

Population-level annotation of lncRNAs in Arabidopsis reveals extensive expression variation associated with transposable element–like silencing

Aleksandra E. Kornienko ^{*}, Viktoria Nizhynska  Almudena Molla Morales , Rahul Pisupati 
and Magnus Nordborg ^{*}

Gregor Mendel Institute, Austrian Academy of Sciences, Vienna Biocenter, Dr. Bohr-gasse 3, Vienna 1030, Austria

^{*}Author for correspondence: aleksandra.kornienko@gmi.oeaw.ac.at (A.E.K.), magnus.nordborg@gmi.oeaw.ac.at (M.N.)

The authors responsible for distribution of materials integral to the findings presented in this article in accordance with the policy described in the Instructions for Authors (<https://academic.oup.com/plcell/pages/General-Instructions>) are: Aleksandra E. Kornienko (aleksandra.kornienko@gmi.oeaw.ac.at) and Magnus Nordborg (magnus.nordborg@gmi.oeaw.ac.at).

Abstract

Long noncoding RNAs (lncRNAs) are understudied and underannotated in plants. In mammals, lncRNA loci are nearly as ubiquitous as protein-coding genes, and their expression is highly variable between individuals of the same species. Using *Arabidopsis thaliana* as a model, we aimed to elucidate the true scope of lncRNA transcription across plants from different regions and study its natural variation. We used transcriptome deep sequencing data sets spanning hundreds of natural accessions and several developmental stages to create a population-wide annotation of lncRNAs, revealing thousands of previously unannotated lncRNA loci. While lncRNA transcription is ubiquitous in the genome, most loci appear to be actively silenced and their expression is extremely variable between natural accessions. This high expression variability is largely caused by the high variability of repressive chromatin levels at lncRNA loci. High variability was particularly common for intergenic lncRNAs (lincRNAs), where pieces of transposable elements (TEs) present in 50% of these lincRNA loci are associated with increased silencing and variation, and such lncRNAs tend to be targeted by the TE silencing machinery. We created a population-wide lncRNA annotation in Arabidopsis and improve our understanding of plant lncRNA genome biology, raising fundamental questions about what causes transcription and silencing across the genome.

Introduction

Long noncoding RNAs (lncRNAs) are a relatively new and still enigmatic class of genes that are increasingly recognized as important regulators participating in nearly every aspect of biology (Statello et al. 2021). There are more lncRNAs than protein-coding genes (PC genes) in the human (*Homo sapiens*) genome (Volders et al. 2019), and they are apparently abundant in the genomes of all eukaryotes (Kapusta and Feschotte 2014; Mattick and Rinn 2015). In human and mouse (*Mus musculus*), lncRNAs have been shown to be involved in various diseases (Wapinski and Chang 2011; Batista and Chang 2013), and medical applications have been

proposed (Wahlestedt 2013). Although many lncRNAs have demonstrated functions, most lncRNAs have not been studied (Leone and Santoro 2016), and many knockouts of seemingly functional candidates showed no phenotypic differences relative to their respective wild types (WTs) (Sauvageau et al. 2013), leading to continuous debate about the functionality and importance of lncRNAs as a gene class (Mattick et al. 2023). Evolutionary studies of lncRNAs have revealed low sequence conservation and highly divergent expression when compared with PC genes (Necsulea and Kaessmann 2014; Nelson et al. 2017), yet some signs of conservation and selection have also been found (Johnsson

Received April 04, 2023. Accepted July 30, 2023. Advance access publication September 8, 2023

© The Author(s) 2023. Published by Oxford University Press on behalf of American Society of Plant Biologists.

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs licence (<https://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial reproduction and distribution of the work, in any medium, provided the original work is not altered or transformed in any way, and that the work is properly cited. For commercial re-use, please contact journals.permissions@oup.com

Open Access

IN A NUTSHELL

Background: Only a small fraction of the genome encodes proteins. We were interested in a special type of gene called long noncoding RNAs (lncRNAs): They are transcribed from the genome but do not encode proteins. lncRNAs can regulate genes or organize cell structures but are largely not studied, and we know very little about lncRNAs as a gene class. For example, we know lncRNAs evolve very quickly and are different between species, but we do not know well how they differ within 1 species and what is responsible for this difference.

Question: We wanted to know how many lncRNAs are present in the model plant *Arabidopsis* (*Arabidopsis thaliana*), how they differ in plants from different regions, and whether the recently reported widespread epigenetic variation in *Arabidopsis* underlies this difference. We used many transcriptome and epigenetic sequencing data sets to answer these questions.

Findings: We discovered that the *Arabidopsis* genome is full of lncRNAs, although most are epigenetically inactivated. Plants from different regions have different sets of active lncRNAs, and epigenetic differences are responsible for much of this difference. Intergenic lncRNAs were particularly variable in their expression levels and contained pieces of transposons, selfish genes that can move and propagate in the genome. Cells fight the spread of transposons with elaborate systems inactivating them and preventing them from harming the genome. We determined that these transposon pieces made lncRNAs look like transposons and become inactivated by the same system.

Next steps: It is unclear what underlies the epigenetic variation causing lncRNA variation. Is it a difference in sequence or the absence of the whole lncRNA gene from the genomes of *Arabidopsis* from certain regions? What contributes the most? Another direction is to understand the nature and origin of transposon pieces inside lncRNAs.

et al. 2014; Mattick et al. 2023). Several studies have looked at how lncRNAs differ between closely related species such as rat (*Rattus norvegicus*) and mouse (Kutter et al. 2012), human and chimp (*Pan troglodytes*) (Necsulea and Kaessmann 2014), or different plant species (Nelson et al. 2017; Zhu et al. 2022), but few have looked at differences within the same species (Melé et al. 2015). It was recently shown that lncRNAs display salient interindividual expression variation in human (Kornienko et al. 2016) and mouse (Andergassen et al. 2017), much higher than that of PC genes, but the meaning, causes, and consequences of this high variability are unknown.

Arabidopsis (*Arabidopsis thaliana*) has higher natural genetic variability than humans (1001 Genomes Consortium 2016) and represents an interesting and convenient model for studying lncRNA variation. As most research on lncRNAs has been performed in human and mouse (Rinn and Chang 2020), relatively little is known about lncRNAs in plants (Liu et al. 2015; Budak et al. 2020). Several studies have identified and annotated lncRNAs in plant species such as *Arabidopsis* (Liu et al. 2012; Palos et al. 2022), wheat (*Triticum aestivum*) (Xin et al. 2011), maize (*Zea mays*) (Li et al. 2014), and strawberry (*Fragaria × vesca*) (Kang and Liu 2015), but, although several databases have been created, the number and comprehensiveness of plant lncRNA annotations are often poorer than those of human and mouse (Szcześniak et al. 2019; Jin et al. 2021; Di Marsico et al. 2022; Zhu et al. 2022). Nonetheless, it is clear that lncRNAs do regulate genes in plants (Liu et al. 2015; Whittaker and Dean 2017; Chen et al. 2023; Gullotta et al. 2023) and that lncRNA expression is particularly responsive to stress and environmental factors (Wang et al. 2017; Budak et al. 2020). Furthermore, this

lncRNA response can be accession-specific to a much greater extent than that of PC genes (Blein et al. 2020). Plant lncRNAs also affect significant crop traits, and their relevance for food security has been highlighted (Gullotta et al. 2023). Understanding the real scope of lncRNA transcription in plants could help identify new candidates for functional studies and shed light on the genome biology of lncRNAs in plants and beyond (Palos et al. 2023).

While many lncRNAs have been shown to participate in epigenetic silencing or activation of PC genes (Statello et al. 2021), much less research exists on the epigenetic regulation of lncRNAs themselves (Yang et al. 2023). In *Arabidopsis*, the epigenetic patterns of some functional lncRNAs have been thoroughly studied (Whittaker and Dean 2017; Yang et al. 2023), but little is known about the epigenetics of lncRNAs on a genome-wide scale. While high epigenetic variation was reported between accessions (Kawakatsu et al. 2016), it is not clear how this variation affects lncRNAs.

lncRNAs are known to sometimes originate from transposable elements (TEs) (Kapusta et al. 2013; Palos et al. 2022; Zhu et al. 2022), yet what implications this origin has for their expression, epigenetics, and variation is not well known. Similarly, while aberrant lncRNA copy number has been connected to disease and other phenotypes in human (*H. sapiens*) (Athie et al. 2020; Xu et al. 2020), general information about lncRNA copy number and its consequences is missing, in particular in plants.

In this study, we aimed to study the extent and natural variability of lncRNA transcription in *Arabidopsis*. We annotated lncRNAs using data from 499 accessions, finding thousands of new lncRNA loci and generating an extended lncRNA annotation. We observed high expression and epigenetic variability

for lncRNAs among accessions, with lncRNAs being generally silenced in any given accession. Epigenetic variability explains expression variation of many lncRNAs. Long intergenic ncRNAs (lincRNAs) showed particularly high variability and can be divided into protein coding–like and TE-like loci that show differences in their epigenetic patterns, copy number, and—most importantly—the presence of pieces of TE sequences. Indeed, such short pieces of TEs were prevalent in intergenic lncRNAs, likely attracting TE-like silencing to these loci. We provide new insights into the biology of lncRNAs in plants, identify a major role for TE-likeness in lncRNA silencing, and provide an extensive annotation and data resource for the Arabidopsis community.

Results

Transcriptome annotation from hundreds of accessions reveals thousands of previously unannotated lncRNAs

To investigate the extent of lncRNA transcription in Arabidopsis, we used newly generated and publicly available (Kawakatsu et al. 2016; Cortijo et al. 2019) polyA⁺ stranded transcriptome deep sequencing (RNA-seq) data sets spanning 5 different tissues/developmental stages (seedlings, 9-leaf rosettes, leaves from 14-leaf rosettes, flowers, and pollen) and 499 accessions (Fig. 1A; see Supplemental Data Set 1 for accession list, Supplemental Data Set 2 for RNA-seq samples, and Supplemental Data Set 3 for RNA-seq mapping statistics). To create a cumulative transcriptome annotation, we mapped the RNA-seq data from all samples onto the TAIR10 genome, assembled transcriptomes from each accession/tissue separately (Supplemental Data Set 4), and then used a series of merging and filtering steps to generate 1 cumulative annotation, which we then classified into several gene classes (Fig. 1B; Materials and methods; Supplemental Fig. S1). We used Araport11 and TAIR10 gene annotations (Cheng et al. 2017) to guide the classification of transcripts corresponding to PC genes, pseudogenes, TE genes and TE fragments, ribosomal RNA (rRNA), and transfer RNA (tRNA) loci and used an additional protein-coding potential filtering step to identify a set of lncRNAs (Supplemental Figs. S1 and S2A). Our transcriptome annotation performed well in assembling known PC genes and known lncRNAs (Supplemental Fig. S2B).

In total, we identified 23,676 PC and 11,295 lncRNA loci, the latter thus representing almost one third (29%) of the cumulative transcriptome annotation (Fig. 1C). The resulting annotation was highly enriched in lncRNAs (Supplemental Fig. S2C) with 10,315 lncRNA loci (91%) being absent from the current public lncRNA annotation by Araport11. Our annotation extended the lncRNA portion of the reference genome from the 2.2% annotated in Araport11 to 10.7%, covering ~13 Mb in total sequence (Supplemental Fig. S2D). Comparison to the recent large-scale lncRNA identification studies in Arabidopsis (Zhao et al. 2018; Kindgren et al. 2020; Ivanov et al. 2021; Corona-Gomez et al. 2022; Palos et al. 2022) that, like Araport 11, were mainly based on the

laboratory accession Columbia 0 (Col-0) showed that we identified 5,954 (53%) new lncRNA loci in the TAIR10 genome. We were also able to detect and annotate many TE genes and TE fragments across accessions/tissues (Fig. 1C), finding spliced isoforms for 579 TE genes previously annotated as single-exon (Supplemental Fig. S2, E to H).

We classified lncRNAs based on their genomic position (Fig. 1D). The largest group (8,195, or 72%) was antisense (AS) lncRNAs that overlapped annotated PC genes in the AS direction (Supplemental Fig. S3) (Ietswaart et al. 2012). We observed that 8,083 Araport11 PC genes have an AS RNA partner, over 5 times more than in the Araport11 reference annotation. Previous studies have reported ubiquitous, unstable AS transcription (Li et al. 2013; Yuan et al. 2015) as well as the activation of AS transcription upon stress in Arabidopsis and other plants (Zhao et al. 2018; Xu et al. 2021); however, our data contained exclusively polyA⁺ RNA-seq data sets of seedlings and plants grown under normal conditions, so we can conclude that relatively stable polyadenylated AS transcripts can be produced over almost a third of PC genes in Arabidopsis across different accessions and tissues.

The second largest class with 2,246 loci (20%) was lincRNAs (Fig. 1, D and E). The third largest class (630, or 6%) consisted of lncRNAs that were AS to TE genes (AS-to-TE lncRNAs). The remaining 3 classes constituted <3% of all lncRNA loci and we will ignore them below.

The genomic distribution of AS lncRNAs and AS-to-TE lncRNAs mirrored the annotation of PC and TE genes, respectively, with the former being enriched in chromosome arms and the latter near centromeres. LincRNAs were also enriched near centromeres but were distributed across the genome (Supplemental Fig. S4).

Analyzing more accessions and tissues reveals more lncRNA loci

We hypothesized that a major reason for our discovering so many previously nonannotated lncRNA loci was that our annotation was based on hundreds of accessions, while most previous studies had only used the reference accession Col-0 (Yuan et al. 2015, 2016; Cheng et al. 2017). If lncRNAs are very variably expressed between individuals, as has been shown in humans (Kornienko et al. 2016), data from a single accession would uncover only the subset expressed in that particular accession. To test this idea, we subsampled the unified rosette RNA-seq data set from the 1001 Genomes Project (Kawakatsu et al. 2016) and ran our annotation pipeline many times (Materials and methods). This saturation analysis showed that the number of annotated lncRNA loci increases 2.5 times by raising the number of accessions analyzed from 10 to 460 (Fig. 1F). Unlike PC genes, the number of lncRNAs strongly depended on the sample size and showed no sign of saturating even with 460 accessions (Supplemental Fig. S5, A and B).

To confirm that the observed increase was not simply due to increased sequencing coverage, we compared these results with very high-coverage RNA-seq data from Col-0 only

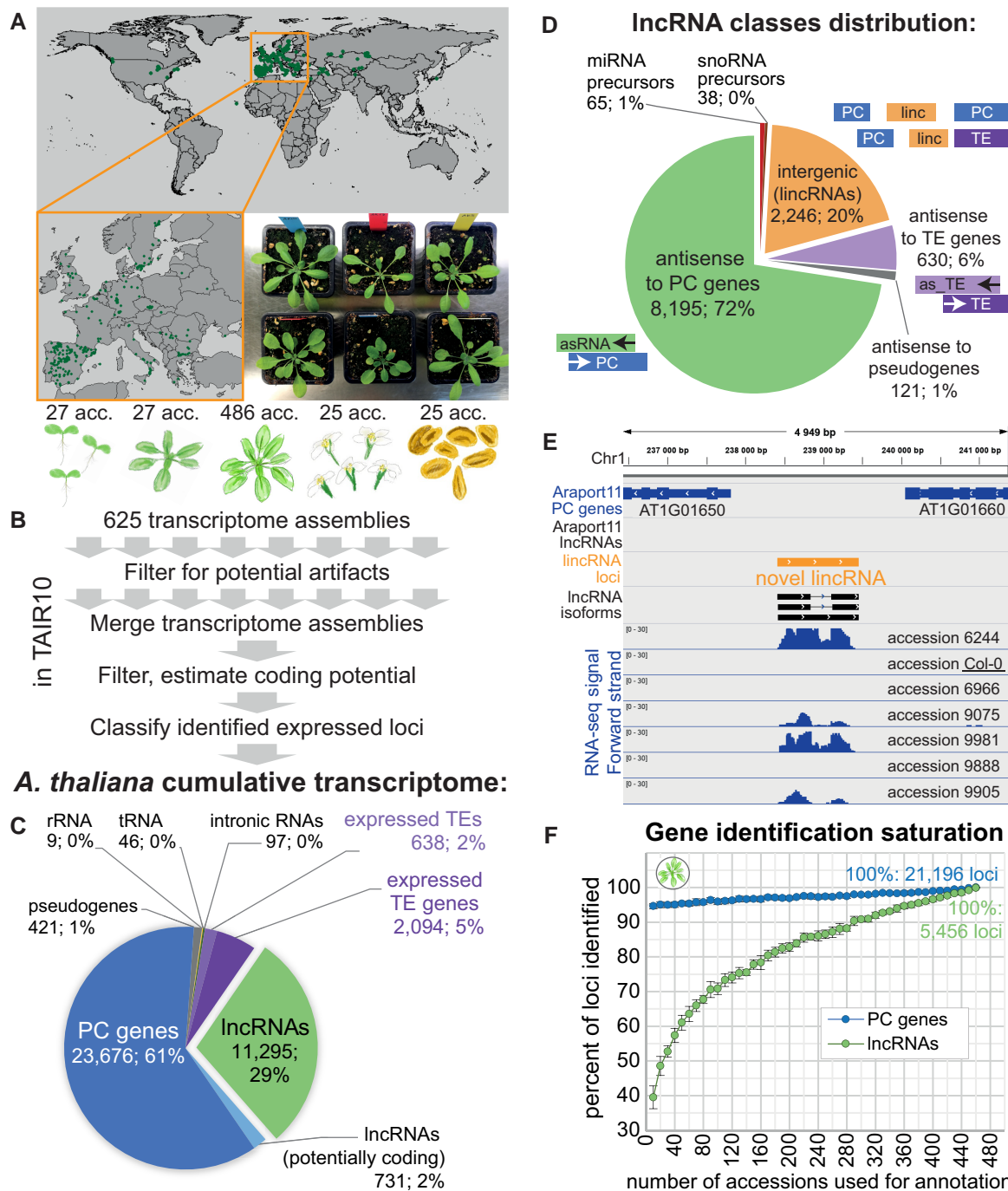


Figure 1. Mapping lincRNA transcription in hundreds of accessions and several tissues reveals thousands of previously unannotated lincRNAs. **A)** Origins of the *Arabidopsis* accessions used for transcriptome annotation and an example photograph of 6 different accessions grown in the growth chamber. **B)** Overview of the pipeline used for cumulative transcriptome annotation. Tissues from left to right: 7-d-old seedlings, 9-leaf rosettes, 14-leaf rosettes, flowers, and pollen. **C)** Distribution of the types of loci in the cumulative annotation. **D)** Distribution of lincRNA positional classes in the genome. **E)** An example of a previously unannotated intergenic lincRNA on chromosome 1. Expression in 7 different *Arabidopsis* accessions is shown. **F)** Number of lincRNA and PC loci identified as a function of the number of accessions used, relative to the number identified using 460 accessions. Random subsampling of accessions was performed in 8 replicates, and the error bars indicate the SD across replicates.

(Cortijo et al. 2019). Applying the same subsampling strategy to these data, we observed a much slower increase that saturated early and could not possibly explain the results in Fig. 1F (Supplemental Fig. S5C).

lincRNAs can be specific to certain tissues and developmental stages (Cabili et al. 2011; Liu et al. 2012; Palos et al. 2022), so it is likely that our use of multiple tissues helped identify more loci. Our final annotation, which was based

on seedlings, rosettes, flowers, and pollen from multiple accessions (Fig. 1B), revealed 11,265 lncRNA loci, while the 460-accession rosette analysis above returned only 5,456 lncRNA loci (Fig. 1F). To better understand the effect of adding different tissues, we performed another saturation analysis where we varied both the number of tissues and accessions (Supplemental Fig. S6). While the number of loci always increased with more accessions, the number of tissues used mattered even more. In particular, adding flowers or pollen to the analysis produced a big jump in the number of genes identified. For example, using data from 20 accessions and 4 tissues allowed the identification of ~3 times more lncRNAs than when using just data from seedlings (Supplemental Fig. S6). Adding flowers alone nearly doubled the number of lncRNAs identified.

