# Evolinc: a tool for the identification and evolutionary comparison of long intergenic non-coding RNAs

Andrew D. Nelson[1*], Upendra K. Devisetty[2], Kyle Palos[1], Asher K. Haug-Baltzell[3], Eric Lyons[4], Mark A. Beilstein[1*]

[1]School of Plant Sciences, University of Arizona School of Plant Sciences, USA, [2]CyVerse, USA, [3]Genetics Graduate Interdisciplinary Group, University of Arizona, USA, [4]School of Plant Sciences, University of Arizona, USA

1 **Evolinc: a tool for the identification and evolutionary comparison of long**
2 **intergenic non-coding RNAs**

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

5

6 **Authors:** [1]School of Plant Sciences, University of Arizona, Tucson, Arizona, 85721, [2]

7 CyVerse, Bio5, University of Arizona, Tucson, Arizona, 85721, [3]Genetics Graduate

8 Interdisciplinary Group, University of Arizona, Tucson, Arizona, 85721

9

10 * These authors contributed equally to this manuscript.

11 † Corresponding Authors

12

13 **Corresponding Authors:**

14 Mark Beilstein, 1140 E. South Campus Drive, 303 Forbes Building, Tucson, Arizona,

15 85721-0036, 520-626-1562, mbeilstein@email.arizona.edu

16 Andrew Nelson, 1140 E. South Campus Drive, 303 Forbes Building, Tucson, Arizona,

17 85721-0036, 520-626-1563, andrewnelson@email.arizona.edu

18

19 **Running Title:** Identification and comparative evolution of lincRNAs with Evolinc

20

21

22

23

24

25

26

**Abstract**

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.

## Introduction

A large, and in some cases predominant, proportion of eukaryotic transcriptomes are composed of long non-coding RNAs (lncRNAs) (Hangauer et al., 2013; Guttman et al., 2009; Cabili et al., 2011; H., Wang et al., 2015; Liu et al., 2012). 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 lncRNA. 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 artefact 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; Necsulea et al., 2014; Nelson et al., 2016; Liu et al., 2012; Li et al., 2014; Zhang et al., 2014; Mohammadin et al., 2015). Importantly, the conserved function of a handful of these lincRNAs have been experimentally verified *in vivo* (Hawkes et al., 2016; Migeon et al., 1999; Quinn et al., 2016).

78    One major factor inhibiting informative comparative genomics analyses of lincRNAs is

79    the lack of robust sampling and user-friendly analytical tools. Here we present Evolinc, a

80    lincRNA identification and comparative analysis pipeline. The goal of Evolinc is to rapidly and

81    reproducibly identify candidate lincRNA loci, and examine their genomic and transcriptomic

82    conservation. Evolinc relies on RNA-seq data to annotate putative lincRNA loci across the target

83    genome. It is designed to utilize cyberinfrastructure such as the CyVerse Discovery

84    Environment (DE), thereby alleviating the computing demands associated with transcriptome

85    assembly (Merchant et al., 2016). The pipeline is divided into two modules. The first module,

86    Evolinc-I, identifies putative lincRNA loci, and provides output files that can be used for analyses

87    of differential expression, as well as visualization of genomic location using the EPIC-CoGe

88    genome browser (Lyons et al., 2014). The second module, Evolinc-II, is a suite of tools that

89    allows users to identify regions of conservation within a candidate lincRNA, assess the extent to

90    which a lincRNA is conserved in the genomes and transcriptomes of related species, and

91    explore patterns of lincRNA evolution. We demonstrate the versatility of Evolinc on both large

92    and small datasets, and explore the evolution of lincRNAs from both plant and animal lineages.

93

94    **Materials and Methods**

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

96    generated by each.

97    ***Evolinc-I: LincRNA identification***

98    Evolinc-I minimally requires the following input data: a set of assembled and merged

99    transcripts from Cuffmerge or Cuffcompare (Trapnell et al., 2010) in gene transfer format (GTF),

100   a reference genome (FASTA), and a reference genome annotation (GTF/GFF/GFF3). From the

101   transcripts provided in the GTF file, only those longer than 200 nt are kept for further analysis.

102   Transcripts with high protein-coding potential are removed using two metrics: 1) open reading

103 frames (ORF) encoding a protein > 100 amino acids, and 2) similarity to the UniProt protein

104 database (based on a 1E-5 threshold). Filtering by these two metrics is carried out by

105 Transdecoder (https://transdecoder.github.io/) with the BLASTp step included. These analyses

106 yield a set of transcripts that fulfill the most basic requirements of lncRNAs. Due to anticipated

107 lack of sequence homology or simple lack of genome data that users may deal with, we did not

108 include ORF conservation as a filtering step within Evolinc-I, but instead suggest users to

109 perform a PhyloCSF or RNAcode (Washietl et al., 2011) step after homology exploration by

110 Evolinc-II.

111        The role of transposable elements (TEs) in the emergence and function of lncRNAs is an

112 active topic of inquiry (Wang et al., 2017; Kapusta et al., 2013). To facilitate these studies,

113 Evolinc allows the user to separate lncRNAs bearing similarity to TEs into a separate FASTA

114 file. This is performed by BLASTn (Camacho et al., 2009; Altschul et al., 1990), with the above

115 lncRNAs as query against a user provided TE database (in FASTA format). Many different TE

116 datasets can be acquired from Repbase (http://www.girinst.org/), PGSB-REdat

117 [http://pgsb.helmholtz-muenchen.de/plant/recat/; (Spannagl et al., 2016)], or DPTEdb: Dioecious

118 Plant Transposable Elements Database (http://genedenovoweb.ticp.net:81/DPTEdb/index.php).

119 We considered lncRNAs that exceeded a bit score value of 200 and an E-value threshold of 1E-

120 20 to be TE-derived. These stringent thresholds remove TE-derived lncRNAs with high similarity

121 to TEs, but allow for retention of lncRNAs with only weak similarity to TEs, perhaps reflecting

122 older TE integration events or TE exaptation events (Johnson and Guigó, 2014). To thoroughly

123 identify TE-derived lncRNAs, we suggest building the TE database from as many closely related

124 and relevant species as possible. The output from these analyses includes a sequence file

125 (FASTA) for each TE-derived lncRNAs, and BED files to permit their visualization via a genome

126 browser. These transcripts are excluded from the file of putative lncRNAs used in downstream

127 analyses by Evolinc-I.

