**frontiers**
in Genetics

# Evolinc: A Tool for the Identification and Evolutionary Comparison of Long Intergenic Non-coding RNAs

Andrew D. L. Nelson[1]*[†], Upendra K. Devisetty[2][†], Kyle Palos[1], Asher K. Haug-Baltzell[3], Eric Lyons[2,3] and Mark A. Beilstein[1]*

[1] Beilstein Lab, School of Plant Sciences, University of Arizona, Tucson, AZ, USA, [2] CyVerse, Bio5, University of Arizona, Tucson, AZ, USA, [3] Lyons Lab, Genetics Graduate Interdisciplinary Group, University of Arizona, Tucson, AZ, USA

Long intergenic non-coding RNAs (lincRNAs) are an abundant and functionally diverse class of eukaryotic transcripts. Reported lincRNA repertoires in mammals vary, but are commonly in the thousands to tens of thousands of transcripts, covering ~90% of the genome. In addition to elucidating function, there is particular interest in understanding the origin and evolution of lincRNAs. Aside from mammals, lincRNA populations have been sparsely sampled, precluding evolutionary analyses focused on their emergence and persistence. Here we present Evolinc, a two-module pipeline designed to facilitate lincRNA discovery and characterize aspects of lincRNA evolution. The first module (Evolinc-I) is a lincRNA identification workflow that also facilitates downstream differential expression analysis and genome browser visualization of identified lincRNAs. The second module (Evolinc-II) is a genomic and transcriptomic comparative analysis workflow that determines the phylogenetic depth to which a lincRNA locus is conserved within a user-defined group of related species. Here we validate lincRNA catalogs generated with Evolinc-I against previously annotated Arabidopsis and human lincRNA data. Evolinc-I recapitulated earlier findings and uncovered an additional 70 Arabidopsis and 43 human lincRNAs. We demonstrate the usefulness of Evolinc-II by examining the evolutionary histories of a public dataset of 5,361 Arabidopsis lincRNAs. We used Evolinc-II to winnow this dataset to 40 lincRNAs conserved across species in Brassicaceae. Finally, we show how Evolinc-II can be used to recover the evolutionary history of a known lincRNA, the human telomerase RNA (TERC). These latter analyses revealed unexpected duplication events as well as the loss and subsequent acquisition of a novel TERC locus in the lineage leading to mice and rats. The Evolinc pipeline is currently integrated in CyVerse's Discovery Environment and is free for use by researchers.

Keywords: lincRNAs, comparative genomics, pipeline, evolution, molecular, comparative transcriptomics

## INTRODUCTION

A large, and in some cases predominant, proportion of eukaryotic transcriptomes are composed of long non-coding RNAs (lncRNAs; Guttman et al., 2009; Cabili et al., 2011; Liu et al., 2012; Hangauer et al., 2013; Wang et al., 2015). LncRNAs are longer than 200 nucleotides (nt) and exhibit low protein-coding potential (non-coding). While some transcripts identified from RNA-seq are

likely the result of aberrant transcription or miss-assembly, others are bona fide lincRNAs with various roles (see Wang and Chang, 2011; Ulitsky and Bartel, 2013; for a review of lncRNA functions). To help factor out transcriptional "noise," additional characteristics are used to delineate lncRNAs. These additional characteristics focus on factors such as reproducible identification between experiments, degree of expression, and number of exons (Derrien et al., 2012). In general, lncRNAs display poor sequence conservation among even closely related species, are expressed at lower levels than protein-coding genes, and lack functional data.

The function of any particular lncRNA is likely to influence its evolution. One means of inferring that a transcript is a functional lncRNA and not an artifact is the degree of conservation we observe at that locus between two or more species. This conservation can be observed at the sequence, positional, and transcriptional level (Ulitsky, 2016). Comparative approaches to identify conserved and potentially functional lncRNAs typically focus on long *intergenic* non-coding RNAs (lincRNAs), since their evolution is not constrained by overlap with protein-coding genes. In vertebrates, lincRNA homologs have been identified in species that diverged some 400 million years ago (MYA), whereas in plants lincRNA homologs are primarily restricted to species that diverged <100 MYA (Ulitsky et al., 2011; Liu et al., 2012; Li et al., 2014; Necsulea et al., 2014; Zhang et al., 2014; Mohammadin et al., 2015; Nelson et al., 2016). Importantly, the conserved function of a handful of these lincRNAs have been experimentally verified *in vivo* (Migeon et al., 1999; Hawkes et al., 2016; Quinn et al., 2016).

One major factor inhibiting informative comparative genomics analyses of lincRNAs is the lack of robust sampling and user-friendly analytical tools. Here we present Evolinc, a lincRNA identification and comparative analysis pipeline. The goal of Evolinc is to rapidly and reproducibly identify candidate lincRNA loci, and examine their genomic and transcriptomic conservation. Evolinc relies on RNA-seq data to annotate putative lincRNA loci across the target genome. It is designed to utilize cyberinfrastructure such as the CyVerse Discovery Environment (DE), thereby alleviating the computing demands associated with transcriptome assembly (Merchant et al., 2016). The pipeline is divided into two modules. The first module, Evolinc-I, identifies putative lincRNA loci, and provides output files that can be used for analyses of differential expression, as well as visualization of genomic location using the EPIC-CoGe genome browser (Lyons et al., 2014). The second module, Evolinc-II, is a suite of tools that allows users to identify regions of conservation within a candidate lincRNA, assess the extent to which a lincRNA is conserved in the genomes and transcriptomes of related species, and explore patterns of lincRNA evolution. We demonstrate the versatility of Evolinc on both large and small datasets, and explore the evolution of lincRNAs from both plant and animal lineages.

## MATERIALS AND METHODS

In this section we describe how the two modules of Evolinc (I and II) work, and explain the data generated by each.

## Evolinc-I: LincRNA Identification

Evolinc-I minimally requires the following input data: a set of assembled and merged transcripts from Cuffmerge or Cuffcompare (Trapnell et al., 2010) in gene transfer format (GTF), a reference genome (FASTA), and a reference genome annotation (GTF/GFF/GFF3). From the transcripts provided in the GTF file, only those longer than 200 nt are kept for further analysis. Transcripts with high protein-coding potential are removed using two metrics: (1) open reading frames (ORF) encoding a protein >100 amino acids, and (2) similarity to the UniProt protein database (based on a 1E-5 threshold). Filtering by these two metrics is carried out by Transdecoder (https://transdecoder.github.io/) with the BLASTp step included. These analyses yield a set of transcripts that fulfill the most basic requirements of lncRNAs. Due to anticipated lack of sequence homology or simple lack of genome data that users may deal with, we did not include ORF conservation as a filtering step within Evolinc-I, but instead suggest users to perform a PhyloCSF or RNAcode (Washietl et al., 2011) step after homology exploration by Evolinc-II.

The role of transposable elements (TEs) in the emergence and function of lncRNAs is an active topic of inquiry (Kapusta et al., 2013; Wang et al., 2017). To facilitate these studies, Evolinc allows the user to separate lncRNAs bearing similarity to TEs into a separate FASTA file. This is performed by BLASTn (Altschul et al., 1990; Camacho et al., 2009), with the above lncRNAs as query against a user provided TE database (in FASTA format). Many different TE datasets can be acquired from Repbase (http://www.girinst.org/), PGSB-REdat [http://pgsb.helmholtz-muenchen.de/plant/recat/; (Spannagl et al., 2016)], or DPTEdb: Dioecious Plant Transposable Elements Database (http://genedenovoweb.ticp.net:81/DPTEdb/index.php). We considered lncRNAs that exceeded a bit score value of 200 and an *E*-value threshold of 1E-20 to be TE-derived. These stringent thresholds remove TE-derived lncRNAs with high similarity to TEs, but allow for retention of lncRNAs with only weak similarity to TEs, perhaps reflecting older TE integration events or TE exaptation events (Johnson and Guigó, 2014). To thoroughly identify TE-derived lncRNAs, we suggest building the TE database from as many closely related and relevant species as possible. The output from these analyses includes a sequence file (FASTA) for each TE-derived lncRNA, and BED files to permit their visualization via a genome browser. These transcripts are excluded from the file of putative lncRNAs used in downstream analyses by Evolinc-I.

Candidate lncRNAs are next compared against reference annotation files using the BEDTools package (Quinlan and Hall, 2010) to determine any overlap with known genes. Some reference annotations distinguish between protein-coding and other genes (lncRNAs, pseudogenes, etc.). If this style of reference annotation is available, we suggest running Evolinc-I twice, once with an annotation file containing only protein-coding genes (generated with a simple grep command) and once with all known genes. This is a simple way to distinguish between the identification of novel putative lncRNAs and known (annotated) lncRNAs. We also recommend using an annotation file that contains 5′ and 3′ UTRs where possible. If this is unknown,

the genome coordinates within the reference annotation file should be manually adjusted to include additional sequence on either end of known genes (i.e., 500 bp). This number can be adjusted to adhere to community-specific length parameters for intergenic space. We provide two simple ways to update genome annotation files, either for the command line: (https://github.com/Evolinc/Accessory-scripts) or an app within the DE (Modify_GFF_Coordinates). Evolinc-I identifies lncRNAs whose coordinates overlap with those of a known gene. These gene-associated lncRNAs are then sorted into groups based on direction of overlap to known genes: sense or antisense-overlapping lncRNA transcripts (SOT or AOT, respectively). In order for these inferences to be made, either strand-specific RNA-sequencing must be performed or the lncRNA must be multi-exonic. Sequence FASTA and BED files for each group of overlapping lncRNAs are generated by Evolinc-I for the user to inspect. Demographic data are also generated for each of these lncRNA types (explained further below).

