

Splicing conservation signals in plant long noncoding RNAs

JOSE ANTONIO CORONA-GOMEZ,¹ IRVING JAIR GARCIA-LOPEZ,¹ PETER F. STADLER,^{2,3,4,5,6,7}
and SELENE L. FERNANDEZ-VALVERDE¹

¹Unidad de Genómica Avanzada, Langebio, Cinvestav, 36821 Irapuato, Guanajuato, Mexico

²Bioinformatics Group, Department of Computer Science, University Leipzig, D-04107 Leipzig, Germany

³Interdisciplinary Center for Bioinformatics, University Leipzig, D-04107 Leipzig, Germany

⁴Max Planck Institute for Mathematics in the Sciences, D-04103 Leipzig, Germany

⁵Department of Theoretical Chemistry, University of Vienna, A-1090 Wien, Austria

⁶Facultad de Ciencias, Universidad Nacional de Colombia, 11001 Sede Bogotá, Colombia

⁷Santa Fe Institute, Santa Fe, New Mexico 87501, USA

ABSTRACT

Long noncoding RNAs (lncRNAs) have recently emerged as prominent regulators of gene expression in eukaryotes. lncRNAs often drive the modification and maintenance of gene activation or gene silencing states via chromatin conformation rearrangements. In plants, lncRNAs have been shown to participate in gene regulation, and are essential to processes such as vernalization and photomorphogenesis. Despite their prominent functions, only over a dozen lncRNAs have been experimentally and functionally characterized. Similar to its animal counterparts, the rates of sequence divergence are much higher in plant lncRNAs than in protein coding mRNAs, making it difficult to identify lncRNA conservation using traditional sequence comparison methods. Beyond this, little is known about the evolutionary patterns of lncRNAs in plants. Here, we characterized the splicing conservation of lncRNAs in Brassicaceae. We generated a whole-genome alignment of 16 Brassica species and used it to identify syntenic lncRNA orthologs. Using a scoring system trained on transcripts from *A. thaliana* and *B. oleracea*, we identified splice sites across the whole alignment and measured their conservation. Our analysis revealed that 17.9% (112/627) of all intergenic lncRNAs display splicing conservation in at least one exon, an estimate that is substantially higher than previous estimates of lncRNA conservation in this group. Our findings agree with similar studies in vertebrates, demonstrating that splicing conservation can be evidence of stabilizing selection. We provide conclusive evidence for the existence of evolutionary deeply conserved lncRNAs in plants and describe a generally applicable computational workflow to identify functional lncRNAs in plants.

Keywords: long noncoding RNAs; lncRNA; splice sites; multiple sequence alignments; evolution; conservation; evolutionary plasticity

INTRODUCTION

Long noncoding RNAs (lncRNAs), by definition, do not code for proteins. Over the last decade, a wide variety of mechanisms have been discovered by which lncRNAs contribute to the regulation of the expression of protein-coding genes and small RNAs (Chekanova 2015; Liu et al. 2015; Ulitsky 2016; Wang and Chekanova 2017; Yamada 2017). Most lncRNAs are found in the nucleus associated with the chromatin, regulating gene expression by recruiting components of the epigenetic machinery to specific genomic locations. Some lncRNAs also influence genome stability and nuclear domain organization. Serving as molecular sponges and decoys, they act both at the transcrip-

tional level, by affecting RNA-directed DNA methylation; in post-transcriptional regulation, by inhibiting the interaction between microRNAs (miRNAs) and their target messenger RNAs (mRNAs); and by controlling alternative splicing due to sequestration of splicing factors (Bardou et al. 2014). Hence, they differ not only in size but also in their biogenesis and molecular mechanisms from small RNAs such as miRNAs and siRNAs (Bánfai et al. 2012). lncRNAs are regulated and processed similar to mRNAs (Mercer and Mattick 2013) and their expression patterns are often very specific to particular tissues or developmental stages (Mercer and Mattick 2013). Recent data suggest

Corresponding authors: studla@bioinf.unileipzig.de, selene.fernandez@cinvestav.mx

Article is online at <http://www.majournal.org/cgi/doi/10.1261/ma.074393.119>.

© 2020 Corona-Gomez et al. This article is distributed exclusively by the RNA Society for the first 12 months after the full-issue publication date (see <http://majournal.cshlp.org/site/misc/terms.xhtml>). After 12 months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

that there appears to be a distinction between highly conserved, constitutively transcribed lncRNAs and tissue-specific lncRNAs with low expression levels (Deng et al. 2018b; Sarropoulos et al. 2019).

Despite their often very poor sequence conservation (Necsulea et al. 2014), the majority of lncRNAs are well-conserved across animals, as evidenced by the conservation of many of their splice sites (Nitsche et al. 2015). While well-conserved as entities, they show much more plasticity in their gene structure and sequence than protein-coding genes. The many lineage-specific differences have implicated lncRNAs as major players in lineage-specific adaptation (Lozada-Chávez et al. 2011): Changes in transcript structure are likely associated with the inclusion or exclusion of sets of protein or miRNA binding sites and hence may have large effects on function and specificity of a particular lncRNA.

The systematic annotation of orthologous lncRNAs is important not only to provide reasonably complete maps of the transcriptome, but also as a means of establishing that a particular lncRNA has a biological function. After all, conservation over long evolutionary timescales is often used as the most important argument for the biological function of an open reading frame in the absence of direct experimental evidence for translation and experimental data characterizing the peptide product. While a large amount of work is available showing that vertebrate genomes contain a large number of secondary elements that are under negative selection (Smith et al. 2013; Hezroni et al. 2015; Nitsche and Stadler 2017; Seemann et al. 2017) and the majority of human lncRNAs are evolutionary old (Nitsche et al. 2015), a much less systematic and complete picture is available for plants.

In fact, detailed studies into the evolution of plant lncRNAs have been rare until very recently. An analysis of lncRNAs in five monocot and five dicot species (Deng et al. 2018b) found that the majority of lncRNAs are poorly conserved at sequence level while a majority is highly divergent but syntenically conserved. These positionally conserved lncRNAs were previously found to be located near telomeres in *A. thaliana* (Mohammadin et al. 2015). Plant lncRNAs have also been shown to display canonical splicing signals (Deng et al. 2018b). Another study in 10 Brassicaceae genomes found 22% conservation of intergenic lncRNA loci (Nelson et al. 2016), as well as little evidence of an impact of whole-genome duplications or transposable element (TE) activity on the emergence of lincRNAs.