128   Candidate lncRNAs are next compared against reference annotation files using the

129   BEDTools package (Quinlan and Hall, 2010) to determine any overlap with known genes. Some

130   reference annotations distinguish between protein-coding and other genes (lncRNAs,

131   pseudogenes, etc). If this style of reference annotation is available, we suggest running Evolinc-

132   I twice, once with an annotation file containing only protein-coding genes (generated with a

133   simple grep command) and once with all known genes. This is a simple way to distinguish

134   between the identification of novel putative lncRNAs and known (annotated) lncRNAs. We also

135   recommend using an annotation file that contains 5' and 3' UTRs where possible. If this is

136   unknown, the genome coordinates within the reference annotation file should be manually

137   adjusted to include additional sequence on either end of known genes (i.e., 500bp). This

138   number can be adjusted to adhere to community-specific length parameters for intergenic

139   space. We provide two simple ways to update genome annotation files, either for the command

140   line: (https://github.com/Evolinc/Accessory-scripts) or an app within the DE

141   (Modify_GFF_Coordinates) Evolinc-I identifies lncRNAs whose coordinates overlap with those

142   of a known gene. These gene-associated lncRNAs are then sorted into groups based on

143   direction of overlap to known genes: sense or antisense-overlapping lncRNA transcripts (SOT

144   or AOT, respectively). Keep in mind that in order for these inferences to be made, either strand-

145   specific RNA-sequencing must be performed or the lncRNA must be multi-exonic. Sequence

146   FASTA and BED files for each group of overlapping lncRNAs are generated by Evolinc-I for the

147   user to inspect. Demographic data are also generated for each of these lncRNA types

148   (explained further below).

149   LncRNAs that do not overlap with known genes and have passed all other filters are

150   considered (putative) lincRNAs. Evolinc-I also deals with optional input data that may increase

151   the confidence in the validity of particular candidate lincRNAs. For example, when users provide

152   transcription start site coordinates (in BED format), Evolinc-I identifies lincRNAs in which the 5'

153    end of the first exon is within 100bp of any transcriptional start site (TSS). LincRNAs with TSS

154    are annotated as "CAGE_PLUS" in the FASTA sequence file (lincRNAs.FASTA), and the

155    identity of such lincRNAs is recorded in the final summary table (Final_summary_table.tsv).

156    Optionally, Evolinc-I identified lincRNAs (termed Evolinc-lincRNAs) can also be tested against a

157    set of user-defined lincRNAs that are not found in the reference annotation (i.e., an in-house set

158    of lincRNAs not included in the genome annotation files). When the coordinates for a set of such

159    lincRNA loci are provided in general feature format (GFF), Evolinc-I will use these data to

160    determine if any putative Evolinc-lincRNAs are overlapping. These loci are appended with

161    "_overlapping_known_lncRNA" in the lincRNA.FASTA file. The identity of the overlapping

162    (known lncRNA) is listed for each Evolinc-lincRNA in the final summary table

163    (Final_summary_table.tsv).

164    ***Output from Evolinc-I***

165    Evolinc-I generates a sequence file and BED file for TE-derived lncRNAs, AOT or SOT

166    lncRNAs, and intergenic lncRNAs (lincRNAs). We highly recommend scanning the FASTA files

167    for the presence of ribosomal and other RNAs against the Rfam database

168    (http://rfam.xfam.org/search#tabview=tab1) and removing these before further analysis. The

169    BED file is useful for direct visualization in a genome browser (Buels et al., 2016) or intersecting

170    with other BED files generated from different Evolinc-I analyses (Quinlan and Hall, 2010). An

171    updated genome annotation file is created, appending only the lincRNA loci to the user-supplied

172    reference annotation file. This file can then be used with differential expression analysis

173    programs such as DESeq2 or edgeR (Anders and Huber, 2010; Robinson and Oshlack, 2010).

174    In addition, two types of demographic outputs are generated. For SOT, AOT, and lincRNAs, a

175    report is created that describes the total number of transcripts identified for each class (isoforms

176    and unique loci), GC content, minimum, maximum, and average length. For lincRNAs only, a

177    final summary table is generated with the length and number of exons for each lincRNA, as well

178    as TSS support and the ID of any overlapping, previously curated lincRNA. The Evolinc-I

179    workflow is shown in Figure 1A.

180    ***Additional Evolinc-I resources***

181    We have also included in the DE and in the GitHub repository

182    (https://github.com/Evolinc/Accessory-scripts/) an assortment of scripts and workflows that will

183    prevent known errors from occurring in transcript assembly and lncRNA identification. For

184    instance, genome FASTA files often have chromosome headers prefaced with lcl| or gi|,

185    whereas the corresponding genome annotation (GFF) file does not. Some tools such as

186    Cuffmerge and Cuffcompare cannot parse genome associated files with non-matching

187    chromosome IDs, resulting in an output file that will not work with Evolinc-I. To address this

188    issue, we have included a short script called "clean_fasta_header.sh" to the GitHub repository

189    and an app with the same name in the DE.

190         We also created an additional workflow to streamline the read mapping and transcript

191    assembly process to generate input for Evolinc-I. This workflow is available as an app in the DE

192    called Hisat2-Cuffcompare v1.0 and as a script in our GitHub repository under

193    Accessory_scripts. Hisat2-Cuffcompare requires one or more SRA IDs, a genome sequence file

194    (FASTA), and a genome reference annotation file (GFF) as input. Hisat2-Cuffcompare uses

195    HISAT2 (Pertea et al., 2016) to map reads, either Cufflinks or StringTie (Trapnell et al., 2010;

196    Pertea et al., 2016) to assemble transcripts, and then Cuffmerge or Cuffcompare to generate

197    the input file for Evolinc-I.

198    ***Identifying lincRNA conservation with Evolinc-II***

199    Evolinc-II minimally requires the following input data: a FASTA file of lincRNA sequences,

200    FASTA file(s) of all genomes to be interrogated, and a single column text file with all species

201    listed in order of phylogenetic relatedness to the query species (example and further elaboration

202   on the species list in File S1). Many of these genomes can be acquired from CoGe

203   ([www.genomevolution.org](www.genomevolution.org)) or the genome_data folder for Evolinc within the DE

204   (/iplant/home/shared/iplantcollaborative/example_data/Evolinc.sample.data), and lincRNA

205   sequences can be obtained from either the output of Evolinc-I or from another source. Genome

206   FASTA files should be cleaned of pipe (|) characters (see above) and lincRNA FASTA files

207   should not include underscores. The number, relationship, and divergence times of the

208   genomes chosen will depend on the hypotheses the user intends to test. We recommend using

209   many closely related species (intra-family), where possible, and then picking species outside of

210   the family of interest depending on quality of genome annotation and number of lincRNAs

211   identified. To determine the transcriptional status of lincRNA homologs across a group of

212   species, Evolinc-II can optionally incorporate genome annotation files (GFF) and known

213   lincRNA datasets from target species in FASTA format. In addition, Evolinc-II can incorporate

214   motif and structure data, in BED format, to highlight any potential overlap between conserved

215   regions and user-supplied locus information.

