reads were mapped to the RefGen_v2 genome using bowtie2 (v2.2.7, parameters: -N 0 -L 20 -p 2) within the Bismark (v0.15.0) software suite (Krueger and Andrews, 2011). We did not allow mismatches and retained only uniquely mapped reads. The number of methylated and unmethylated reads at each cytosine in the genome was calculated using bismark methylation extractor. Positions with >2 reads were retained for further analysis. Bisulfite conversion error rates, or false methylation rates (FMR), were estimated from reads that mapped to the chloroplast genome. A binomial test incorporating the estimated FMR (p < 0.05 after Benjamini-Yekutieli false discovery rate correction) was then used to identify methylated cytosines (Lister et al., 2008). For each locus we retrieved the number of covered and methylated CG, CHG, and CHH sites and calculated methylation levels for each context with ≥10 covered sites by dividing the number of methylated to covered cytosines (Takuno et al., 2016). Methylation BS-seq data for the leaf tissue in sorghum were retrieved from Seymour and Gaut (2020).

Finally, we downloaded gene expression data from the ATLAS Expression database (www.ebi.ac.uk/gxa/) for ear (E-GEOD-50191), leaf (E-MTAB-4342), and various tissues of the maize kernel (E-GEOD-62778). Only genes with >0.1 TPM are included in the ATLAS database, hence we classified all other genes as having no expression.

Statistical Analyses of Donor Genes for siRNA Mapping, Expression, and Methylation

We used a one-sided binomial test to test whether crosstalk siRNAs map to donor genes more often than expected by chance. The number of successes is the number of crosstalk siRNAs, the number of trials is the total number of siRNAs that mapped to the donor gene, and the probability of success is the proportion of the total exonic length that has been captured by all TEs. If we assume a random distribution of siRNA across the donor gene, the expected probability of mapping of any siRNA to the captured area is the length of the captured area divided by total gene length.

Binomial exact test *p* values were corrected for multiple testing using Benjamini and Hochberg (1995).

To study the link between methylation levels of regions of donor genes and the number of crosstalk siRNAs, the Imer function of the R package Ime4 (Bates et al., 2015) was used to write a linear model with mixed effects. The r.squaredGLMM function of the R package MuMIn (Barton, 2009) was used to compute the marginal R^2 (the variance explained by the fixed effects, here the number of crosstalk siRNAs). The proportion of methylated cytosines was log transformed, and the gene was set as a random factor (each gene had one measurement for leaf and one for ear). The analysis was repeated separately for the three methylation contexts and each siRNA length for captured and non-captured regions of donor genes:

log(proportion of methylated cytosines + 1) - number of crosstalk siRNAs + (1|gene)

Age of Capture Events

To estimate the age of gene capture events, we estimated dS between donor genes and the captured fragments within TEs. The RefGen_v2 genome GFF file was used to split sequences into coding and noncoding (since in v2 UTRs are included in the first/last exons). The coding parts of donor genes and captured fragments were aligned using MACSE v2 (Ranwez et al., 2018). In cases where stop codons were found in the captured gene fragment, they were replaced by "NNN" to compute dS using the yn00 program in the paml package (Yang, 2007). To obtain capture age, dS values were divided by 2 × (1.3 × 10^{-8}) (Ma and Bennetzen, 2004). When there were multiple captured fragments in a single TE, we used the oldest capture event (i.e. the maximum dS) as the age estimate.