Nevertheless, there are some plant lncRNAs whose regulatory functions have been studied extensively and are understood at a level of detail comparable to most proteins (Rai et al. 2019): *COOLAIR* in Brassicaceae has a crucial role in the vernalization process (Hawkes et al. 2016) and its transcription accelerates epigenetic silencing of the flowering locus C (*FLC*) (Rosa et al. 2016). The

lncRNA *HID1* is a key component in promoting photomorphogenesis in response to different levels of red light (Wang et al. 2014). *HID1* is highly conserved and acts through binding to chromatin in *trans* to act upon the *PIF3* promoter. A similar trans-acting lncRNA is *ELENA1*, which functions in plant immunity (Mach 2017). Competing endogenous RNAs (ceRNAs) acts as “sponges” for miRNAs. In plants, ceRNAs are a large class of lncRNAs (Yuan et al. 2017; Paschoal et al. 2018) and form extensive regulatory networks (Meng et al. 2018; Zhang et al. 2018). The paradigmatic example in *A. thaliana* is *IPS1*, which sequesters miR399, resulting in changes in phosphate homeostasis (Franco-Zorrilla et al. 2007).

Although the functional characterization of plant lncRNAs is confined to a small number of cases, plant lncRNAs are being reported at a rapidly increasing pace (Nelson et al. 2016). As in the case of animals, it is important therefore to amass evidence for the functionality of individual transcripts. Differential expression, or correlations with important regulatory proteins or pathways alone do not provide sufficient evidence to decide whether a transcript has a causal effect or whether its expression pattern is a coincidental downstream effect. As a first step toward prioritizing candidates for functional characterization, we advocate for the use of unexpected deep conservation of the gene structure as an indicator of biological function. While logically this still does not inform about function in a specific context, it is much less likely that changes in expression patterns of a conserved and thus presumably functional molecule are without biological consequence.

The much higher level of plasticity in plant genomes, compared to animal genomes, potentially makes it more difficult to trace the evolution of lncRNAs. We therefore concentrate here on a phylogenetically relatively narrow group, the Brassicaceae, with genomes that are largely alignable with each other. We track the conservation of functional elements, in particular splice junctions, through the entire data set. This provides direct evidence also in cases where transcriptome data is not available in sufficient coverage and or sufficient diversity of tissues and/or developmental stages. As a final result, we provide a list of homologous lncRNAs in Brassicaceae as well as a detailed map of the conservation of splice sites in this clade.

RESULTS

Identification of splice sites and lncRNAs

To build a *A. thaliana* splice junction reference, we identified about 125,000 introns using the transcriptomes of Liu et al. (2012) compared with 175,000 introns annotated in TAIR10 (Release 38) (Berardini et al. 2015). The smaller number was expected as (i) only introns with convincing coverage by uniquely mapping reads were considered

and (ii) not all *A. thaliana* genes are expressed in these four transcriptomes. Consistent with previous reports (Brown et al. 1996; Hebsgaard 1996), the vast majority of the detected splice junctions have the canonical GT/AG motif required for inclusion into our splice site map. In total, we identified 222,772 individual sites in *A. thaliana* (117,644 donor and 121,002 acceptor sites). 55% of all donors but only 13% of the acceptors have aligned sequences in the WGA (Supplemental Fig. S3). In addition, many splice sites have evidence of expression in transcriptome data from other species (Supplemental Table S3).

To characterize splicing conservation in lincRNAs, we focused solely on intergenic long noncoding RNAs (lincRNAs). Conservation of splice sites in lincRNAs overlapping with coding genes may be confounded by the coding gene conservation signal, resulting in false positives. The lincRNAs described by Liu et al. (2012) comprise 595 lincRNAs with predicted introns, with only 18 with confirmed introns as annotated in Araport9 (Liu et al. 2012), while in Araport11 (Cheng et al. 2017) 288 annotated lincRNAs out of 2444 have introns. We also used an additional set of lincRNAs expressed in *A. thaliana* cotyledons and hypocotyls in Col-0 plants in normal light or shade conditions (Kohnen et al. 2016). These libraries were stranded, and had three replicates as well as sufficient depth to produce a high confidence lincRNA annotation. As these transcriptomes are only derived from two experimental conditions (shadow and light) (Kohnen et al. 2016), they encompass only a fraction of the lincRNAs expressed throughout the *A. thaliana* life cycle. We identified 2375 lincRNA transcripts, 1465 of which overlapped with protein coding RNAs, while 808 were found in intergenic regions and were thus considered bona fide lincRNAs. In our analysis, we found 159 lincRNAs that were included in neither Araport11 nor TAIR10 (Berardini et al. 2015; Cheng et al. 2017). Furthermore, we excluded all lincRNAs that had any overlap with other annotated ncRNAs thus depleting our set of lincRNAs that may be microRNA or snoRNA precursors, as small RNAs are generally conserved. All 808 lincRNAs transcripts aggregated in 627 lincRNA genes, of which 58 have multiple isoforms. In contrast to the situation in animals, lincRNAs are therefore mostly mono-exonic in *A. thaliana*. Of the 627 lincRNA genes, only 173 had at least one intron and thus were used to test splice site conservation in lincRNAs; of these 173, only 35 were previously annotated in the Araport11 database.

Conservation of lincRNAs

To identify conserved elements by position, we extracted aligned sequences corresponding to different annotations sets from the WGA. Between 69.6% to 44.2% of the *A. thaliana* genome was aligned with other Brassicaceae species. For the protein-coding genes annotated in Araport11 (Cheng et al. 2017), the alignment recovery rate ranges

from 95.3% (26,153/27,445) (*A. lyrata*) to 86.9% (23,856/27,445) (*Aethionema arabicum*). As expected, the values are substantially lower for the Araport11 lincRNAs, where we recover between 77.1% (1885/2444) in *A. lyrata* and 50.8% (1243/2444) in *A. arabicum*. Using our own annotation, we recover between 62.0% (389/627) in *A. lyrata* and 38.1% (239/627) in *A. arabicum*, i.e., values comparable to the overall coverage of the genome. This reflects the fact that lincRNA sequences experience very little constraint on their sequence. Conservation (as measured by alignability) is summarized in Figure 1 for different types of RNA elements. These values are comparable to a previous estimate of ~22% of the lincRNA loci are at least partially conserved at the sequence level in the last common ancestor of Brassicaceae (Nelson et al. 2016).

Conservation of splice sites is a strong indication for the functionality of the transcript. In order to evaluate splice site conservation quantitatively, we constructed a splicing map that identifies for every experimentally determined splice site the homologous position in the other genomes and evaluates them using the MES (see Materials and Methods for details). Figure 2 shows the splicing map for the lincRNA TCONS00053212-00053217 as an illustrative example. Despite the unusually complex transcript structure and the conservation throughout the Brassicaceae, so far nothing is known about the function of this lincRNA. While not all splice sites are represented in all species in the WGA, almost all MES values in this lincRNA are well above the threshold of $MES > 0$. This contrasts with a random sampling of splice sites in coding and noncoding regions in all genomes in the WGA (Supplemental Fig. S4). Indeed, the probability of identifying a random splice site with an MES value greater than 0 in *A. thaliana* is 0.0237 (acceptor) and 0.0165 (donor) for coding genes, and 0.0225 (acceptor) and 0.0168 (donor) for lincRNAs (Supplemental Fig. S4). Most of this lincRNA isoforms therefore can be expected to be present throughout the Brassicaceae, even though the locus is not annotated in Ensembl Plants (release 42) for *B. oleracea*, *B. rapa*, and *A. lyrata*. Only the short first exon and the 5' most acceptor of the last exon are poorly conserved by sequence even in close relatives of *A. thaliana*.