216       Evolinc-II starts by performing a series of reciprocal BLASTn (Camacho et al., 2009)

217   searches against provided target genomes, using a user-defined set of lincRNAs as query and

218   user chosen E-value cutoff. We suggest starting at an E-value cutoff of 1E-20 because we

219   found that across 10 Brassicaceae genomes, and independently among human, orangutan and

220   mouse, this value was optimal for recovering reciprocal and syntenic sequence homologs

221   (Nelson et al., 2016). While 1E-20 represents a starting point for these analyses, lincRNA

222   homolog recovery relies on a variety of factors (i.e., background mutation rate, genome stability,

223   evolutionary distance of species / taxa being analyzed, and genome size) that could affect the

224   E-value cutoff most likely to return homologous loci among related genomes. Thus, we

225   recommend "calibrating" Evolinc-II using varying E-values with at least three genomes (two

226   genomes aside from the query) of varying evolutionary distances from the query species before

227     including a larger (> 3) set of genomes. If few sequence homologs are recovered for distantly

228     related species, the user should try lowering the E-value. For command-line users examining

229     transposable element derived lncRNAs identified by Evolinc-I, it might be useful to replace all

230     instances of "blastn" within "Building_Families.sh" with "rmblastn". RmBLASTn is a version of

231     BLASTn with Repeat Masker extensions, which will provide more sensitivity when examining

232     conservation of this set of lncRNAs (www.repeatmasker.org). After BLASTn (or RmBLASTn),

233     the top blast hit (TBH) to the query lincRNA is identified for each additional genome included.

234     Multiple, non-redundant hits falling within the same genomic region, which is likely to occur

235     when the query lincRNA is multi-exonic, are merged as a single TBH. Sequence for all TBHs

236     are then used as query in reciprocal BLAST searches (see below). For researchers interested in

237     inferring orthology versus paralogy of a sequence homolog in a particular subject species, the

238     coordinates of all BLAST hits that passed the E-value cutoff are retained in the file:

239     Homology_search/Subject_species.out.merged.gff. However, to reduce computing time,

240     subsequent analyses are confined to TBHs. Query lincRNAs for which a TBH is not identified in

241     the first iteration (i.e., did not pass the E-value cutoff), are subdivided into non-overlapping

242     segments of 200 nt and each segment is used as query in a second set of BLAST searches

243     using similar parameters as the initial search. This reiterative step can be useful in finding short

244     regions of sequence similarity in long query lincRNAs.

245          TBHs from each species included in the analysis are then used as query sequences in a

246     reciprocal BLAST against the genome of the species whose lincRNA library was used in the

247     original query. For a locus to be considered homologous to the original query lincRNA locus,

248     both loci must be identified as the TBH to each other. This is especially useful when performing

249     searches using a low E-value cutoff, as it reduces the chance of random sequence being

250     returned as a sequence homolog. TBHs that pass the reciprocity test are appended with

251     "Homolog" in the final FASTA sequence alignment file ("query_lincRNA_1"_alignment.FASTA).

252    As TBHs from each target genome are identified, they are scanned for overlap against

253    optional genome reference annotation datasets (GFF) and known lincRNA files (FASTA). The

254    identifier number (ID) of all TBHs with overlap against these two datasets is appended with

255    either "Known_gene" or "Known_lincRNA". The identity of the overlapping gene is retained in

256    the final summary table (final_summary_table.tsv) as well as in each FASTA sequence

257    alignment file (see below). Many genes and almost all lincRNAs are annotated based on

258    transcriptional evidence. Thus, this is a simple way of determining if a query lincRNA

259    corresponds to a locus with evidence of transcription in another species. In addition, when

260    working with a poorly annotated genome, comparing against well-annotated species can

261    provide additional levels of information about the putative function of query lincRNAs. For

262    example, if the homologous locus of a query lincRNA overlaps a protein-coding gene in that

263    species, it could indicate that the query lincRNA is a protein-coding gene, or a pseudogene.

264    All TBH sequences for a given query lincRNA are clustered into a family. For example,

265    an Evolinc-II analysis that queries ten lincRNAs across a set of target genomes will result in ten

266    lincRNA families, populated with the TBH from each target genome. Genomes that do not return

267    a TBH at the specified E-value cutoff (from either full-length or segmented searches), or whose

268    TBH does not pass the reciprocity test, will not be represented in the family. These lincRNA

269    families are then batch aligned using MAFFT under default settings with 1000 iterations (Katoh

270    and Standley, 2013). Command-line users wishing to modify the MAFFT parameters can do so

271    on line 27 of the Batch_MAFFT script available in our GitHub repository (below). The alignment

272    file for each lincRNA family can be downloaded into a sequence viewer. Evolinc-II will also infer

273    phylogeny from the sequence alignment using RAxML v8.2.9 (Stamatakis, 2014) under the

274    GTRGAMMA model, with rapid bootstrap analysis of 1000 bootstrap datasets. Parameters for

275    RAxML are viewable and modifiable in the Batch_RAxML file. Gene trees are reconciled with a

276    user-provided species tree, in Newick format, using Notung (Durand et al., 2006). This latter

277 analysis pinpoints duplication and loss events that may have occurred during the evolution of

278 the lincRNA locus. Bootstrap support of 70 is required for Notung to choose the gene tree model

279 over the species tree. The Notung reconciled tree is available to view in PNG format within the

280 CyVerse DE. Duplication and loss events are denoted by a red D or L, respectively (Example in

281 Figure S4). The Evolinc-II workflow is shown in Figure 2A.

282 **Output from Evolinc-II**

283 Evolinc-II generates sequence files containing lincRNA families with all identified sequence

284 homologs from the user-defined target genomes. In addition, a summary statistics table of

285 identified lincRNA loci based on depth of conservation and overlapping features (e.g., genes,

286 lincRNAs, or other user defined annotations) is generated. The identity of overlapping features

287 (e.g., gene, known lincRNAs) in each genome for which a sequence homolog was identified is

288 listed (Shown for the Liu-lincRNAs in File S3). To visualize conserved regions of all query

289 lincRNAs, a query-centric BED file is generated that is ready for import into any genome

290 browser. An example using the genome browser embedded within CoGe (Tang and Lyons,

291 2012) is shown below (Figure 2C). Following phylogenetic analysis, a reconciled gene tree is

292 produced with predicted duplication and loss events indicated. Lastly, to provide the user with a

293 broad picture of lincRNA conservation within their sample set, a bar graph is produced that

294 indicates the number and percent of recovered sequence homologs in each species (Figure

295 S2A).

296 **Data and software availability**

297 All genomes used in this work, including version and source, are listed in File S1. The accession

298 number of all short read archive files (SRA) used in this work, including project ID, TopHat (Kim

299 et al., 2013) read mapping rate, and total reads mapped for each SRA are shown in File S1.