LncRNAs that do not overlap with known genes and have passed all other filters are considered (putative) lincRNAs. Evolinc-I also deals with optional input data that may increase the confidence in the validity of particular candidate lincRNAs. For example, when users provide transcription start site coordinates (in BED format), Evolinc-I identifies lincRNAs in which the 5′ end of the first exon is within 100 bp of any transcriptional start site (TSS). LincRNAs with TSS are annotated as "CAGE_PLUS" in the FASTA sequence file (lincRNAs.FASTA), and the identity of such lincRNAs is recorded in the final summary table (Final_summary_table.tsv). Optionally, Evolinc-I identified lincRNAs (termed Evolinc-lincRNAs) can also be tested against a set of user-defined lincRNAs that are not found in the reference annotation (i.e., an in-house set of lincRNAs not included in the genome annotation files). When the coordinates for a set of such lincRNA loci are provided in general feature format (GFF), Evolinc-I will use these data to determine if any putative Evolinc-lincRNAs are overlapping. These loci are appended with "_overlapping_known_lncRNA" in the lincRNA.FASTA file. The identity of the overlapping (known lncRNA) is listed for each Evolinc-lincRNA in the final summary table (Final_summary_table.tsv).

## Output from Evolinc-I

Evolinc-I generates a sequence file and BED file for TE-derived lncRNAs, AOT, or SOT lncRNAs, and intergenic lncRNAs (lincRNAs). We highly recommend scanning the FASTA files for the presence of ribosomal and other RNAs against the Rfam database (http://rfam.xfam.org/search#tabview=tab1) and removing these before further analysis. The BED file is useful for direct visualization in a genome browser (Buels et al., 2016) or intersecting with other BED files generated from different Evolinc-I analyses (Quinlan and Hall, 2010). An updated genome annotation file is created, appending only the lincRNA loci to the user-supplied reference annotation file. This file can then be used with differential expression analysis programs such as DESeq2 or edgeR (Anders and Huber, 2010; Robinson and Oshlack, 2010). In addition, two types of demographic outputs are generated. For SOT, AOT, and lincRNAs, a report is created that describes

the total number of transcripts identified for each class (isoforms and unique loci), GC content, minimum, maximum, and average length. For lincRNAs only, a final summary table is generated with the length and number of exons for each lincRNA, as well as TSS support and the ID of any overlapping, previously curated lincRNAs. The Evolinc-I workflow is shown in **Figure 1A**.

## Additional Evolinc-I Resources

We have also included in the DE and in the GitHub repository (https://github.com/Evolinc/Accessory-scripts/) an assortment of scripts and workflows that will prevent known errors from occurring in transcript assembly and lncRNA identification. For instance, genome FASTA files often have chromosome headers prefaced with lcl| or gi|, whereas the corresponding genome annotation (GFF) file does not. Some tools such as Cuffmerge and Cuffcompare cannot parse genome associated files with non-matching chromosome IDs, resulting in an output file that will not work with Evolinc-I. To address this issue, we have included a short script called "clean_fasta_header.sh" to the GitHub repository and an app with the same name in the DE.

We also created an additional workflow to streamline the read mapping and transcript assembly process to generate input for Evolinc-I. This workflow is available as an app in the DE called Hisat2-Cuffcompare v1.0 and as a script in our GitHub repository under Accessory_scripts. Hisat2-Cuffcompare requires one or more SRA IDs, a genome sequence file (FASTA), and a genome reference annotation file (GFF) as input. Hisat2-Cuffcompare uses HISAT2 (Pertea et al., 2016) to map reads, either Cufflinks or StringTie (Trapnell et al., 2010; Pertea et al., 2016) to assemble transcripts, and then Cuffmerge or Cuffcompare to generate the input file for Evolinc-I.

## Identifying lincRNA Conservation with Evolinc-II

Evolinc-II minimally requires the following input data: a FASTA file of lincRNA sequences, FASTA file(s) of all genomes to be interrogated, and a single column text file with all species listed in order of phylogenetic relatedness to the query species (example and further elaboration on the species list in **File S1**). Many of these genomes can be acquired from CoGe (www.genomevolution.org) or the genome_data folder for Evolinc within the DE (/iplant/home/shared/iplantcollaborative/example_data/Evolinc.sample.data), and lincRNA sequences can be obtained from either the output of Evolinc-I or from another source. Genome FASTA files should be cleaned of pipe (|) characters (see above) and lincRNA FASTA files should not include underscores. The number, relationship, and divergence times of the genomes chosen will depend on the hypotheses the user intends to test. We recommend using many closely related species (intra-family), where possible, and then picking species outside of the family of interest depending on quality of genome annotation and number of lincRNAs identified. To determine the transcriptional status of lincRNA homologs across a group of species, Evolinc-II can optionally incorporate genome annotation files (GFF) and known lincRNA datasets from target species in FASTA format. In addition, Evolinc-II can incorporate motif and structure
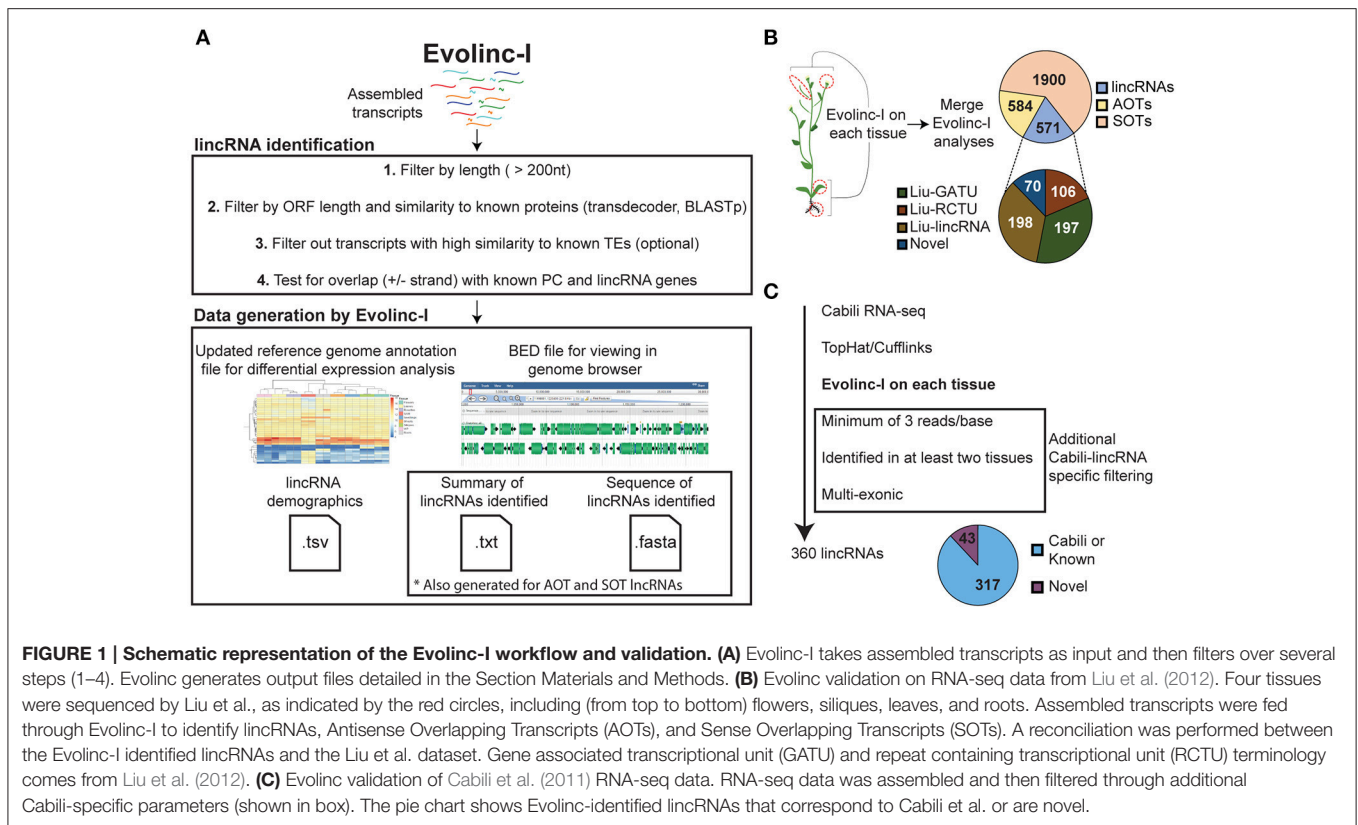
**FIGURE 1 | Schematic representation of the Evolinc-I workflow and validation. (A)** Evolinc-I takes assembled transcripts as input and then filters over several steps (1–4). Evolinc generates output files detailed in the Section Materials and Methods. **(B)** Evolinc validation on RNA-seq data from Liu et al. (2012). Four tissues were sequenced by Liu et al., as indicated by the red circles, including (from top to bottom) flowers, siliques, leaves, and roots. Assembled transcripts were fed through Evolinc-I to identify lincRNAs, Antisense Overlapping Transcripts (AOTs), and Sense Overlapping Transcripts (SOTs). A reconciliation was performed between the Evolinc-I identified lincRNAs and the Liu et al. dataset. Gene associated transcriptional unit (GATU) and repeat containing transcriptional unit (RCTU) terminology comes from Liu et al. (2012). **(C)** Evolinc validation of Cabili et al. (2011) RNA-seq data. RNA-seq data was assembled and then filtered through additional Cabili-specific parameters (shown in box). The pie chart shows Evolinc-identified lincRNAs that correspond to Cabili et al. or are novel.

data, in BED format, to highlight any potential overlap between conserved regions and user-supplied locus information.