In order to validate the predicted lincRNA splice sites, we investigated publicly available RNA-seq data from eight of the species included in this study (Supplemental Table S2). The depth of these data varied considerably. We therefore compared the fraction of recovered lincRNA predictions with the fraction of mRNAs that were detectable in the same RNA-seq data (Supplemental Table S4). As expected, we observed that the relative validation rate increases with the depth of data, presumably owing to the fact that lincRNAs are on average less highly expressed and more specifically expressed than mRNAs. Nevertheless, the validation rate in our data of lincRNAs is on average 33.3% in the eight species used for validation

	own		Araport 11					
	lincRNA	lincRNA	NAT	mRNAs	Pseudocoding	miRNA	TE	snoRNA
<i>Arabidopsis thaliana</i>	627	2,444	1,037	27,445	941	325	3,897	287
<i>Arabidopsis lyrata</i>	389	1,885	1,010	26,153	672	244	1,402	278
<i>Arabidopsis halleri</i>	346	1,662	957	25,333	611	208	974	232
<i>Camelina sativa</i>	335	1,700	1,009	25,592	593	205	915	270
<i>Capsella rubella</i>	300	1,538	996	25,008	516	191	703	266
<i>Boechera stricta</i>	325	1,659	1,004	25,382	548	203	730	258
<i>Leavenworthia alabamica</i>	253	1,350	973	24,300	481	158	411	221
<i>Arabis alpina</i>	249	1,360	972	24,301	489	173	506	211
<i>Sisymbrium irio</i>	244	1,383	987	24,352	490	161	478	209
<i>Eutrema salsugineum</i>	252	1,375	979	24,330	485	155	488	217
<i>Thellungiella parvula</i>	249	1,379	987	24,338	488	158	404	216
<i>Raphanus sativus</i>	230	1,342	977	24,177	469	148	462	208
<i>Brassica rapa</i>	229	1,323	965	23,972	459	145	391	193
<i>Brassica napus</i>	234	1,332	970	24,137	469	145	393	202
<i>Brassica oleracea</i>	232	1,311	964	23,997	457	144	375	201
<i>Aethionema arabicum</i>	239	1,243	964	23,856	457	147	450	237

FIGURE 1. Conservation of genes by position in WGA. Own: lincRNAs genes expressed in shade experiments (Kohnen et al. 2016). *Araport11* database annotations (Cheng et al. 2017): lincRNAs (long intergenic noncoding RNAs), NAT (Natural antisense transcripts), Coding genes (messenger RNAs), miRNA (microRNAs), Pseudocoding (Pseudocoding genes), TE (Transposable elements), snoRNA (Small nucleolar RNAs)

and 10.7% for lincRNAs in Araport11 (Cheng et al. 2017), while for coding genes it is 57.5%.

On a genome-wide scale, the conservation of splice sites in lincRNAs provides a lower bound on the fraction of lincRNAs that are under selective constraint as a transcript. We find that 112 of the 173 spliced *A. thaliana* lincRNAs have at least one conserved splice site in another species (Fig. 3).

As expected, we find that splice sites in lincRNAs are much less well conserved than splice sites in protein coding genes (Fig. 4). In total, we identified 39 lincRNAs conserved between the most distant species and *A. thaliana* and 26 lincRNAs with conservation in at least one splice site in the 16 species included in the WGA. These numbers are much lower than for coding genes. Albeit this is expected, given the high conservation of protein coding genes, one has to keep in mind that coding genes on average have at least six introns (Deng et al. 2018b), hence it is much more likely to observe conservation of at least one splice site and in lincRNAs with only one or two introns (see Fig. 3).

The potential incompleteness of annotated lncRNAs, for example, due to low expression levels, is of concern in this context. It has little influence on our conclusions, however, since incomplete or fragmented annotation only causes us to underestimate the depth of conservation: We might occasionally miss the best-conserved splice junction and we might count fragments as independent, less conserved lncRNAs. Unrecognized overlap with known short ncRNAs is of little concern because the latter are almost never spliced. The only exception are the “splice-site-overlapping” SO-microRNAs (Mattioli et al. 2014), which, however, are almost exclusively found in coding genes (Pianigiani et al. 2018) and thus removed by our filters. We therefore assume that such artefacts have a very minor impact in our analysis.

In comparison to vertebrates, we observe a much lower level of conservation as measured by gene structure. For instance, 35.2% of the transcripts are conserved between human and mouse (Nitsche et al. 2015), while between *A. thaliana* and *A. arabicum* we only find splice site conservation in 6.2% (39/627) of our own lincRNAs and 1.3% (32/2444) of lincRNAs annotated in Araport11. This difference is even more striking given the fact that the evolutionary distance between human and mouse (~75 Mya) (Waterston et al. 2002) is larger than between *A. thaliana* and *A. arabicum* (~54 Mya) (Beilstein et al. 2010).

Transposable elements (TEs) are important factors in lncRNA origin (Kapusta et al. 2013). To explore if conserved lincRNAs may be related to TEs, we compared our 627 lincRNAs with the genomic positions of TEs described in Araport11 database. We find only 149 of 627 lincRNAs overlap with TEs and these lincRNAs display significantly lower positional conservation than other lincRNAs in the WGA. Indeed, only 11 were found to be positionally conserved between *A. thaliana* and *B. rapa*. The number of TEs which coincide with lincRNAs with conserved splice sites is even smaller; of the 173 lincRNAs with introns only 11 overlapped with TEs. From all 3897 TEs in the Araport11 database, only 450 are conserved by position in the WGA between *A. thaliana* and *A. arabicum*. This represents only 11.5% of the TEs, that is, less than the percentage of the lincRNAs conserved by genomic position (Fig. 1).

DISCUSSION

In this work we explore the conservation of lncRNAs in the Brassicaceae plant family and we find conservation at different levels: From 627 lincRNAs identified we have 38.1% (239/627) conserved by genomic position as determined by the presence of alignable sequence. A small