300 Genomic coordinates for lincRNAs identified by Evolinc-I are listed by species in BED/GFF

301  format in File S2. LincRNAs were scanned for the presence of ribosomal and other known

302  RNAs by batch searching against the Rfam database

303  (http://rfam.xfam.org/search#tabview=tab1). Novel lincRNAs have also been deposited within

304  the CoGe environment as tracks for genome browsing (Links found in File S2). Evolinc is

305  available as two apps (Evolinc-I and Evolinc-II) in CyVerse's DE (https://de.cyverse.org/de/), for

306  which a tutorial and sample data are available

307  (https://wiki.cyverse.org/wiki/display/TUT/Evolinc+in+the+Discovery+Environment). Evolinc is

308  also available as self-contained Docker images (https://hub.docker.com/r/evolinc/evolinc-i/ and

309  https://hub.docker.com/r/evolinc/evolinc-ii/) for use in a Linux or Mac OSX command-line

310  environment. The code for Evolinc is available to download/edit as a GitHub repository

311  (https://github.com/Evolinc). Information for installation of the Docker image in a command-line

312  environment, as well as FAQs associated with this process are available in the Evolinc GitHub

313  repository readme file. Both Evolinc tools make use of several open source tools, such as

314  BLAST for sequence comparisons (Altschul et al., 1990; Camacho et al., 2009), Cufflinks

315  (Trapnell et al., 2010) for GFF to FASTA conversion, Bedtools (Quinlan and Hall, 2010) for

316  sequence intersect comparisons, MAFFT (Katoh and Standley, 2013) for sequence alignment,

317  RAxML (Stamatakis, 2014) for inferring phylogeny, Notung (Durand et al., 2006) for reconciling

318  gene and species trees, and python, perl, and R for file manipulation and data reporting.

319  ***RNA-seq read mapping and transcript assembly***

320  SRA files were uploaded directly into CyVerse DE from (http://www.ncbi.nlm.nih.gov/sra) by

321  using the "Import from URL" option. All further read processing was performed using

322  applications within DE. Briefly, uncompressed paired end reads were trimmed (5 nt from 5' end

323  and 10 nt from 3' end) using FASTX trimmer, whereas single end read files were filtered with

324  the FASTX quality filter so that only reads where ≥ 70% of bases with a minimum quality score

325  of 25 were retained (http://hannonlab.cshl.edu/fastx_toolkit/index.html). Reads were mapped to

326 their corresponding genomes using TopHat2 version 2.0.9 (Kim et al., 2013). TopHat2 settings

327 varied based on organism and SRA, and are listed in File S1. Transcripts were assembled using

328 the Cufflinks2 app version 2.1.1 under settings listed in File S1 (Trapnell et al., 2010). TopHat2

329 and Cufflinks2 were executed on reads from each SRA file independently.

330 ***Validation of lincRNA expression in vivo***

331 RNA was extracted from 2-week old seedlings and flower buds from 4-week old Arabidopsis

332 Col-0 using Trizol (ThermoFisher Life Sciences catalog # 15596018). These tissues and age at

333 extraction most closely matched the experiments from which the RNA-seq data was obtained

334 (Liu et al., 2012). cDNA was synthesized using SuperScript III (ThermoFisher Life Sciences

335 catalog # 18080051) and 2μg of RNA as input. Primers were first validated by performing PCR

336 with genomic DNA as template using GoTaq Green polymerase master mix (Promega catalog

337 #M712) with 95°C for 3' to denature, followed by 35 cycles of 95°C for 15", 55°C for 30" and

338 72°C for 30" and a final extension step of 5' at 72°C. Primers used are listed in File S2.

339 **Results**

340 ***An overview of lincRNA identification with Evolinc-I***

341 *Evolinc-I validation*

342 After establishing a workflow using the most commonly accepted parameters for defining a

343 lincRNA (detailed in Materials and Methods), we wanted to evaluate its efficiency at

344 distinguishing between unknown or novel protein-coding genes and non-coding loci. For this, we

345 used a random set of 5,000 protein-coding transcripts selected from the TAIR10 annotation to

346 determine Evolinc-I's false discovery rate (FDR) (i.e., protein-coding transcripts erroneously

347 classified as lincRNAs). ORFs for this test dataset ranged in length from 303 to 4182 nts, with

348 an average ORF of 1131 nts (File S3). Because Evolinc is designed to automatically remove

349 transcripts that map back to known genes, we removed these 5,000 genes from the reference

350  genome annotation file, and then generated a transcript assembly file from RNA-seq data where

351  these 5,000 genes were known to be expressed. We fed the transcript assembly file to Evolinc-

352  I. Out of 5,000 protein-coding genes, only 11 were categorized as non-coding by Evolinc-I

353  (0.22% FDR; File S3). Further investigation of the 11 loci revealed that they were predominantly

354  low coverage transcripts with ORFs capable of producing polypeptides greater than 90, but less

355  than 100 amino acids (aa). Moreover, low read coverage for these transcripts led to incomplete

356  transcript assembly. Together these factors were responsible for the miss-annotation of these

357  loci as non-coding. Importantly, our results indicate that read depth and transcript assembly

358  settings impact lincRNA identification, a finding also noted by Cabilli et al. (2011). Therefore,

359  exploring transcript assembly parameters may be necessary prior to running Evolinc-I. In sum,

360  Evolinc-I has a low FDR that can be further reduced by increasing read per base coverage

361  thresholds during transcript assembly as performed in Cabilli et al. (2011).

362       We determined the overlap of Evolinc predicted lincRNAs with previously published

363  datasets from humans and Arabidopsis, following as closely as possible the methods published

364  for each dataset. We first used Evolinc-I to identify lincRNAs from an RNA-seq dataset

365  generated by Liu et al. (2012) in Arabidopsis (File S1). From nearly one billion reads generated

366  from four different tissues (siliques, flowers, leaves, and roots), Liu et al. (2012) identified 278

367  lincRNAs (based on the TAIR9 reference genome annotation). Using the Liu et al. (2012) SRA

368  data, we mapped RNA-seq reads and assembled transcripts with Tophat2 and Cufflinks2 in the

369  DE. From these transcripts, Evolinc-I, identified 571 lincRNAs. We then reconciled the lincRNAs

370  identified in Liu et al. (Liu-lincRNAs) with those from Evolinc-I (Evolinc-lincRNAs), by identifying

371  overlapping genomic coordinates for lincRNAs from the two datasets using the Bedtools suite

372  (Quinlan and Hall, 2010). Of the 278 Liu-lincRNAs, 261 were also recovered by Evolinc-I (Table

373  S1). Cufflinks failed to assemble the 17 unrecovered Liu-lincRNAs, due to low coverage, and

374    thus differences in recovery for these loci reflect differences in the Cufflinks parameters

375    employed.