Evolinc-II starts by performing a series of reciprocal BLASTn (Camacho et al., 2009) searches against provided target genomes, using a user-defined set of lincRNAs as query and user chosen $E$-value cutoff. We suggest starting at an $E$-value cutoff of 1E-20 because we found that across 10 Brassicaceae genomes, and independently among human, orangutan, and mouse, this value was optimal for recovering reciprocal and syntenic sequence homologs (Nelson et al., 2016). While 1E-20 represents a starting point for these analyses, lincRNA homolog recovery relies on a variety of factors (i.e., background mutation rate, genome stability, evolutionary distance of species/taxa being analyzed, and genome size) that could affect the $E$-value cutoff most likely to return homologous loci among related genomes. Thus, we recommend "calibrating" Evolinc-II using varying $E$-values with at least three genomes (two genomes aside from the query) of varying evolutionary distances from the query species before including a larger (>3) set of genomes. If few sequence homologs are recovered for distantly related species, the user should try lowering the $E$-value. For command-line users examining transposable element derived lncRNAs identified by Evolinc-I, it might be useful to replace all instances of "blastn" within "Building_Families.sh" with "rmblastn." RmBLASTn is a version of BLASTn with Repeat Masker extensions, which will provide more sensitivity when examining conservation of this set of lncRNAs (www.repeatmasker.org). After BLASTn (or RmBLASTn), the top blast hit (TBH) to the query lincRNA is

identified for each additional genome included. Multiple, non-redundant hits falling within the same genomic region, which is likely to occur when the query lincRNA is multi-exonic, are merged as a single TBH. Sequence for all TBHs are then used as query in reciprocal BLAST searches (see below). For researchers interested in inferring orthology vs. paralogy of a sequence homolog in a particular subject species, the coordinates of all BLAST hits that passed the $E$-value cutoff are retained in the file: Homology_search/Subject_species.out.merged.gff. However, to reduce computing time, subsequent analyses are confined to TBHs. Query lincRNAs for which a TBH is not identified in the first iteration (i.e., did not pass the $E$-value cutoff), are subdivided into non-overlapping segments of 200 nt and each segment is used as query in a second set of BLAST searches using similar parameters as the initial search. This reiterative step can be useful in finding short regions of sequence similarity in long query lincRNAs.

TBHs from each species included in the analysis are then used as query sequences in a reciprocal BLAST against the genome of the species whose lincRNA library was used in the original query. For a locus to be considered homologous to the original query lincRNA locus, both loci must be identified as the TBH to each other. This is especially useful when performing searches using a low $E$-value cutoff, as it reduces the chance of random sequence being returned as a sequence homolog. TBHs that pass the reciprocity test are appended with "Homolog" in the final FASTA sequence alignment file ("query_lincRNA_1"_alignment.FASTA).

As TBHs from each target genome are identified, they are scanned for overlap against optional genome reference annotation datasets (GFF) and known lincRNA files (FASTA). The identifier number (ID) of all TBHs with overlap against these two datasets is appended with either "Known_gene" or "Known_lincRNA." The identity of the overlapping gene is retained in the final summary table (final_summary_table.tsv) as well as in each FASTA sequence alignment file (see below). Many genes and almost all lincRNAs are annotated based on transcriptional evidence. Thus, this is a simple way of determining if a query lincRNA corresponds to a locus with evidence of transcription in another species. In addition, when working with a poorly annotated genome, comparing against well-annotated species can provide additional levels of information about the putative function of query lincRNAs. For example, if the homologous locus of a query lincRNA overlaps a protein-coding gene in that species, it could indicate that the query lincRNA is a protein-coding gene, or a pseudogene.

All TBH sequences for a given query lincRNA are clustered into a family. For example, an Evolinc-II analysis that queries 10 lincRNAs across a set of target genomes will result in 10 lincRNA families, populated with the TBH from each target genome. Genomes that do not return a TBH at the specified *E*-value cutoff (from either full-length or segmented searches), or whose TBH does not pass the reciprocity test, will not be represented in the family. These lincRNA families are then batch aligned using MAFFT under default settings with 1,000 iterations (Katoh and Standley, 2013). Command-line users wishing to modify the MAFFT parameters can do so on line 27 of the Batch_MAFFT script available in our GitHub repository (below). The alignment file for each lincRNA family can be downloaded into a sequence viewer. Evolinc-II will also infer phylogeny from the sequence alignment using RAxML v8.2.9 (Stamatakis, 2014) under the GTRGAMMA model, with rapid bootstrap analysis of 1,000 bootstrap datasets. Parameters for RAxML are viewable and modifiable in the Batch_RAxML file. Gene trees are reconciled with a user-provided species tree, in Newick format, using Notung (Durand et al., 2006). This latter analysis pinpoints duplication and loss events that may have occurred during the evolution of the lincRNA locus. Bootstrap support of 70 is required for Notung to choose the gene tree model over the species tree. The Notung reconciled tree is available to view in PNG format within the CyVerse DE. Duplication and loss events are denoted by a red D or L, respectively (Example in **Figure S4**). The Evolinc-II workflow is shown in **Figure 2A**.

## Output from Evolinc-II

Evolinc-II generates sequence files containing lincRNA families with all identified sequence homologs from the user-defined target genomes. In addition, a summary statistics table of identified lincRNA loci based on depth of conservation and overlapping features (e.g., genes, lincRNAs, or other user defined annotations) is generated. The identity of overlapping features (e.g., gene, known lincRNAs) in each genome for which a sequence homolog was identified is listed (Shown for the Liu-lincRNAs in **File S3**). To visualize conserved regions of all query lincRNAs, a query-centric BED file is generated that is ready for import into any genome browser. An example using the genome browser embedded within CoGe (Tang and Lyons, 2012) is shown below (**Figure 2C**). Following phylogenetic analysis, a reconciled gene tree is produced with predicted duplication and loss events indicated. Lastly, to provide the user with a broad picture of lincRNA conservation within their sample set, a bar graph is produced that indicates the number and percent of recovered sequence homologs in each species (**Figure S2A**).

## Data and Software Availability

All genomes used in this work, including version and source, are listed in **File S1**. The accession number of all short read archive files (SRA) used in this work, including project ID, TopHat (Kim et al., 2013) read mapping rate, and total reads mapped for each SRA are shown in **File S1**. Genomic coordinates for lincRNAs identified by Evolinc-I are listed by species in BED/GFF format in **File S2**. LincRNAs were scanned for the presence of ribosomal and other known RNAs by batch searching against the Rfam database (http://rfam.xfam.org/search#tabview=tab1). Novel lincRNAs have also been deposited within the CoGe environment as tracks for genome browsing (links found in **File S2**). Evolinc is available as two apps (Evolinc-I and Evolinc-II) in CyVerse's DE (https://de.cyverse.org/de/), for which a tutorial and sample data are available (https://wiki.cyverse.org/wiki/display/TUT/Evolinc+in+the+Discovery+Environment). Evolinc is also available as self-contained Docker images (https://hub.docker.com/r/evolinc/evolinc-i/ and https://hub.docker.com/r/evolinc/evolinc-ii/) for use in a Linux or Mac OSX command-line environment. The code for Evolinc is available to download/edit as a GitHub repository (https://github.com/Evolinc). Information for installation of the Docker image in a command-line environment, as well as FAQs associated with this process are available in the Evolinc GitHub repository readme file. Both Evolinc workflows make use of several open source tools, such as BLAST for sequence comparisons (Altschul et al., 1990; Camacho et al., 2009), Cufflinks (Trapnell et al., 2010) for GFF to FASTA conversion, Bedtools (Quinlan and Hall, 2010) for sequence intersect comparisons, MAFFT (Katoh and Standley, 2013) for sequence alignment, RAxML (Stamatakis, 2014) for inferring phylogeny, Notung (Durand et al., 2006) for reconciling gene and species trees, and python, perl, and R for file manipulation and data reporting.