The Imer function of the R package Ime4 (Bates et al., 2015) was used to write a generalized linear model with mixed effects to study the link between capture age and the number of crosstalk siRNAs. The r.squaredGLMM function of the R package MuMIn (Barton, 2009) was used to compute the marginal R² (the variance explained by the fixed effects, here capture age). The Poisson family was used and gene was set as a random factor:

crosstalk siRNA number - capture age + (1|gene)

Inference of Selective Constraint

Selective constraint was determined for each gene using the dN/dS ratio between maize and sorghum. Orthology for syntelog pairs was extracted from Springer et al. (2018). Orthologs of maize translocated genes in sorghum were inferred by their best blast hit. CDS sequences of translocated genes from the maize RefGen_v2 were blasted against the CDS sequences of sorghum genome reference v3.1.1 (extracted from Phytozome V13) using the NCBI blastn tool version 2.10.1 with an E value cutoff of 1 \times 10⁻⁵. The best hits were selected with lowest E value and highest bit score. Hits with alignment length under 200 bp were discarded. This approach resulted in 1580 out of 2685 maize translocated genes with an identified sorghum ortholog. All maize-sorghum ortholog pairs were aligned with the MACSE v2.03 alignSequences program (Ranwez et al., 2018). Internal stop codons were removed before running the yn00 program from the paml package version 4.9j (Yang, 2007) to infer dS and dN. To eliminate ancient paralogs among our inferred orthologous pairs, genes with maize-sorghum dS \geq 0.9 were excluded. This left 1095 translocated genes and 16 989 syntelogs with a dN/dS value (Supplemental Table 8). We discriminated between constrained (dS/dN > 0.4) and less-constrained genes (dS/dN < 0.4) based on the distribution of dS/dN values across all genic categories (Supplemental Figure 11). We chose this threshold because donor translocated genes had a median dN/dS value of ~0.4, which allowed us to contrast this category into two equal sets of constrained and less-constrained genes.

Gene capture by TEs triggers epigenetic conflict

Statistical Analyses of TEs for siRNA Mapping and Methylation

For siRNA mapping, the glmer function of the R package Ime4 (Bates et al., 2015) was used to write an exponential model with mixed effects to study the effect of TE type (with or without gene capture) and TE age. The r.squaredGLMM function of the R package MuMIn (Barton, 2009) was used to compute the marginal R² (the variance explained by the fixed effects, here TE type and age). The Ismeans function from the R package Ismeans (Length, 2016) was used to compute the contrast between TEs with and without gene capture. The TE was set as a random factor (tissue was the repetition) and the number of siRNAs per kb was log transformed:

log(siRNA per kb+1) - TE type + TE age + (1|TE)

A simpler model was also used:

log(siRNA per kb+1) - TE type + (1|TE)

Similarly, a generalized linear model with mixed effects was used to study the effects of TE type and TE age on TE methylation. The binomial family was used, and TE was set as a random factor (tissue was the repetition). The analysis was repeated separately for the three methylation contexts:

proportion of methylated cytosines - TE type + TE age + (1|TE)

SUPPLEMENTAL INFORMATION

Supplemental Information can be found online at Molecular Plant Online.

FUNDING

A.M. is supported by an EMBO Postdoctoral Fellowship ALTF 775-2017 and by HFSPO Fellowship LT000496/2018-L. D.K.S. is supported by a Postdoctoral Fellowship from the National Science Foundation (NSF) Plant Genome Research Program (1609024). B.S.G. is supported by an NSF grant 1655808. A.B. is supported by The Royal Society (award nos. UF160222 and RGF/R1/180006).

AUTHOR CONTRIBUTIONS

A.B. and B.S.G. designed the research. A.M., N.D., and A.B. performed the research. A.M., D.S., E.P., and A.B. analyzed the data. A.B. and B.S.G. wrote the paper.

ACKNOWLEDGMENTS

No conflict of interest is declared.

Received: July 6, 2020 Revised: October 15, 2020 Accepted: November 5, 2020 Published: November 7, 2020

REFERENCES

Barrett, L.W., Fletcher, S., and Wilton, S.D. (2012). Regulation of eukaryotic gene expression by the untranslated gene regions and other non-coding elements. Cell. Mol. Life Sci. 69:3613–3634.

Barton, K. (2009). MuMIn: multi-model inference. http://r-forge.r-project.org/projects/mumin/.