376         The Arabidopsis genome reference has been updated since Liu et al. (2012), from

377    TAIR9 to TAIR10 (Lamesch et al., 2012). We also ran Evolinc-I with the TAIR10 annotation and

378    found that only 198 of the 261 Liu-lincRNAs were still considered intergenic (Figure 1B). The

379    remaining 63 were reclassified as overlapping a known gene (either sense overlapping

380    transcript, SOT, or antisense overlapping transcript, AOT). This highlights an important aspect

381    of Evolinc-I. While Evolinc-I is able to identify long non-coding RNAs without a genome

382    annotation, genome annotation quality can impact whether an lncRNA is considered intergenic

383    versus AOT or SOT. In sum, 198 of the 571 lincRNAs identified by Evolinc-I correspond to a

384    previously identified Liu-lincRNA (Figure 1B).

385         Of the 571 lincRNAs identified by Evolinc-I, 373 were not classified as lincRNAs by Liu

386    et al. (2012). Evolinc-I removes transcripts that overlap with the 5' and 3' UTRs of a known

387    gene, whereas Liu et al. (2012) removed transcripts that were within 500 bp of a known gene

388    (Liu et al., 2012). This difference in the operational definition of intergenic space accounts for

389    the omission of 197 Evolinc-lincRNAs from the Liu et al. (2012) lincRNA catalog. In addition,

390    Evolinc-I removes transcripts with high similarity to transposable elements, but not tandem di- or

391    trinucleotide repeats. We could see no biological reason for excluding these simple repeat

392    containing transcripts, and in fact, transcripts with simple tandem repeats have been attributed

393    to disease phenotypes and therefore might be of particular interest (Usdin, 2008). The inclusion

394    of these transcripts accounts for 106 of the unique Evolinc-lincRNAs.

395         Finally, 70 of the 571 Evolinc-lincRNAs were entirely novel, and did not correspond to

396    any known Liu-lincRNA or gene within the TAIR10 genome annotation. To determine whether

397    these represented *bona fide* transcripts, we tested expression of a subset ($n$ = 20) of single and

398    multi-exon putative lincRNAs by RT-PCR using RNA extracted from two different tissues

399    (seedlings and flowers, Figure S1A). We considered expression to be positive if we recovered a

400    band in two different tissues or in the same tissue but from different biological replicates. We

401    recovered evidence of expression for 18 of these putative lincRNAs out of 20 tested. Based on

402    these data we conclude that a majority of the 70 novel lincRNAs identified by Evolinc-I for

403    Arabidopsis are likely to reflect *bona fide* transcripts, and thus valid lincRNA candidates.

404            We next compared Evolinc-I against a well-annotated set of human lincRNAs

405    characterized by Cabili et al. (2011). Cabili et al. (2011) used RNA-seq data from 24 different

406    tissues and cell types, along with multiple selection criteria to identify a "gold standard"

407    reference set of 4,662 lincRNAs. We assembled transcripts from RNA-seq data for seven of

408    these tissues (File S1) using Cufflinks under the assembly parameters and read-per-base

409    coverage cut-offs of Cabili et al. (2011) (see Materials and Methods). We then fed these

410    transcripts to Evolinc-I. To directly compare Evolinc-I identified lincRNAs with the Cabili et al.

411    (2011) reference dataset (Cabili-lincRNAs), we used the BED files generated by Evolinc-I to

412    identify a subset of 360 multi-exon putative lincRNAs that were observed in at least two tissues

413    (consistent with criteria employed in Cabili et al. [2011] when using a single transcript

414    assembler). We then asked whether these 360 Evolinc-I lincRNAs were found in either the

415    Cabili-lincRNAs, or the hg19 human reference annotation (UCSC). A total of 317 (88%) of the

416    Evolinc-I lincRNAs matched known lincRNAs from the two annotation sources (Figure 1C). The

417    remaining 43 transcripts (12% of the 360 tested) passed all other criteria laid out by Cabili et. al.

418    (2011) and therefore may be *bona fide* lincRNAs, but will require further testing.

419

420    **Evolution of lincRNA loci with Evolinc-II**

421    *Evolinc-II validation*

422      Evolinc-II is an automated and improved version of a workflow we previously used to determine

423      the depth to which Liu-lincRNAs (Liu et al., 2012) were conserved in other species of the

424      Brassicaceae (A., D., L., Nelson et al., 2016). The Evolinc-II workflow is outlined in Figure 2A.

425      While most Liu-lincRNAs were restricted to Arabidopsis, or shared only by Arabidopsis and *A.*

426      *lyrata*, 3% were conserved across the family, indicating that the lincRNA-encoding locus was

427      present in the common ancestor of all Brassicaceae ~54 MYA (Beilstein et al., 2010). We used

428      Evolinc-II to recapitulate our previous analysis in three ways. First, to provide replicates for

429      statistical analysis, we randomly divided the 5,361 Liu-lincRNAs into 200-sequence groups prior

430      to Evolinc-II analysis (*n* = 27; Figure 2B and Figure S2B). Second, we performed a separate

431      comparison by dividing the Liu-lincRNAs based upon chromosomal location (*n* = 5). Lastly, we

432      used Evolinc-II to search for sequence homologs using the complete Liu-lincRNA dataset but

433      querying with varying E-value cutoffs (E-20, E-15, E-10, E-05, and E-01). This analysis allowed

434      us to test the impact of the requirement for reciprocity on the recovery of putative homologs

435      under different E-value criteria (Figure 2B and Figure S2D). The number of sequence homologs

436      increased for each decrement in BLAST stringency (Figure S2D), indicating that a significant

437      number of putative homologs fulfill the reciprocity requirement even as sequence similarity

438      decreases. The percentage of sequence homologs retrieved by Evolinc-II was statistically

439      indistinguishable for lincRNAs assigned to groups, chromosomes, or the average from all E-

440      value cutoffs (Figure 2B and Figure S2C). Thus, Evolinc-II is a robust method to identify sets of

441      lincRNAs that are conserved across a user-defined set of species, such as the Brassicaceae.