## RNA-Seq Read Mapping and Transcript Assembly

SRA files were uploaded directly into CyVerse DE from (http://www.ncbi.nlm.nih.gov/sra) by using the "Import from URL" option. All further read processing was performed using applications within DE. Briefly, uncompressed paired end reads were trimmed (5 nt from 5′ end and 10 nt from 3′ end) using FASTX trimmer, whereas single end read files were filtered with the FASTX quality filter so that only reads where ≥70% of bases with a minimum quality score of 25 were retained (http://hannonlab.cshl.edu/fastx_toolkit/index.html). Reads were mapped to their corresponding genomes using TopHat2 version 2.0.9 (Kim et al., 2013). TopHat2 settings
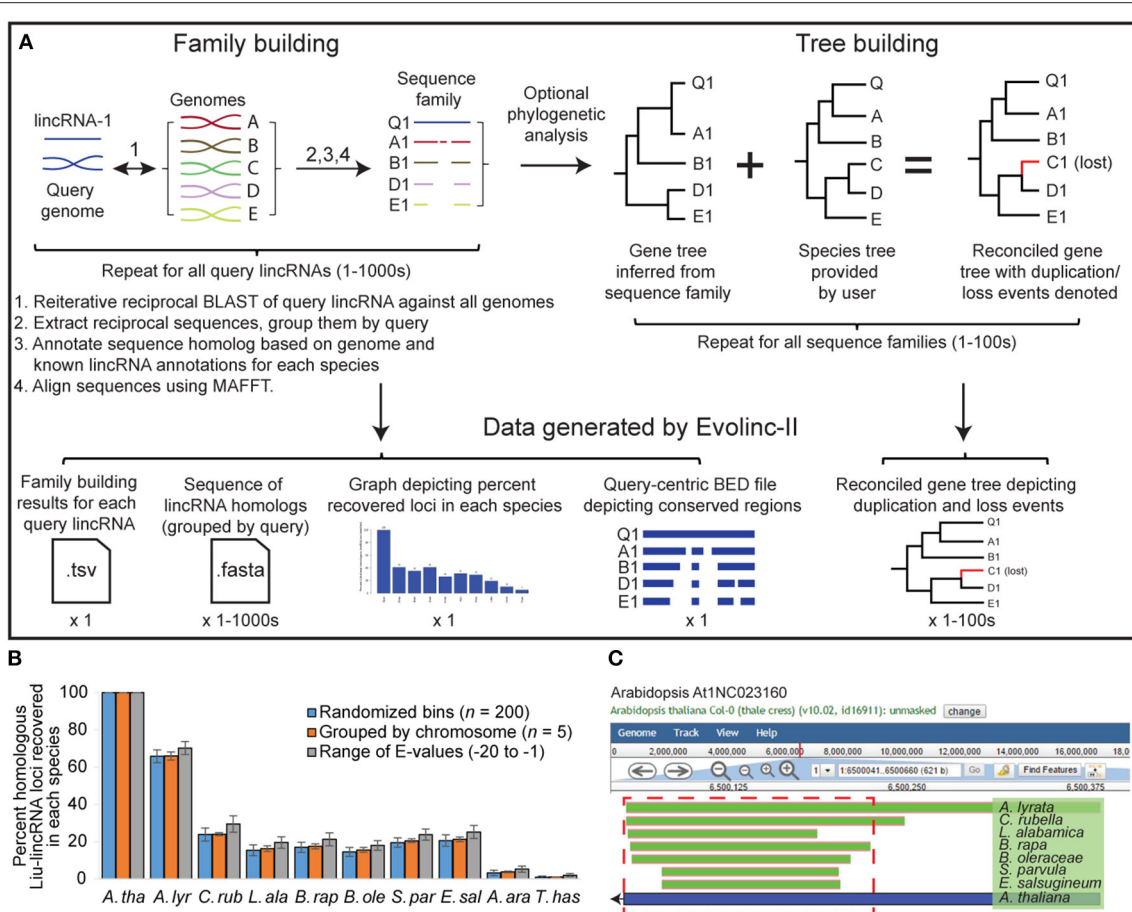
**FIGURE 2 | Schematic representation of the Evolinc-II workflow and validation of Liu-lincRNA and Evolinc-identified lincRNAs. (A)** Evolinc-II uses lincRNAs as a query in reciprocal BLAST analyses against any number of genomes. Sequences that match the filters (see Section Materials and Methods) are grouped into families of sequences based on the query lincRNA. Each sequence homolog is classified using user-defined data or annotations, such as expression or overlap with known gene or lincRNA. Sequences are aligned to highlight conserved regions and to infer phylogeny. These steps can be performed on thousands to tens of thousands of query lincRNAs. Gene trees are inferred for each sequence family using RAxML. The resulting trees are reconciled with the known species tree using Notung 2.0. Notung delineates gene loss and duplication events by marking the output tree with a D (duplication) and blue branch, or L (loss) and red branch. Phylogenetic inference is computationally intensive, and thus we suggest limiting the number of sequence families for which the analysis is performed. Data files generated by Evolinc-II are described in the Section Materials and Methods. **(B)** Validation of Evolinc-II by repeating the Liu-lincRNA dataset in three different ways. The ∼5,400 Liu-lincRNAs were randomly divided into 200 sequence bins (blue bar), each bin was run through Evolinc-II (total number of runs = 27), and then the results were averaged, with standard deviation denoted. In the second analysis, the Liu-lincRNAs were divided based on chromosome, and then each set of Liu-lincRNAs (five groups) were run through Evolinc-II separately. Lastly, all Liu-lincRNAs were run through Evolinc using different BLAST *E*-value cutoffs (E-1, -5, -10, -15, -20), and the results averaged. Bars represent the percent of Liu-lincRNAs for which sequence homologs were identified. *A. tha*, *Arabidopsis thaliana*; *A. lyr*, *Arabidopsis lyrata*; *C. rub*, *Capsella rubella*; *L. ala*, *Leavenworthia alabamica*; *B. rap*, *Brassica rapa*; *B. ole*, *Brassica oleracea*; *S. par*, *Schrenkiella parvula*; *E. sal*, *Eutrema salsugineum*; *A. ara*, *Aethionema arabicum*; *T. has*, *Tarenaya hassleriana*. **(C)** Genome browser visualization of the At1NC023160 locus and its conservation in other Brassicaceae. Regions of the Arabidopsis locus that Evolinc-II identified to be conserved are shown in green, with species of origin listed to the right. The blue bar indicates the length of the locus in Arabidopsis, with the arrow indicating direction of transcription. The region of the locus selected for structural prediction is shown in the red dashed box.

varied based on organism and SRA, and are listed in **File S1**. Transcripts were assembled using the Cufflinks2 app version 2.1.1 under settings listed in **File S1** (Trapnell et al., 2010). TopHat2 and Cufflinks2 were executed on reads from each SRA file independently.

## Validation of lincRNA Expression *In vivo*

RNA was extracted from 2-week old seedlings and flower buds from 4-week old Arabidopsis Col-0 using Trizol (ThermoFisher Life Sciences catalog # 15596018). These tissues and age at extraction most closely matched the experiments from which the RNA-seq data was obtained (Liu et al., 2012). cDNA was synthesized using SuperScript III (ThermoFisher Life Sciences catalog # 18080051) and 2 μg of RNA as input. Primers were first validated by performing PCR with genomic DNA as template using GoTaq Green polymerase master mix (Promega catalog #M712) with 95°C for 3′ to denature, followed by 35 cycles of 95°C for 15″, 55°C for 30″, and 72°C for 30″ and a final extension step of 5′ at 72°C. Primers used are listed in **File S2**.

# RESULTS

## An Overview of lincRNA Identification with Evolinc-I

### Evolinc-I Validation

After establishing a workflow using the most commonly accepted parameters for defining a lincRNA (detailed in Section Materials and Methods), we wanted to evaluate its efficiency at distinguishing between unknown or novel protein-coding genes and non-coding loci. For this, we used a random set of 5,000 protein-coding transcripts selected from the TAIR10 annotation to determine Evolinc-I's false discovery rate (FDR; i.e., protein-coding transcripts erroneously classified as lincRNAs). ORFs for this test dataset ranged in length from 303 to 4,182 nts, with an average ORF of 1,131 nts (**File S3**). Because Evolinc is designed to automatically remove transcripts that map back to known genes, we removed these 5,000 genes from the reference genome annotation file, and then generated a transcript assembly file from RNA-seq data where these 5,000 genes were known to be expressed. We fed the transcript assembly file to Evolinc-I. Out of 5,000 protein-coding genes, only 11 were categorized as non-coding by Evolinc-I (0.22% FDR; **File S3**). Further investigation of the 11 loci revealed that they were predominantly low coverage transcripts with ORFs capable of producing polypeptides >90, but <100 amino acids (aa). Moreover, low read coverage for these transcripts led to incomplete transcript assembly. Together these factors were responsible for the miss-annotation of these loci as non-coding. Importantly, our results indicate that read depth and transcript assembly settings impact lincRNA identification, a finding also noted by Cabili et al. (2011). Therefore, exploring transcript assembly parameters may be necessary prior to running Evolinc-I. In sum, Evolinc-I has a low FDR that can be further reduced by increasing read per base coverage thresholds during transcript assembly as performed in Cabili et al. (2011).