Bates, D., Machler, M., Bolker, B., and Walker, S. (2015). Fitting linear mixed-effects models using Ime4. J. Stat. Softw.. https://www. jstatsoft.org/article/view/v067i01

Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. J. R. Stat. Soc. Series B. Stat. Methodol. 57:289–300.

Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics **30**:2114–2120.

Bousios, A., and Darzentas, N. (2013). Sirevirus LTR retrotransposons: phylogenetic misconceptions in the plant world. Mobile DNA 4:9.

- Bousios, A., Diez, C.M., Takuno, S., Bystry, V., Darzentas, N., and Gaut, B.S. (2016). A role for palindromic structures in the cis-region of maize Sirevirus LTRs in transposable element evolution and host epigenetic response. Genome Res. 26:226-237.
- Bousios, A., Gaut, B.S., and Darzentas, N. (2017). Considerations and complications of mapping small RNA high-throughput data to transposable elements. Mobile DNA 8:3.
- Bousios, A., Kourmpetis, Y.A.I., Pavlidis, P., Minga, E., Tsaftaris, A., and Darzentas, N. (2012a). The turbulent life of Sirevirus retrotransposons and the evolution of the maize genome: more than ten thousand elements tell the story. Plant J. 69:475-488.
- Bousios, A., Minga, E., Kalitsou, N., Pantermali, M., Tsaballa, A., and Darzentas, N. (2012b). MASiVEdb: the Sirevirus plant retrotransposon database. BMC Genomics 13:158.
- Catoni, M., Jonesman, T., Cerruti, E., and Paszkowski, J. (2019). Mobilization of Pack-CACTA transposons in Arabidopsis suggests the mechanism of gene shuffling. Nucleic Acids Res. 47:1311–1320.
- Cerbin, S., Wai, C.M., VanBuren, R., and Jiang, N. (2019). GingerRoot: a novel DNA transposon encoding integrase-related transposase in plants and animals. Genome Biol. Evol. 11:3181-3193.
- Cho, J. (2018). Transposon-derived non-coding RNAs and their function in plants. Front. Plant Sci. 9:600.
- Cuerda-Gil, D., and Slotkin, R.K. (2016). Non-canonical RNA-directed DNA methylation. Nat. Plants 2:16163.
- Darzentas, N., Bousios, A., Apostolidou, V., and Tsaftaris, A.S. (2010). MASiVE: mapping and analysis of SireVirus elements in plant genome sequences. Bioinformatics 26:2452-2454.
- Diez, C.M., Meca, E., Tenaillon, M.I., and Gaut, B.S. (2014). Three groups of transposable elements with contrasting copy number dynamics and host responses in the maize (Zea mays ssp. mays) genome. PLoS Genet. 10:e1004298.
- Du, C., Fefelova, N., Caronna, J., He, L., and Dooner, H.K. (2009). The polychromatic Helitron landscape of the maize genome. Proc. Natl. Acad. Sci. U S A 106:19916-19921.
- Dvir, S., Velten, L., Sharon, E., Zeevi, D., Carey, L.B., Weinberger, A., and Segal, E. (2013). Deciphering the rules by which 5 '-UTR sequences affect protein expression in yeast. Proc. Natl. Acad. Sci. U S A **110**:E2792-E2801.
- El Baidouri, M., Do Kim, K., Abernathy, B., Li, Y.H., Qiu, L.J., and Jackson, S.A. (2018). Genic C-methylation in soybean is associated with gene paralogs relocated to transposable element-rich pericentromeres. Mol. Plant 11:485-495.
- Gent, J.I., Ellis, N.A., Guo, L., Harkess, A.E., Yao, Y., Zhang, X., and Dawe, R.K. (2013). CHH islands: de novo DNA methylation in neargene chromatin regulation in maize. Genome Res. 23:628-637.
- Grabundzija, I., Messing, S.A., Thomas, J., Cosby, R.L., Bilic, I., Miskey, C., Gogol-Doring, A., Kapitonov, V., Diem, T., Dalda, A., et al. (2016). A Helitron transposon reconstructed from bats reveals a novel mechanism of genome shuffling in eukaryotes. Nat. Commun. **7**:10716.
- Hanada, K., Vallejo, V., Nobuta, K., Slotkin, R.K., Lisch, D., Meyers, B.C., Shiu, S.H., and Jiang, N. (2009). The functional role of Pack-MULEs in rice inferred from purifying selection and expression profile. Plant Cell 21:25-38.
- Holligan, D., Zhang, X.Y., Jiang, N., Pritham, E.J., and Wessler, S.R. (2006). The transposable element landscape of the model legume Lotus japonicus. Genetics 174:2215-2228.
- Hollister, J.D., and Gaut, B.S. (2009). Epigenetic silencing of transposable elements: a trade-off between reduced transposition and deleterious effects on neighboring gene expression. Genome Res. 19:1419-1428.