442      In addition to identifying sets of conserved lincRNAs, Evolinc-II also highlights conserved

443      regions within each query lincRNA. To demonstrate these features, we scanned through the

444      Liu-lincRNA Evolinc-II summary statistics file (at 1E-10; File S4) to identify a conserved

445      lincRNA. At1NC023160 is conserved as a single copy locus in eight of the ten species we

446      examined. It was identified by Liu et al. (2012) based on both RNA-seq and tiling array data, as

447  well as validated by Evolinc-I. During the comparative analyses, Evolinc-II generates a query-

448  centric coordinate file that allows the user to visualize within a genome browser (e.g., JBrowse;

449  [Buels et al., 2016]) what regions of the query lincRNA are most conserved. Using this query-

450  centric coordinate file, we examined the 332 nt At1NC023160 locus in the CoGe genome

451  browser and determined that the 3' end was most highly conserved (Figure 2C). We used the

452  MAFFT multiple sequence alignment generated by Evolinc-II for At1NC023160 to perform

453  structure prediction with RNAalifold (Figure S3A; (Lorenz et al., 2011)). The structural prediction

454  based on the multiple sequence alignment had a greater base pair probability score and lower

455  minimum free energy than the structure inferred from the Arabidopsis lincRNA alone (Figure

456  S3B and S3C). Conserved regions of a lincRNA serve as potential targets for disruption via

457  genome editing techniques, thereby facilitating its functional dissection.

458

459  *Using Evolinc-II to infer the evolution of the human telomerase RNA locus TERC*

460  In addition to exploring the evolutionary history of a lincRNA catalog, Evolinc-II is an effective

461  tool to infer the evolution of individual lincRNA loci. To showcase the insights Evolinc-II can

462  provide for datasets comprised of a small number of lincRNAs, we focused on the well-

463  characterized human lincRNA, TERC. TERC is the RNA subunit of the ribonucleoprotein

464  complex telomerase that is essential for chromosome end maintenance in stem cells, germ-line

465  cells, and single-cell eukaryotes (Theimer and Feigon, 2006; Zhang et al., 2011; Blackburn and

466  Collins, 2011). TERC is functionally conserved across almost all eukarya, but is highly

467  sequence divergent. Building on work performed by Chen et al. (2000) we used Evolinc-II to

468  examine the evolutionary history of the human TERC locus in 26 mammalian species that last

469  shared a common ancestor between 100-130 MYA (Figure 3) (Glazko, 2003; Arnason et al.,

470  2008).

471     Evolinc-II identified a human TERC sequence homolog in 23 of the 26 species examined

472     (Figure 3; raw output shown in Figure S4). We were unable to identify a human TERC homolog

473     in *Ornithoryhnchus anatinus* (platypus), representing the earliest diverging lineage within class

474     Mammalia, using our search criteria. In addition, *Mus musculus* (mouse) and *Rattus norvegicus*

475     (rat) were also lacking a human TERC homolog. However, close relatives of mouse and rat,

476     such as *Ictidomys tridecemlineatus* (squirrel) and *Oryctolagus cuniculus* (rabbit) retained clear

477     human TERC sequence homologs, suggesting that loss of the human TERC-like locus is

478     restricted to the Muridae (mouse/rat family). This is in agreement with the previous identification

479     of the mouse TERC, which exhibits much lower sequence similarity with the human TERC than

480     do other mammals (Chen et al, 2000). All identified human TERC homologs also share synteny,

481     suggesting similar evolutionary origins for this locus throughout mammals (Figure 3). Evolinc-II

482     also identified lineage-specific duplication events for the human TERC-like locus in the

483     orangutan, lemur, and galago genomes (Figure 3), similar to previous observations in pig and

484     cow (Chen et al., 2000). In sum, Evolinc-II can be applied to both large and small datasets to

485     uncover patterns of duplication, loss, and conservation across large phylogenetic distances.

486

487     **Discussion**

488     ***Rapid identification of lincRNAs using Evolinc-I***

489     With Evolinc-I our goal was to develop an automated and simple pipeline for rapid lincRNA

490     discovery from RNA-seq data. In addition to identification, Evolinc-I generates output files that

491     put downstream analyses and data visualization into the hands of biologists, making it simpler

492     for researchers to discover and explore lincRNAs. Evolinc-I makes use of standard lincRNA

493     discovery criteria, and packages each step into easy-to-use applications within the CyVerse DE

494     or for command-line use via a Docker image with all dependencies pre-installed. We

495    recommend the DE-version of Evolinc-I for novice users, whereas the command-line version of

496    Evolinc-I is useful for knowledgeable users wishing to tweak parameters to fit their system or

497    question. By using Evolinc-I within the DE, the user can take advantage of the

498    cyberinfrastructure support of CyVerse (Merchant et al., 2016). One of the key advantages of

499    combining Evolinc-I with cyberinfrastructure such as the CyVerse's DE is the ability to combine

500    various applications together in one streamlined workflow, and making the workflow easier to

501    implement by interested researchers. For instance, a user can download an RNA-seq SRA file

502    into their DE account, quickly process and map reads, assemble transcripts, and execute

503    Evolinc-I. All of this occurs within the DE without downloading a single file or installing a

504    program on a desktop computer.

505       We demonstrated the ability of Evolinc-I to identify lincRNAs from previously curated

506    catalogs for plants and mammals. Note that we were able to account for all differences between

507    results from Evolinc-I and the published studies, indicating that our pipeline is operating under

508    definitions and filters currently used by the community. Moreover, because we have formalized

509    the process by which annotations of genome data can be incorporated into the search strategy,

510    Evolinc-I gives researchers the ability to easily explore the contributions of TEs, repetitive

511    elements, or other user defined features to the prediction of lincRNA loci. Finally, we stress that

512    this tool permits experiments to be repeated by researchers to compare the contribution of

513    recently released annotations, or to repeat experiments from other groups. This latter point

514    cannot be overemphasized as interest in lincRNAs grows.

515

516    ***Examining evolutionary history and patterns of conservation of lincRNA loci using***

517    ***Evolinc-II***

518     Evolinc-II is designed to perform a series of comparative genomic and transcriptomic analyses

519     across an evolutionary timescale of the user's choosing and on any number (1-1000s) of query

520     lincRNAs. Similar to the lncRNA discovery and evolutionary analysis tool Slncky (J., Chen et al.,

521     2016), the analyses performed by Evolinc-II highlight conserved lincRNA loci, conserved

522     regions within those loci, and overlap with transcripts in other species. To develop an

523     informative evolutionary profile, we recommend users incorporate as many genomes as

524     possible for closely related species and then choose more distantly related species based on

525     the level of genome annotation, genome quality, and quantity of lncRNAs identified for those

526     species. The computationally intensive nature of these analyses is ameliorated by taking

527     advantage of a high-performance computing cluster such as CyVerse. While sequence

528     conservation is certainly not the only filtering mechanism to identify functional lncRNAs, we

529     believe that is a critical first step. In the future, as more becomes known about structural

530     conservation within lncRNAs, this aspect of lncRNA evolution will be added as an additional

531     filter. We envision Evolinc-II being useful for both scientists attempting to identify functional

532     regions of a lincRNA as well as those wanting to understand the pressures impacting lincRNA

533     evolution.