We determined the overlap of Evolinc predicted lincRNAs with previously published datasets from humans and Arabidopsis, following as closely as possible the methods published for each dataset. We first used Evolinc-I to identify lincRNAs from an RNA-seq dataset generated by Liu et al. (2012) in Arabidopsis (**File S1**). From nearly one billion reads generated from four different tissues (siliques, flowers, leaves, and roots), Liu et al. (2012) identified 278 lincRNAs (based on the TAIR9 reference genome annotation). Using the Liu et al. (2012) SRA data, we mapped RNA-seq reads and assembled transcripts with Tophat2 and Cufflinks2 in the DE. From these transcripts, Evolinc-I, identified 571 lincRNAs. We then reconciled the lincRNAs identified in Liu et al. (Liu-lincRNAs) with those from Evolinc-I (Evolinc-lincRNAs), by identifying overlapping genomic coordinates for lincRNAs from the two datasets using the Bedtools suite (Quinlan and Hall, 2010). Of the 278 Liu-lincRNAs, 261 were also recovered by Evolinc-I (**Table S1**). Cufflinks failed to assemble the 17 unrecovered Liu-lincRNAs, due to low coverage, and thus differences in recovery for these loci reflect differences in the Cufflinks parameters employed.

The Arabidopsis genome reference has been updated since Liu et al. (2012), from TAIR9 to TAIR10 (Lamesch et al., 2012). We also ran Evolinc-I with the TAIR10 annotation and found

that only 198 of the 261 Liu-lincRNAs were still considered intergenic (**Figure 1B**). The remaining 63 were reclassified as overlapping a known gene (either sense overlapping transcript, SOT, or antisense overlapping transcript, AOT). This highlights an important aspect of Evolinc-I. While Evolinc-I is able to identify long non-coding RNAs without a genome annotation, genome annotation quality can impact whether an lncRNA is considered intergenic vs. AOT or SOT. In sum, 198 of the 571 lincRNAs identified by Evolinc-I correspond to a previously identified Liu-lincRNA (**Figure 1B**).

Of the 571 lincRNAs identified by Evolinc-I, 373 were not classified as lincRNAs by Liu et al. (2012). Evolinc-I removes transcripts that overlap with the 5′ and 3′ UTRs of a known gene, whereas Liu et al. (2012) removed transcripts that were within 500 bp of a known gene (Liu et al., 2012). This difference in the operational definition of intergenic space accounts for the omission of 197 Evolinc-lincRNAs from the Liu et al. (2012) lincRNA catalog. In addition, Evolinc-I removes transcripts with high similarity to transposable elements, but not tandem di- or tri-nucleotide repeats. We could see no biological reason for excluding these simple repeat containing transcripts, and in fact, transcripts with simple tandem repeats have been attributed to disease phenotypes and therefore might be of particular interest (Usdin, 2008). The inclusion of these transcripts accounts for 106 of the unique Evolinc-lincRNAs.

Finally, 70 of the 571 Evolinc-lincRNAs were entirely novel, and did not correspond to any known Liu-lincRNA or gene within the TAIR10 genome annotation. To determine whether these represented *bona fide* transcripts, we tested expression of a subset ($n = 20$) of single and multi-exon putative lincRNAs by RT-PCR using RNA extracted from two different tissues (seedlings and flowers, **Figure S1**). We considered expression to be positive if we recovered a band in two different tissues or in the same tissue but from different biological replicates. We recovered evidence of expression for 18 of these putative lincRNAs out of 20 tested. Based on these data we conclude that a majority of the 70 novel lincRNAs identified by Evolinc-I for Arabidopsis are likely to reflect *bona fide* transcripts, and thus valid lincRNA candidates.

We next compared Evolinc-I against a well-annotated set of human lincRNAs characterized by Cabili et al. (2011). Cabili et al. (2011) used RNA-seq data from 24 different tissues and cell types, along with multiple selection criteria to identify a "gold standard" reference set of 4,662 lincRNAs. We assembled transcripts from RNA-seq data for seven of these tissues (**File S1**) using Cufflinks under the assembly parameters and read-per-base coverage cut-offs of Cabili et al. (2011) (see Section Materials and Methods). We then fed these transcripts to Evolinc-I. To directly compare Evolinc-I identified lincRNAs with the Cabili et al. (2011) reference dataset (Cabili-lincRNAs), we used the BED files generated by Evolinc-I to identify a subset of 360 multi-exon putative lincRNAs that were observed in at least two tissues (consistent with criteria employed in Cabili et al. (2011) when using a single transcript assembler). We then asked whether these 360 Evolinc-I lincRNAs were found in either the Cabili-lincRNAs, or the hg19 human reference annotation (UCSC). A total of 317 (88%) of the Evolinc-I lincRNAs matched known

lincRNAs from the two annotation sources (**Figure 1C**). The remaining 43 transcripts (12% of the 360 tested) passed all other criteria laid out by Cabili et al. (2011) and therefore may be *bona fide* lincRNAs, but will require further testing.

## Evolution of lincRNA Loci with Evolinc-II
### Evolinc-II Validation
Evolinc-II is an automated and improved version of a workflow we previously used to determine the depth to which Liu-lincRNAs (Liu et al., 2012) were conserved in other species of the Brassicaceae (Nelson et al., 2016). The Evolinc-II workflow is outlined in **Figure 2A**. While most Liu-lincRNAs were restricted to Arabidopsis, or shared only by Arabidopsis and *A. lyrata*, 3% were conserved across the family, indicating that the lincRNA-encoding locus was present in the common ancestor of all Brassicaceae ∼54 MYA (Beilstein et al., 2010). We used Evolinc-II to recapitulate our previous analysis in three ways. First, to provide replicates for statistical analysis, we randomly divided the 5,361 Liu-lincRNAs into 200-sequence groups prior to Evolinc-II analysis ($n = 27$; **Figure 2B** and **Figure S2B**). Second, we performed a separate comparison by dividing the Liu-lincRNAs based upon chromosomal location ($n = 5$). Lastly, we used Evolinc-II to search for sequence homologs using the complete Liu-lincRNA dataset but querying with varying E-value cutoffs (E-20, E-15, E-10, E-05, and E-01). This analysis allowed us to test the impact of the requirement for reciprocity on the recovery of putative homologs under different E-value criteria (**Figure 2B** and **Figure S2D**). The number of sequence homologs increased for each decrement in BLAST stringency (**Figure S2D**), indicating that a significant number of putative homologs fulfill the reciprocity requirement even as sequence similarity decreases. The percentage of sequence homologs retrieved by Evolinc-II was statistically indistinguishable for lincRNAs assigned to groups, chromosomes, or the average from all E-value cutoffs (**Figure 2B** and **Figure S2C**). Thus, Evolinc-II is a robust method to identify sets of lincRNAs that are conserved across a user-defined set of species, such as the Brassicaceae.

In addition to identifying sets of conserved lincRNAs, Evolinc-II also highlights conserved regions within each query lincRNA. To demonstrate these features, we scanned through the Liu-lincRNA Evolinc-II summary statistics file (at 1E-10; **File S4**) to identify a conserved lincRNA. At1NC023160 is conserved as a single copy locus in eight of the 10 species we examined. It was identified by Liu et al. (2012) based on both RNA-seq and tiling array data, as well as validated by Evolinc-I. During the comparative analyses, Evolinc-II generates a query-centric coordinate file that allows the user to visualize within a genome browser (e.g., JBrowse; Buels et al., 2016) what regions of the query lincRNA are most conserved. Using this query-centric coordinate file, we examined the 332 nt At1NC023160 locus in the CoGe genome browser and determined that the 3′ end was most highly conserved (**Figure 2C**). We used the MAFFT multiple sequence alignment generated by Evolinc-II for At1NC023160 to perform structure prediction with RNAalifold (**Figure S3A**; Lorenz et al., 2011). The structural prediction based on the multiple sequence alignment had a greater base pair probability score and lower minimum free energy than

the structure inferred from the Arabidopsis lincRNA alone (**Figures S3B,C**). Conserved regions of a lincRNA serve as potential targets for disruption via genome editing techniques, thereby facilitating its functional dissection.