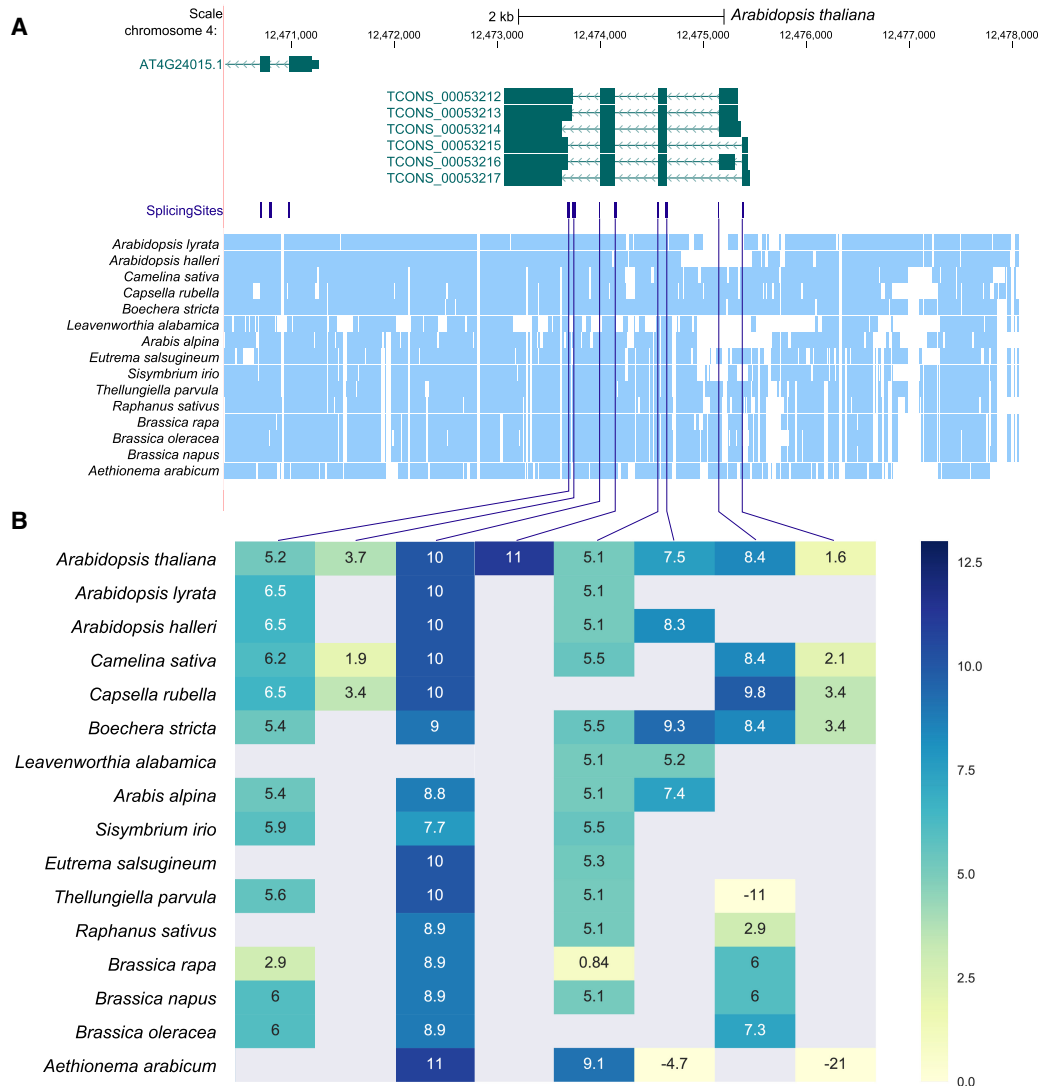


FIGURE 2. Splicing conservation map of lincRNA locus *TCONS00053212-TCONS00053217*. (A) UCSC Genome browser screenshot of the *TCONS00053212-TCONS00053217* locus; blocks denote exons, and line with arrows, introns. The arrow direction indicates direction of transcription. Splicing sites are shown in purple. Light blue blocks represent aligned regions as identified by Cactus. (B) Heatmap of *TCONS00053212-TCONS00053217* MES in each splice site (columns) in each species (rows), linked to its position in A with a purple line. MES are shown from more negative (light yellow) to more positive (dark blue). MES values >0 were used to identify conserved splice sites.

fraction (27.6%) of these lincRNAs contain introns. Only 19.1% of spliced lincRNAs are conserved between *A. thaliana* and *B. oleracea*, the species with the lowest level of conservation in our data set. While sequence conservation may be a consequence of selective constraints on DNA elements, conservation of splice sites directly indicates selective constraints at the transcript level, and thus can be interpreted as evidence for an (unknown) functional role of the lincRNA. The 112 lincRNAs with conserved splice sites are therefore attractive candidates for studies into lincRNA function.

In spite of the small number of spliced lincRNAs analyzed, we find most of them (nearly 65%) have at least one conserved splice site. This is substantially higher

than estimates of conservation by sequence of about 22% amongst Brassica species (Nelson et al. 2016). Thus there is a stronger evolutionary constraint on plant lincRNA processing as measured by splice site conservation than by sequence. This is similar to what was previously found in placental mammals (Nitsche et al. 2015), where ~70% of the lincRNAs have splice site conservation. However, this level of conservation should be considered lower given the divergence time between placental mammals is larger than the divergence times between the Brassicaceae analyzed in this study (52.6 Mya [Kagale et al. 2014]). At least in part this difference is the consequence of the prevalence of single-exon lincRNAs in this clade and the small number of splice sites in those

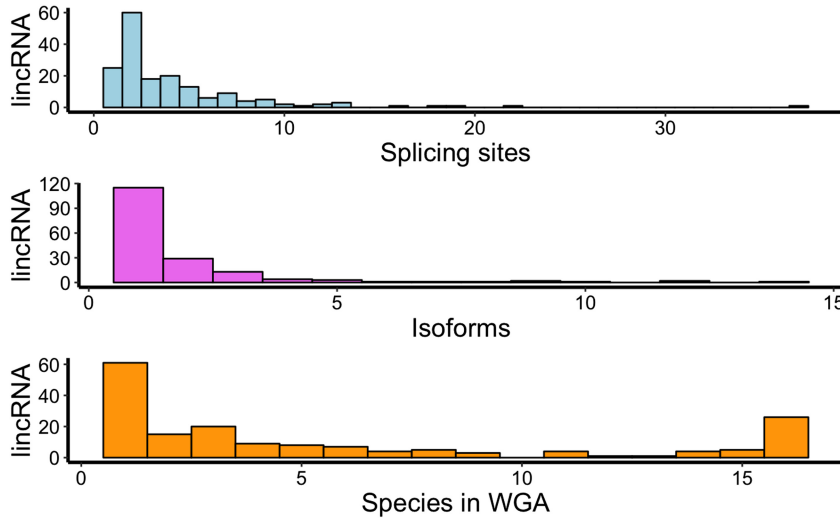


FIGURE 3. Histograms showing number of splicing sites, isoforms, and conservation in WGA of the 173 lincRNAs genes with introns in Own data set. Blue bars indicate the number of splicing sites per lincRNA gene with introns. Purple bars visualize the number of isoforms by lincRNA gene. Orange bars refer to the number of species in which lincRNA genes are conserved.

lincRNAs that contain introns. This reduced the power of the method we used to detect splice site conservation, and hints at a reduced importance of introns in the small genomes of the Brassicaceae. The apparent lower conservation of splice sites may also result from our decision to use *A. thaliana* as a reference which, in addition to having a drastically reduced genome, may have also been subjected to clade-specific intron-loss. Transcriptomes of other Brassicas and other plant families that have not undergone drastic genome reduction will help clarify the actual prevalence on monoexonic and intron-gain -loss in plant lncRNAs.