534        In addition to highlighting large-scale lincRNA patterns of conservation, we also

535     demonstrated how Evolinc-II can be used to examine the detailed evolutionary history of a

536     single lincRNA, using the human TERC as a test-case. We performed an Evolinc-II analysis

537     with human TERC on 26 genomes in the class Mammalia, 14 of which had not been included in

538     previous studies (Chen et al., 2000). As expected, we recovered a human TERC-like locus in

539     most mammals, as well as three previously unrecorded lineage-specific duplication events.

540     Whether these duplicate TERC loci are expressed and interact with telomerase is unknown; if

541     so they may represent potential regulatory molecules, similar to TER2 in Arabidopsis (Xu et al.,

542     2015; A., D., L., Nelson and Shippen, 2015). We also determined that the human TERC-like

543 locus was lost (or experienced an accelerated mutation rate relative to other mammals) in the

544 common ancestor of mouse and rat. The conservation of the TERC locus across mammals,

545 characterized by rare evolutionary transitions such as that in mouse and rat, stands in stark

546 contrast to the evolution of the telomerase RNA in Brassicaceae (Beilstein et al., 2012), despite

547 the fact that other telomere components are highly conserved (Nelson et al., 2014).

548 Interestingly, mammalian TERCs appear to evolve more slowly than their plant counterparts,

549 similar to the protein components of telomerase (Wyatt et al., 2010). These discoveries highlight

550 the novel insights that can be uncovered using Evolinc-II on even well studied lincRNAs.

551 In summary, Evolinc streamlines lincRNA identification and evolutionary analysis. Given

552 the wealth of RNA-seq data being uploaded on a daily basis to NCBI's SRA, and the increased

553 availability of high performance computing resources, we believe that Evolinc will prove to be

554 tremendously useful. Combining these resources, Evolinc can uncover broad and fine-scale

555 patterns in the way that lincRNAs evolve and ultimately help in linking lincRNAs to their function.

556

568

569  **Figure Legends**

570  **Figure 1. Schematic representation of the Evolinc-I workflow and validation. (A)**

571  Evolinc-I takes assembled transcripts as input and then filters over several steps (1-4).

572  Evolinc generates output files detailed in the materials and methods. **(B)** Evolinc

573  validation on RNA-seq data from Liu et al. (2012). Four tissues were sequenced by Liu

574  et al., as indicated by the red circles, including (from top to bottom) flowers, siliques,

575  leaves, and roots. Assembled transcripts were fed through Evolinc-I to identify

576  lincRNAs, Antisense Overlapping Transcripts (AOTs), and Sense Overlapping

577  Transcripts (SOTs). A reconciliation was performed between the Evolinc-I identified

578  lincRNAs and the Liu et al. dataset. Gene associated transcriptional unit (GATU) and

579  repeat containing transcriptional unit (RCTU) terminology comes from Liu et al. (2012).

580  **(C)** Evolinc validation of Cabili et al. (2011) RNA-seq data. RNA-seq data was

581  assembled and then filtered through additional Cabili-specific parameters (shown in

582  box). The pie chart shows Evolinc-identified lincRNAs that correspond to Cabili et al. or

583  are novel.

584  **Figure 2. Schematic representation of the Evolinc-II workflow and validation of**

585  **Liu-lincRNA and Evolinc-identified lincRNAs. (A)** Evolinc-II uses lincRNAs as a

586  query in reciprocal BLAST analyses against any number of genomes. Sequences that

587  match the filters (see Materials and Methods) are grouped into families of sequences

588  based on the query lincRNA. Each sequence homolog is classified using user-defined

589    data or annotations, such as expression or overlap with known gene or lincRNA.

590    Sequences are aligned to highlight conserved regions and to infer phylogeny. These

591    steps can be performed on thousands to tens of thousands of query lincRNAs. Gene

592    trees are inferred for each sequence family using RAxML. The resulting trees are

593    reconciled with the known species tree using Notung 2.0. Notung delineates gene loss

594    and duplication events by marking the output tree with a D (duplication) and blue

595    branch, or L (loss) and red branch. Phylogenetic inference is computationally intensive,

596    and thus we suggest limiting the number of sequence families for which the analysis is

597    performed. Data files generated by Evolinc-II are described in the Materials and

598    Methods. **(B)** Validation of Evolinc-II by repeating the Liu-lincRNA dataset in three

599    different ways. The ~5400 Liu-lincRNAs were randomly divided into 200 sequence bins

600    (blue bar), each bin was run through Evolinc-II (total number of runs = 27), and then the

601    results were averaged, with standard deviation denoted. In the second analysis, the Liu-

602    lincRNAs were divided based on chromosome, and then each set of Liu-lincRNAs (five

603    groups) were run through Evolinc-II separately. Lastly, all Liu-lincRNAs were run

604    through Evolinc using different BLAST E-value cutoffs (E-1, -5, -10, -15, -20), and the

605    results averaged. Bars represent the percent of Liu-lincRNAs for which sequence

606    homologs were identified. *A. tha = Arabidopsis thaliana*, *A. lyr = Arabidopsis lyrata*, *C.*

607    *rub = Capsella rubella*, *L. ala = Leavenworthia alabamica*, *B. rap = Brassica rapa*, *B. ole*

608    *= Brassica oleracea*, *S. par = Schrenkiella parvula*, *E. sal = Eutrema salsugineum*, *A.*

609    *ara = Aethionema arabicum*, and *T. has = Tarenaya hassleriana*. **(C)** Genome browser

610    visualization of the At1NC023160 locus and its conservation in other Brassicaceae.

611    Regions of the Arabidopsis locus that Evolinc-II identified to be conserved are shown in

612   green, with species of origin listed to the right. The blue bar indicates the length of the

613   locus in Arabidopsis, with the arrow indicating direction of transcription. The region of

614   the locus selected for structural prediction is shown in the red dashed box.