## Using Evolinc-II to Infer the Evolution of the Human Telomerase RNA Locus TERC
In addition to exploring the evolutionary history of a lincRNA catalog, Evolinc-II is an effective tool to infer the evolution of individual lincRNA loci. To showcase the insights Evolinc-II can provide for datasets composed of a small number of lincRNAs, we focused on the well-characterized human lincRNA, TERC. TERC is the RNA subunit of the ribonucleoprotein complex telomerase that is essential for chromosome end maintenance in stem cells, germ-line cells, and single-cell eukaryotes (Theimer and Feigon, 2006; Blackburn and Collins, 2011; Zhang et al., 2011). TERC is functionally conserved across almost all eukarya, but is highly sequence divergent. Building on work performed by Chen et al. (2000) we used Evolinc-II to examine the evolutionary history of the human TERC locus in 26 mammalian species that last shared a common ancestor between 100 and 130 MYA (**Figure 3**; Glazko, 2003; Arnason et al., 2008).

Evolinc-II identified a human TERC sequence homolog in 23 of the 26 species examined (**Figure 3**; raw output shown in **Figure S4**). We were unable to identify a human TERC homolog in *Ornithoryhnchus anatinus* (platypus), representing the earliest diverging lineage within class Mammalia, using our search criteria. In addition, *Mus musculus* (mouse) and *Rattus norvegicus* (rat) were also lacking a human TERC homolog. However, close relatives of mouse and rat, such as *Ictidomys tridecemlineatus* (squirrel) and *Oryctolagus cuniculus* (rabbit) retained clear human TERC sequence homologs, suggesting that loss of the human TERC-like locus is restricted to the Muridae (mouse/rat family). This is in agreement with the previous identification of the mouse TERC, which exhibits much lower sequence similarity with the human TERC than do other mammals (Chen et al., 2000). All identified human TERC homologs also share synteny, suggesting similar evolutionary origins for this locus throughout mammals (**Figure 3**). Evolinc-II also identified lineage-specific duplication events for the human TERC-like locus in the orangutan, lemur, and galago genomes (**Figure 3**), similar to previous observations in pig and cow (Chen et al., 2000). In sum, Evolinc-II can be applied to both large and small datasets to uncover patterns of duplication, loss, and conservation across large phylogenetic distances.

## DISCUSSION

## Rapid Identification of lincRNAs Using Evolinc-I
With Evolinc-I our goal was to develop an automated and simple pipeline for rapid lincRNA discovery from RNA-seq data. In addition to identification, Evolinc-I generates output files that put downstream analyses and data visualization into the hands of biologists, making it simpler for researchers to discover and explore lincRNAs. Evolinc-I makes use of standard
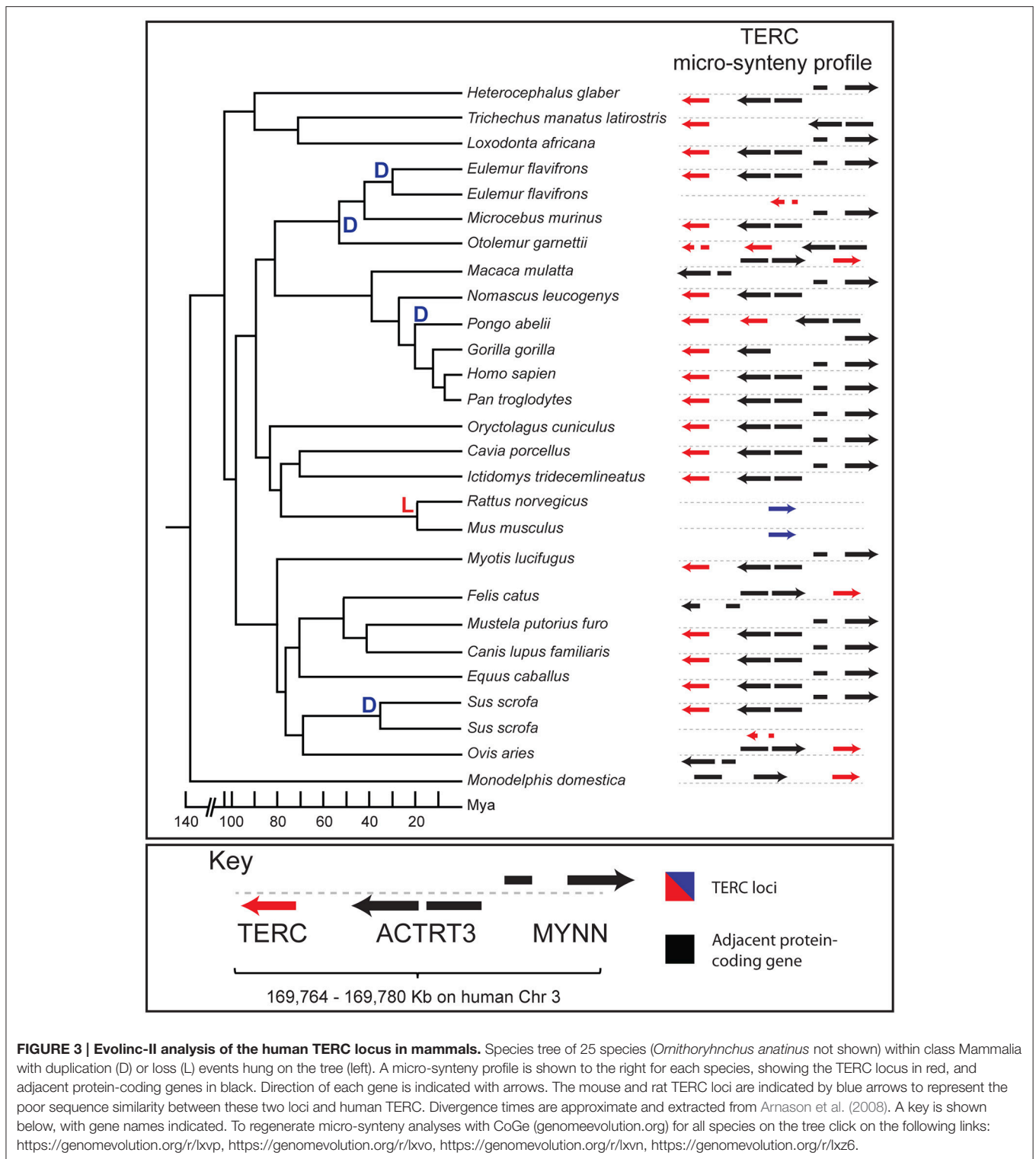
**FIGURE 3 | Evolinc-II analysis of the human TERC locus in mammals.** Species tree of 25 species (*Ornithoryhnchus anatinus* not shown) within class Mammalia with duplication (D) or loss (L) events hung on the tree (left). A micro-synteny profile is shown to the right for each species, showing the TERC locus in red, and adjacent protein-coding genes in black. Direction of each gene is indicated with arrows. The mouse and rat TERC loci are indicated by blue arrows to represent the poor sequence similarity between these two loci and human TERC. Divergence times are approximate and extracted from Arnason et al. (2008). A key is shown below, with gene names indicated. To regenerate micro-synteny analyses with CoGe (genomeevolution.org) for all species on the tree click on the following links: https://genomevolution.org/r/lxvp, https://genomevolution.org/r/lxvo, https://genomevolution.org/r/lxvn, https://genomevolution.org/r/lxz6.

lincRNA discovery criteria, and packages each step into easy-to-use applications within the CyVerse DE or for command-line use via a Docker image with all dependencies pre-installed. We recommend the DE-version of Evolinc-I for novice users, whereas the command-line version of Evolinc-I is useful for

knowledgeable users wishing to tweak parameters to fit their system or question. By using Evolinc-I within the DE, the user can take advantage of the cyberinfrastructure support of CyVerse (Merchant et al., 2016). One of the key advantages of combining Evolinc-I with cyberinfrastructure such as the

CyVerse's DE is the ability to combine various applications together in one streamlined workflow, and making the workflow easier to implement by interested researchers. For instance, a user can download an RNA-seq SRA file into their DE account, quickly process and map reads, assemble transcripts, and execute Evolinc-I. All of this occurs within the DE without downloading a single file or installing a program on a desktop computer.

We demonstrated the ability of Evolinc-I to identify lincRNAs from previously curated catalogs for plants and mammals. Note that we were able to account for all differences between results from Evolinc-I and the published studies, indicating that our pipeline is operating under definitions and filters currently used by the community. Moreover, because we have formalized the process by which annotations of genome data can be incorporated into the search strategy, Evolinc-I gives researchers the ability to easily explore the contributions of TEs, repetitive elements, or other user defined features to the prediction of lincRNA loci. Finally, we stress that this tool permits experiments to be repeated by researchers to compare the contribution of recently released annotations, or to repeat experiments from other groups. This latter point cannot be overemphasized as interest in lincRNAs grows.