When comparing with other plant families, for example Poaceae, we find that ~20% of maize and rice lincRNAs are conserved by position (Wang et al. 2015), while we find 38.1% (239/627) of lincRNAs conserved in Brassicaceae. These numbers are roughly comparable given that the divergence times of the two families are similar: Brassicaceae, 52.6 Mya (Kagale et al. 2014); Poaceae, 60 Mya (Charles et al. 2009). The lower conservation observed in Poaceae may be explained by the much larger genome size, and thus higher content of repetitive and unconstrained sequences, leaving conserved sequence regions more “concentrated”—and therefore easier to align—in the small genomes of the Brassicaceae. Consistent with previous findings (Nelson et al. 2016), we find that only a small fraction of our

lincRNAs associated with TEs, compared to a much stronger association in Poaceae (Wang et al. 2017). We interpret this to be a consequence of the substantial reduction of genome size in Brassicas. More detailed comparisons of lincRNA conservation with other plant families will have to await better assembled and annotated genomes to construct adequate WGs.

A limitation of our work is the restriction to intergenic lncRNAs, caused by the need to avoid potential overlaps of the splice sites with other constrained elements. High quality transcriptomes from diverse tissues for most species could alleviate this shortcoming, allowing us to

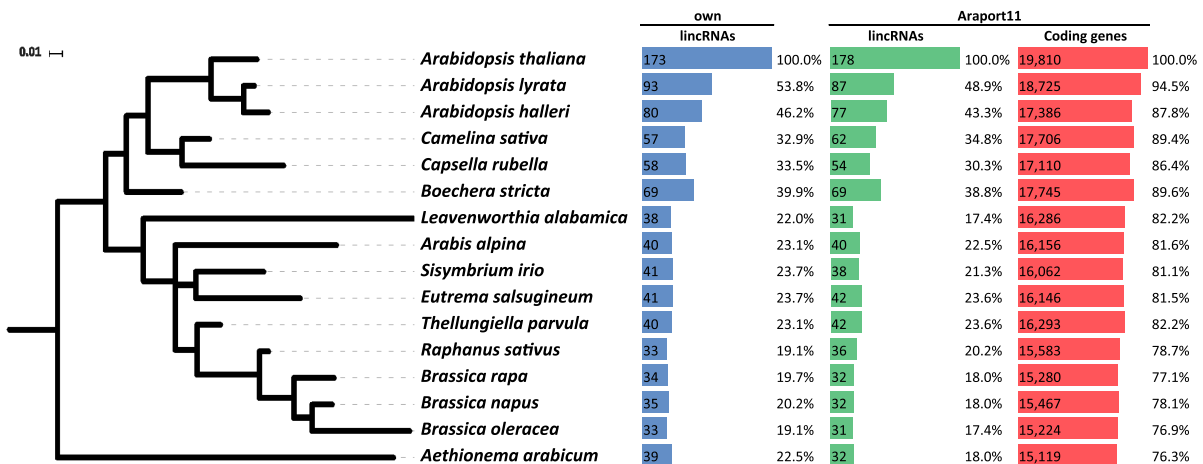


FIGURE 4. Conservation genes in the Brassicaceae family measured by the conservation of splice sites. (Blue) Own lincRNA set (627); (green) lincRNAs in Araport11 (2444); (red) coding RNA genes (27,445). Only genes with at least one intron are shown. Phylogenetic tree scale is in changes per site.

construct splicing maps using only experimental evidence. Spurious sequence conservation would then no longer influence the results. This is of particular relevance in Brassicaceae, since ~70% of transcripts have antisense lncRNAs (Wang et al. 2014). These had to be excluded from our analysis even though at least some of them, for example, *COOLAIR* (Hawkes et al. 2016), are known to have important biological functions. Complementarily to the analysis of splice site conservation, conserved RNA secondary structure can serve as evidence of selection constraints at the RNA level (Washietl et al. 2005). Moreover, structural analysis can be applied to both spliced and monoexonic transcripts. So far, no genome-wide assessment of conservation of RNA secondary structure has been reported for plants. However, recent transcriptome sequence data indicates that RNA structure is also under selection at the genome-wide level in plants (Deng et al. 2018a).

In summary we showed here that higher plants contain at least dozens and most likely hundreds of well-conserved—and with near certainty functional—long noncoding RNAs. We provide an initial catalog of candidates for more detailed exploration, in many cases supported by direct evidence for expression in several species. We furthermore contribute a generic workflow that can be used to uncover conserved lncRNAs in other groups of plants. Given the rapidly expanding collection of publicly available RNA-seq data sets, we suggest that a comparative analysis of lncRNA conservation can complement standard procedures for genome annotation and thus eventually lead to a comprehensive picture of lncRNA diversity and evolution in plants.

MATERIALS AND METHODS

Whole-genome alignment

We selected sixteen plant genomes from those available for the Brassicaceae family in NCBI, Phytozome, and Ensembl-Plants (Supplemental Table S1) based on the quality of assembly, as measured by the number of contigs/scaffolds. All genomes were downloaded in fasta format. Mitochondrial and chloroplast sequences were excluded based on annotation.

The genomes were aligned using Cactus v0 (Paten et al. 2011). Like other whole-genome alignments (WGA) methods, Cactus v0 uses small regions with very high sequence similarity as anchors. To resolve conflicts at this level, Cactus v0 uses a specialized graph data structure that produces better overall alignments than other WGA approaches (Earl et al. 2014). The final WGA result was stored in HAL format (Hickey et al. 2013) for further processing.

Transcriptome data and assembly

We used four previously published base-line transcriptomes for *A. thaliana* (Liu et al. 2012) (GEO accession number GSE38612), as

well as transcriptomes of shade response experiments from Kohnen et al. (2016) (GEO accession number GSE81202). For *Brassica oleracea* we used transcriptomes from Yu et al. (2014) (Expression Atlas accession number E-GEO4-42891). To validate predicted lncRNAs, we used the publicly available transcriptome data sets listed in Supplemental Table S2. All transcriptomes were downloaded as raw reads in fastq format.

We generated our own lncRNA annotation using all single-end stranded sequencing libraries from Kohnen et al. (2016). Libraries were quality-filtered using Trimmomatic v0.32 (Bolger et al. 2014), and mapped to the TAIR10 genome (Berardini et al. 2015) using TopHat v2.1.1 (Trapnell et al. 2009) with parameters: `-l 20 -l 1000 -read-edit-dist 3 -read-realign-edit-dist 0 -library-type fr-firststrand -g 1`. Transcripts were assembled with Cufflinks v2.2.1 (Trapnell et al. 2010) with parameters: `-overlap-radius 1 -p 8 -l 1000 -min-intron-length 20 -g TAIR10_GFF3.gff -library-type fr-firststrand` and subsequently merged into a single reference transcriptome using Cuffmerge v2.2.1.

lncRNA annotation

lncRNAs in the (Kohnen et al. 2016) data set were annotated using two independent methods. First, coding and noncoding transcripts were identified with CPC v0.9.r2 (Coding Potential Calculator) (Kong et al. 2007), a support vector machine classifier. Additionally, we used a strict stepwise annotation workflow (Cabili et al. 2011) on all transcripts. Specifically, we removed transcripts <200 nt in length and identified ORFs 75 aminoacids or longer. Identified ORFs were compared against the NCBI non redundant (nr) database using blastx v2.2.31 and blastp v2.2.31 (Altschul et al. 1990) with *E-value* and cutoff of <10 for a sequence to be considered potentially coding. In addition, we used HMMER v3.1b2 (Wheeler and Eddy 2013) to search for Pfam protein domains, signalP v4.1 (Petersen et al. 2011) to identify signal peptides, and tmhmm v2.0 (Krogh et al. 2001) for transmembrane helices. Only sequences that had no similarity with proteins in nr and no identifiable protein domains, signal peptides or transmembrane domains were annotated as bona fide lncRNAs.