To summarize, by combining RNA-seq data from hundreds of accessions and 4 developmental stages, we provide a massively expanded lncRNA annotation for Arabidopsis, identifying thousands of previously missed loci. Our extended lncRNA annotation is available in the supplement (Supplemental Data Set 5).

High lncRNA expression variability between accessions

We showed above that including more accessions allowed the identification of more lncRNA loci (Fig. 1F) and hypothesized that the reason was that not every lncRNA was expressed in every accession. Indeed, this appears to be the case: while most PC genes were expressed in nearly all accessions, most lncRNAs, as well as most TE genes and fragments, were expressed in fewer than 5% of all accessions (Fig. 2A) (throughout this article, we use “expressed” as in “detected” to refer to loci with transcripts per million mapped reads [TPM] > 0.5 in a given data set, aware of our inability to detect unstable or nonpolyadenylated transcription using polyA⁺ data). Our analysis of the expression frequency (ON/OFF state) of the 4 main types of loci showed that while about 50% of PC loci are expressed in every accession, the same was true for no more than 1% of all AS lncRNAs, lincRNAs, and TE genes (Supplemental Fig. S7).

To quantify the natural variability in expression of lncRNAs and other gene types, we calculated the coefficient of variance using rosette RNA-seq data across 461 accessions (Kawakatsu et al. 2016). Similarly to human (Kornienko et al. 2016) and mouse (Andergassen et al. 2017), the expression of both AS lncRNAs and lincRNAs was significantly more variable than that of PC genes (Fig. 2B). In particular, lincRNAs showed expression variability almost to the level of TE genes and fragments. The variability in expression of lncRNAs that were AS to TE genes was similar to that of TE genes and fragments (Fig. 2B). Analyzing expression variability in other rosette RNA-seq data sets (Supplemental Fig. S8, A and B) and other tissues (Supplemental Fig. S8, C to E) confirmed these results.

Two factors crucially affect expression variability values and must be controlled for when comparing lncRNAs with

PC genes: gene length and absolute expression level. lncRNAs are known to be shorter and have lower expression than PC genes (Cabili et al. 2011), which we confirmed in our data (Supplemental Fig. S9, A and B). Both gene length and absolute expression level were negatively correlated with the coefficient of variance, while this observation held true for all gene types, the anticorrelation slopes were different (Supplemental Fig. S9, C and D). When we controlled for expression level or gene length alone, the trend shown in Fig. 2B was preserved (Supplemental Fig. S10, A and B). When we controlled for both expression and gene length, the trend was preserved for lincRNAs, while AS lncRNAs were similar to PC genes (Supplemental Fig. S9, E and F), which might be explained by particularly high variability in the expression of short PC genes (Cortijo et al. 2019).

As the 14-leaf rosette data set produced in this study contained 2 to 4 repeats for each accession, we could assess the level of intraaccession expression variation. For all classes of genes except AS lncRNAs, the intraaccession expression variation was significantly lower than the interaccession variation; the difference between the classes of genes mirrored interaccession variation (Supplemental Fig. S10C). This result suggests that compared with PC genes, the expression of lincRNAs and TEs is more unstable and prone to be affected by the precise conditions or noise, while much of the AS lncRNA expression variation between accessions might be defined by generally unstable expression. To estimate the true noise in lncRNA expression, we analyzed the RNA-seq data consisting of Arabidopsis seedlings collected every 2 h over 24 h, with 14 technical replicates per time point (Cortijo et al. 2019). We determined that lncRNAs have significantly noisier expression (Fig. 2C) as well as higher circadian expression variability (Supplemental Fig. S10D). Interestingly, while TE genes showed higher variability between accessions (Fig. 2B), both lincRNAs and AS lncRNAs were noisier than TE genes (Fig. 2C).

To illustrate the extent of lncRNA expression variation, we plotted expression across 4 tissues in 3 accessions as a heatmap for different types of genes (Fig. 2D). While PC gene expression levels clustered the samples according to tissue, lincRNA and TE expression clustered the samples according to accession (AS lncRNA was similar to PC gene but noisier). While pollen samples always clustered separately due to the particular transcriptome of pollen (Slotkin et al. 2009), the expression of lincRNAs and TEs was strikingly different between accessions. It was also notable that particularly many lincRNAs and TE genes were expressed in pollen, while flowers appeared to have higher expression for all 4 gene types. In general, Fig. 2D illustrates how few lncRNAs are expressed in each accession and how striking the interaccession variation is. Two randomly chosen accessions effectively express the same PC genes, whereas they share only about half of their expressed lncRNAs (Fig. 2E). Furthermore, while ~70% of PC genes were expressed in both seedlings and rosettes of any given accession, only 7% of AS lncRNAs and 4% of lincRNAs were expressed in these tissues (Fig. 2F).

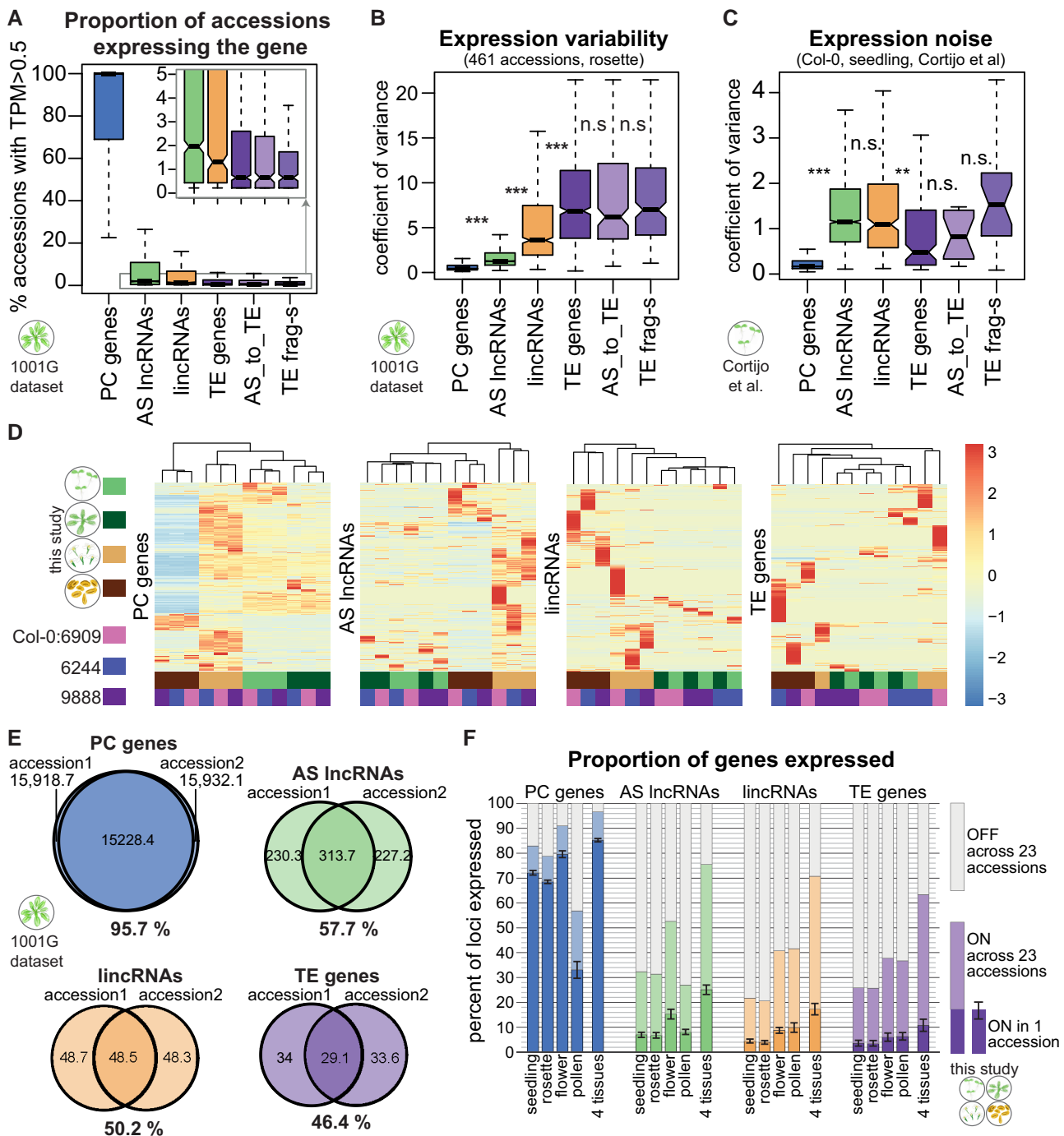


Figure 2. lncRNAs display extensive variability in their expression across accessions and appear to be largely silent. **A)** Proportion of accessions in the 1001 Genomes data set (Kawakatsu et al. 2016) where the indicated type of gene is expressed (TPM > 0.5). Only genes that are expressed in at least 1 accession are plotted. **B)** Coefficient of variance of expression in 461 accessions from the 1001 Genomes data set (Kawakatsu et al. 2016). Only genes with TPM > 1 in at least 1 accession are plotted. **C)** Expression noise, calculated from 14 technical replicates of Col-0 seedlings; expression noise values averaged across 12 samples are shown (Cortijo et al. 2019). Only genes with TPM > 1 in at least 1 sample are plotted. Boxplots: outliers are not shown, and *P*-values were calculated using a Mann–Whitney test on equalized sample sizes: ****P* < 10⁻¹⁰, ***P* < 10⁻⁵, **P* < 0.01, n.s. *P* > 0.01. **D)** Gene expression levels for different types of genes in 4 tissues for the reference accession Col-0 (accession number 6909) and 2 randomly picked accessions. Heatmaps were built using “pheatmap” in R with scaling by row. Only genes expressed in at least 1 sample are plotted. Clustering trees for rows not shown. **E)** Average number of genes expressed in an accession and its randomly selected partner accession from the 1001 Genomes data set and the number of genes expressed (TPM > 0.5) in both accessions. Percentages indicate the extent of overlap between accessions. **F)** Proportion of genes expressed in 1 accession in seedlings, 9-leaf rosettes, flowers, pollen, or all 4 tissues combined (lower part of the bars). The error bars show SD across the 23 accessions. The light part of the bars indicates the additional proportion of genes that can be detected as expressed when all 23 accessions are considered.

Strikingly, almost twice as many AS lncRNA loci were expressed in flowers, but not in pollen, while twice as many lincRNAs were expressed in both flowers and pollen (Fig. 2F). Like lincRNAs, TE genes showed increased expression in flowers and pollen, in agreement with a previous report (Slotkin et al. 2009). Across all 4 tissues in 23 accessions, 96% of PC genes, 76% of AS lncRNAs, 71% of lincRNAs, and 63% of TE genes were expressed (Fig. 2F), thus covering many more lncRNAs, in line with the saturation analysis (Fig. 1F) and the increased individual and tissue specificity of lncRNAs (Fig. 2, B and D).

In summary, lncRNA expression shows high variability between accessions, between tissues, and also between replicates. LincRNAs differ from AS lncRNAs in that they show higher expression variability and increased expression in pollen, but both classes are predominantly silent in any given sample.

The epigenetic landscape of lncRNA loci suggests ubiquitous silencing

To characterize the epigenetic patterns of lncRNAs in *Arabidopsis* and investigate their apparently ubiquitous silencing, we performed chromatin immunoprecipitation followed by deep sequencing (ChIP-seq) and bisulfite sequencing in multiple accessions using leaves from rosettes at the 14-leaf stage (Materials and methods; Supplemental Fig. S11A and Data Sets 6 and 7). For the ChIP experiments, we chose 2 active marks (H3K4me3 and H3K36me3) and 3 repressive marks associated with different types of silencing (histone H1, H3K9me2, and H3K27me3) (Supplemental Fig. S11B). Histone H1 was shown to be involved in silencing TEs, but also AS transcription (Choi et al. 2020), and H3K9me2 is a common heterochromatin mark and is known to silence TEs (Zemach et al. 2013), while H3K27me3 is commonly associated with polycomb repressive complex 2 (PRC2)-mediated silencing and is mostly deposited on PC genes (Feng and Jacobsen 2011). H3K27me3 can also present found on some TEs when the normal silencing machinery is inactivated (Déléris et al. 2021; Zhao et al. 2022).

We first analyzed the ChIP-seq data focusing on the reference accession Col-0. The gene body profiles of the different histone modifications were distinct for the different types of genes (Fig. 3A; Supplemental Fig. S11C). For example, while AS lncRNAs and PC genes showed similar levels of chromatin modifications—which is expected given their overlapping positions—their profiles differed, with PC genes showing a characteristic drop in H1 and H3K9me2 levels at their transcription start site (TSS) and an increase toward the transcription end site (TES), whereas AS lncRNAs showed an even distribution across the gene body (Fig. 3A). LincRNAs and TE genes showed increased levels for the heterochromatic marks H3K9me2 and H1 and lower levels for the active marks H3K36me3 and H3K4me3 (Fig. 3A). Calculating normalized and replicate-averaged ChIP-seq coverage over the entire locus (Fig. 3B; Supplemental Fig. S11D) and the

promoter region (Supplemental Fig. S11E) confirmed the above observations. Overall, AS lncRNAs were similar to PC genes, while lincRNAs were intermediate between PC genes and TE genes in their heterochromatic marks and the lowest in their active marks (Fig. 3B; Supplemental Fig. S11, D and E).

We then analyzed the bisulfite sequencing data, quantifying DNA methylation in 3 different contexts: CG, CHG, and CHH, where H stands for A, C, or T (Materials and methods). CHG and CHH methylation are common in plants and are involved in TE silencing (Fultz et al. 2015). While PC genes and AS lncRNAs displayed low levels of CG methylation and no CHG or CHH methylation, as expected, lincRNAs exhibited a very significant methylation increase in all 3 contexts (Fig. 3C; Supplemental Fig. S12, A and B). Interestingly, the distribution of CG methylation over lincRNAs was bimodal (Fig. 3C, right; Supplemental Fig. S12, C and D), with some loci looking like PC genes, while others looked like TE genes.

As TE genes and lincRNAs are enriched next to the centromeres while PC genes and AS RNAs are not (Supplemental Fig. S4), we checked if the observed epigenetic differences held true when controlling for chromosomal position. While pericentromeric genes within 2 Mb of the centromeres all showed more heterochromatic patterns, the observed trends for histone marks and DNA methylation held true, especially for the genes further than 2 Mb from the centromeres (Supplemental Figs. S13 and S14).

To confirm that the repressive marks at lncRNA loci are associated with silencing, we checked for epigenetic differences between expressed and silent genes (the same samples were used for ChIP-seq, bisulfite sequencing, and RNA-seq; Fig. 3D). The repressive marks H3K9me2 (Fig. 3E; Supplemental Fig. S15A) and H1 (Supplemental Fig. S15, B and C) were significantly more abundant on silent genes. While this result was true for all gene categories, the H3K9me2 difference for TE genes was particularly high, underscoring the fact that TEs are normally silenced by H3K9me2 deposition (Feng and Jacobsen 2011). Another repressive mark, H3K27me3, also showed significantly higher levels on silent genes of all categories, but here, PC genes showed a striking increase, while TE genes were minimally different (Fig. 3F; Supplemental Fig. S15D). We also found that silent AS lncRNAs show increased CG methylation and that silent lincRNAs have strikingly increased CG and CHH methylation levels, although less so than TE genes (Fig. 3G; Supplemental Fig. S16). Expressed PC genes had higher CG gene body methylation levels than silent ones, which is a known phenomenon of as yet unclear function (Bewick and Schmitz 2017).

Since lincRNAs showed both CHH methylation and H3K9me2, both characteristic of TEs and absent from PC genes (Fultz et al. 2015), we performed small RNA (sRNA) sequencing of flowers (Fig. 3H; Supplemental Data Set 8) to look for evidence of targeting by the 24-nucleotide (nt) sRNAs that are normally involved in TE silencing by the RNA-directed DNA methylation (RdDM) pathway

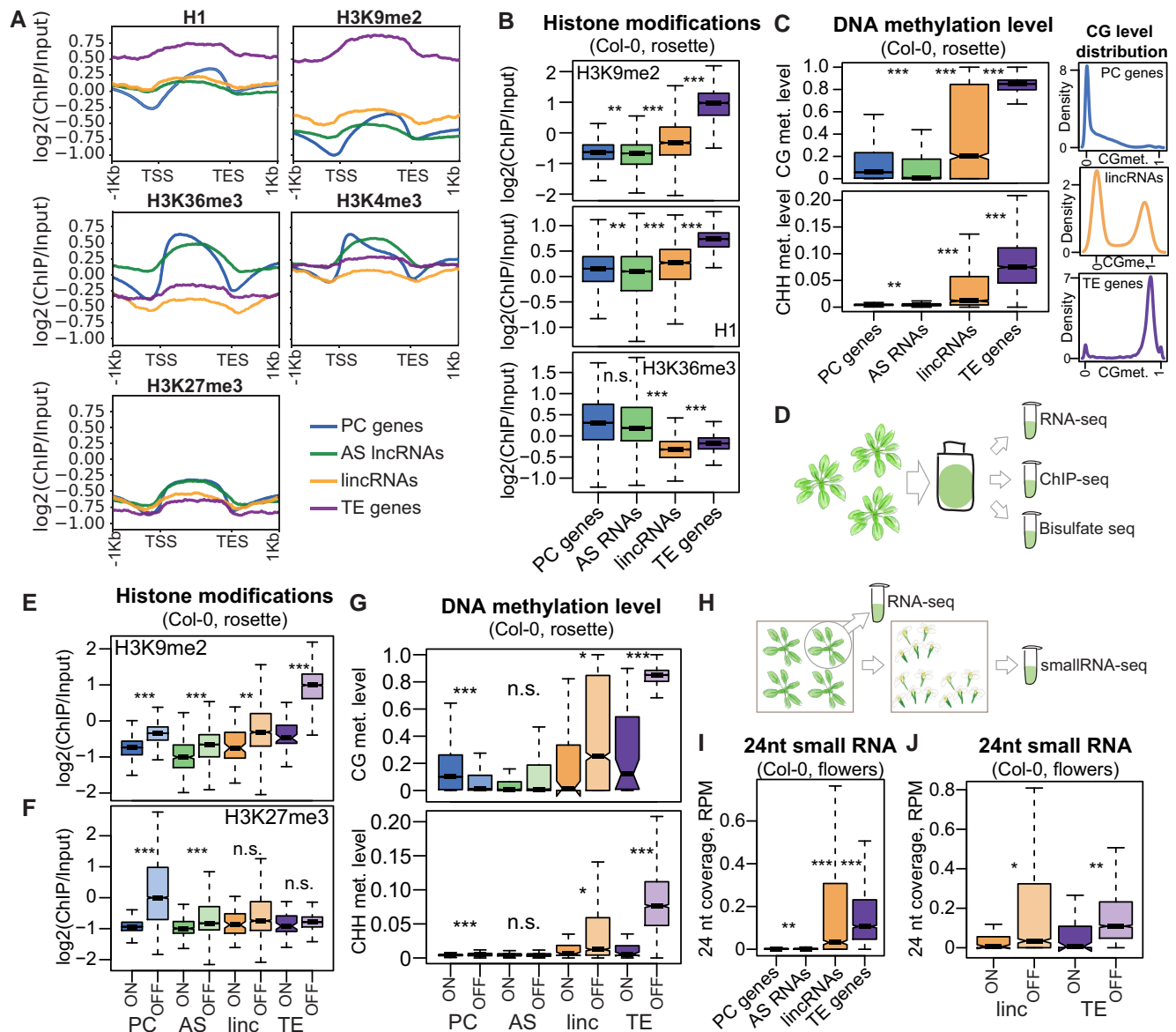


Figure 3. Epigenetic patterns of lincRNAs in Arabidopsis indicate their ubiquitous silencing. **A**) Averaged profiles of the input-normalized ChIP-seq signal for the epigenetic marks histone H1, H3K9me2, H3K36me3, H3K4me3, and H3K27me3 over 4 gene types from our cumulative transcriptome annotation. The plots show data from Col-0 rosettes, replicate 2. All genes, expressed and silent in Col-0, are used for the analysis. Metaplot profiles were built using plotProfile from deeptools (Ramírez et al. 2016). **B**) H3K9me2, H1, and H3K36me3 histone modifications in Col-0 rosette. The log₂ of the gene body coverage normalized by input and averaged between the 2 replicates is plotted. **C**) Left: CG and CHH DNA methylation levels in Col-0 rosettes. Right: density of CG methylation levels for PC genes, lincRNA loci, and TE genes. Methylation level is calculated as the ratio between the number of methylated and unmethylated reads over all Cs in the respective context (CG and CHH) in the gene body and averaged over 4 replicates. **D**) Diagram of the experiment: the same tissue from 14-leaf rosettes was used for RNA-seq, ChIP-seq, and bisulfite-seq in this study. **E**) H3K9me2 input-normalized coverage, plotted separately for expressed (ON, TPM > 0.5) and silent (OFF, TPM < 0.5) genes. **F**) H3K27me3 normalized coverage, plotted separately for expressed (ON, TPM > 0.5) and silent (OFF, TPM < 0.5) genes. **G**) Methylation levels for expressed (ON, TPM > 0.5) and silent (OFF, TPM < 0.5) genes. Expression was calculated in the corresponding 14-leaf rosette samples. **H**) Diagram of the experiment: 1 9-leaf rosette individual from the batch was used for RNA-seq, and flowers for small RNA-Seq were collected from the remaining individuals at a later point. **I**) Coverage of 24-nt small RNAs in the gene body, calculated as the number of 24-nt reads mapping to the locus and divided by the total number of reads and locus length. **J**) Coverage of 21–22-nt small RNA in Col-0 flowers, plotted separately for expressed (ON, TPM > 0.5) and silent (OFF, TPM < 0.5) genes. Expression was calculated in the corresponding 9-leaf rosette samples. *P*-values were calculated using a Mann–Whitney test on equalized sample sizes: ****P* < 10^{−10}, ***P* < 10^{−5}, **P* < 0.01, n.s. *P* > 0.01. Outliers in the boxplots are not shown.