- Hollister, J.D., Smith, L.M., Guo, Y.-L., Ott, F., Weigel, D., and Gaut, B.S. (2011). 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.
- Jiang, N., Bao, Z.R., Zhang, X.Y., Eddy, S.R., and Wessler, S.R. (2004). Pack-MULE transposable elements mediate gene evolution in plants. Nature 431:569-573.
- Jiang, N., Ferguson, A.A., Slotkin, R.K., and Lisch, D. (2011). Pack-Mutator-like transposable elements (Pack-MULEs) induce directional modification of genes through biased insertion and DNA acquisition. Proc. Natl. Acad. Sci. U S A 108:1537-1542.
- Jiao, Y.P., Peluso, P., Shi, J.H., Liang, T., Stitzer, M.C., Wang, B., Campbell, M.S., Stein, J.C., Wei, X.H., Chin, C.S., et al. (2017). Improved maize reference genome with single-molecule technologies. Nature 546:524.
- Juretic, N., Hoen, D.R., Huynh, M.L., Harrison, P.M., and Bureau, T.E. (2005). The evolutionary fate of MULE-mediated duplications of host gene fragments in rice. Genome Res. 15:1292-1297.
- Krueger, F., and Andrews, S.R. (2011). Bismark: a flexible aligner and methylation caller for Bisulfite-Seq applications. Bioinformatics **27**:1571-1572.
- Lee, Y.C.G., and Karpen, G.H. (2017). Pervasive epigenetic effects of Drosophila euchromatic transposable elements impact their evolution. Elife 6:e25762.
- Length, R.V. (2016). Least-squares means: the R package Ismeans. J. Stat. Softw. 69:1-33.
- Li, H., and Durbin, R. (2010). Fast and accurate long-read alignment with Burrows-Wheeler transform. Bioinformatics 26:589-595.
- Li, Q., Eichten, S.R., Hermanson, P.J., Zaunbrecher, V.M., Song, J.W., Wendt, J., Rosenbaum, H., Madzima, T.F., Sloan, A.E., Huang, J., et al. (2014). Genetic perturbation of the maize methylome. Plant Cell **26**:4602–4616.
- Lisch, D. (2009). Epigenetic regulation of transposable elements in plants. Annu. Rev. Plant Biol. 43-66.
- Lisch, D. (2013). How important are transposons for plant evolution? Nat. Rev. Genet. 14:49-61.
- Lister, R., O'Malley, R.C., Tonti-Filippini, J., Gregory, B.D., Berry, C.C., Millar, A.H., and Ecker, J.R. (2008). Highly integrated single-base resolution maps of the epigenome in Arabidopsis. Cell 133:523-536.
- Lockton, S., and Gaut, B.S. (2009). The contribution of transposable elements to expressed coding sequence in Arabidopsis thaliana. J. Mol. Evol. 68:80-89.
- Lu, L., Chen, J.F., Robb, S.M.C., Okumoto, Y., Stajich, J.E., and Wessler, S.R. (2017). Tracking the genome-wide outcomes of a transposable element burst over decades of amplification. Proc. Natl. Acad. Sci. U S A 114:E10550-E10559.
- Luo, R.S., Bai, C.L., Yang, L., Zheng, Z., Su, G.H., Gao, G.Q., Wei, Z.Y., Zuo, Y.C., and Li, G.P. (2018). DNA methylation subpatterns at distinct regulatory regions in human early embryos. Open Biol. 8:180131.
- Ma, J.X., and Bennetzen, J.L. (2004). Rapid recent growth and divergence of rice nuclear genomes. Proc. Natl. Acad. Sci. U S A 101:12404-12410.
- Matzke, M.A., and Mosher, R.A. (2014). RNA-directed DNA methylation: an epigenetic pathway of increasing complexity. Nat. Rev. Genet. **15**:394–408.
- Maumus, F., and Quesneville, H. (2014). Ancestral repeats have shaped epigenome and genome composition for millions of years in Arabidopsis thaliana. Nat. Commun. 5:4104.
- Morgante, M., Brunner, S., Pea, G., Fengler, K., Zuccolo, A., and Rafalski, A. (2005). Gene duplication and exon shuffling by helitron-