## Examining Evolutionary History and Patterns of Conservation of lincRNA Loci Using Evolinc-II

Evolinc-II is designed to perform a series of comparative genomic and transcriptomic analyses across an evolutionary timescale of the user's choosing and on any number (1–1,000 s) of query lincRNAs. Similar to the lncRNA discovery and evolutionary analysis tool Slncky (Chen et al., 2016), the analyses performed by Evolinc-II highlight conserved lincRNA loci, conserved regions within those loci, and overlap with transcripts in other species. To develop an informative evolutionary profile, we recommend users incorporate as many genomes as possible for closely related species and then choose more distantly related species based on the level of genome annotation, genome quality, and quantity of lncRNAs identified for those species. The computationally intensive nature of these analyses is ameliorated by taking advantage of a high-performance computing cluster such as CyVerse. While sequence conservation is certainly not the only filtering mechanism to identify functional lncRNAs, we believe that it is a critical first step. In the future, as more becomes known about structural conservation within lncRNAs, this aspect of lncRNA evolution will be added as an additional filter. We envision Evolinc-II being useful for both scientists attempting to identify functional regions of a lincRNA as well as those wanting to understand the pressures impacting lincRNA evolution.

In addition to highlighting large-scale lincRNA patterns of conservation, we also demonstrated how Evolinc-II can be used to examine the detailed evolutionary history of a single lincRNA, using the human TERC as a test-case. We performed an Evolinc-II analysis with human TERC on 26 genomes in the class Mammalia, 14 of which had not been included in previous studies (Chen et al., 2000). As expected, we recovered a human TERC-like locus in most mammals, as well as three previously unrecorded lineage-specific duplication events. Whether these duplicate TERC loci are expressed and interact with telomerase is unknown; if so they may represent potential regulatory molecules, similar to TER2 in Arabidopsis (Nelson and Shippen, 2015; Xu et al., 2015). We also determined that the human TERC-like locus was lost (or experienced an accelerated mutation rate relative to other mammals) in the common ancestor of mouse and rat. The conservation of the TERC locus across mammals, characterized by rare evolutionary transitions such as that in mouse and rat, stands in stark contrast to the evolution of the telomerase RNA in Brassicaceae (Beilstein et al., 2012), despite the fact that other telomere components are highly conserved (Nelson et al., 2014). Interestingly, mammalian TERCs appear to evolve more slowly than their plant counterparts, similar to the protein components of telomerase (Wyatt et al., 2010). These discoveries highlight the novel insights that can be uncovered using Evolinc-II on even well-studied lincRNAs.

In summary, Evolinc streamlines lincRNA identification and evolutionary analysis. Given the wealth of RNA-seq data being uploaded on a daily basis to NCBI's SRA, and the increased availability of high performance computing resources, we believe that Evolinc will prove to be tremendously useful. Combining these resources, Evolinc can uncover broad and fine-scale patterns in the way that lincRNAs evolve and ultimately help in linking lincRNAs to their function.

## AUTHOR CONTRIBUTIONS

AN, UD, EL, and MB designed the outline of Evolinc. AN, UD, and AH wrote the code. AN, UD, and KP validated and tested the code. AN, UD, EL, and MB wrote the manuscript.

## ACKNOWLEDGMENTS

## SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: http://journal.frontiersin.org/article/10.3389/fgene.2017.00052/full#supplementary-material

**Figure S1 | RT-PCR validation of lincRNAs identified in Arabidopsis by Evolinc-I.** LincRNA IDs match those found in **File S2**. G, genomic DNA positive control; F, flower cDNA; S, seedling cDNA.

**Figure S2 | Examining conservation of Liu-lincRNAs in multiple ways with Evolinc-II. (A)** Example of the type of bar graph produced by Evolinc-II, in this case for the Liu-lincRNAs at 1E-20. **(B)** Bar graph of level of lincRNA conservation observed when dividing the Liu-lincRNAs into 27 random bins of 200 lincRNAs each. Standard deviation is based on the difference seen between the 27 bins. **(C)** Bar graph depicting the level of lincRNA conservation seen when dividing the Liu-lincRNAs by Arabidopsis chromosome (E-cutoff value of 1E-20). **(D)** Bar graph demonstrating the level of conservation of the Liu-lincRNAs throughout Brassicaceae at different E-cutoff values.

**Figure S3 | Using At1NC023160 to highlight the structural information that can be gleaned from Evolinc-II. (A)** Multiple sequence alignment, generated by MAFFT and visualized within Geneious v7.1 (Kearse et al., 2012). Similar sequences are highlighted, with the consensus sequence shown on top. Nucleotide identity is shown below the consensus sequence, with green representing 100% identity across all sequences. **(B)** RNAalifold (Lorenz et al., 2011) consensus secondary structure prediction based on multiple sequence alignment in **(A)**. Base-pair probabilities are shown, with red being more probable and blue least probable. **(C)** RNAfold structure prediction based on the same region as in **(B)**, but limited to just the Arabidopsis sequence. Base-pair probabilities are shown as in **(B)**.

**Figure S4 | Raw phylogenetic output from Evolinc-II for TERC. (A)** A gene tree for the TERC sequence homologs identified in each of the species shown. Sequences without "TBH" indicate paralogs. **(B)** Notung (Durand et al., 2006) reconciliation of the gene tree shown in **(A)** to the known species tree. Duplication (red "D") and loss events (gray "LOST") are shown. Support for duplication or loss events are indicated by the green numbers at the nodes that represent the predicted origin of those events.

**Table S1 | Percent similarity between transcripts identified following transcript assembly and lincRNA identification.**

**File S1 | List of publically available genome and sequence files used, as well as conditions and results from TopHat and Cufflinks for each assembly.**

**File S2 | Evolinc-I output for all species from which lincRNAs were identified, as well as bed files for genome browser viewing, and primers used in RT-PCR verification of transcription of novel Arabidopsis lincRNAs.** Also contains CoGe genome browser links to the novel lincRNAs identified.

**File S3 | False-positive testing of Evolinc-I with Arabidopsis protein-coding genes.**

**File S4 | Evolinc-II output summary table for Liu-lincRNAs at different E-values.**

# REFERENCES

Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. (1990). Basic local alignment search tool. *J. Mol. Biol.* 215, 403–410. doi: 10.1016/S0022-2836(05)80360-2

Anders, S., and Huber, W. (2010). Differential expression analysis for sequence count data. *Genome Biol.* 11:R106. doi: 10.1186/gb-2010-11-10-r106

Arnason, U., Adegoke, J. A., Gullberg, A., Harley, E. H., Janke, A., and Kullberg, M. (2008). Mitogenomic relationships of placental mammals and molecular estimates of their divergences. *Gene* 421, 37–51. doi: 10.1016/j.gene.2008.05.024

Beilstein, M. A., Brinegar, A. E., and Shippen, D. E. (2012). Evolution of the Arabidopsis telomerase RNA. *Front. Genet.* 3:188. doi: 10.3389/fgene.2012.00188

Beilstein, M. A., Nagalingum, N. S., Clements, M. D., Manchester, S. R., and Mathews, S. (2010). Dated molecular phylogenies indicate a Miocene origin for *Arabidopsis thaliana*. *Proc. Natl. Acad. Sci. U.S.A.* 107, 18724–18728. doi: 10.1073/pnas.0909766107

Blackburn, E. H., and Collins, K. (2011). Telomerase: an RNP enzyme synthesizes DNA. *Cold Spring Harb. Perspect. Biol.* 3:a003558. doi: 10.1101/cshperspect.a003558

Buels, R., Yao, E., Diesh, C. M., Hayes, R. D., Munoz-Torres, M., Helt, G., et al. (2016). JBrowse: a dynamic web platform for genome visualization and analysis. *Genome Biol.* 17:66. doi: 10.1186/s13059-016-0924-1

Cabili, M. N., Trapnell, C., Goff, L., Koziol, M., Tazon-Vega, B., Regev, A., et al. (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

Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., et al. (2009). BLAST+: architecture and applications. *BMC Bioinformatics* 10:421. doi: 10.1186/1471-2105-10-421

Chen, J., Shishkin, A. A., Zhu, X., Kadri, S., Maza, I., Guttman, M., et al. (2016). Evolutionary analysis across mammals reveals distinct classes of long non-coding RNAs. *Genome Biol.* 17:19. doi: 10.1186/s13059-016-0880-9

Chen, J.-L., Blasco, M. A., and Greider, C. W. (2000). Secondary structure of vertebrate telomerase RNA. *Cell* 100, 503–514. doi: 10.1016/S0092-8674(00)80687-X

Derrien, T., Johnson, R., Bussotti, G., Tanzer, A., Djebali, S., Tilgner, H., et al. (2012). The GENCODE v7 catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and expression. *Genome Res.* 22, 1775–1789. doi: 10.1101/gr.132159.111

Durand, D., Halldórsson, B. V., and Vernot, B. (2006). A hybrid micro-macroevolutionary approach to gene tree reconstruction. *J. Comput. Biol.* 13, 320–335. doi: 10.1089/cmb.2006.13.320

Glazko, G. V. (2003). Estimation of divergence times for major lineages of primate species. *Mol. Biol. Evol.* 20, 424–434. doi: 10.1093/molbev/msg050

Guttman, M., Amit, I., Garber, M., French, C., Lin, M. F., Feldser, D., et al. (2009). Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in mammals. *Nature* 458, 223–227. doi: 10.1038/nature07672

Hangauer, M. J., Vaughn, I. W., and McManus, M. T. (2013). Pervasive transcription of the human genome produces thousands of previously unidentified long intergenic noncoding RNAs. *PLoS Genet.* 9:e1003569. doi: 10.1371/journal.pgen.1003569