(Matzke and Mosher 2014). This analysis demonstrated that sRNAs indeed target many lincRNA loci (Fig. 3I) (1,131 [50.4%] loci with RPM > 0.03) and were associated with silencing for both lincRNAs and TE genes (Fig. 3J). Analyzing published sRNA-seq data (Papareddy et al. 2020) from leaves and a very early embryonic stage known as “early heart” where RdDM-mediated silencing is particularly active (Papareddy et al. 2020) confirmed targeting of lincRNAs by 24-nt sRNAs, with samples at the early heart showing very high levels of sRNAs (Supplemental Fig. S17, A and B). An analysis of shorter 21–22-nt sRNAs, reported to also participate in TE silencing (Pontier et al. 2012), showed that lincRNAs also show higher levels of targeting by this type of siRNAs (Supplemental Fig. S17, C to E); however, there was no clear association between the levels of these shorter siRNAs and the lack of expression of their corresponding lincRNA (Supplemental Fig. S17F).

Epigenetic variation explains expression variation of many lncRNAs

In the previous section, we described the epigenetic patterns only in the reference accession, Col-0. As we produced ChIP-seq, bisulfite-seq, and sRNA-seq data for several accessions (Supplemental Data Sets 6 to 8), we were able to confirm that the epigenetic patterns we observed in Col-0 were similar in other accessions (Supplemental Fig. S18). However, while the overall patterns were similar, the variability between accessions at particular loci was very high for both lncRNA (especially lincRNA) and TE genes (Fig. 4, A and B; Supplemental Fig. S19).

To test whether epigenetic variation can explain this variation in expression, we analyzed methylation patterns and expression data from rosettes collected from 444 accessions (Kawakatsu et al. 2016). We determined that for 454 lincRNAs and 509 AS lncRNAs, expression across accessions is indeed explained by the level of CG or CHH methylation at their gene body or promoter (Fig. 4C; Supplemental Fig. S20A; Materials and methods). While these numbers correspond to only 20.2% and 6.2% of all lincRNAs and AS lncRNAs, respectively, this analysis could only be performed on a limited number of informative loci with sufficiently high methylation variation and expression frequency (Supplemental Fig. S20B; Materials and methods). Among these informative loci, we could explain expression variation by variation in DNA methylation for 50.7% of lincRNAs and 21.5% of AS lncRNAs.

An example of such an lncRNA with high variation between accessions is displayed in Fig. 4D: the accession that expresses the lincRNA lacks CG and CHH methylation in the locus as well as 24-nt sRNAs, while the accession where the lincRNA is silent has both CG and CHH methylation and 24-nt sRNAs in flowers. The epigenetic variation at this locus was extensive and quite binary with very strong association

between the presence of methylation and the lack of expression and vice versa (Fig. 4E). For the accessions with available data, the repressive histone modifications H1 (Supplemental Fig. S20C) and H3K9me2 (Fig. 4F), as well as 24-nt sRNA coverage (Fig. 4G), were also anticorrelated with expression across accessions.

In summary, we establish that lncRNAs display distinctive epigenetic patterns, consistent with the above observation of the lack of expression and suggesting ubiquitous silencing. Compared with PC genes, lncRNAs display increased epigenetic variation between accessions that explained the variation in expression of ~51% of lincRNAs and ~22% of AS lncRNAs. Many lincRNAs show TE-like epigenetic status that is associated with silencing, as well as being targeted by 24-nt siRNAs, which are characteristic of the RdDM pathway for silencing TEs. Because of the interesting patterns and the outstanding variation observed for lincRNAs, we focused on them below.

lincRNAs are enriched for TE pieces

Several similarities between lincRNAs and TE genes were apparent. First, both showed increased expression in flowers and pollen (Fig. 2F). Second, lincRNA expression was dramatically more variable than that of PC genes and AS lncRNAs—almost at the level of TE genes (Fig. 2A). Third, our survey of the epigenetic landscape showed that lincRNAs display TE-like characteristics, although to a lesser extent (Fig. 3, A, B, and F). Fourth, similarly to expression variation, levels of repressive chromatin at lincRNA loci were more variable than at PC genes and AS lncRNAs, trending toward the pattern seen in TE genes (Fig. 4, A and B). Furthermore, lncRNAs can originate from TEs and contain parts of their sequences, both in plants and animals (Kapusta et al. 2013; Corona-Gomez et al. 2022; Palos et al. 2022); moreover, TE domains within lncRNAs can play significant roles in lncRNA biology (Johnson and Guigó 2014), such as their nuclear export or retention (Lubelsky and Ulitsky 2018), or even have a crucial role in their function (Colognori et al. 2020). Accordingly, we asked if TE sequences contributed to lincRNA loci in Arabidopsis and affected their expression, epigenetics, and variability.

To this end, we used a BLAST-based analysis to identify sequences similar to TAIR10-annotated TEs inside loci and their borders (Fig. 5A; Materials and methods). We called each match a “TE piece” and merged overlapping same-direction TE pieces. We further refer to each TE-like region of a locus as a “TE patch,” no matter whether it was constituted by a single TE piece or several merged ones (Fig. 5A). lincRNA loci were clearly enriched in TE patches compared with AS lncRNAs and PC genes, as well as randomly picked intergenic regions of corresponding length (Fig. 5B; Materials and methods). We determined that 52% of lincRNAs but only 27% of

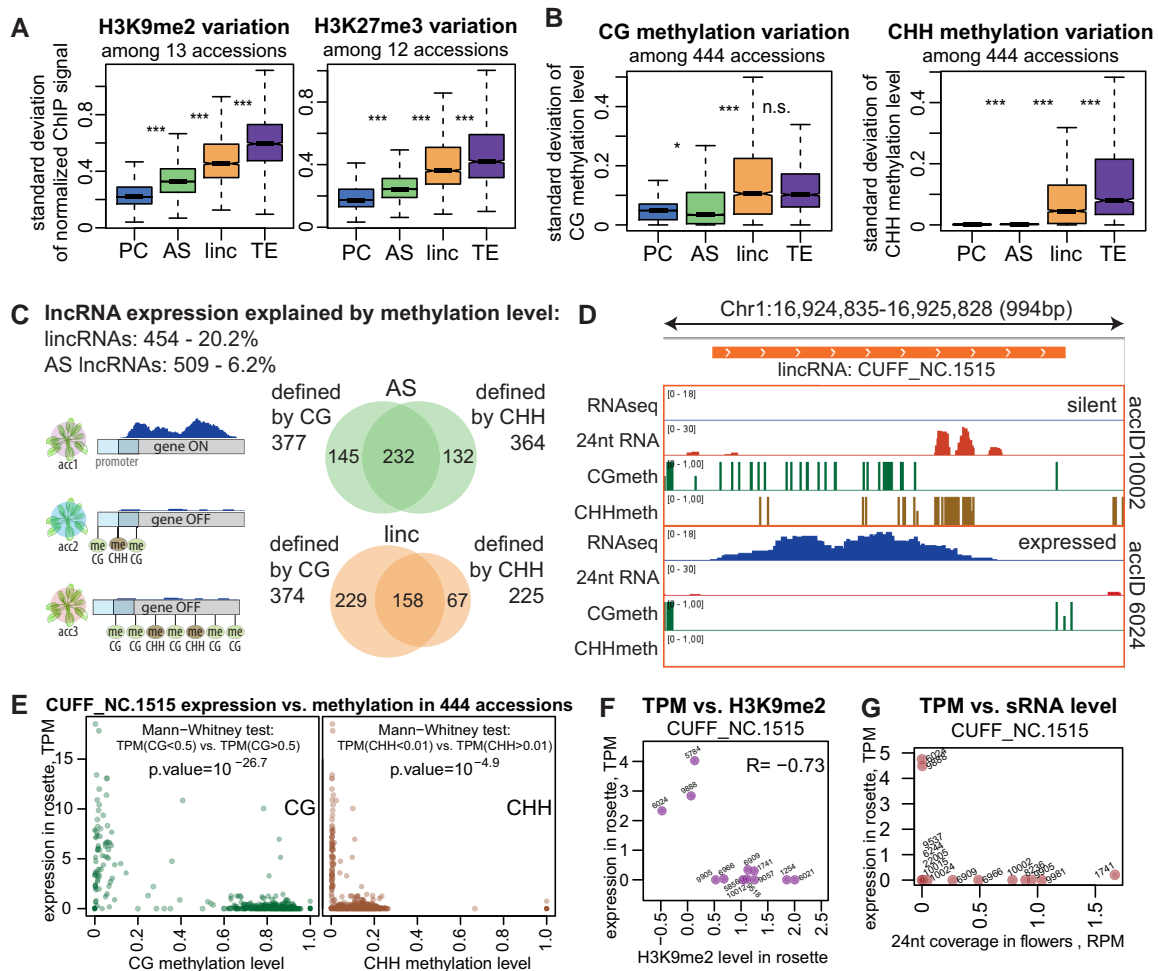


Figure 4. lncRNAs display increased epigenetic variation that explains the variation in expression of many lncRNAs. **A**) sd of input and quantile-normalized coverage (see Materials and methods) of H3K9me2 (left) and H3K27me3 (right) signals in rosettes across 13 or 12 accessions, respectively. **B**) sd of CG (left) and CHH (right) methylation levels across 444 accessions (rosettes, 1001 Genomes data set (Kawakatsu et al. 2016)). *P*-values were calculated using Mann–Whitney test on equalized sample sizes. ****P* < 10⁻¹⁰, ***P* < 10⁻⁵, **P* < 0.01, n.s. *P* > 0.01. Outliers in the boxplots are not shown. **C**) Summary of lncRNAs for which expression can be explained by methylation (Supplemental Fig. S20B). The Venn diagrams show the overlap between loci for AS lncRNAs (green) and lincRNAs (orange) that were found to be defined by CG or CHH methylation levels. **D**) An example of a lincRNA defined by CG and CHH methylation, showing RNA-seq signal (forward strand), CG and CHH methylation levels in rosettes, and the 24-nt sRNA signal in flowers in 2 accessions. **E**) Expression level as a function of CG (left) or CHH (right) methylation for the example lincRNA across 444 accessions (Kawakatsu et al. 2016). The results of the Mann–Whitney tests used for defining the explanatory power of CG/CHH methylation are shown. **F**) Expression in rosettes as a function of H3K9me2 level in rosettes of the example lincRNA in 13 accessions. **G**) Expression in rosettes as a function of 24-nt sRNA coverage in flowers of the example lincRNA in 14 accessions.

matching intergenic controls contain a TE patch. We also observed an enrichment for TE patches in upstream and downstream lincRNA border regions compared with matching controls (Fig. 5B). On a per kb basis, lincRNA borders showed the highest density of TE patches, even higher than for lincRNA loci, and the difference between lincRNAs and other gene types became even more prominent (Fig. 5C). TE genes had fewer TE patches per 1 kb than lincRNAs, presumably because TE genes usually contain 1 large TE patch corresponding to the full TE, while lincRNAs contained several smaller patches (Supplemental Fig. S21A).

It is important to note that lincRNAs are not simply expressed TEs. Our lincRNA annotation pipeline required that

an lncRNA did not overlap with any TE genes and allowed for a maximum of 60% same-strand exonic overlap with annotated TE fragments (Supplemental Fig. S1). While only 486 (21.6%) of all lincRNA loci had a same-strand exonic overlap with a TAIR10-annotated TE fragment, 1,176 (52.4%) contained a TE patch (Fig. 5B), both sense and AS to the lincRNA direction (Fig. 5A; Supplemental Fig. S21B). In addition, 136 (12%) of TE-containing lincRNA loci fully (>90%) overlapped with annotated TE fragments but were transcribed in the direction AS to those (Supplemental Fig. S21, C and D). While 893 (76%) of the TE patch-containing lincRNA loci overlapped with an annotated TE fragment (Supplemental Fig. S21, E and F), 472 in the sense and 605

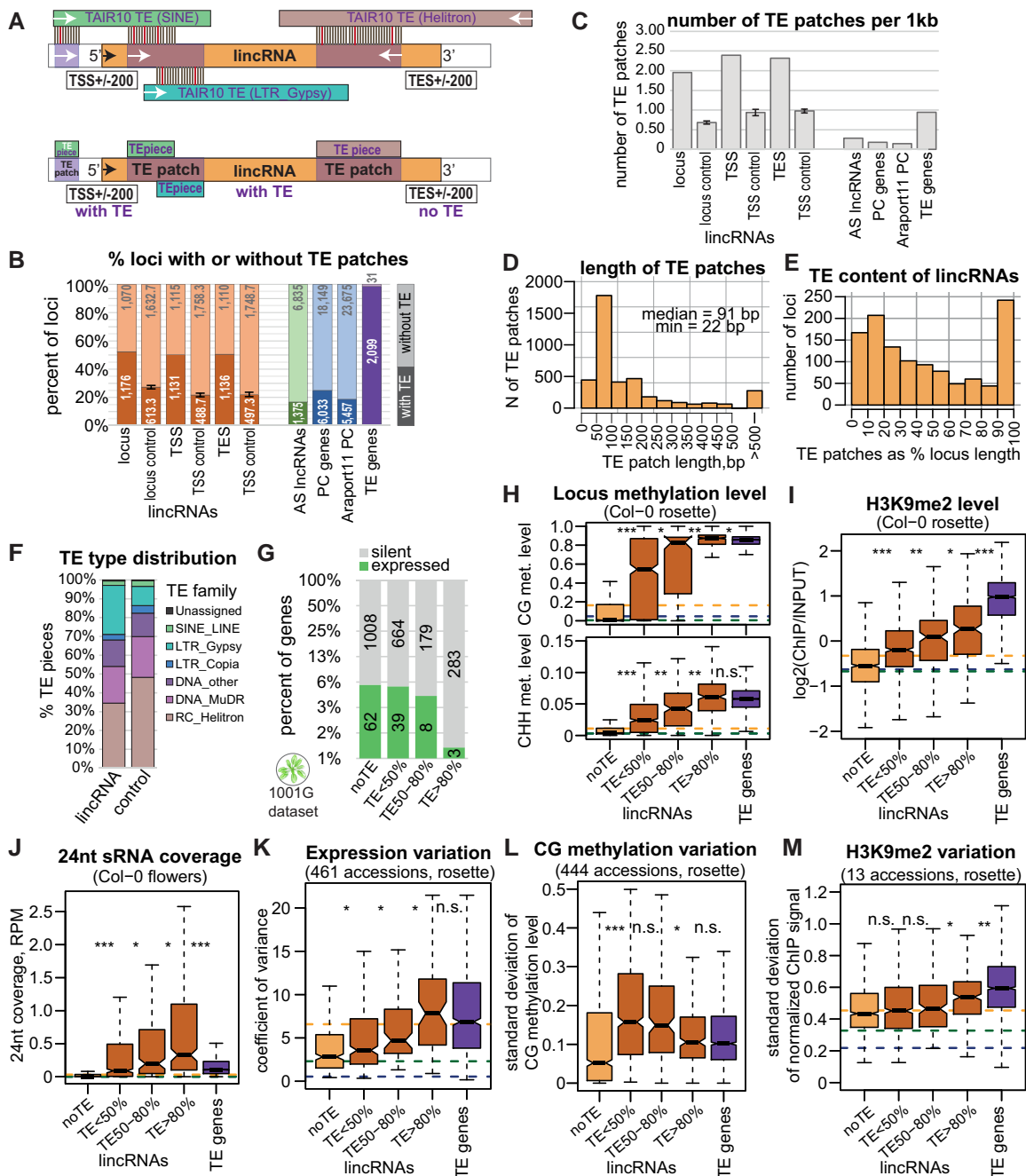


Figure 5. Many lincRNAs contain pieces of TEs that affect their silencing and expression variation. **A)** Outline of TE content analysis. Top: TAIR10-annotated TEs were compared with the sequences of lincRNAs (and other loci) by BLAST. Bottom: the mapped pieces of different TEs overlapping in the same direction were merged into “TE patches.” The upstream and downstream “borders” of genes were analyzed in the same way. **B)** Proportion of loci containing a TE piece. The intergenic controls for lincRNAs, lincRNA TSS \pm 200 bp, and TES \pm 200 bp were obtained by shuffling the corresponding loci within intergenic regions (lincRNAs excluded) 3 times and averaging the results. The error bars on controls represent the SD between the 3 shuffling replicates. **C)** Number of TE patches per 1 kb. **D)** Distribution of the length of TE patches (any relative direction) within lincRNA loci. **E)** TE content distribution among lincRNAs. TE patches were considered in any relative direction. The loci with large TE content are those where the TE patches map AS to the lincRNA locus. **F)** Proportion of TE pieces of different TE families within different types of loci. **G)** Proportion of expressed lincRNAs as a function of their TE content. The y axis is displayed in log₂ scale. **H)** to **J)** Levels of methylation (**H**), H3K9me2 (**I**), and 24-nt sRNAs (**J**) for lincRNA loci as a function of their TE content, with TE genes for comparison. CG and CHH methylation data displayed are from Col-0 rosettes (Kawakatsu et al. 2016). **K)** to **M)** Expression variability between 461 accessions (Kawakatsu et al. 2016) (**K**), SD of CG methylation levels across 444 accessions (Kawakatsu et al. 2016) (**L**), and SD of quantile- and input-normalized H3K9me2 levels in rosettes across 13 accessions (**M**) of lincRNA loci as a function of their TE content, with TE genes for comparison. *P*-values were calculated using Mann-Whitney tests: ****p* < 10⁻¹⁰, ***p* < 10⁻⁵, **p* < 0.01, n.s. *P* > 0.01. Outliers in the boxplots are not shown.

in the AS direction, 283 (24%) did not overlap with any annotated TEs but did contain a TE patch (Supplemental Fig. S21, G and H), that is, contained pieces of sequences bearing resemblance (see Materials and methods) to TE sequences. TE patches within lincRNAs (and other genes) were generally short with a median length of 91 bp and a minimal length of 22 bp (Fig. 5D), thus much shorter than TAIR10-annotated TE fragments or the TE patches that our analysis identified in TE genes (Supplemental Fig. S21I). Moreover, 74% of lincRNAs with multiple TE pieces contained pieces of TEs from different families and 24% pieces of TEs from both Class I and Class II (Supplemental Fig. S21, F, G, and H). The protein-coding potentials of lincRNAs and TE genes were also very different (Supplemental Fig. S2A). Collectively, these results suggest that lincRNAs can be considered as a gene category separate from expressed TE fragments.