- like transposons generate intraspecies diversity in maize. Nat. Genet. **37**:997–1002.
- Niederhuth, C.E., Bewick, A.J., Ji, L.X., Alabady, M.S., Kim, K.D., Li, Q., Rohr, N.A., Rambani, A., Burke, J.M., Udall, J.A., et al. (2016). Widespread natural variation of DNA methylation within angiosperms. Genome Biol. 17:194.
- Nobuta, K., Lu, C., Shrivastava, R., Pillay, M., De Paoli, E., Accerbi, M., Arteaga-Vazquez, M., Sidorenko, L., Jeong, D.H., Yen, Y., et al. (2008). Distinct size distribution of endogenous siRNAs in maize: evidence from deep sequencing in the mop1-1 mutant. Proc. Natl. Acad. Sci. U S A 105:14958–14963.
- Panda, K., and Slotkin, R.K. (2020). Long-read cDNA sequencing enables a "gene-like" transcript annotation of transposable elements. Plant Cell 32:2687–2698.
- Ranwez, V., Douzery, E.J.P., Cambon, C., Chantret, N., and Delsuc, F. (2018). MACSE v2: toolkit for the alignment of coding sequences accounting for frameshifts and stop codons. Mol. Biol. Evol. 35:2582–2584.
- Schnable, J.C. (2015). Genome evolution in maize: from genomes back to genes. S.S. Merchant, ed. Vol. 66:329–343.
- Schnable, J.C. (2019). Genes and gene models, an important distinction. New Phytol. https://doi.org/10.1111/nph.16011.
- Schnable, P.S., Ware, D., Fulton, R.S., Stein, J.C., Wei, F., Pasternak, S., Liang, C., Zhang, J., Fulton, L., Graves, T.A., et al. (2009). The B73 maize genome: complexity, diversity, and dynamics. Science 326:1112–1115.
- **Seymour, D.K., and Gaut, B.S.** (2020). Phylogenetic shifts in gene body methylation correlate with gene expression and reflect trait conservation. Mol. Biol. Evol. **37**:31–43.
- Sigman, M.J., and Slotkin, R.K. (2016). The first rule of plant transposable element silencing: location, location, location. Plant Cell 28:304–313.
- Springer, N.M., Anderson, S.N., Andorf, C.M., Ahern, K.R., Bai, F., Barad, O., Barbazuk, W.B., Bass, H.W., Baruch, K., Ben-Zvi, G., et al. (2018). The maize W22 genome provides a foundation for functional genomics and transposon biology. Nat. Genet. 50:1282.
- Stitzer, M.C., Anderson, S.N., Springer, N.M., and Ross-Ibarra, J. (2019). The genomic ecosystem of transposable elements in maize. bioRxiv https://doi.org/10.1101/559922.
- Su, W.J., Gu, X., and Peterson, T. (2019). TIR-learner, a new ensemble method for TIR transposable element annotation, provides evidence for abundant new transposable elements in the maize genome. Mol. Plant 12:447–460.
- Sultana, T., Zamborlini, A., Cristofari, G., and Lesage, P. (2017). Integration site selection by retroviruses and transposable elements in eukaryotes. Nat. Rev. Genet. 18:292–308.