Hawkes, E. J., Hennelly, S. P., Novikova, I. V., Irwin, J. A., Dean, C., Sanbonmatsu, K. Y., et al. (2016). COOLAIR antisense RNAs form evolutionarily conserved elaborate secondary structures. *Cell Rep.* 16, 3087–3096. doi: 10.1016/j.celrep.2016.08.045

Johnson, R., and Guigó, R. (2014). The RIDL hypothesis: transposable elements as functional domains of long noncoding RNAs. *RNA* 20, 959–976. doi: 10.1261/rna.044560.114

Kapusta, A., Kronenberg, Z., Lynch, V. J., Zhuo, X., Ramsay, L., Bourque, G., et al. (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

Katoh, K., and Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* 30, 772–780. doi: 10.1093/molbev/mst010

Kearse, M., Moir, R., Wilson, A., Stones-Havas, S., Cheung, M., Sturrock, S., et al. (2012). Geneious Basic: an integrated and extendable desktop software platform for the organization and analysis of sequence data. *Bioinformatics* 28, 1647–1649. doi: 10.1093/bioinformatics/bts199

Kim, D., Pertea, G., Trapnell, C., Pimentel, H., Kelley, R., Salzberg, S. L., et al. (2013). TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol.* 14:R36. doi: 10.1186/gb-2013-14-4-r36

Lamesch, P., Berardini, T. Z., Li, D., Swarbreck, D., Wilks, C., Sasidharan, R., et al. (2012). The Arabidopsis Information Resource (TAIR): improved gene annotation and new tools. *Nucleic Acids Res.* 40, D1202–D1210. doi: 10.1093/nar/gkr1090

Li, L., Eichten, S. R., Shimizu, R., Petsch, K., Yeh, C.-T., Wu, W., et al. (2014). Genome-wide discovery and characterization of maize long non-coding RNAs. *Genome Biol.* 15:R40. doi: 10.1186/gb-2014-15-2-r40

Liu, J., Jung, C., Xu, J., Wang, H., Deng, S., Bernad, L., et al. (2012). Genome-wide analysis uncovers regulation of long intergenic noncoding RNAs in Arabidopsis. *Plant Cell* 24, 4333–4345. doi: 10.1105/tpc.112.102855

Lorenz, R., Bernhart, S. H., Höner Zu Siederdissen, C., Tafer, H., Flamm, C., Stadler, P. F., et al. (2011). ViennaRNA Package 2.0. *Algorithms Mol. Biol.* 6:26. doi: 10.1186/1748-7188-6-26

Lyons, E., Bomhoff, M., Li, F., and Gregory, B. D. (2014). *EPIC-CoGe: Functional and Diversity Comparative Genomics*. Available online at: https://pag.confex.com/pag/xxii/webprogram/Paper10615.html

Merchant, N., Lyons, E., Goff, S., Vaughn, M., Ware, D., Micklos, D., et al. (2016). The iPlant collaborative: cyberinfrastructure for enabling data to discovery for the life sciences. *PLoS Biol.* 14:e1002342. doi: 10.1371/journal.pbio.1002342

Migeon, B. R., Kazi, E., Haisley-Royster, C., Hu, J., Reeves, R., Call, L., et al. (1999). Human X inactivation center induces random X chromosome inactivation in male transgenic mice. *Genomics* 59, 113–121. doi: 10.1006/geno.1999.5861

Mohammadin, S., Edger, P. P., Pires, J. C., and Schranz, M. E. (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., et al. (2014). The evolution of lncRNA repertoires and expression patterns in tetrapods. *Nature* 505, 635–640. doi: 10.1038/nature12943

Nelson, A. D., and Shippen, D. E. (2015). Evolution of TERT-interacting lncRNAs: expanding the regulatory landscape of telomerase. *Front. Genet.* 6:277. doi: 10.3389/fgene.2015.00277

Nelson, A. D., Forsythe, E. S., Devisetty, U. K., Clausen, D. S., Haug-Batzell, A. K., Meldrum, A. M. R., et al. (2016). A genomic analysis of factors driving lincRNA diversification: lessons from plants. *G3* 6, 2881–2891. doi: 10.1534/g3.116.030338

Nelson, A. D., Forsythe, E. S., Gan, X., Tsiantis, M., and Beilstein, M. A. (2014). Extending the model of Arabidopsis telomere length and composition across Brassicaceae. *Chromosom. Res.* 22, 153–166. doi: 10.1007/s10577-014-9423-y

Pertea, M., Kim, D., Pertea, G. M., Leek, J. T., and Salzberg, S. L. (2016). Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat Protoc.* 11, 1650–1667. doi: 10.1038/nprot.2016.095

Quinlan, A. R., and Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 26, 841–842. doi: 10.1093/bioinformatics/btq033

Quinn, J. J., Zhang, Q. C., Georgiev, P., Ilik, I. A., Akhtar, A., and Chang, H. Y. (2016). Rapid evolutionary turnover underlies conserved lncRNA-genome interactions. *Genes Dev.* 30, 191–207. doi: 10.1101/gad.272187.115

Robinson, M. D., and Oshlack, A. (2010). A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol.* 11:R25. doi: 10.1186/gb-2010-11-3-r25

Spannagl, M., Nussbaumer, T., Bader, K. C., Martis, M. M., Seidel, M., Kugler, K. G., et al. (2016). PGSB plantsDB: updates to the database framework for comparative plant genome research. *Nucleic Acids Res.* 44, D1141–D1147. doi: 10.1093/nar/gkv1130

Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30, 1312–1313. doi: 10.1093/bioinformatics/btu033

Tang, H., and Lyons, E. (2012). Unleashing the genome of brassica rapa. *Front. Plant Sci.* 3:172. doi: 10.3389/fpls.2012.00172

Theimer, C. A., and Feigon, J. (2006). Structure and function of telomerase RNA. *Curr. Opin. Struct. Biol.* 16, 307–318. doi: 10.1016/j.sbi.2006.05.005

Trapnell, C., Williams, B. A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M. J., et al. (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. Genet.* 17, 601–614. doi: 10.1038/nrg.2016.85

Ulitsky, I., and Bartel, D. P. (2013). LincRNAs: genomics, evolution, and mechanisms. *Cell* 154, 26–46. doi: 10.1016/j.cell.2013.06.020

Ulitsky, I., Shkumatava, A., Jan, C. H., Sive, H., and Bartel, D. P. (2011). Conserved function of lincRNAs in vertebrate embryonic development despite rapid sequence evolution. *Cell* 147, 1537–1550. doi: 10.1016/j.cell.2011.11.055

Usdin, K. (2008). The biological effects of simple tandem repeats: lessons from the repeat expansion diseases. *Genome Res.* 18, 1011–1019. doi: 10.1101/gr.070409.107

Wang, D., Qu, Z., Yang, L., Zhang, Q., Liu, Z.-H., Do, T., et al. (2017). Transposable elements (TEs) contribute to stress-related long intergenic noncoding RNAs in plants. *Plant J.* 90, 133–146. doi: 10.1111/tpj.13481

Wang, H., Niu, Q.-W., Wu, H.-W., Liu, J., Ye, J., Yu, N., et al. (2015). Analysis of non-coding transcriptome in rice and maize uncovers roles of conserved lincRNAs associated with agriculture traits. *Plant J.* 84, 404–416. doi: 10.1111/tpj.13018

Wang, K. C., and Chang, H. Y. (2011). Molecular mechanisms of long noncoding RNAs. *Mol. Cell* 43, 904–914. doi: 10.1016/j.molcel.2011.08.018

Washietl, S., Findeiss, S., Müller, S. A., Kalkhof, S., von Bergen, M., Hofacker, I. L., et al. (2011). RNAcode: robust discrimination of coding and noncoding regions in comparative sequence data. *RNA* 17, 578–594. doi: 10.1261/rna.2536111

Wyatt, H. D., West, S. C., and Beattie, T. L. (2010). In*TERT*preting telomerase structure and function. *Nucleic Acids Res.* 38, 5609–5622. doi: 10.1093/nar/gkq370

Xu, H., Nelson, A. D., and Shippen, D. E. (2015). A transposable element within the non-canonical telomerase RNA of *Arabidopsis thaliana* modulates telomerase activity in response to DNA damage. *PLoS Genet.* 11:e1005281. doi: 10.1371/journal.pgen.1005281

Zhang, Q., Kim, N.-K., and Feigon, J. (2011). Architecture of human telomerase RNA. *Proc. Natl. Acad. Sci. U.S.A.* 108, 20325–20332. doi: 10.1073/pnas.1100279108

Zhang, Y.-C., Liao, J.-Y., Li, Z.-Y., Yu, Y., Zhang, J.-P., Li, Q.-F., et al. (2014). Genome-wide screening and functional analysis identify a large number of long noncoding RNAs involved in the sexual reproduction of rice. *Genome Biol.* 15:512. doi: 10.1186/s13059-014-0512-1