The relative TE content of TE-containing lincRNA loci differed greatly, with a few lincRNAs fully covered by TE patches, corresponding to lincRNAs that are AS to a TE fragment (Fig. 5E; Supplemental Fig. S22A). On average, lincRNAs contained 309 bp of same-strand TE-like sequences and 436 bp of AS TE-like sequences per kb. LincRNA loci were particularly enriched in TE sequences from the long terminal repeat (LTR) type TE *Gypsy* compared with matching intergenic controls and other gene types (Fig. 5F; Supplemental Fig. S22B), and this enrichment was particularly pronounced in lincRNAs with AS TE patches (Supplemental Fig. S22, C and D). We did not, however, observe any particular localization for TE pieces from different families within the lincRNA loci (Supplemental Fig. S22, E and F).

The TE content of lincRNAs affects their expression and epigenetics

We asked if the TE content of a lincRNA affected its expression, epigenetic characteristics, and their variability. Indeed, when binned based on relative TE sequence content, lincRNAs with higher TE content were less often expressed (Fig. 5G) and showed higher levels of CG and CHH methylation (Fig. 5H), H3K9me2 (Fig. 5I), and 24-nt siRNAs (Fig. 5J). The expression variation (Fig. 5K; Supplemental Fig. S23) and epigenetic variation (Fig. 5, L and M; Supplemental Fig. S24) of lincRNAs also depended on their TE content, but not as strongly.

LincRNAs are enriched in the pericentromeric regions (Supplemental Fig. S4) that are naturally enriched in TEs and heterochromatin, which might confound our TE piece (Fig. 5, B and C) and epigenetic analyses (Fig. 5, H to J). Controlling for the proximity to centromeres, we first discovered that while all gene types have higher TE content closer to centromeres, the increased TE content of lincRNAs observed in Fig. 5B was preserved (Supplemental Fig. S25A). Second, while all pericentromeric lincRNAs, even those without TE patches, showed high repressive chromatin, the level of heterochromatic marks at lincRNA loci further from

centromeres strongly depended on their TE content (Supplemental Fig. S25, B to D). Furthermore, while 24-nt siRNA coverage was generally low near centromeres (consistent with previous findings, see (Sigman and Slotkin 2016)), it strongly depended on the TE content in chromosome arms (Supplemental Fig. S25E). Thus, the presence and the relative size of TE sequences inside lincRNA loci are indeed associated with a more repressive chromatin state irrespective of chromosomal location.

In summary, we showed that intergenic lincRNAs are highly enriched for short pieces of TEs. About half of all lincRNAs have a TE sequence within them, and higher TE content is associated with more repressive epigenetic marks when comparing different lincRNAs in the genome.

Copy number of lincRNAs affects their expression variability and epigenetic patterns

Apart from expression variability, epigenetic patterns, and TE sequence content, another classical TE feature was evident for lincRNAs: lincRNAs were often present in multiple copies (Supplemental Fig. S26A). We decided to investigate this pattern further and see whether it affects their epigenetic patterns and expression.

We used a BLAST-based approach to look for multiple gene copies in TAIR10 (Materials and methods) and found that lincRNAs are much more commonly multiplied than PC genes and AS lincRNAs, with 28% being present in more than 1 copy and 8% in more than 10 copies (Fig. 6A). Again, lincRNAs were intermediate between PC genes and TE genes. We split lincRNAs into 4 categories: single- or multicopy lincRNAs, with or without TE patches (Fig. 6B, top). Similarly to the overall lincRNA distribution (Fig. 5B), about half of all single-copy lincRNAs contained a TE patch, while most multicopy lincRNAs did (Fig. 6B). LincRNAs with higher copy numbers also showed higher TE sequence content (Supplemental Fig. S26B). Although *Helitrons* have the highest copy number in the Arabidopsis genome (Quesneville 2020), lincRNA with pieces of *Gypsy* elements showed the highest copy number (Supplemental Fig. S26C), even when the TE sequence content was no more than 20% of the locus (Supplemental Fig. S26D).

We analyzed the features of all 4 categories of lincRNAs and observed that increased copy number is associated with lower expression (Fig. 6C) and increased repressive chromatin marks (Fig. 6, D to F; Supplemental Fig. S27, A to F), as well as expression and epigenetic variability (Fig. 6, G to I; Supplemental Fig. S27, G to I). The presence of a TE patch within multicopy lincRNA loci was associated with strikingly increased CG and CHH methylation levels (Fig. 6D) and targeting by 24-nt siRNAs (Fig. 6F) but did not appear to affect the level of H3K9me2 or H1 or their variability (Fig. 6, E and I; Supplemental Fig. S27, D and I).

In summary, we show that many lincRNAs are present in multiple copies and that increased copy number is associated with increased silencing and variability in expression and epigenetic marks. This effect comes in addition to the effect of

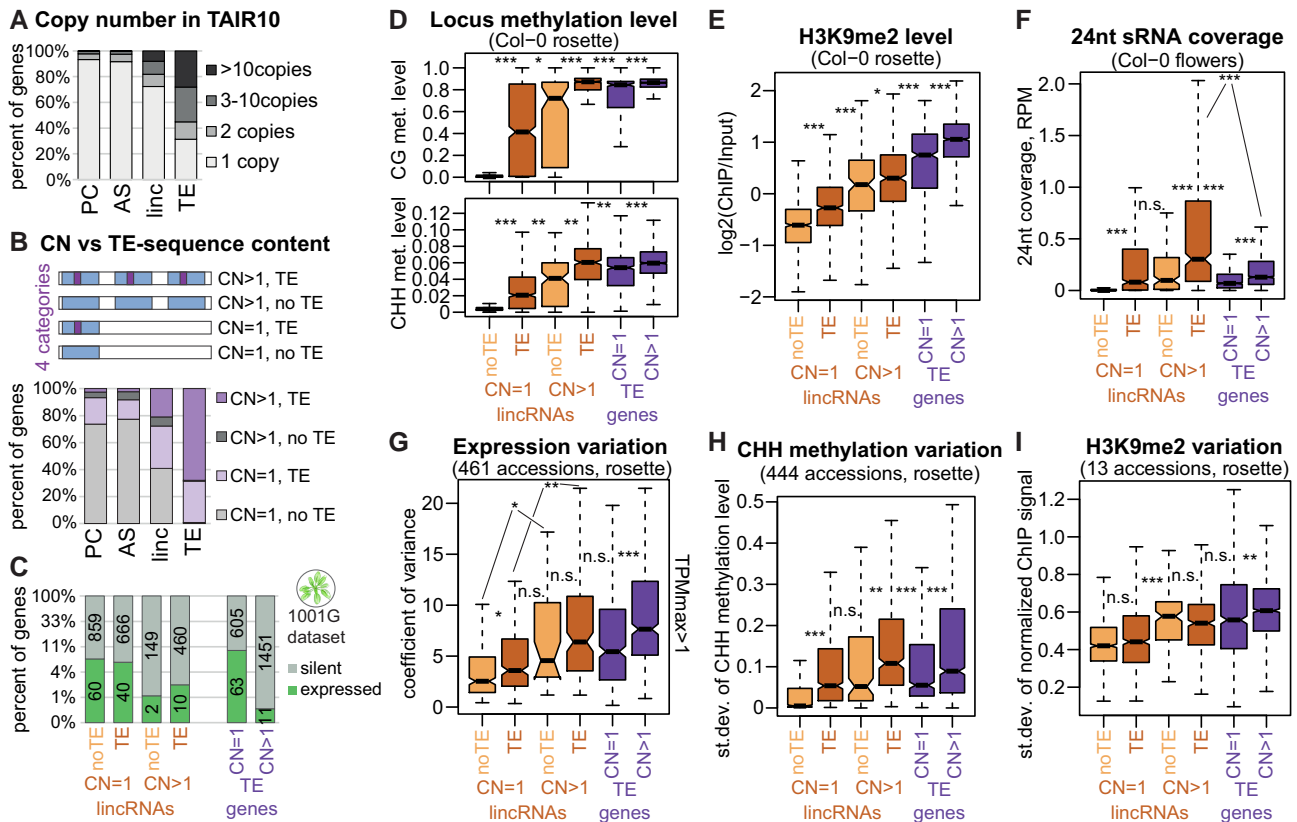


Figure 6. Copy number of lincRNAs affects their epigenetic patterns and variability. **A**) Distribution in copy number for PC genes, AS lincRNAs, lincRNAs, and TE genes from the cumulative transcriptome annotation in the TAIR10 genome. **B**) Top, diagram of the 4 types of loci; bottom, the distribution of copy number of the 4 types of loci: 1 copy with no TE patch, 1 copy with a TE patch, multiple copies with no TE patch in the original locus, and multiple copies with a TE patch in the original locus. **C**) Proportion of the 4 types of lincRNAs and 2 types of TE genes: expressed (TPM > 0.5, green) or silent (TPM < 0.5, gray) in Col-0 rosettes. The y axis is displayed in log3 scale. **D**) to **F**) CG and CHH methylation levels (**D**), H3K9me2 levels (**E**) in Col-0 rosettes, and 24-nt sRNA coverage in Col-0 flowers (**F**) for the 4 types of lincRNAs and 2 types of TE genes. CG and CHH methylation data displayed are from Col-0 rosettes (Kawakatsu et al. 2016). **G**) to **I**) Boxplots showing expression (**G**), CG methylation (**H**), and H3K9me2 (**I**) variability for the 4 types of lincRNAs and the 2 types of TE genes. *P*-values in the boxplots are calculated using a Mann–Whitney test: ****P* < 10^{−10}, ***P* < 10^{−5}, **P* < 0.01, n.s. *P* > 0.01. Outliers in the boxplots are not plotted.

TE patches, in that lincRNAs with both multiple copies and TE patches (i.e. most TE-like) show the highest level of silencing.

lincRNAs are silenced by TE-like and PC-like mechanisms

We saw that lincRNAs are ubiquitously silenced, with very few lincRNAs being expressed in any particular accession (Fig. 2F) and with very few accessions expressing any particular lincRNA (Fig. 2A). We also observed that TE pieces within lincRNAs were associated with heterochromatin and siRNA targeting (Fig. 4, F to H), at least when comparing lincRNA loci within a single genome. We investigated these patterns in greater detail, connecting them to known silencing pathways.

First, we observed that silent lincRNAs show a binary behavior when it comes to which silencing mark—H3K9me2 or H3K27me3—covers the locus (Fig. 7A). We observed the same pattern across all gene types, but whereas almost all TE genes showed H3K9me2 silencing and most PC genes

and AS lincRNAs were covered with H3K27me3, lincRNAs were split into 2 large categories (Fig. 7B; Supplemental Fig. S28). We thus defined 2 nonoverlapping classes of lincRNAs based on which epigenetic mark they presented: H3K9me2 lincRNAs and H3K27me3 lincRNAs or K9 lincRNAs and K27 lincRNAs for short (Fig. 7A). K27 lincRNAs were almost free of TE patches, were present as 1 copy, and showed low DNA methylation and targeting by 24-nt siRNAs, while K9 lincRNAs tended to have higher TE content, were present as multiple copies, and showed strikingly more DNA methylation and targeting by sRNA (Fig. 7, C to G). We concluded that K27 lincRNAs and K9 lincRNAs are PC-like and TE-like, respectively. The bimodality in the epigenetic features we observed before (Fig. 3) can thus be explained by lincRNAs being a heterogeneous group of PC-like and TE-like lincRNAs with different features.

PC-like lincRNAs are likely silenced by PRC2, which establishes the H3K27me3 repressive mark (Hansen et al. 2008); however, we did not try to confirm this hypothesis.

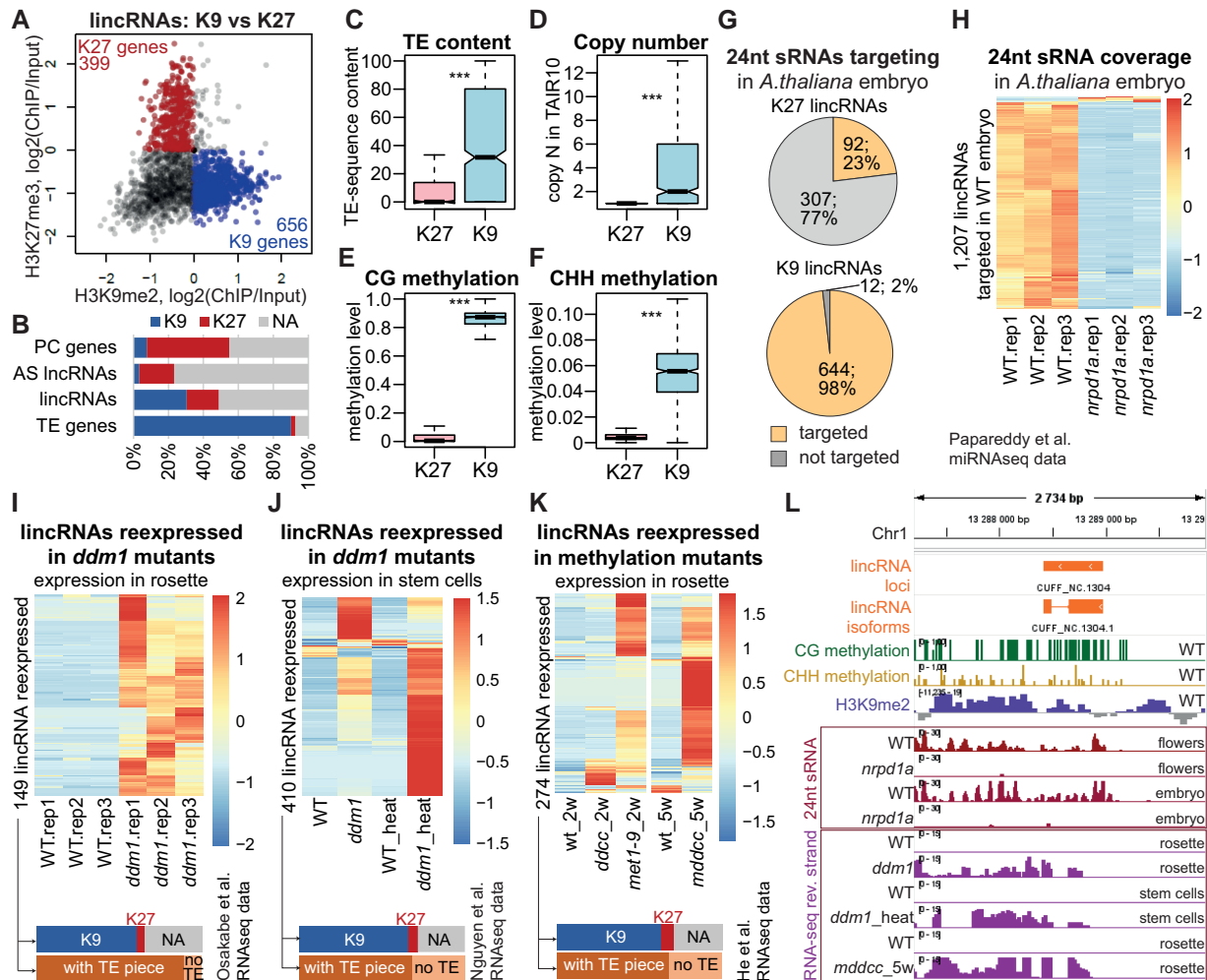


Figure 7. lincRNAs are silenced by PC-like and TE-like mechanisms. **A)** H3K27me3 levels as a function of H3K9me2 levels over lincRNA loci in Col-0 14-leaf rosettes (average of 2 replicates). K27 genes, red, K27 signal > 0, K9 signal < 0; K9 genes, blue, K27 signal < 0, K9 signal > 0. **B)** Proportion of K9 (blue) and K27 (red) genes among the 4 gene types. NA, genes with neither mark (gray, K27 signal < 0, K9 signal < 0). **C) to F)** Boxplots showing the relative TE sequence content (**C**), copy number (**D**), CG methylation level (**E**), and CHH methylation level (**F**) of lincRNA loci classified as K27 or K9 genes. Outliers not plotted. *P*-values were calculated using Mann–Whitney tests: ****P* < 10^{−10}. CG and CHH methylation data displayed are from Col-0 rosettes (Kawakatsu et al. 2016). **G)** Distribution of K27 and K9 lincRNAs targeted by 24-nt sRNAs (reads per million [RPM] > 0.03) in Arabidopsis embryos (“early heart” stage) (Papareddy et al. 2020). The sRNA coverage was averaged across 3 replicates. **H)** 24-nt sRNA coverage in Arabidopsis embryos (“early heart” stage) in the WT (Col-0) and in Pol IV-deficient mutants (*nrpd1a*, Col-0 background) (Papareddy et al. 2020). The 1,207 lincRNAs that are targeted (RPM > 0.03, average of 3 replicates) by 24-nt sRNAs in the WT are plotted. **I)** Expression level of the 149 lincRNAs reexpressed in rosettes of the *ddm1* mutant in the Col-0 background (Osakabe et al. 2021). The bars at the bottom show the distribution of K9 (blue), K27 (red), and TE-containing (dark orange) or TE-free (light orange) loci among the reexpressed lincRNAs (same for **J** and **K**). **J)** Expression level of the 410 lincRNAs reexpressed in shoot stem cells of the *ddm1* mutant in the Col-0 background under mock conditions or with heat stress treatment (Nguyen et al. 2023). **K)** Expression level of lincRNAs reexpressed in the rosettes of DNA methylase mutants *met1-9*, *ddcc*, and *mdcc* (all in the Col-0 background) (He et al. 2022) (see Materials and methods). Heatmaps were built using “pheatmap” in R with scaling by row. No column clustering and row clustering dendrograms not displayed. **L)** An example of a lincRNA epigenetically silenced in Col-0 WT but expressed in the silencing mutants.

Instead, we focused on the TE-like lincRNAs and hypothesized that the TE-like epigenetic patterns we observed were due to TE silencing pathways.

TEs in plants are thought to be silenced by 2 main mechanisms. First, in the RNA-directed DNA methylation mechanism known as RdDM (Onodera et al. 2005), RNA polymerase IV (Pol IV)-transcribed RNA from TE loci is turned into 24-nt sRNAs that guide the DNA methylation

machinery to the locus being transcribed as well as to all homologous loci, allowing this mechanism to recognize and silence newly inserted TEs as well (Fultz et al. 2015). The second mechanism known to maintain TE silencing involves DECREASED DNA METHYLATION 1 (DDM1), METHYLTRANSFERASE 1 (MET1), CHROMOMETHYLASE 2 (CMT2), and CMT3, working together to establish the repressive H3K9me2 histone mark and DNA methylation at TE loci

(Sigman and Slotkin 2016; Osakabe et al. 2021). To test whether lincRNAs are also actively silenced by these mechanisms, we made use of publicly available RNA-seq data from mutants in components of the TE silencing machinery in Arabidopsis.

First, we analyzed the effect of inactivating the RdDM pathway. We observed above that 24-nt siRNA targeted ~50% of all lincRNA loci in flowers. Analysis of the sRNA data from Papareddy et al. (Papareddy et al. 2020) showed that 54% of lincRNA loci are targeted in early embryos; this targeting was highly specific to K9 lincRNAs (Fig. 7G). Knocking out *NUCLEAR RNA POLYMERASE D1* (*NRPD1*), encoding the largest subunit of PolIV, caused a dramatic loss of 24-nt sRNA coverage over 98% of those lincRNAs in embryos (Fig. 7H) as well as flowers (Supplemental Fig. S29A) (Papareddy et al. 2020). As 21–22-nt sRNA were also shown to trigger RdDM-mediated TE silencing (Nuthikattu et al. 2013) and be produced by Pol IV (Pontier et al. 2012; Panda et al. 2020), we analyzed the 21–22-nt sRNA levels at lincRNA loci. We determined that, similarly to 24-nt sRNAs, increased levels of 21–22-nt sRNAs in early embryos were associated with silencing in TE-containing lincRNAs but not in TE-free lincRNAs (Supplemental Fig. S29B). Moreover, targeting by 21–22-nt sRNAs was specific to K9 lincRNAs (Supplemental Fig. S29C) and knocking out *NRPD1* sharply reduced the level of 21–22-nt small RNAs at lincRNA loci that are normally targeted in the WT (Supplemental Fig. S29, D and E).