- Sun, S.L., Zhou, Y.S., Chen, J., Shi, J.P., Zhao, H.M., Zhao, H.N., Song, W.B., Zhang, M., Cui, Y., Dong, X.M., et al. (2018). Extensive intraspecific gene order and gene structural variations between Mo17 and other maize genomes. Nat. Genet. 50:1289.
- **Takuno, S., and Gaut, B.S.** (2013). Gene body methylation is conserved between plant orthologs and is of evolutionary consequence. Proc. Natl. Acad. Sci. U S A **110**:1797–1802.
- Takuno, S., Ran, J.H., and Gaut, B.S. (2016). Evolutionary patterns of genic DNA methylation vary across land plants. Nat. Plants 2:15222.. https://doi.org/10.1038/nplants.2015.222.
- Thomas, J., and Pritham, E.J. (2015). Helitrons, the eukaryotic rolling-circle transposable elements. Microbiol. Spectr. 3. https://doi.org/10.1128/microbiolspec.MDNA3-0049-2014.
- Wang, J., Yu, Y., Tao, F., Zhang, J.W., Copetti, D., Kudrna, D., Talag, J., Lee, S., Wing, R.A., and Fan, C.Z. (2016). DNA methylation changes facilitated evolution of genes derived from Mutator-like transposable elements. Genome Biol. 17:92.
- Wendte, J.M., Zhang, Y.W., Ji, L.X., Shi, X.L., Hazarika, R.R., Shahryary, Y., Johannes, F., and Schmitz, R.J. (2019). Epimutations are associated with CHROMOMETHYLASE 3-induced de novo DNA methylation. Elife 8:e47891.
- Wicker, T., Buchmann, J.P., and Keller, B. (2010). Patching gaps in plant genomes results in gene movement and erosion of colinearity. Genome Res. 20:1229–1237.
- Xiong, W.W., He, L.M., Lai, J.S., Dooner, H.K., and Du, C.G. (2014). HelitronScanner uncovers a large overlooked cache of Helitron transposons in many plant genomes. Proc. Natl. Acad. Sci. U S A 111:10263–10268.
- Yang, L.X., and Bennetzen, J.L. (2009). Distribution, diversity, evolution, and survival of Helitrons in the maize genome. Proc. Natl. Acad. Sci. U S A 106:19922–19927.
- Yang, Z.H. (2007). Paml 4: phylogenetic analysis by maximum likelihood. Mol. Biol. Evol. 24:1586–1591.
- Zhang, L., Chia, J.-M., Kumari, S., Stein, J.C., Liu, Z., Narechania, A., Maher, C.A., Guill, K., McMullen, M.D., and Ware, D. (2009). A genome-wide characterization of microRNA genes in maize. PLoS Genet. 5:e1000716.
- Zhang, X., Zhao, M., McCarty, D.R., and Lisch, D. (2020). Transposable elements employ distinct integration strategies with respect to transcriptional landscapes in eukaryotic genomes. Nucleic Acids Res. 48:6685–6698.
- Zhao, D.Y., Hamilton, J.P., Vaillancourt, B., Zhang, W.L., Eizenga, G.C., Cui, Y.H., Jiang, J.M., Buell, C.R., and Jiang, N. (2018). The unique epigenetic features of Pack-MULEs and their impact on chromosomal base composition and expression spectrum. Nucleic Acids Res. 46:2380–2397.