To characterize the genomic context of identified lncRNAs, we used bedtools v2.25.0 (Quinlan and Hall 2010) and compared the lncRNA annotation with the protein coding gene annotation in Araport11 (Cheng et al. 2017). All lncRNA candidates that overlapped a coding sequence or some other ncRNA (miRNA, snoRNA, snRNA) by at least 1 nt were discarded.

Splicing map

The construction of splicing maps requires a seed set of experimentally determined splice sites in at least one species as well as a statistical model to assess the conservation of splice donors and splice acceptors whenever no direct experimental evidence is available.

To obtain these data for Brassicaceae, we mapped the reference transcriptomes to the corresponding reference genome using STAR v2.4.0.1 (Dobin et al. 2013) with default parameters. The table of splice junctions produced by STAR v2.4.0.1 for each data set were concatenated. Only splice junctions that (a) had at least 10 uniquely mapped reads crossing the junction, and (b) showed the canonical GT/AG dinucleotides delimiting the intron (c) within

an intron of size between 59 bp and 999 bp were retained for subsequent analyses. Since some of the transcriptome data sets were not strand-specific, we included CT/AC delimiters, interpreting these as reverse-complements. The same procedure was used for splice site validation in other species, where each transcriptome was mapped against their respective genomes prior to splice junction identification. See Supplemental Table S1 for accessions.

For each identified splice site in *A. thaliana*, we used the HalTools v2.1 liftover tool (Hickey et al. 2013) to determine the corresponding orthologous positions in all other genome sequences in the Cactus v0 generated WGA. For each of the retained splice sites, we extracted the genomic sequences surrounding the donor and acceptor sites. If more than one homolog per species is contained in the WGA, we retained the candidate with the highest sequence similarity to *A. thaliana*. For each known splice site and their orthologous position, the MaxEntScan v0 splice-site score (MES) (Yeo and Burge 2004) was computed with either the donor or acceptor model provided the region contained neither gaps nor ambiguous nucleotides (Supplemental Fig. S1). Otherwise, the regions were treated as nonconserved. MaxEntScan v0 models sequence motifs with a probabilistic model based on the Maximum Entropy Principle, which considers adjacent and nonadjacent dependencies between positions. Several works have verified that the MES is an informative score to measure splice site conservation (Eng et al. 2004; Nitsche et al. 2015). A MaxEntScan v0 splice-site score cut-off of 0 was used (Supplemental Fig. S2). This cut-off value was estimated from the distribution of the MES values obtained from *A. thaliana* and *B. oleracea* transcriptome data (Supplemental Fig. S2). To estimate the rate of false positives, we calculated the probability of finding random splice sites in coding genes and in lncRNAs. For this, we sampled 10,000 random splice positions for both acceptor and donor splicing motifs. In addition to this we calculate the MES values of the same random positions conserved in all WGA species, verifying that they follow the same distribution as in *A. thaliana*. All positively predicted splice-sites, that is, those with MES > 0, were added to the splicing map. The pipeline implementing this analysis is available at: bitbucket.org/JoseAntonioCorona/splicing_map_plants.

DATA DEPOSITION

TrackHubs for all data sets and lncRNAs used in this study as well as WGA are available here: www.bioinf.uni-leipzig.de/Publications/SUPPLEMENTS/19-001/BrassicaceaeWGA/hub.txt. Additional information and machine readable intermediate results are provided at <http://www.bioinf.uni-leipzig.de/Publications/SUPPLEMENTS/19-001>.

SUPPLEMENTAL MATERIAL

Supplemental material is available for this article.

ACKNOWLEDGMENTS

This work was funded in part by Consejo Nacional de Ciencia y Tecnologia (CONACYT PhD Scholarship 338379 [J.A.C.-G.],

CONACYT Research Fellowship 2015-72223 [S.L.F.-V.]), the Deutsche Forschungsgemeinschaft (DFG) grant no. STA850/19-1 to P.F.S., and by a Royal Society Newton Advanced Fellowship (NAF\R1\180303) awarded to S.L.F.-V. We are grateful to Thomas Gatter for advice on transcriptome assembly.

Received December 17, 2019; accepted March 28, 2020.

REFERENCES

- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol* **215**: 403–410. doi:10.1016/S0022-2836(05)80360-2
- Bardou F, Ariel F, Simpson CG, Romero-Barrios N, Laporte P, Balzergue S, Brown JWS, Crespi M. 2014. Long noncoding RNA modulates alternative splicing regulators in *Arabidopsis*. *Dev Cell* **30**: 166–176. doi:10.1016/j.devcel.2014.06.017
- Bánfai B, Jia H, Khatun J, Wood E, Risk B, Gundling WE, Kundaje A, Gunawardena HP, Yu Y, Xie L, et al. 2012. Long noncoding RNAs are rarely translated in two human cell lines. *Genome Res* **22**: 1646–1657. doi:10.1101/gr.134767.111
- Beilstein MA, Nagalingum NS, Clements MD, Manchester SR, Mathews S. 2010. Dated molecular phylogenies indicate a Miocene origin for *Arabidopsis thaliana*. *Proc Natl Acad Sci* **107**: 18724–18728. doi:10.1073/pnas.0909766107
- Berardini TZ, Reiser L, Li D, Mezheritsky Y, Muller R, Strait E, Huala E. 2015. The *Arabidopsis* information resource: making and mining the 'gold standard' annotated reference plant genome. *Genesis* **53**: 474–485. doi:10.1002/dvg.22877
- Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**: 2114–2120. doi:10.1093/bioinformatics/btu170
- Brown JWS, Smith P, Simpson CG. 1996. *Arabidopsis* consensus intron sequences. *Plant Mol Biol* **32**: 531–535. doi:10.1007/BF00019105
- Cabili MN, Trapnell C, Goff L, Koziol M, Tazon-Vega B, Regev A, Rinn JL. 2011. Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses. *Genes Dev* **25**: 1915–1927. doi:10.1101/gad.17446611
- Charles M, Tang H, Belcram H, Paterson A, Gornicki P, Chalhoub B. 2009. Sixty million years in evolution of soft grain trait in grasses: emergence of the softness locus in the common ancestor of pooidae and ehrhartoideae, after their divergence from panicoideae. *Mol Biol Evol* **26**: 1651–1661. doi:10.1093/molbev/msp076
- Chekanova JA. 2015. Long non-coding RNAs and their functions in plants. *Curr Opin Plant Biol* **27**: 207–216. doi:10.1016/j.pbi.2015.08.003
- Cheng CY, Krishnakumar V, Chan AP, Thibaud-Nissen F, Schobel S, Town CD. 2017. Araport11: a complete reannotation of the *Arabidopsis thaliana* reference genome. *Plant J* **89**: 789–804. doi:10.1111/tpj.13415
- Deng H, Cheema J, Zhang H, Woolfenden H, Norris M, Liu Z, Liu Q, Yang X, Yang M, Deng X, et al. 2018a. Rice in vivo RNA structure reveals RNA secondary structure conservation and divergence in plants. *Mol Plant* **11**: 607–622. doi:10.1016/j.molp.2018.01.008
- Deng P, Liu S, Nie X, Weining S, Wu L. 2018b. Conservation analysis of long non-coding RNAs in plants. *Sci China Life Sci* **61**: 190–198. doi:10.1007/s11427-017-9174-9
- Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15–21. doi:10.1093/bioinformatics/bts635