Next, we checked for the effect of removing *DDM1*, a key factor in TE silencing (Osakabe et al. 2021). Using the *ddm1* mutant in the Col-0 background, we observed that 149 of our lincRNAs become reexpressed in rosette leaves in this mutant compared with Col-0 (Fig. 7I) and 410 lincRNAs were reexpressed in *ddm1* stem cells (Fig. 7J). Heat stress combined with knocking out *ddm1* was particularly beneficial for the reactivation of lincRNA, which is similar to TE behavior (Nguyen et al. 2023) (Supplemental Fig. S30). The removal of CG and non-CG DNA methylation in Arabidopsis also allowed reexpression of many lincRNAs in rosette leaves (Fig. 7K). The reexpressed lincRNAs were again predominantly K9 lincRNAs and thus mostly TE-containing (Fig. 7, I to K, bottom; Supplemental Data Set 9).

The reexpression of lincRNAs in the *nrpd1* and *ddm1* mutants underscores 2 important points. First, that lincRNAs, predominantly the TE-like lincRNAs, are indeed silenced by the TE silencing machinery. Second, while we see most lincRNAs as being silent in any given accession, many retain the potential to be expressed and must therefore be actively silenced rather than having been inactivated by mutations. In fact, an analysis spanning across the different tissues and mutants with deactivated TE silencing pathways in Col-0 showed that over 50% of our annotated lincRNAs can be expressed in Col-0 (Supplemental Fig. S31), in contrast to 4% to 10% normally expressed in 1 sample (Fig. 2F). Thus, it appears that any genome is capable of expressing a large fraction of the numerous lincRNAs it harbors, but they are actively silenced, presumably largely via TE silencing pathways.

TE pieces within lincRNA loci appear to attract silencing to them

The presence of TE pieces within lincRNA loci was associated with increased epigenetic silencing (Fig. 5); in addition, TE silencing pathways predominantly affected lincRNAs with TE pieces (Fig. 7, I to K). We hypothesized that TE pieces might be decisive for TE-like silencing of lincRNAs by attracting the silencing machinery to the locus. To investigate this idea, we made use of the fact that different TE types show different silencing patterns. In particular, the RdDM pathway is more prevalent for DNA elements (Class II TEs), which are heavily targeted by 24-nt siRNAs, while retrotransposons (Class I TEs) such as LTR elements are more affected by the DDM1/CMT2 pathway showing heterochromatic patterns with high H3K9me2 levels (Sigman and Slotkin 2016; Sasaki et al. 2019). If TE pieces inside lincRNA loci are decisive for their silencing, we would expect that lincRNAs with pieces of different types of TEs would show silencing patterns resembling that of the corresponding TEs. Our analysis confirmed this hypothesis: lincRNA loci with pieces of DNA TEs, especially those derived from *Mutator Don Robertson* (*MuDR*) elements, showed significantly increased levels of 24-nt sRNAs (Fig. 8A), and lincRNAs with pieces of LTRs, especially those of *Gypsy* elements, showed significantly increased H3K9me2 levels (Fig. 8B). Class I TEs are more prevalent in the chromosome arms, and LTRs are enriched closer to the centromeres (Quesneville 2020); these trends were preserved when controlling for chromosomal position (Supplemental Fig. S32).

Although TE patches usually constitute only a portion of a lincRNA locus (Fig. 5E), they are associated with silencing over the full length of the locus (Fig. 5G). We analyzed repressive chromatin marks on the TE patch and TE patch-free parts of lincRNA loci and determined that while TE patches show higher repressive chromatin epigenetic modification, there was a very significant increase in repressive chromatin also outside of TE patches (Fig. 8, C and D; Supplemental Fig. S33), consistent with spreading of silencing (Sigman and Slotkin 2016). While H3K9me2 generally covered the whole TE-containing locus, 24-nt sRNAs and DNA methylation were more restricted to the TE pieces inside loci (Fig. 8, C to E; Supplemental Fig. S33).

Finally, we noticed that the numbers of lincRNAs and TE genes expressed in a given accession were quite well correlated (Supplemental Fig. S34A). TE silencing can vary across accessions, and indeed, we observed that the number of TE genes expressed across accessions varies nearly 3-fold. The correlation was much stronger for lincRNAs that contained TE pieces (Fig. 7N) supporting our hypothesis of shared silencing mechanisms. The number of TEs and lincRNAs expressed was correlated in every tissue (Supplemental Fig. S34B), indicating the organism-wide success or failure of silencing. Interestingly, while the number of expressed loci correlated well between accessions, the correlation between the mean expression levels across expressed lincRNAs and TE genes was much lower (Supplemental Fig. S34, C and D),

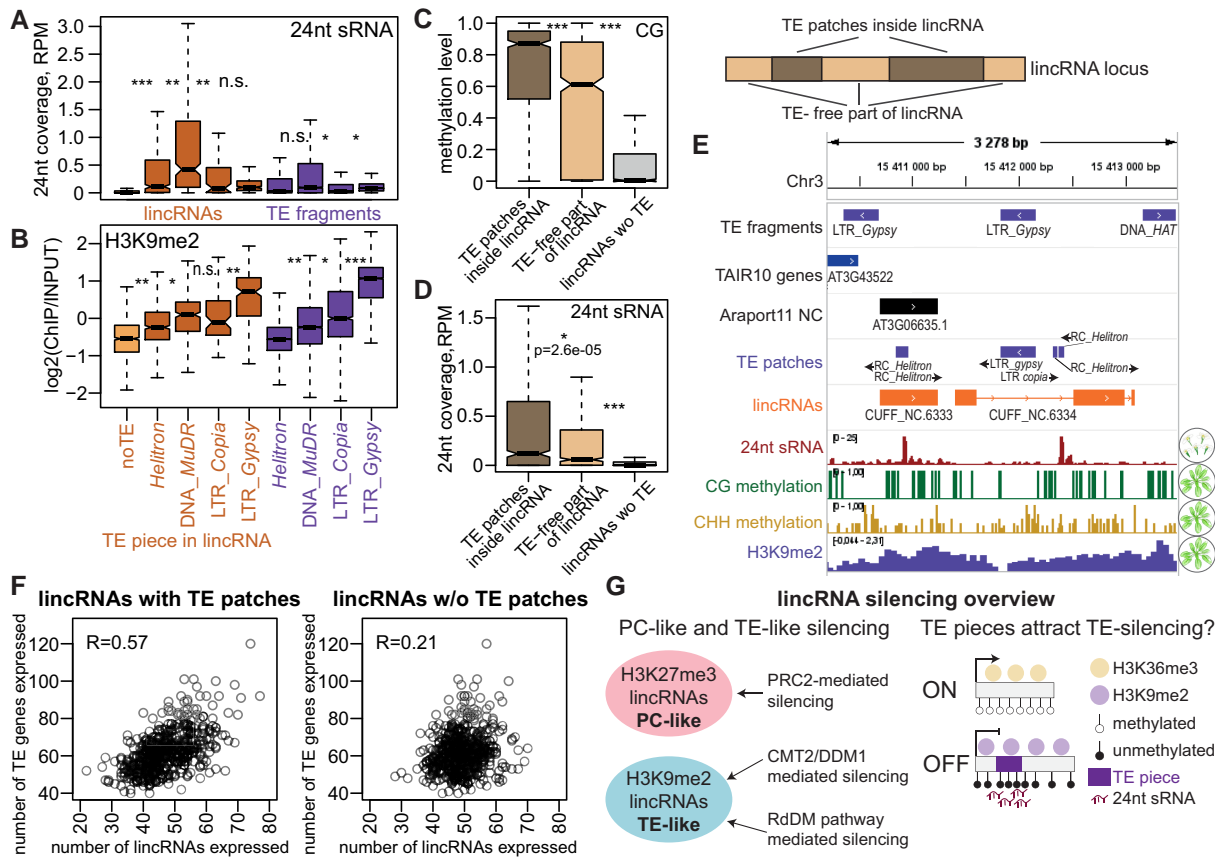


Figure 8. TE pieces appear to attract silencing to lincRNA loci. **A**) and **B**) 24-nt sRNA levels in Col-0 flowers (**A**) and H3K9me2 levels in rosettes (**B**) for lincRNAs with pieces of TEs from 4 superfamilies and TAIR10 TE fragments from the same superfamilies. Only lincRNAs with TE pieces from 1 superfamily are plotted. The light orange boxplot indicates lincRNAs without TE pieces (noTE). **C**) and **D**) Boxplots showing CG methylation level (**C**) and 24-nt sRNA coverage (**D**) for TE patches inside lincRNAs, TE patch-free parts of TE-containing lincRNA loci, and lincRNA loci without TE patches. Outliers not plotted. *P*-values were calculated using Mann–Whitney tests: ****P* < 10^{−10}, **P* < 0.01. **E**) Integrative Genomics Viewer (IGV) screenshot showing an example of lincRNAs with TE patches that have higher levels of CG methylation and 24-nt sRNA coverage over TE patches than over the rest of the locus. **F**) Scatterplot showing the number of TE genes expressed in rosettes of 460 different accessions (Kawakatsu et al. 2016) as a function of the number of lincRNAs with TE pieces (left) and without TE pieces (right) expressed in the same accession. Pearson’s correlation coefficient is displayed. **G**) Summary of lincRNA silencing pathways. PC-like lincRNAs that show H3K27me3 repressive histone marks are likely silenced by PRC2, while TE-like lincRNAs that display H3K9me2 are silenced by CMT2/DDM1 and RdDM pathways. TE piece presence likely attracts TE silencing and repressive chromatin to the lincRNA locus.

indicating that the 2 types of loci might share the same silencing machinery but likely not the general transcription apparatus and factors. We tried to identify genetic factors associated with the number of TE genes and lincRNAs expressed using genome-wide association study (GWAS) but could not see any clear association, except for 1 nearly significant peak on chromosome 2 near the *XERICO* gene (At2g04240), encoding a protein with a zinc finger domain (Supplemental Fig. S35), which is interesting as proteins with such domains are thought to participate in TE silencing (Yang et al. 2017).

In sum, lincRNAs display 2 distinct silencing mechanisms (Fig. 8G): PC-like silencing via H3K27me3 that is normally deposited by PRC2 (Hansen et al. 2008) and TE-like silencing, achieved via DDM1–CMT2 and RdDM silencing pathways (Fultz et al. 2015). The presence of TE pieces within lincRNAs appears to induce their TE-like silencing (Fig. 8G).

Discussion

An extended Arabidopsis lincRNA annotation

Unlike annotations based only on the reference accession Col-0, we used almost 500 Arabidopsis accessions and several developmental stages and identified several thousand previously unannotated lincRNA loci in the TAIR 10 reference genome. We conclude that over 10% of the genome can express lincRNAs but that most are not expressed in any particular accession or tissue, preventing a comprehensive lincRNA identification from few accessions or tissues. Analyzing more accessions allows identification of more lincRNA loci, with little evidence of saturation even when using data from several hundred accessions (Fig. 1F; Supplemental Fig. S5A). We provide an extended lincRNA annotation (Supplemental Data Set 5) as a resource for the Arabidopsis research community. Our results also suggest that lincRNA

annotations in other plant species could similarly be extended by population-wide studies.

In our study, we annotated lncRNAs using polyA⁺ RNA-seq data. While affordable and less prone to transcriptome assembly artifacts, polyA⁺ RNA-seq can miss nonpolyadenylated and/or unstable lncRNAs. Other methods, such as Global run-on sequencing (GRO-seq) (Hetzel et al. 2016), plant native elongating transcripts sequencing (plaNET-seq) (Kindgren et al. 2020), or exosome depletion (Thomas et al. 2020), can successfully detect an extended set of transcripts in Arabidopsis. Characterizing nascent and non-polyA transcription in multiple accessions and tissues would help capture the truly full scope of possible transcription and extend our understanding of transcription variation, allowing the distinction between variability in stability and in transcription initiation.

The largest part of our lncRNA transcriptome annotation consists of lncRNAs that are AS to PC genes. Apart from the general problem of natural variation impeding lncRNA identification described above, identifying AS lncRNAs crucially depends on having high-quality stranded RNA-seq data and a careful analysis to avoid artifacts (Supplemental Fig. S1). We were able to annotate almost 9,000 AS lncRNAs with nearly 30% of all PC genes having an AS partner, which greatly extends the scope of AS transcription. This is an important finding, since most functional lncRNAs reported in Arabidopsis, such as *COOLAIR* (Csorba et al. 2014), *antisense DELAY OF GERMINATION 1 (asDOG1)* (Fedak et al. 2016), *SVALKA* (Kindgren et al. 2018), and recently *SERRATE antisense intragenic RNA A (SEAIRa)* (Chen et al. 2023), are AS lncRNAs, and the massive extension of AS lncRNA annotation reported here thus opens a broad field for functional studies. A deeper investigation into AS lncRNAs and their function is beyond this study, but we provide a list of 14 AS lncRNAs that show striking negative correlation in expression with their partner PC gene (Supplemental Data Set 10 and Fig. S36, A and B) and thus are excellent candidates for being regulatory.

The second largest class of lncRNAs was intergenic lncRNAs that do not overlap with any PC genes, and these are the main focus of this article. This type of lncRNAs is very actively studied in mammals, with many functional examples reported (Rinn and Chang 2020). Arabidopsis lncRNA loci we annotate in this study are enriched for previously reported interesting genetic associations (Supplemental Fig. S36C) (Togninalli et al. 2020) as 157 lncRNA loci contained top GWAS hits associated with 65 different, mostly climate-related, phenotypes in Arabidopsis (Supplemental Data Set 11) (Togninalli et al. 2020). In this study, we focused on lncRNAs because they showed extreme expression variability and an interesting position intermediate between PC genes and TE genes in terms of expression, epigenetic features, and variation (Figs. 2 and 3). Their bimodal distribution in CG methylation levels was particularly striking (Fig. 3C). We also observed a clear dichotomy between H3K27me3 and H3K9me2 silencing (Fig. 5A) that

was recently also reported by Zhao et al. (Zhao et al. 2022). K27-silenced and K9-silenced lncRNAs were distinct in many features, most strikingly TE piece content, which made them similar to PC genes and TE genes, respectively, thus allowing us to distinguish 2 lncRNA subclasses: PC-like and TE-like lncRNAs. TE-likeness was conferred by the presence of TE pieces within the lncRNA locus.

TE pieces in lncRNAs

We showed that about half of all Arabidopsis lncRNA loci contained sequences similar to TAIR10-annotated TEs, which we refer to as TE pieces or patches when they held similarity to more than 1 TE superfamily. Strikingly, TE pieces were nearly 20 times more common within lncRNA loci than within PC genes and about 3 times more common than in random intergenic regions (Fig. 5C). It is unclear why lncRNA loci are so dramatically enriched in TE pieces. While this enrichment over PC genes is understandable, as TE insertions can be more deleterious for PC genes than for lncRNAs, the enrichment over random intergenic regions is very interesting. As lncRNAs are simply expressed intergenic regions of the genome without protein-coding capacity, the enrichment suggests that having a TE piece within the locus increases the probability of transcription. While our analyses suggest that TE pieces are associated with silencing, they might also provide the ability to be expressed when silencing fails. TEs are known to be the source of novel promoters in various organisms (Sundaram and Wysocka 2020). Thus, we can hypothesize that for many lncRNAs, TE pieces within the locus provide the potential for being transcribed, as well as contribute to it being silenced, albeit imperfectly, leading to our ability to detect these loci in our population-wide annotation. This hypothesis would go along with the extreme expression variability of TE-containing lncRNAs (Fig. 5K) and the very high variability in the overall level of TE gene and lncRNA expression (Fig. 7O) that indicates high TE silencing variability. Alternatively (and arguably more obscurely), the enrichment of TE pieces within lncRNAs may be caused by their transcriptional activity, if actively transcribed loci are more attractive for insertions compared with nontranscribed intergenic regions.

One very interesting group of TE-containing lncRNAs are lncRNAs with AS *Gypsy* elements. LncRNAs showed significant enrichment in *Gypsy* pieces or often full elements in the AS direction (Supplemental Fig. S22C). Why *Gypsy* elements show AS transcription more commonly than other elements remains to be investigated. We speculate that these elements might increase their mobilization chances by being transcribed from another strand, as strandedness does not matter for the transposition of retroelements, since it involves a double-stranded DNA step.

The nature of TE pieces

Another major topic raised by our results is the nature and origin of the TE pieces we identified in lncRNA loci. Some

of these TE pieces are simply parts of intact TE fragments that are overlapped by the lincRNA locus (Supplemental Fig. S21, E and F), and some are full TE fragments in the direction AS to the direction of lincRNA transcription (Supplemental Fig. S21, C and D). In these cases, the nature of the TE sequence inside lincRNA is clear, but the question of what came first—the expression or the TE—remains. Most intriguing are the many cases of short, and sometimes very short, independent pieces of TEs within lincRNA loci, the nature of which is puzzling. First, these TE pieces might represent insertions into the loci. However, their small size (Fig. 5D) raises the question of how they were able to mobilize and get inserted into the lincRNA loci. Nonautonomous TEs (Quesneville 2020), in particular, DNA-TE-derived MITEs (miniature inverted-repeat TEs) (Okamoto et al. 2008) and LTR-TE-derived SMARTs (small LTR retrotransposons), have been studied (Mhiri et al. 2022), yet those still have a length of a few hundred bp, while our pieces are often around 100 bp or shorter. It has also been suggested that small nonautonomous TEs can transpose with a piece of a nearby genomic sequence, thus shuffling it around, but there is little understanding of how this might work (Quesneville 2020).

As many lincRNAs are known to originate from TEs (Kapusta et al. 2013), such as the famous *X-inactive specific transcript* (*XIST*) lincRNA (Colognori et al. 2020), it is also possible that the TE pieces we find within lincRNAs are not insertions but rather remnants of decaying TEs. One approach to distinguishing between the 2 possibilities would be to study the structural variation of TE pieces: variability of the presence of that precise piece would clearly indicate insertion/excision rather than the decay of a larger TE. What we could assess within the scope of this study is whether multiple TE pieces within 1 lincRNA locus resemble 1 or multiple TE families. If a locus contains pieces of different TEs, this would be evidence against the TE decay hypothesis. Among lincRNAs with more than 1 TE piece in TAIR10, 74% have TE pieces from different superfamilies and 24% from both Class I and Class II TEs. Further research and an analysis of full genomes from multiple accessions are crucial for understanding the nature, evolutionary history, and population dynamics of TE pieces inside lincRNAs.

Silencing

We discovered that the Arabidopsis genome has a large potential for lincRNA expression that is massively repressed by silencing. While many lincRNAs are repressed by PC-like H3K27me3-mediated silencing, about as many are repressed by TE silencing, which is associated with having a TE piece within the locus. The presence of a TE piece was correlated with repressive chromatin marks and silencing, with higher TE contents in a locus being correlated with stronger silencing (Fig. 5, E to H). We also showed that inactivating TE silencing pathways in the reference accession Col-0 allowed expression of many TE-like lincRNAs that are normally completely silent in this accession. TE pieces appear to attract silencing to the locus, as we observed that lincRNAs seem to be

preferentially silenced by RdDM- or CMT2-silencing pathways depending on which TE family the TE piece within the lincRNA locus came from (Fig. 7E). Interestingly, TE pieces and multiple copy number were associated with the same patterns of silencing in both AS lincRNAs and PC genes (Supplemental Fig. S37), although the relative number of such TE-like genes was much smaller (Fig. 4B). This observation suggests that a genome-wide mechanism for suppression of TE-like loci exists (Sigman and Slotkin 2016).

The mechanism by which short TE pieces attract TE-like silencing to a lincRNA locus is unclear. It is known that full-length TEs can induce the silencing of nearby genes by the spreading of repressive chromatin (Sigman and Slotkin 2016), and we hypothesize that TE pieces are capable of this as well. However, how they themselves obtain repressive chromatin is unclear. One possibility is that 24-nt siRNAs produced at TE loci find the TE pieces by homology and initiate silencing at this “TE-like” locus (Fultz et al. 2015). They likely initially target only the TE piece and not the full locus, and we do see that the 24-nt siRNA and CG/CHH methylation signal is highest at TE patches (Fig. 8, C and D; Supplemental Fig. S33). However, we also observed a significant increase of 24-nt siRNA as well as CG and CHH methylation levels outside of TE patches, which may suggest that the spreading might include sRNAs starting to be produced at the locus. It is also possible that many lincRNA loci with TE patches initiate their own silencing through the RdDM pathway, thus producing their own Pol IV-dependent sRNAs.

It is also unclear what causes the failure of silencing of certain lincRNAs in certain accessions. It is possible that the silencing machinery varies in efficiency, and we see some evidence for this in the 3-fold range of variation in the number of TE genes and TE-containing lincRNAs expressed across accessions (Fig. 7N). However, we could not find any gene expression level or single-nucleotide polymorphisms (SNPs) that were clearly associated with the overall extent of lincRNA or TE transcription. It is also unclear how variation in silencing efficiency could account for such a strong lincRNA landscape variability across accessions with similar overall lincRNA transcription (Supplemental Fig. S38). This variation may reflect the presence of particular TE loci producing the appropriate siRNAs for TE pieces within particular lincRNA loci.

Further studies are clearly needed. In this study, we focused on the reference genome, demonstrating that TE pieces within lincRNA loci are important for silencing. Direct experiments, like inserting a TE piece into a TE-free lincRNA locus and assessing the resulting expression change, are outside the scope for this study. Similarly, an analysis of the full genomes of multiple accessions, including variation for TE and TE fragment content, would be informative, and such an analysis is underway.

The distinction between lincRNAs and TEs

As lincRNAs with TE pieces showed many similarities to TEs, including similar epigenetic patterns, silencing pathways, and

increased copy number, the question might arise as to whether these TE-containing lincRNAs are distinct from TEs. By definition, an lincRNA is a transcript longer than 200 nt without protein-coding potential and thus technically any non-protein coding transcript arising from a TE can be considered a lincRNA. However, our annotation pipeline did distinguish lincRNA loci from expressed TE fragments by a <60% sense exonic overlap with annotated TE fragments, as well as by applying a protein-coding capacity cutoff to lincRNAs but not TEs. The 60% cutoff we applied is admittedly arbitrary, although common in the lincRNA field, and strongly depends on the TE annotation used. The TE annotation in Arabidopsis is far from complete, and studies analyzing recent genetic mobility in Arabidopsis and annotation of new TEs are underway. A largely extended TE annotation could affect the TE overlap filtering step we used in our annotation pipeline classifying some of our lincRNAs as “expressed TEs” (see Supplemental Fig. S1). However, many of the TE patch-containing lincRNA loci showed only a minor overlap with annotated TE fragments, while quite a few (24%) had no overlap at all. TEs furthermore have a direction, and many lincRNAs were transcribed AS to the TE fragment they overlapped with, indicating that these are separate transcriptional units, even though the inherently non-strand-specific epigenetic marks were shared between the 2. Moreover, most of the lincRNAs contained a mix of TE pieces from different families, which is a strong indication that these lincRNAs are unlikely to be intact TEs. Thus, we think that most of our lincRNAs are distinct from intact TEs and that the effect TE patches have on the expression and silencing of lincRNAs, and other loci they occur in, is a very interesting phenomenon deserving future research.

Nevertheless, some of the lincRNAs we detected might actually be previously unannotated active or recently active TEs. The presence of a patch similar to annotated TEs might resemble the occasional sequence likeness between annotated TEs of different families. We identified 58 lincRNA loci with a sense TE patch that were present in more than 10 copies; they represent the most likely candidates for unannotated TEs. We also found 39 lincRNAs that, while having no TE patches, were also present in more than 10 copies, highlighting that all TE sequences are unlikely to have been annotated. However, to definitively conclude that a lincRNA locus is in fact a TE, we would need evidence of mobilization between accessions or species and evidence of the TE piece within the lincRNA locus being an integral part of it rather than an insertion—thus not showing variability between accessions. These analyses represent future directions and are outside the scope of this study.

lincRNA expression variation and future directions

Our study initially had 2 major goals: to create a population-wide map of lincRNA transcription in Arabidopsis and characterize its natural variation. We discovered that the extent of lincRNA transcription in Arabidopsis is much larger than previously thought and that lincRNA expression patterns are largely

variable between accessions with half of all lincRNAs being expressed in 1 accession while being off in another (Fig. 2E). In this study, we characterized the expression variability of lincRNAs in Arabidopsis, but we only accessed the epigenetic patterns among the factors that could explain the expression variation across accessions. We showed that lincRNAs display extensive epigenetic variation (Fig. 4, A and B), and this variation can explain the expression of ~50% of informative lincRNAs and ~20% of informative AS lincRNAs (Supplemental Fig. S20B). While purely epigenetic variation is well known (Xu et al. 2019; Rajpal et al. 2022), our analysis did not distinguish between this and when the epigenetic variation that defines expression variation is itself defined by an underlying genetic or structural variation. We showed that 2 structural features of lincRNA loci—their TE content and their copy number—are associated with silencing and increased expression and epigenetic variation (Figs. 5 and 6), and it is clear that variation in these 2 features might be responsible for the variation in expression that we observed between accessions. In this study, we constrained our analysis to the reference genome, because an analysis of structural variation in copy number or TE piece presence requires full-genome assemblies of nonreference accessions. We will perform these analyses in an upcoming study that investigates the determinants of lincRNA expression across accessions in greater depth.

In conclusion, analyzing transcriptomes from multiple accessions and tissues of Arabidopsis accessions allowed us to drastically extend its lincRNA annotation and study the natural variation of lincRNA expression. We established that 10% of the Arabidopsis genome is covered with almost 12,000 lincRNA loci; however, most of them are silent in any given sample. lincRNAs, particularly long intergenic ncRNAs, show very high expression and epigenetic variation. The silencing of lincRNAs is achieved via PC-like and TE-like mechanisms, with the latter being defined by the pieces of TEs present in about half of all lincRNAs. We produced a multiaccession transcriptome and epigenetic resource, as well as an extended lincRNA annotation useful for the Arabidopsis community and provide new insights into the genome biology and composition of lincRNAs.

Materials and methods

Sample collection

Arabidopsis thaliana seeds were surface sterilized with chlorine gas for ~1 h, stratified at 4 °C for ~5 d to induce germination before being sown onto soil (3 parts Peat moss Gramoflor professional mixture [Gramoflor GmbH] mixed with 1 part Gramoflor Premium Perlite 2-6 [Gramoflor GmbH]). Plants were grown in growth chambers at 21 °C under long-day conditions (16 h light/8 h dark) with a light intensity of 130 to 150 $\mu\text{mol}/\text{m}^2/\text{s}$ (HLG-240H-30A LED lamps, Mean Well Enterprises Co., LTD). For each tissue type, all accessions were grown and processed in parallel at all stages to avoid non-accession-related variation. The “14-leaf rosette”

(or mature leaves) samples were collected at the 13- to 16-leaf stage before plants started to bolt. Approximately 8 leaves (avoiding the oldest and the youngest leaves) were collected from 2 to 3 individuals of the same accession into 20-mL Polyvial bottles (Zinsser Analytic) with metal beads inside, snap-frozen in liquid nitrogen, and stored at -70°C . Tissue was ground while frozen, producing 1 to 2 mL of tissue powder that was used for preparation of RNA-seq, bisulfate-seq, and ChIP-seq libraries. The 14-leaf rosette samples had 2 to 4 replicates per accession (Supplemental Data Set 2): accessions were grown in the growth chamber, followed by tissue harvesting 2 to 4 times, with a gap of several weeks between each batch. For each accession, the samples collected in each batch are referred to as replicates. For the “9-leaf rosette” samples, the full rosette at the 9-leaf stage was collected, with 1 plant harvested per sample. Seedlings were collected at 7 d postsowing (~ 5 d postgermination). Full seedlings with the root were harvested with ~ 10 seedlings harvested per sample. For the “flower” samples, flowers and flower buds were collected from ~ 5 individuals per sample. Accessions 1741, 6024, 6244, 9075, 9543, 9638, 9728, 9764, 9888, 9905, 22003, 22004, 22005, 22006, and 22007 were vernalized by taking them out from the 21°C growth chamber at the age of ~ 3 wks and placing them into 10°C growth chambers (under long-day conditions) for ~ 4 wks to induce flowering. Accessions 6069 and 6124 did not flower even after the cold treatment. Pollen was collected using the method described in (Johnson-Brousseau and McCormick 2004) that uses vacuum suction and a series of filters to harvest dry pollen from flowering plants. Polyester mesh filters of 3 sizes were used (150, 60, and 10 mm) for sample collection. Collected pollen was snap-frozen in liquid nitrogen, stored at -70°C , and ground for total RNA ahead of RNA-seq using ~ 0.5 mL of 0.5-mm-diameter glass beads (Scientific Industries, Inc.). All tissue grinding was performed using the Retsch Oscillating Mill MM400 (Retsch GmbH) with 1/30 frequency for 90 s using custom-made metal adapters that were precooled in liquid nitrogen to keep the samples frozen while being ground.

RNA sequencing and analysis

Total RNA was isolated and treated with DNase I (NEB) using a KingFisher Robot with an in-house magnetic RNA isolation kit. Total extracted RNA was diluted in nuclease-free water (Ambion) and stored at -80°C . Libraries for RNA-seq were prepared using a TruSeq Stranded mRNA kit (Illumina) following the manufacturer’s protocol with a 4 min RNA fragmentation time and 12 PCR cycles for library amplification. RNA-seq was performed at the Vienna Bio Center (VBC) NGS facility on an Illumina HiSeq 2500 machine in paired-end read mode of 150 and 125 bp. Raw RNA-seq data were aligned to the TAIR10 genome using STAR (Dobin et al. 2013) with the following options: `--alignIntronMax 6000 --alignMates GapMax 6000 --outFilterIntronMotifs RemoveNoncanonical --outFilterMismatchNoverReadLmax 0.1 --outFilterMismatchNoverLmax 0.3 --outFilterMultimapNmax 10 --align`

`SjoverhangMin 8 --outSAMattributes NH HI AS nM NM MD jM jI XS`. Gene expression levels were calculated using featurecounts from the Subread package with `-t exon` option and an exonic SAF file as an annotation (Liao et al. 2014).

Transcriptome assembly and lncRNA annotation

Transcriptome assembly was performed in several steps as described in Supplemental Fig. S1; the scripts are provided at https://github.com/aleksandrakornienko/Kornienko_et_al_lncRNA_expression_variation_and_silencing. In brief, the following RNA-seq data sets were used for transcriptome assembly: 14-leaf rosette data from 461 accessions (100 bp single-end reads) (Kawakatsu et al. 2016), seedling data from Cortijo et al. (Cortijo et al. 2019) (75 bp paired-end, 14 replicates for each of the 12 samples were pooled), our 14-leaf rosette data from 28 accessions (2 to 4 replicates each) (125 bp paired-end), our seedling and 9-leaf rosette data from 27 accessions (150 bp paired-end), and our flower and pollen data from 25 accessions (150 bp paired-end). First, the transcriptomes of each sample were assembled separately using Stringtie v.2.1.5 (Pertea et al. 2015) with options: `-c 2 -m 150 -j 2.5 -a 15` guided by the TAIR10 gene annotation (-G). The transcriptome assemblies of the same tissue and data types were then merged using Cuffmerge (Cufflinks v.2.2.1) (Trapnell et al. 2010) with `--min-isoform-fraction 0` before performing a second merging of the resulting 7 transcriptomes to obtain the cumulative transcriptome annotation. A series of filtering steps were applied, including a transcript length cutoff of 200 nt for multiexon genes and 400 nt for single-exon genes, and then, the genes were split into (i) PC genes by exonic overlap with TAIR10- or Araport11-annotated PC genes, (ii) TE genes based on exonic overlap with Araport11-annotated TE genes, (iii) TE fragments with $>60\%$ same-strand exonic overlap with TE fragments annotated in Araport11, (iv) pseudogenes by exonic overlap with Araport11-annotated pseudogenes, and (v) initial lncRNAs showing no overlap with PC genes, TE genes, or pseudogenes and with $<60\%$ same-strand exonic overlap with TE fragments annotated in Araport11. Then, lncRNA transcripts (and the corresponding loci containing those transcripts) with protein-coding capacity as tested by CPC2 (Kang et al. 2017) were removed; rRNA, tRNA, sn/snoRNA, and miRNA precursor lncRNAs were classified based on overlap with the appropriate annotations. The remaining lncRNAs were classified into (i) AS lncRNAs by AS overlap with TAIR10 or Araport11-annotated PC genes, (ii) lncRNAs AS to pseudogenes, (iii) lncRNAs AS to Araport11 TE genes (AS_to_TE), and (iv) intergenic lncRNAs (lincRNAs) with no overlap with PC genes, TE genes, or pseudogenes. LincRNAs were additionally filtered against loci that started <100 bp downstream from annotated genes to avoid read-through transcripts. The number of Araport11 PC genes with an AS transcript was calculated using Araport11 noncoding and novel_transcribed_region annotations filtered for genes longer than 200 bp.

Gene saturation curve

To create the gene saturation curves for accession and tissue number, the annotation pipeline was automated and run many times with different numbers of accessions and tissues. The accession saturation curve was generated by inputting 10 to 460 transcriptome assemblies (1 assembly being 1 accession) obtained from the 1001 Arabidopsis genome data set (Kawakatsu et al. 2016) into the same annotation pipeline used for the main gene annotation defined above. Subsampling of accessions was done randomly with 8 iterations for each number of accessions. The curve fitting and prediction of the saturation curve behavior with up to 1,000 accessions was done by fitting a linear model using the `lm` function in R with the command line: `model <- lm(y ~ x + l(log2(x)))` (Supplemental Fig. S5, A and B). The control for the accession saturation curve was done using the data from Cortijo et al. (Cortijo et al. 2019), from which 1 to 12 transcriptome assemblies (corresponding to 12 samples with 14 replicates per sample pooled into 1 BAM file pre-assembly) were randomly chosen and fed into our standard annotation pipeline, counting the number of loci identified as an output. The procedure was performed 8 times for each assembly number. Then, the number of reads was calculated and juxtaposed to the number of reads in the multi-accession saturation curve (Supplemental Fig. S5C). As different data sets had different read modes, the results were aligned by calculating the total read length and multiplying it by the total read number. The tissue saturation curve analysis was performed on 23 accessions that had data from all 4 tissues. Random sampling of accessions was performed with 8 iterations as replicates for each number of accessions. Tissues were assessed in this particular order: (i) seedling, (ii) rosette, (iii) flowers, and (iv) pollen, without random sampling (Supplemental Fig. S6).

Expression variation analysis

Interaccession variability was calculated as coefficient of variance (SD divided by mean) of TPM of the locus across accessions in the data set. When multiple replicates were available, the average between the replicates was taken for the coefficient of variation calculation. Intraaccession variability was calculated using our 14-leaf multireplicate data set as follows: for each accession, the coefficient of variance for TPMs across replicates was calculated, and then, the coefficients of variance for each accession were averaged. Expression noise was calculated using the seedling data set from Cortijo et al. (2019) that contained 14 replicates for 12 time points: coefficient of variance across 14 replicates was calculated for each of the 12 time points, and then, these 12 coefficients of variance were averaged to produce the resulting “noise” level. Circadian expression variation was also calculated using the seedling data set from Cortijo et al.: for each time point, TPMs from 14 replicates were averaged, and then, coefficient of variance was calculated across the 12 time points. The 12 time points represent the samples collected every 2 h within

a 24 h period of time (12 h light, 12 h dark), and by “circadian expression variation,” we mean expression variation during a 24 h period (across the 12 time points).

ChIP sequencing

Chromatin immunoprecipitation was performed with a protocol adapted from Yelagandula et al. (2014) (full protocol is available at https://github.com/aleksandrakornienko/Kornienko_et_al_lncRNA_expression_variation_and_silencing). Briefly, 1 to 2 g of ground frozen leaf tissue was fixed with 1% formaldehyde at 4 °C for 5 min, and then, nuclei were isolated and lysed using a series of lysis and centrifugation steps; chromatin was fragmented in 1-mL Covaris milliTUBEs using a Covaris E220 Focused-ultrasonicator for 15 min at 4 °C with the following settings: duty factor of 5.0, peak incident power of 140, and 200 cycles per burst. An aliquot was then taken out as input sample and frozen at –20 °C; the remaining supernatant was split in 5 tubes, 1 for each antibody, and was processed together. The antibodies used were against H1 (Agrisera), 3 µg per reaction; H3K4me3 (Abcam), 3 µg per reaction; H3K9me2 (Abcam), 4 µg per reaction; H3K27me3 (Millipore), 4 µg per reaction; and H3K36me3 (Abcam), 4 µg per reaction. The immunoprecipitation was performed using prewashed Dynabeads Protein A magnetic beads (Invitrogen) and incubated at 4 °C overnight. Afterwards, the samples were washed, followed by elution, overnight reverse-crosslinking, treated with RNase A (Fermentas) for 30 min at room temperature, and the DNA isolated using a QIAquick PCR purification kit (Qiagen) with 0.3 M sodium acetate. Next, ChIP-seq libraries were prepared from half of the resulting sample (due to very low amounts, we did not measure the DNA concentrations) with a NEBNext Ultra II DNA kit (New England Biolabs) according to the manufacturer’s protocol and sequenced as 100 bp single-end reads on an Illumina NovaSeq 6000 instrument.

ChIP-seq analysis

Raw ChIP-seq reads were mapped using STAR (Dobin et al. 2013) adjusted for ChIP-seq with the following options: `--alignIntronMax 5 --outFilterMismatchNmax 10 --outFilterMultimapNmax 1 --alignEndsType EndToEnd`. Only samples with >1 million unique nonduplicated reads were used for analysis. Aligned BAM files from each ChIP sample were then normalized by the corresponding input samples using `bamCompare` from `deeptools` (Ramírez et al. 2016) with the following options: `--operation log2 --ignoreDuplications --effectiveGenomeSize 119481543`; bigwig and bedgraph files were created. Read coverage over loci and promoters was estimated using `bedtools map` on the bedgraph files with the “mean” operation. To estimate the variation in histone modification levels, the ChIP-seq coverage values were normalized again to achieve the same range of values across accessions, applying quantile normalization, setting the 20% and 80% quantile values for each sample to the same value across samples with the function in R: `quantile_minmax <- function(x) {(x-quantile(x,.20)) / (quantile(x,.80) - quantile(x,.20))}`. Histone modification variation was then calculated as the

SD of quantile-normalized levels averaged across replicates for each accession.

Bisulfite sequencing

Bisulfite sequencing was performed as described in [Pisupati et al. \(2022\)](#). Briefly, DNA was extracted from frozen leaf tissue (14-leaf rosettes) using a Nuclear Mag Plant kit (Machery-Nagel) and the bisulfate sequencing libraries were prepared using a tagmentation method described in [Wang et al. \(2013\)](#) using an in-house Tn5 transposase (IMBA-IMP-GMI Molecular Biology Services) and an EZ-96 DNA Methylation-Gold Mag Prep kit (Zymo Research) for bisulfite conversion. Bisulfate sequencing libraries were sequenced on the Illumina NovaSeq 6000 instrument in the 100 bp paired-end mode.

DNA methylation analysis

Bisulfite sequencing data were used to call methylation in 3 contexts (CG, CHG, and CHH where H stands for A, C, or T) using the method described in [Pisupati et al. \(2022\)](#). The methylation level per locus for each context was determined by dividing the number of methylated reads by the total number of reads covering the cytosines in the CG, CHG, or CHH context. Thus, the values of methylation of each locus range from 0 to 1 and roughly correspond to the ratio of methylated to total cytosines in the locus (we did not take the average of the ratios for each cytosine to avoid high error rates caused by low read coverage).