- Earl D, Nguyen N, Hickey G, Harris RS, Fitzgerald S, Beal K, Seledtsov I, Molodtsov V, Raney BJ, Clawson H, et al. 2014. Alignathon: a competitive assessment of whole-genome alignment methods. *Genome Res* **24**: 2077–2089. doi:10.1101/gr.174920.114
- Eng L, Coutinho G, Nahas S, Yeo G, Tanouye R, Babaei M, Dörk T, Burge C, Gatti RA. 2004. Nonclassical splicing mutations in the coding and noncoding regions of the ATM gene: maximum entropy estimates of splice junction strengths. *Hum Mutat* **23**: 67–76. doi:10.1002/humu.10295
- Franco-Zorrilla JM, Valli A, Todesco M, Mateos I, Puga MI, Rubio-Somoza I, Leyva A, Weigel D, Garcia JA, Paz-Ares J. 2007. Target mimicry provides a new mechanism for regulation of microRNA activity. *Nat Genet* **39**: 1033–1037. doi:10.1038/ng2079
- Hawkes EJ, Hennelly SP, Novikova IV, Irwin JA, Dean C, Sanbonmatsu KY. 2016. COOLAIR antisense RNAs form evolutionarily conserved elaborate secondary structures. *Cell Rep* **16**: 3087–3096. doi:10.1016/j.celrep.2016.08.045
- Hebsgaard S. 1996. Splice site prediction in *Arabidopsis thaliana* pre-mRNA by combining local and global sequence information. *Nucleic Acids Res* **24**: 3439–3452. doi:10.1093/nar/24.17.3439
- Hezroni H, Koppstein D, Schwartz MG, Avrutin A, Bartel DP, Ulitsky I. 2015. Principles of long noncoding RNA evolution derived from direct comparison of transcriptomes in 17 species. *Cell Rep* **11**: 1110–1122. doi:10.1016/j.celrep.2015.04.023
- Hickey G, Paten B, Earl D, Zerbino D, Haussler D. 2013. HAL: a hierarchical format for storing and analyzing multiple genome alignments. *Bioinformatics* **29**: 1341–1342. doi:10.1093/bioinformatics/btt128
- Kagale S, Robinson SJ, Nixon J, Xiao R, Huebert T, Condie J, Kessler D, Clarke WE, Edger PP, Links MG, et al. 2014. Polyploid evolution of the Brassicaceae during the cenozoic era. *Plant Cell* **26**: 2777–2791. doi:10.1105/tpc.114.126391
- Kapusta A, Kronenberg Z, Lynch VJ, Zhuo X, Ramsay LA, Bourque G, Yandell M, Feschotte C. 2013. Transposable elements are major contributors to the origin, diversification, and regulation of vertebrate long noncoding RNAs. *PLoS Genet* **9**: e1003470. doi:10.1371/journal.pgen.1003470
- Kohnen MV, Schmid-Siegert E, Trevisan M, Petrolati LA, Sénéchal F, Müller-Moulé P, Maloof J, Xenarios I, Fankhauser C. 2016. Neighbor detection induces organ-specific transcriptomes, revealing patterns underlying hypocotyl-specific growth. *Plant Cell* **28**: 2889–2904. doi:10.1105/tpc.16.00463
- Kong L, Zhang Y, Ye ZQ, Liu XQ, Zhao SQ, Wei L, Gao G. 2007. CPC: assess the protein-coding potential of transcripts using sequence features and support vector machine. *Nucleic Acids Res* **35**: 345–349. doi:10.1093/nar/gkm391
- Krogh A, Larsson B, Von Heijne G, Sonnhammer ELL. 2001. Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *J Mol Biol* **305**: 567–580. doi:10.1006/jmbi.2000.4315
- Liu J, Jung C, Xu J, Wang H, Deng S, Bernad L, Arenas-Huertero C, Chua NH. 2012. Genome-wide analysis uncovers regulation of long intergenic noncoding RNAs in *Arabidopsis*. *Plant Cell* **24**: 4333–4345. doi:10.1105/tpc.112.102855
- Liu J, Wang H, Chua NH. 2015. Long noncoding RNA transcriptome of plants. *Plant Biotechnol J* **13**: 319–328. doi:10.1111/pbi.12336
- Lozada-Chávez I, Stadler PF, Prohaska SJ. 2011. “Hypothesis for the modern RNA world”: a pervasive non-coding RNA-based genetic regulation is a prerequisite for the emergence of multicellular complexity. *Orig Life Evol Biosph* **41**: 587–607. doi:10.1007/s11084-011-9262-1
- Mach J. 2017. The long-noncoding RNA *Elena1* functions in plant immunity. *Plant Cell* **29**: 916. doi:10.1105/tpc.17.00343.
- Mattioli C, Pianigiani G, Pagani F. 2014. Cross talk between spliceosome and microprocessor defines the fate of pre-mRNA. *Wiley Interdiscip Rev RNA* **5**: 647–658. doi:10.1002/wrna.1236
- Meng X, Zhang P, Chen Q, Wang J, Chen M. 2018. Identification and characterization of ncRNA-associated ceRNA networks in *Arabidopsis* leaf development. *BMC Genomics* **19**: 607. doi:10.1186/s12864-018-4993-2
- Mercer TR, Mattick JS. 2013. Structure and function of long noncoding RNAs in epigenetic regulation. *Nature* **20**: 300–307. doi:10.1038/nsmb.2480
- Mohammadin S, Edger PP, Pires JC, Schranz ME. 2015. Positionally-conserved but sequence-diverged: identification of long non-coding RNAs in the Brassicaceae and Cleomaceae. *BMC Plant Biol* **15**: 217. doi:10.1186/s12870-015-0603-5
- Necsulea A, Soumillon M, Warnefors M, Liechti A, Daish T, Zeller U, Baker JC, Grützner F, Kaessmann H. 2014. The evolution of lncRNA repertoires and expression patterns in tetrapods. *Nature* **505**: 635–640. doi:10.1038/nature12943
- Nelson AD, Forsythe ES, Devisetty UK, Clausen DS, Haug-Batzell AK, Meldrum AM, Frank MR, Lyons E, Beilstein MA. 2016. A genomic analysis of factors driving lincRNA diversification: lessons from plants. *G3(Bethesda)* **6**: 2881–2891. doi:10.1534/g3.116.030338
- Nitsche A, Stadler PF. 2017. Evolutionary clues in lncRNAs. *Wiley Interdiscip Rev RNA* **8**: 14–17. doi:10.1002/wrna.1376
- Nitsche A, Rose D, Fasold M, Reiche K, Stadler PF. 2015. Comparison of splice sites reveals that long noncoding RNAs are evolutionarily well conserved. *RNA* **21**: 801–812. doi:10.1261/rna.046342.114
- Paschoal AR, Lozada-Chávez I, Domingues DS, Stadler PF. 2018. ceRNAs in plants: computational approaches and associated challenges for target mimics research. *Brief Bioinformatics* **19**: 1273–1289. doi:10.1093/bib/bbx058
- Paten B, Earl D, Nguyen N, Diekhans M, Zerbino D, Haussler D. 2011. Cactus: algorithms for genome multiple sequence alignment. *Genome Res* **21**: 1512–1528. doi:10.1101/gr.123356.111
- Petersen TN, Brunak S, von Heijne G, Nielsen H. 2011. SignalP 4.0: discriminating signal peptides from transmembrane regions. *Nat Methods* **8**: 785–786. doi:10.1038/nmeth.1701
- Pianigiani G, Licastro D, Fortugno P, Castiglia D, Petrovic I, Pagani F. 2018. Microprocessor-dependent processing of splice site overlapping microRNA exons does not result in changes in alternative splicing. *RNA* **24**: 1158–1171. doi:10.1261/rna.063438.117
- Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**: 841–842. doi:10.1093/bioinformatics/btq033
- Rai MI, Maheen A, Lightfoot DA, Gurha P, Afzal AJ. 2019. Classification and experimental identification of plant long non-coding RNAs. *Genomics* **111**: 997–1005. doi:10.1016/j.ygeno.2018.04.014
- Rosa S, Duncan S, Dean C. 2016. Mutually exclusive sense-antisense transcription at *FLC* facilitates environmentally induced gene repression. *Nat Commun* **7**: 13031. doi:10.1038/ncomms13031
- Sarropoulos I, Marin R, Cardoso-moreira M, Kaessmann H. 2019. Developmental dynamics of lncRNAs across mammalian organs and species. *Nature* **571**: 510–514. doi:10.1038/s41586-019-1341-x
- Seemann SE, Mirza AH, Hansen C, Bang-Berthelsen CH, Garde C, Christensen-Dalsgaard M, Torarinsson E, Yao Z, Workman CT, Pociot F, et al. 2017. The identification and functional annotation of RNA structures conserved in vertebrates. *Genome Res* **27**: 1371–1383. doi:10.1101/gr.208652.116
- Smith MA, Gesell T, Stadler PF, Mattick JS. 2013. Widespread purifying selection on RNA structure in mammals. *Nucleic Acids Res* **41**: 8220–8236. doi:10.1093/nar/gkt596
- Trapnell C, Pachter L, Salzberg SL. 2009. TopHat: discovering splice junctions with RNA-seq. *Bioinformatics* **25**: 1105–1111. doi:10.1093/bioinformatics/btp120

- Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, Van Baren MJ, Salzberg SL, Wold BJ, Pachter L. 2010. Transcript assembly and quantification by RNA-seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* **28**: 511–515. doi:10.1038/nbt.1621
- Ulitsky I. 2016. Evolution to the rescue: using comparative genomics to understand long non-coding RNAs. *Nat Rev Genets* **17**: 601–614. doi:10.1038/nrg.2016.85
- Wang HV, Chekanova JA. 2017. Long noncoding RNAs in plants. *Adv Exp Med Biol* **1008**: 133–154. doi:10.1007/978-981-10-5203-3_5
- Wang Y, Fan X, Lin F, He G, Terzaghi W, Zhu D, Deng XW. 2014. *Arabidopsis* noncoding RNA mediates control of photomorphogenesis by red light. *Proc Natl Acad Sci* **111**: 10359–10364. doi:10.1073/pnas.1409457111
- Wang H, Niu QW, Wu HW, Liu J, Ye J, Yu N, Chua NH. 2015. Analysis of non-coding transcriptome in rice and maize uncovers roles of conserved lncRNAs associated with agriculture traits. *Plant J* **84**: 404–416. doi:10.1111/tpj.13018
- Wang D, Qu Z, Yang L, Zhang Q, Liu ZH, Do T, Adelson DL, Wang ZY, Searle I, Zhu JK. 2017. Transposable elements (TEs) contribute to stress-related long intergenic noncoding RNAs in plants. *Plant J* **90**: 133–146. doi:10.1111/tpj.13481
- Washietl S, Hofacker IL, Stadler PF. 2005. Fast and reliable prediction of noncoding RNAs. *Proc Natl Acad Sci* **102**: 2454–2459. doi:10.1073/pnas.0409169102
- Waterston RH, Lindblad-Toh K, Birney E, Rogers J, Abril JF, Agarwal P, Agarwala R, Ainscough R, Alexandersson M, An P, et al. 2002. Initial sequencing and comparative analysis of the mouse genome. *Nature* **420**: 520–562. doi:10.1038/nature01262
- Wheeler TJ, Eddy SR. 2013. nhmmer: DNA homology search with profile HMMs. *Bioinformatics* **29**: 2487–2489. doi:10.1093/bioinformatics/btt403
- Yamada M. 2017. Functions of long intergenic non-coding (Linc) RNAs in plants. *J Plant Res* **130**: 67–73. doi:10.1007/s10265-016-0894-0
- Yeo G, Burge CB. 2004. Maximum entropy modeling of short sequence motifs with applications to RNA splicing signals. *J Comput Biol* **11**: 377–394. doi:10.1089/1066527041410418
- Yu J, Tehrim S, Zhang F, Tong C, Huang J, Cheng X, Dong C, Zhou Y, Qin R, Hua W, et al. 2014. Genome-wide comparative analysis of NBS-encoding genes between Brassica species and *Arabidopsis thaliana*. *BMC Genomics* **15**: 3. doi:10.1186/1471-2164-15-3
- Yuan C, Meng X, Li X, Illing N, Ingle RA, Wang J, Chen M. 2017. PceRBase: a database of plant competing endogenous RNA. *Nucleic Acids Res* **45**: D1009–D1014. doi:10.1093/nar/gkw916
- Zhang J, Wei L, Jiang J, Mason AS, Li H, Cui C, Chai L, Zheng B, Zhu Y, Xia Q, et al. 2018. Genome-wide identification, putative functionality and interactions between lncRNAs and miRNAs in Brassica species. *Sci Rep* **8**: 4960. doi:10.1038/s41598-018-23334-1