Small RNA sequencing and analysis

Small RNA was isolated from frozen and ground flower samples using a NucleoSpin miRNA kit (Macherey-Nagel) following the manufacturer's protocol. Small RNA-seq libraries were prepared using a QIAseq miRNA Library Kit (Qiagen). Raw fastq files were trimmed with cutadapt (v.1.18) ([Martin 2011](#)) using `-a AACTGTAGGCACCATCAAT` and `--minimum-length 18` options. Trimmed reads were aligned to the TAIR10 genome using STAR (v.2.9.6) adjusted for sRNA-seq, allowing 10 multimappers and 2 mismatches. Reads that were 24 nt long, and likewise 21–22-nt-long reads, were extracted, and read coverage for each bp of the genome was calculated using genomeCoverageBed (bedtools v.2.27.1) and normalized by dividing by the unique number of reads in the sample. Final sRNA coverage was calculated by mapping the normalized read coverage per bp over the loci of interest and calculating the average coverage across all bp of each locus. Raw sRNA-seq data from [Papareddy et al. \(2020\)](#) were processed using the same pipeline. The cutoff for calling a locus as being targeted by 24-nt or 21–22-nt sRNAs was set to $RPM = 0.03$.

Explaining expression variation by DNA methylation variation levels

To determine if DNA methylation can explain expression variation, matching RNA-seq and DNA methylation data

from the 1001 Genomes Project was used ([Kawakatsu et al. 2016](#)). Out of 461 RNA-seq samples, 444 had matching bisulfite sequencing data; thus, data for these 444 accessions were used for this analysis. For each lncRNA in our annotation, we asked if accessions where an lncRNA was highly methylated showed significantly lower expression (Mann–Whitney test, $P < 0.01$) than accessions with low methylation levels at this lncRNA locus ([Supplemental Fig. S20, A and B](#)). This analysis was performed for exonic gene body TPM calculated as described above using 4 estimates of methylation level: CG and CHH methylation level of gene body and CG and CHH methylation level of promoters ($TSS \pm 200$ bp). We set “low” CG methylation level as “ <0.5 ” and “high” CG as “ ≥ 0.5 ” and “low” CHH methylation level as “ <0.01 ” and “high” CHH as “ ≥ 0.01 .” These cutoffs were defined based on the distribution of CG and CHH methylation levels of PC genes (average gene body methylation level across 444 accessions): 90% of PC genes in our transcriptome annotation had accession-wide mean CG methylation levels below 0.37 and CHH methylation levels below 0.01.

TE piece analysis

To find TE pieces in various loci, 31,189 annotated TE sequences from TAIR10 were compared with the sequence of each locus using BLASTN (BLAST+ v2.8.1) with options `-word_size 10 -strand both -evalue 1e-7`. We required $>80\%$ sequence identity and did not restrict the length. We then merged same-strand overlapping TE pieces into TE patches. For all of our analyses, we grouped TE families into 7 superfamilies: DNA_other, DNA_MuDR, SINE_LINE, RC_Helitron, LTR_Gypsy, LTR_Copia, and Unassigned_NA.

Copy number analysis

Copy number was estimated by extracting the sequence of the locus from the TAIR10 genome and using it as a query against the TAIR10 genome using BLASTN (BLAST+ v2.8.1) with options `-word_size 10 -strand both -outfmt 7 -evalue 1e-7`. We allowed for copies to be disrupted by insertions of no more than 1.5 kb and applied a cutoff of $>80\%$ on sequence identity and $>80\%$ on length to all regions identified by BLASTN.

lncRNA reexpression analysis

Raw RNA-seq data from the silencing mutants from [Osakabe et al. \(2021\)](#), [He et al. \(2022\)](#), and [Nguyen et al. \(2023\)](#) were processed using the same RNA-seq pipeline as described above. The reexpression in the mutants was generally defined as the lack of expression in the WT ($TPM < 0.5$, averaged from all available replicates), the presence of expression in the mutant ($TPM > 0.5$), and additionally a 3-fold difference between the expression in the mutant and the WT ($MUT > 3*WT$). For the *ddm1* knockout in stem cells, the mock *ddm1* mutant sample was matched with the mock WT, and the heat-treated *ddm1* mutant was matched with the heat-treated WT sample (heat treatment as described in [Nguyen et al. \(2023\)](#)). For the methylation mutants, *ddcc*

and *met1-9* mutants were 2-wk-old seedlings and matched with the 2-wk-old WT control, while the *mddcc* mutant was matched with the 5-wk-old WT control. The *met1-9* mutant corresponds to a *MET1* knockout with loss of CG methylation; the *ddcc* mutant is the quadruple mutant for *DOMAINS REARRANGED METHYLASE 1* (*DRM1*), *DRM2*, *CMT2*, and *CMT3* with a loss of CHG and CHH methylation; the *mddcc* mutant is a quintuple mutant for *MET1*, *DRM1*, *DRM2*, *CMT2*, and *CMT3* with a nearly full loss of all methylation (He et al. 2022).

Use of public data sets

The summary of the public data sets used in our study and the corresponding mapping statistics are available in [Supplemental Data Sets 2, 3, and 12](#). The public data sets were downloaded from NCBI GEO using the specified GEO accession numbers below:

- 1) RNA-seq and bisulfite-seq from mature leaves of 14-leaf rosettes from the 1001 Genomes Project (Kawakatsu et al. 2016): **GSE80744** and **GSE43857**.
- 2) RNA-seq data from Col-0 seedlings from 12 time points with 14 technical replicates each (Cortijo et al. 2019): GEO accession number **GSE115583**.
- 3) Early embryo and flower bud sRNA-seq data from *nprpd1* knockouts (Papareddy et al. 2020): **GSE152971**.
- 4) Rosette RNA-seq data from *ddm1* knockouts (Osakabe et al. 2021): **GSE150436**.
- 5) Stem cell RNA-seq data from *ddm1* knockouts with and without heat stress (Nguyen et al. 2023): **GSE223915**.
- 6) Rosette RNA-seq data from DNA methylation-free mutants (He et al. 2022): **GSE169497**.

Statistical analysis

Statistical analyses were performed as described in each figure legend. Statistical data are provided in [Supplemental Data Set 13](#).

Accession numbers

The sequencing data produced in this study are available at the NCBI Gene Expression Omnibus (<https://www.ncbi.nlm.nih.gov/geo/>) as SuperSeries GSE224761. The gene annotations and the 14-leaf rosette RNA-seq dataset from 28 Arabidopsis accessions are available under the accession number GSE224760 (note that gene annotations are also available in the supplement as [Supplemental Data Set S1](#)); the corresponding bisulfite sequencing data are under the GEO accession number GSE226560. The 14-leaf rosette ChIP-seq data from 14 accessions are available under accession number GSE226682. The RNA-seq data from seedlings, 9-leaf rosettes, flowers, and pollen from 25 to 27 accessions are available under the GEO accession number GSE226691. The flower sRNA-seq data from 14 Arabidopsis accessions are available under the GEO accession number GSE224571.

The code used for the analyses and figures is available at https://github.com/aleksandrakornienko/Kornienko_et_al_lncRNA_expression_variation_and_silencing.

Acknowledgments

We would like to thank the Vienna Biocenter Core Facilities GmbH (VBCF) Next-Generation Sequencing for NGS services; VBCF Plant Sciences for the excellent growth chambers; IMBA-IMP-GMI Molecular Biology Services for providing the access to instruments, molecular biology reagents, and support; and the VBC Ethics, Health and Safety team for their support during the COVID-19 pandemic. We would like to thank Mirjam Bissmeier for help with preliminary ChIP-seq analyses, Michael Schon for lncRNA discussions and advice on transcriptome assembly, and Bhagyshree Jamge, Vu Nguyen, Ramesh Yelagandula, Zdravko Lorkovic, Nathalie Durut, and Ortrun Mittelsten Scheid for their advice with experiments and data and fruitful discussions. We would like to greatly thank Detlef Weigel for helping to secure funding for the “1001 Genomes Plus” project and for his helpful comments on the manuscript.

Author contributions

A.E.K. and M.N. designed the study. A.E.K. performed most of the experiments and data analyses. A.M.M. helped with sample collection. R.P. performed the bisulfite sequencing data processing. V.N. prepared all libraries for sequencing. A.E.K. and M.N. wrote the paper.

Supplemental data

The following materials are available in the online version of this article.

Supplemental Figure S1. Cumulative transcriptome annotation pipeline.

Supplemental Figure S2. Transcriptome annotation supplement.

Supplemental Figure S3. AS lncRNA supplement.

Supplemental Figure S4. Genomic distribution of annotated loci.

Supplemental Figure S5. Gene identification saturation analysis.

Supplemental Figure S6. Gene identification saturation analysis: tissues and accessions.

Supplemental Figure S7. Expression frequency.

Supplemental Figure S8. Interaccession expression variation in different tissues.

Supplemental Figure S9. Expression variability controls: absolute expression level and gene length.

Supplemental Figure S10. Inter- and intraaccession expression variability with replicates.

Supplemental Figure S11. Histone mark profiling.

Supplemental Figure S12. DNA methylation level supplement.

Supplemental Figure S13. DNA methylation vs distance to the centromere.

Supplemental Figure S14. Heterochromatic histone marks vs distance to the centromere.

Supplemental Figure S15. Histone modifications of silent and expressed genes: supplement.

Supplemental Figure S16. DNA methylation of silent and expressed genes: supplement.

Supplemental Figure S17. Coverage of 24-nt and 21–22-nt sRNA in early embryo and leaves.

Supplemental Figure S18. Epigenetic patterns in nonreference accessions.

Supplemental Figure S19. Epigenetic variation supplement 1.

Supplemental Figure S20. Epigenetic variation supplement 2.

Supplemental Figure S21. TE patches in lincRNAs.

Supplemental Figure S22. Sense and AS TE pieces in lincRNAs: content, family, and position.

Supplemental Figure S23. Expression variation vs TE content: expression control.

Supplemental Figure S24. lincRNA methylation variation vs. TE content: supplement.

Supplemental Figure S25. TE pieces affect epigenetics when controlled for chromosomal location.

Supplemental Figure S26. Copy number supplement.

Supplemental Figure S27. Copy number affects lincRNA epigenetic pattern: supplement.

Supplemental Figure S28. H3K27me3 and H3K9me2 dichotomy.

Supplemental Figure S29. Loss of 24 nt targeting in *nprp1a* mutants.

Supplemental Figure S30. TE silencing mutant supplement.

Supplemental Figure S31. The scope of lincRNA expression potential in Col-0.

Supplemental Figure S32. Genomic position and epigenetics of lincRNAs with pieces of Class I and II TEs.

Supplemental Figure S33. Spreading of silencing from TE patches.

Supplemental Figure S34. Variability of the number of TE genes and lincRNAs expressed.

Supplemental Figure S35. GWAS on expressed TE gene number.

Supplemental Figure S36. lincRNA candidates.

Supplemental Figure S37. TE pieces affect expression and epigenetic patterns of AS lincRNAs and PC genes.

Supplemental Figure S38. High expression variability despite similar expressed gene numbers.

Supplemental Data Set 1. Accession overview.

Supplemental Data Set 2. RNA-seq data overview.

Supplemental Data Set 3. RNA-seq data statistics.

Supplemental Data Set 4. Overview of transcriptome assemblies.

Supplemental Data Set 5. Compiled annotations for all types of loci produced in this study (BED6 for loci and BED12 for isoforms).

Supplemental Data Set 6. ChIP-seq summary and statistics.

Supplemental Data Set 7. Summary of bisulfite-seq samples and read number.

Supplemental Data Set 8. Small RNA-seq summary and statistics.

Supplemental Data Set 9. lincRNAs reexpressed in TE silencing mutants.

Supplemental Data Set 10. AS lincRNA candidates with anticorrelated expression from their partner PC gene.

Supplemental Data Set 11. lincRNAs with AraGWAS hits.

Supplemental Data Set 12. Public RNA-seq data set summary and statistics.

Supplemental Data Set 13. Statistical tests for Figs. 1 to 8.

Funding

This study was funded by the Austrian Science Fund FWF (A.E.K. is the recipient of the Hertha Firnberg Postdoctoral Fellowship, project T-1018 “Role of long non-coding RNA variation in *Arabidopsis thaliana*”) and the European Research Area Network for Coordinating Action in Plant Sciences (ERA-CAPS project: “1001 Genomes Plus”).

Conflict of interest statement. Authors declare no conflict of interest.

References

- 1001 Genomes Consortium.** 1,135 genomes reveal the global pattern of polymorphism in *Arabidopsis thaliana*. *Cell* 2016;**166**(2):481–491. <https://doi.org/10.1016/j.cell.2016.05.063>
- Andergassen D, Dotter CP, Wenzel D, Sigl V, Bammer PC, Muckenhuber M, Mayer D, Kulinski TM, Theussl H-C, Penninger JM, et al.** Mapping the mouse Allelome reveals tissue-specific regulation of allelic expression. *eLife* 2017;**6**:e25125. <https://doi.org/10.7554/eLife.25125>
- Athie A, Marchese FP, González J, Lozano T, Raimondi I, Juvvuna PK, Abad A, Marin-Bejar O, Serizay J, Martínez D, et al.** Analysis of copy number alterations reveals the lincRNA ALAL-1 as a regulator of lung cancer immune evasion. *J Cell Biol.* 2020;**219**(9): e201908078. <https://doi.org/10.1083/jcb.201908078>
- Batista PJ, Chang HY.** Long noncoding RNAs: cellular address codes in development and disease. *Cell* 2013;**152**(6):1298–1307. <https://doi.org/10.1016/j.cell.2013.02.012>
- Bewick AJ, Schmitz RJ.** Gene body DNA methylation in plants. *Curr Opin Plant Biol.* 2017;**36**:103–110. <https://doi.org/10.1016/j.pbi.2016.12.007>
- Blein T, Balzergue C, Roulé T, Gabriel M, Scalisi L, François T, Sorin C, Christ A, Godon C, Delannoy E, et al.** Landscape of the non-coding transcriptome response of two *Arabidopsis* ecotypes to phosphate starvation. *Plant Physiol.* 2020;**183**(3):1058–1072. <https://doi.org/10.1104/pp.20.00446>
- Budak H, Kaya SB, Cagirici HB.** Long non-coding RNA in plants in the era of reference sequences. *Front Plant Sci.* 2020;**11**:276. <https://doi.org/10.3389/fpls.2020.00276>
- Cabili MN, Trapnell C, Goff L, Koziol M, Tazon-Vega B, Regev A, Rinn JL.** Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses. *Genes Dev.* 2011;**25**(18):1915–1927. <https://doi.org/10.1101/gad.17446611>
- Chen W, Zhu T, Shi Y, Chen Y, Li WJ, Chan RJ, Chen D, Zhang W, Yuan YA, Wang X, et al.** An antisense intragenic lincRNA SEAIRa mediates transcriptional and epigenetic repression of SERRATE in

- Arabidopsis. Proc Natl Acad Sci U S A. 2023;120(10):e2216062120. <https://doi.org/10.1073/pnas.2216062120>
- Cheng C-Y, Krishnakumar V, Chan AP, Thibaud-Nissen F, Schobel S, Town CD.** Araport11: a complete reannotation of the *Arabidopsis thaliana* reference genome. Plant J. 2017;89(4):789–804. <https://doi.org/10.1111/tpj.13415>
- Choi J, Lyons DB, Kim MY, Moore JD, Zilberman D.** DNA Methylation and histone H1 jointly repress transposable elements and aberrant intragenic transcripts. Mol Cell. 2020;77(2):310–323.e7. <https://doi.org/10.1016/j.molcel.2019.10.011>
- Colognori D, Sunwoo H, Wang D, Wang C-Y, Lee JT.** Xist repeats A and B account for two distinct phases of X inactivation establishment. Dev Cell. 2020;54(1):21–32.e5. <https://doi.org/10.1016/j.devcel.2020.05.021>
- Corona-Gomez JA, Coss-Navarrete EL, Garcia-Lopez IJ, Klapproth C, Pérez-Patiño JA, Fernandez-Valverde SL.** Transcriptome-guided annotation and functional classification of long non-coding RNAs in *Arabidopsis thaliana*. Sci Rep. 2022;12(1):14063. <https://doi.org/10.1038/s41598-022-18254-0>
- Cortijo S, Aydin Z, Ahnert S, Locke JC.** Widespread inter-individual gene expression variability in *Arabidopsis thaliana*. Mol Syst Biol. 2019;15(1):e8591. <https://doi.org/10.15252/msb.20188591>
- Csorba T, Questa JI, Sun Q, Dean C.** Antisense COOLAIR mediates the coordinated switching of chromatin states at FLC during vernalization. Proc Natl Acad Sci U S A. 2014;111(45):16160–16165. <https://doi.org/10.1073/pnas.1419030111>
- Délérís A, Berger F, Duharcourt S.** Role of polycomb in the control of transposable elements. Trends Genet. 2021;37(10):882–889. <https://doi.org/10.1016/j.tig.2021.06.003>
- Di Marsico M, Paytuyi Gallart A, Sanseverino W, Aiese Cigliano R.** GreeNC 2.0: a comprehensive database of plant long non-coding RNAs. Nucleic Acids Res. 2022;50(D1):D1442–D1447. <https://doi.org/10.1093/nar/gkab1014>
- Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR.** STAR: ultrafast universal RNA-seq aligner. Bioinformatics. 2013;29(1):15–21. <https://doi.org/10.1093/bioinformatics/bts635>
- Fedak H, Palusinska M, Krzyczmonik K, Brzezniak L, Yatusevich R, Pietras Z, Kaczanowski S, Swiezewski S.** Control of seed dormancy in Arabidopsis by a cis-acting noncoding antisense transcript. Proc Natl Acad Sci U S A. 2016;113(48):E7846–E7855. <https://doi.org/10.1073/pnas.1608827113>
- Feng S, Jacobsen SE.** Epigenetic modifications in plants: an evolutionary perspective. Curr Opin Plant Biol. 2011;14(2):179–186. <https://doi.org/10.1016/j.pbi.2010.12.002>
- Fultz D, Choudury SG, Slotkin RK.** Silencing of active transposable elements in plants. Curr Opin Plant Biol. 2015;27:67–76. <https://doi.org/10.1016/j.pbi.2015.05.027>
- Gullotta G, Korte A, Marquardt S.** Functional variation in the non-coding genome: molecular implications for food security. J Exp Bot. 2023;74(7):2338–2351. <https://doi.org/10.1093/jxb/erac395>
- Hansen KH, Bracken AP, Pasini D, Dietrich N, Gehani SS, Monrad A, Rappsilber J, Lerdrup M, Helin K.** A model for transmission of the H3K27me3 epigenetic mark. Nat Cell Biol. 2008;10(11):1291–1300. <https://doi.org/10.1038/ncb1787>
- He L, Huang H, Bradai M, Zhao C, You Y, Ma J, Zhao L, Lozano-Durán R, Zhu J-K.** DNA methylation-free Arabidopsis reveals crucial roles of DNA methylation in regulating gene expression and development. Nat Commun. 2022;13(1):1335. <https://doi.org/10.1038/s41467-022-28940-2>
- Hetzl J, Duttke SH, Benner C, Chory J.** Nascent RNA sequencing reveals distinct features in plant transcription. Proc Natl Acad Sci U S A. 2016;113(43):12316–12321. <https://doi.org/10.1073/pnas.1603217113>
- Ietswaart R, Wu Z, Dean C.** Flowering time control: another window to the connection between antisense RNA and chromatin. Trends Genet. 2012;28(9):445–453. <https://doi.org/10.1016/j.tig.2012.06.002>
- Ivanov M, Sandelin A, Marquardt S.** TranscriptomeReconstructoR: data-driven annotation of complex transcriptomes. BMC Bioinformatics. 2021;22:290.
- Jin J, Lu P, Xu Y, Li Z, Yu S, Liu J, Wang H, Chua N-H, Cao P.** PLncDB V2.0: a comprehensive encyclopedia of plant long noncoding RNAs. Nucleic Acids Res. 2021;49(D1):D1489–D1495. <https://doi.org/10.1093/nar/gkaa910>
- Johnson-Brousseau SA, McCormick S.** A compendium of methods useful for characterizing Arabidopsis pollen mutants and gametophytically-expressed genes. Plant J. 2004;39(5):761–775. <https://doi.org/10.1111/j.1365-313X.2004.02147.x>
- Johnson R, Guigó R.** The RIDL hypothesis: transposable elements as functional domains of long noncoding RNAs. RNA. 2014;20(7):959–976. <https://doi.org/10.1261/rna.044560.114>
- Johnsson P, Lipovich L, Grandér D, Morris KV.** Evolutionary conservation of long non-coding RNAs; sequence, structure, function. Biochim Biophys Acta. 2014;1840(3):1063–1071. <https://doi.org/10.1016/j.bbagen.2013.10.035>
- Kang C, Liu Z.** Global identification and analysis of long non-coding RNAs in diploid strawberry *Fragaria vesca* during flower and fruit development. BMC Genomics. 2015;16:815. <https://doi.org/10.1186/s12864-015-2014-2>
- Kang Y-J, Yang D-C, Kong L, Hou M, Meng Y-Q, Wei L, Gao G.** CPC2: a fast and accurate coding potential calculator based on sequence intrinsic features. Nucleic Acids Res. 2017;45(W1):W12–W16. <https://doi.org/10.1093/nar/gkx428>
- Kapusta A, Feschotte C.** Volatile evolution of long noncoding RNA repertoires: mechanisms and biological implications. Trends Genet. 2014;30(10):439–452. <https://doi.org/10.1016/j.tig.2014.08.004>
- Kapusta A, Kronenberg Z, Lynch VJ, Zhuo X, Ramsay L, Bourque G, Yandell M, Feschotte C.** Transposable elements are major contributors to the origin, diversification, and regulation of vertebrate long noncoding RNAs. PLoS Genet. 2013;9(4):e1003470. <https://doi.org/10.1371/journal.pgen.1003470>
- Kawakatsu T, Huang SC, Jupe F, Sasaki E, Schmitz RJ, Urlich MA, Castanon R, Nery JR, Barragan C, He Y, et al.** Epigenomic diversity in a global collection of *Arabidopsis thaliana* accessions. Cell. 2016;166(2):492–505. <https://doi.org/10.1016/j.cell.2016.06.044>
- Kindgren P, Ard R, Ivanov M, Marquardt S.** Transcriptional read-through of the long non-coding RNA SVALKKA governs plant cold acclimation. Nat Commun. 2018;9(1):4561. <https://doi.org/10.1038/s41467-018-07010-6>
- Kindgren P, Ivanov M, Marquardt S.** Native elongation transcript sequencing reveals temperature dependent dynamics of nascent RNAPII transcription in Arabidopsis. Nucleic Acids Res. 2020;48(5):2332–2347. <https://doi.org/10.1093/nar/gkz1189>
- Kornienko AE, Dotter CP, Guenzl PM, Gisslinger H, Gisslinger B, Cleary C, Kralovics R, Pauler FM, Barlow DP.** Long non-coding RNAs display higher natural expression variation than protein-coding genes in healthy humans. Genome Biol. 2016;17:14. <https://doi.org/10.1186/s13059-016-0873-8>
- Kutter C, Watt S, Stefflova K, Wilson MD, Goncalves A, Ponting CP, Odom DT, Marques AC.** Rapid turnover of long noncoding RNAs and the evolution of gene expression. PLoS Genet. 2012;8(7):e1002841. <https://doi.org/10.1371/journal.pgen.1002841>
- Leone S, Santoro R.** Challenges in the analysis of long noncoding RNA functionality. FEBS Lett. 2016;590(15):2342–2353. <https://doi.org/10.1002/1873-3468.12308>
- Li L, Eichten SR, Shimizu R, Petsch K, Yeh C-T, Wu W, Chetoor AM, Givan SA, Cole RA, Fowler JE, et al.** Genome-wide discovery and characterization of maize long non-coding RNAs. Genome Biol. 2014;15(2):R40. <https://doi.org/10.1186/gb-2014-15-2-r40>
- Li S, Liberman LM, Mukherjee N, Benfey PN, Ohler U.** Integrated detection of natural antisense transcripts using strand-specific RNA sequencing data. Genome Res. 2013;23(10):1730–1739. <https://doi.org/10.1101/gr.149310.112>

- Liao Y, Smyth GK, Shi W.** Featurecounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics*. 2014;**30**(7):923–930. <https://doi.org/10.1093/bioinformatics/btt656>
- Liu X, Hao L, Li D, Zhu L, Hu S.** Long non-coding RNAs and their biological roles in plants. *Genomics Proteomics Bioinformatics*. 2015;**13**(3):137–147. <https://doi.org/10.1016/j.gpb.2015.02.003>
- Liu J, Jung C, Xu J, Wang H, Deng S, Bernad L, Arenas-Huertero C, Chua N-H.** Genome-wide analysis uncovers regulation of long intergenic noncoding RNAs in *Arabidopsis*. *Plant Cell*. 2012;**24**(11):4333–4345. <https://doi.org/10.1105/tpc.112.102855>
- Lubelsky Y, Ulitsky I.** Sequences enriched in Alu repeats drive nuclear localization of long RNAs in human cells. *Nature* 2018;**555**(7694):107–111. <https://doi.org/10.1038/nature25757>
- Martin M.** Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal*. 2011;**17**(1):10–12. <https://doi.org/10.14806/ej.17.1.200>
- Mattick JS, Amaral PP, Carninci P, Carpenter S, Chang HY, Chen L-L, Chen R, Dean C, Dinger ME, Fitzgerald KA, et al.** Long non-coding RNAs: definitions, functions, challenges and recommendations. *Nat Rev Mol Cell Biol*. 2023;**24**(6):430–447. <https://doi.org/10.1038/s41580-022-00566-8>
- Mattick JS, Rinn JL.** Discovery and annotation of long noncoding RNAs. *Nat Struct Mol Biol*. 2015;**22**(1):5–7. <https://doi.org/10.1038/nsmb.2942>
- Matzke MA, Mosher RA.** RNA-directed DNA methylation: an epigenetic pathway of increasing complexity. *Nat Rev Genet*. 2014;**15**(6):394–408. <https://doi.org/10.1038/nrg3683>
- Melé M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, Young TR, Goldmann JM, Pevouchine DD, Sullivan TJ, et al.** The human transcriptome across tissues and individuals. *Science* 2015;**348**(6235):660–665. <https://doi.org/10.1126/science.aaa0355>
- Mhiri C, Borges F, Grandbastien M-A.** Specificities and dynamics of transposable elements in land plants. *Biology (Basel)*. 2022;**11**(4):488. <https://doi.org/10.3390/biology11040488>
- Necsulea A, Kaessmann H.** Evolutionary dynamics of coding and non-coding transcriptomes. *Nat Rev Genet*. 2014;**15**(11):734–748. <https://doi.org/10.1038/nrg3802>
- Nelson ADL, Devisetty UK, Palos K, Haug-Baltzell AK, Lyons E, Beilstein MA.** Evolinc: a tool for the identification and evolutionary comparison of long intergenic non-coding RNAs. *Front Genet*. 2017;**8**:52. <https://doi.org/10.3389/fgene.2017.00052>
- Nguyen VH, Scheid OM, Gutzat R.** Heat stress response and transposon control in plant shoot stem cells. *bioRxiv* 2023.02.24.529891. <https://doi.org/10.1101/2023.02.24.529891>, 26 February 2023, preprint: not peer reviewed.
- Nuthikattu S, McCue AD, Panda K, Fultz D, DeFraia C, Thomas EN, Slotkin RK.** The initiation of epigenetic silencing of active transposable elements is triggered by RDR6 and 21-22 nucleotide small interfering RNAs. *Plant Physiol*. 2013;**162**(1):116–131. <https://doi.org/10.1104/pp.113.216481>
- Oki N, Yano K, Okumoto Y, Tsukiyama T, Teraishi M, Tanisaka T.** A genome-wide view of miniature inverted-repeat transposable elements (MITEs) in rice, *Oryza sativa* ssp. *japonica*. *Genes Genet Syst*. 2008;**83**(4):321–329. <https://doi.org/10.1266/ggs.83.321>
- Onodera Y, Haag JR, Ream T, Costa Nunes P, Pontes O, Pikaard CS.** Plant nuclear RNA polymerase IV mediates siRNA and DNA methylation-dependent heterochromatin formation. *Cell* 2005;**120**(5):613–622. <https://doi.org/10.1016/j.cell.2005.02.007>
- Osakabe A, Jamge B, Axelsson E, Montgomery SA, Akimcheva S, Kuehn AL, Pisupati R, Lorković ZJ, Yelagandula R, Kakutani T, et al.** The chromatin remodeler DDM1 prevents transposon mobility through deposition of histone variant H2A.W. *Nat Cell Biol*. 2021;**23**(4):391–400. <https://doi.org/10.1038/s41556-021-00658-1>
- Palos K, Nelson Dittrich AC, Yu L, Brock JR, Railey CE, Wu H-YL, Sokolowska E, Skirycz A, Hsu PY, Gregory BD, et al.** Identification and functional annotation of long intergenic non-coding RNAs in Brassicaceae. *Plant Cell*. 2022;**34**(9):3233–3260. <https://doi.org/10.1093/plcell/koac166>
- Palos K, Yu L, Railey CE, Nelson Dittrich AC, Nelson ADL.** Linking discoveries, mechanisms, and technologies to develop a clearer perspective on plant long noncoding RNAs. *Plant Cell*. 2023;**35**(6):1762–1786. <https://doi.org/10.1093/plcell/koad027>
- Panda K, McCue AD, Slotkin RK.** *Arabidopsis* RNA polymerase IV generates 21-22 nucleotide small RNAs that can participate in RNA-directed DNA methylation and may regulate genes. *Philos Trans R Soc Lond B Biol Sci*. 2020;**375**(1795):20190417. <https://doi.org/10.1098/rstb.2019.0417>
- Papareddy RK, Páldi K, Paulraj S, Kao P, Lutzmayr S, Nodine MD.** Chromatin regulates expression of small RNAs to help maintain transposon methylome homeostasis in *Arabidopsis*. *Genome Biol*. 2020;**21**(1):251. <https://doi.org/10.1186/s13059-020-02163-4>
- Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL.** Stringtie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*. 2015;**33**(3):290–295. <https://doi.org/10.1038/nbt.3122>
- Pisupati R, Nizhynska V, Morales AM, Nordborg M.** On the causes of gene-body methylation variation in *Arabidopsis thaliana*. *bioRxiv* 2022.12.04.519028. <https://doi.org/10.1101/2022.12.04.519028>, 04 December 2022, preprint: not peer reviewed.
- Pontier D, Picart C, Roudier F, Garcia D, Lahmy S, Azevedo J, Alart E, Laudí M, Karłowski WM, Cooke R, et al.** NERD, a plant-specific GW protein, defines an additional RNAi-dependent chromatin-based pathway in *Arabidopsis*. *Mol Cell*. 2012;**48**(1):121–132. <https://doi.org/10.1016/j.molcel.2012.07.027>
- Quesneville H.** Twenty years of transposable element analysis in the *Arabidopsis thaliana* genome. *Mob DNA*. 2020;**11**:28. <https://doi.org/10.1186/s13100-020-00223-x>
- Rajpal VR, Rathore P, Mehta S, Wadhwa N, Yadav P, Berry E, Goel S, Bhat V, Raina SN.** Epigenetic variation: a major player in facilitating plant fitness under changing environmental conditions. *Front Cell Dev Biol*. 2022;**10**:1020958. <https://doi.org/10.3389/fcell.2022.1020958>
- Ramírez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, Heyne S, Dündar F, Manke T.** DeepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Res*. 2016;**44**(W1):W160–W165. <https://doi.org/10.1093/nar/gkw257>
- Rinn JL, Chang HY.** Long noncoding RNAs: molecular modalities to organismal functions. *Annu Rev Biochem*. 2020;**89**:283–308. <https://doi.org/10.1146/annurev-biochem-062917-012708>
- Sasaki E, Kawakatsu T, Ecker JR, Nordborg M.** Common alleles of CMT2 and NRPE1 are major determinants of CHH methylation variation in *Arabidopsis thaliana*. *PLoS Genet*. 2019;**15**(12):e1008492. <https://doi.org/10.1371/journal.pgen.1008492>
- Sauvageau M, Goff LA, Lodato S, Bonev B, Groff AF, Gerhardinger C, Sanchez-Gomez DB, Hacisuleyman E, Li E, Spence M, et al.** Multiple knockout mouse models reveal lincRNAs are required for life and brain development. *eLife*. 2013;**2**:e01749. <https://doi.org/10.7554/eLife.01749>
- Sigman MJ, Slotkin RK.** The first rule of plant transposable element silencing: location, location, location. *Plant Cell*. 2016;**28**(2):304–313. <https://doi.org/10.1105/tpc.15.00869>
- Slotkin RK, Vaughn M, Borges F, Tanurdzić M, Becker JD, Feijó JA, Martienssen RA.** Epigenetic reprogramming and small RNA silencing of transposable elements in pollen. *Cell* 2009;**136**(3):461–472. <https://doi.org/10.1016/j.cell.2008.12.038>
- Statello L, Guo C-J, Chen L-L, Huarte M.** Gene regulation by long non-coding RNAs and its biological functions. *Nat Rev Mol Cell Biol*. 2021;**22**(2):96–118. <https://doi.org/10.1038/s41580-020-00315-9>
- Sundaram V, Wysocka J.** Transposable elements as a potent source of diverse cis-regulatory sequences in mammalian genomes. *Philos Trans R Soc Lond B Biol Sci*. 2020;**375**(1795):20190347. <https://doi.org/10.1098/rstb.2019.0347>
- Szczeniak MW, Bryzgalov O, Ciombrowska-Basheer J, Makalowska I.** CANTATAdb 2.0: expanding the collection of plant long noncoding

- RNAs. *Methods Mol Biol.* 2019;**1933**:415–429. https://doi.org/10.1007/978-1-4939-9045-0_26
- Thomas QA, Ard R, Liu J, Li B, Wang J, Pelechano V, Marquardt S.** Transcript isoform sequencing reveals widespread promoter-proximal transcriptional termination in *Arabidopsis*. *Nat Commun.* 2020;**11**(1):2589. <https://doi.org/10.1038/s41467-020-16390-7>
- Togninalli M, Seren Ü, Freudenthal JA, Monroe JG, Meng D, Nordborg M, Weigel D, Borgwardt K, Korte A, Grimm DG.** Arapheno and the AraGWAS Catalog 2020: a major database update including RNA-Seq and knockout mutation data for *Arabidopsis thaliana*. *Nucleic Acids Res.* 2020;**48**(D1):D1063–D1068. <https://doi.org/10.1093/nar/gkz925>
- Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ, Pachter L.** Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol.* 2010;**28**(5):511–515. <https://doi.org/10.1038/nbt.1621>
- Volders P-J, Anckaert J, Verheggen K, Nuytens J, Martens L, Mestdagh P, Vandesompele J.** LNCipedia 5: towards a reference set of human long non-coding RNAs. *Nucleic Acids Res.* 2019;**47**(D1):D135–D139. <https://doi.org/10.1093/nar/gky1031>
- Wahlestedt C.** Targeting long non-coding RNA to therapeutically up-regulate gene expression. *Nat Rev Drug Discov.* 2013;**12**(6):433–446. <https://doi.org/10.1038/nrd4018>
- Wang Q, Gu L, Adey A, Radlwimmer B, Wang W, Hovestadt V, Bähr M, Wolf S, Shendure J, Eils R, et al.** Tagmentation-based whole-genome bisulfite sequencing. *Nat Protoc.* 2013;**8**(10):2022–2032. <https://doi.org/10.1038/nprot.2013.118>
- Wang J, Meng X, Dobrovolskaya OB, Orlov YL, Chen M.** Non-coding RNAs and their roles in stress response in plants. *Genomics Proteomics Bioinformatics.* 2017;**15**(5):301–312. <https://doi.org/10.1016/j.gpb.2017.01.007>
- Wapinski O, Chang HY.** Long noncoding RNAs and human disease. *Trends Cell Biol.* 2011;**21**(6):354–361. <https://doi.org/10.1016/j.tcb.2011.04.001>
- Whittaker C, Dean C.** The FLC locus: a platform for discoveries in epigenetics and adaptation. *Annu Rev Cell Dev Biol.* 2017;**33**:555–575. <https://doi.org/10.1146/annurev-cellbio-100616-060546>
- Xin M, Wang Y, Yao Y, Song N, Hu Z, Qin D, Xie C, Peng H, Ni Z, Sun Q.** Identification and characterization of wheat long non-protein coding RNAs responsive to powdery mildew infection and heat stress by using microarray analysis and SBS sequencing. *BMC Plant Biol.* 2011;**11**:61. <https://doi.org/10.1186/1471-2229-11-61>
- Xu J, Chen G, Hermanson PJ, Xu Q, Sun C, Chen W, Kan Q, Li M, Crisp PA, Yan J, et al.** Population-level analysis reveals the widespread occurrence and phenotypic consequence of DNA methylation variation not tagged by genetic variation in maize. *Genome Biol.* 2019;**20**(1):243. <https://doi.org/10.1186/s13059-019-1859-0>
- Xu Y, Wu T, Li F, Dong Q, Wang J, Shang D, Xu Y, Zhang C, Dou Y, Hu C, et al.** Identification and comprehensive characterization of lncRNAs with copy number variations and their driving transcriptional perturbed subpathways reveal functional significance for cancer. *Brief Bioinform.* 2020;**21**(6):2153–2166. <https://doi.org/10.1093/bib/bbz113>
- Xu Y-C, Zhang J, Zhang D-Y, Nan Y-H, Ge S, Guo Y-L.** Identification of long noncoding natural antisense transcripts (lncNATs) correlated with drought stress response in wild rice (*Oryza nivara*). *BMC Genomics.* 2021;**22**(1):424. <https://doi.org/10.1186/s12864-021-07754-4>
- Yang W, Bai Q, Li Y, Chen J, Liu C.** Epigenetic modifications: allusive clues of lncRNA functions in plants. *Comput Struct Biotechnol J.* 2023;**21**:1989–1994. <https://doi.org/10.1016/j.csbj.2023.03.008>
- Yang P, Wang Y, Macfarlan TS.** The role of KRAB-ZFPs in transposable element repression and mammalian evolution. *Trends Genet.* 2017;**33**(11):871–881. <https://doi.org/10.1016/j.tig.2017.08.006>
- Yelagandula R, Stroud H, Holec S, Zhou K, Feng S, Zhong X, Muthurajan UM, Nie X, Kawashima T, Groth M, et al.** The histone variant H2A.W defines heterochromatin and promotes chromatin condensation in *Arabidopsis*. *Cell* 2014;**158**(1):98–109. <https://doi.org/10.1016/j.cell.2014.06.006>
- Yuan C, Wang J, Harrison AP, Meng X, Chen D, Chen M.** Genome-wide view of natural antisense transcripts in *Arabidopsis thaliana*. *DNA Res.* 2015;**22**(3):233–243. <https://doi.org/10.1093/dnares/dsv008>
- Yuan J, Zhang Y, Dong J, Sun Y, Lim BL, Liu D, Lu ZJ.** Systematic characterization of novel lncRNAs responding to phosphate starvation in *Arabidopsis thaliana*. *BMC Genomics.* 2016;**17**:655. <https://doi.org/10.1186/s12864-016-2929-2>
- Zemach A, Kim MY, Hsieh P-H, Coleman-Derr D, Eshed-Williams L, Thao K, Harmer SL, Zilberman D.** The *Arabidopsis* nucleosome remodeler DDM1 allows DNA methyltransferases to access H1-containing heterochromatin. *Cell* 2013;**153**(1):193–205. <https://doi.org/10.1016/j.cell.2013.02.033>
- Zhao X, Li J, Lian B, Gu H, Li Y, Qi Y.** Global identification of *Arabidopsis* lncRNAs reveals the regulation of MAF4 by a natural antisense RNA. *Nat Commun.* 2018;**9**(1):5056. <https://doi.org/10.1038/s41467-018-07500-7>
- Zhao L, Zhou Q, He L, Deng L, Lozano-Duran R, Li G, Zhu J-K.** DNA methylation underpins the epigenomic landscape regulating genome transcription in *Arabidopsis*. *Genome Biol.* 2022;**23**(1):197. <https://doi.org/10.1186/s13059-022-02768-x>
- Zhu Y, Chen L, Hong X, Shi H, Li X.** Revealing the novel complexity of plant long non-coding RNA by strand-specific and whole transcriptome sequencing for evolutionarily representative plant species. *BMC Genomics.* 2022;**23**(S4):381. <https://doi.org/10.1186/s12864-022-08602-9>