615   **Figure 3. Evolinc-II analysis of the human TERC locus in mammals.** Species tree of

616   26 species within class Mammalia with duplication (D) or loss (L) events hung on the

617   tree (left). A micro-synteny profile is shown to the right for each species, showing the

618   TERC locus in red, and adjacent protein-coding genes in black. Direction of each gene

619   is indicated with arrows. The mouse and rat TERC loci are indicated by blue arrows to

620   represent the poor sequence similarity between these two loci and human TERC.

621   Divergence times are approximate and extracted from Arnason et al. (2008). A key is

622   shown below, with gene names indicated. All pertinent links are shown below to

623   regenerate micro-synteny analyses with CoGe (genomeevolution.org) for all species on

624   the tree.

625   **File S1** List of publically available genome and sequence files used, as well as

626   conditions and results from TopHat and Cufflinks for each assembly.

627   **File S2** Evolinc-I output for all species from which lincRNAs were identified, as well as

628   bed files for genome browser viewing, and primers used in RT-PCR verification of

629   transcription of novel Arabidopsis lincRNAs. Also contains CoGe genome browser links

630   to the novel lincRNAs identified.

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

632 **Figure S1** RT-PCR validation of lincRNAs identified in Arabidopsis by Evolinc-I.

633 LincRNA IDs match those found in File S2. G = genomic DNA positive control. F =

634 flower cDNA, S = seedling cDNA.

635 **Figure S2** Examining conservation of Liu-lincRNAs in multiple ways with Evolinc-II. **(A)**

636 Example of the type of bar graph produced by Evolinc-II, in this case for the Liu-

637 lincRNAs at 1E-20. **(B)** Bar graph of level of lincRNA conservation observed when

638 dividing the Liu-lincRNAs into 27 random bins of 200 lincRNAs each. Standard deviation

639 is based on the difference seen between the 27 bins. **(C)** Bar graph depicting the level

640 of lincRNA conservation seen when dividing the Liu-lincRNAs by Arabidopsis

641 chromosome (E-cutoff value of 1E-20). **(D)** Bar graph demonstrating the level of

642 conservation of the Liu-lincRNAs throughout Brassicaceae at different E-cutoff values.

643 **Figure S3** Using At1NC023160 to highlight the structural information that can be

644 gleaned from Evolinc-II. (A) Multiple sequence alignment, generated by MAFFT and

645 visualized within Geneious v7.1 (Kearse et al., 2012). Similar sequences are

646 highlighted, with the consensus sequence shown on top. Nucleotide identity is shown

647 below the consensus sequence, with green representing 100% identity across all

648 sequences. **(B)** RNAalifold (Lorenz et al., 2011) consensus secondary structure

649 prediction based on multiple sequence alignment in **(A)**. Base-pair probabilities are

650 shown, with red being more probable and blue least probable. **(C)** RNAfold structure

651 prediction based on the same region as in **(B)**, but limited to just the Arabidopsis

652 sequence. Base-pair probabilities are shown as in **(B)**.

653 **Figure S4** Raw phylogenetic output from Evolinc-II for TERC. **(A)** A gene tree for the

654 TERC sequence homologs identified in each of the species shown. Sequences without

655 "TBH" indicate paralogs. **(B)** Notung (Durand et al., 2006) reconciliation of the gene tree

656 shown in (A) to the known species tree. Duplication (red "D") and loss events (grey

657 "LOST") are shown. Support for duplication or loss events are indicated by the green

658 numbers at the nodes that represent the predicted origin of those events.

659 **Table S1** Percent similarity between transcripts identified following transcript assembly

660 and lincRNA identification.

661

662

663 **References:**

664 Altschul, Stephen F., Warren Gish., Webb Miller., Eugene W. Myers., and David J.

665 Lipman. 1990. Basic Local Alignment Search Tool. *Journal of Molecular Biology*

666 215 (3). Academic Press: 403–10. doi:10.1016/S0022-2836(05)80360-2.

667 Anders, Simon., and Wolfgang Huber. 2010. Differential Expression Analysis for

668 Sequence Count Data. *Genome Biology* 11 (10). BioMed Central: R106.

669 doi:10.1186/gb-2010-11-10-r106.

670 Arnason, Ulfur., Joseph A Adegoke., Anette Gullberg., Eric H Harley., Axel Janke., and

671 Morgan Kullberg. 2008. Mitogenomic Relationships of Placental Mammals and

672 Molecular Estimates of Their Divergences. *Gene* 421 (1–2): 37–51.

673 doi:10.1016/j.gene.2008.05.024.

674 Beilstein, Mark a., Amy E. Brinegar., and Dorothy E. Shippen. 2012. Evolution of the

675 Arabidopsis Telomerase RNA. *Frontiers in Genetics* 3 (SEP). Frontiers: 1–8.

676  doi:10.3389/fgene.2012.00188.

677 Beilstein, Mark A., Nathalie S Nagalingum., Mark D Clements., Steven R Manchester.,

678  and Sarah Mathews. 2010. Dated Molecular Phylogenies Indicate a Miocene Origin

679  for Arabidopsis Thaliana. *Proceedings of the National Academy of Sciences of the*

680  *United States of America* 107 (43): 18724–28. doi:10.1073/pnas.0909766107.

681 Blackburn, Elizabeth H., and Kathleen Collins. 2011. Telomerase: An RNP Enzyme

682  Synthesizes DNA. *Cold Spring Harbor Perspectives in Biology* 3 (5): a003558-.

683  doi:10.1101/cshperspect.a003558.

684 Buels, Robert., Eric Yao., Colin M. Diesh., Richard D. Hayes., Monica Munoz-Torres.,

685  Gregg Helt., et al. 2016. JBrowse: A Dynamic Web Platform for Genome

686  Visualization and Analysis. *Genome Biology* 17 (1). BioMed Central: 66.

687  doi:10.1186/s13059-016-0924-1.

688 Cabili, Moran N., Cole Trapnell., Loyal Goff., Magdalena Koziol., Barbara Tazon-Vega.,

689  Aviv Regev., et al. 2011. Integrative Annotation of Human Large Intergenic

690  Noncoding RNAs Reveals Global Properties and Specific Subclasses. *Genes &*

691  *Development* 25 (18). Cold Spring Harbor Laboratory Press: 1915–27.

692  doi:10.1101/gad.17446611.

693 Camacho, C., G Coulouris., V Avagyan., N Ma., J Papadopoulos., K Bealer., et al.

694  2009. BLAST+: Architecture and Applications. *BMC Bioinformatics* 10: 421.

695  doi:10.1186/1471-2105-10-421.

696 Chen, Jenny., Alexander A. Shishkin., Xiaopeng Zhu., Sabah Kadri., Itay Maza.,

697  Mitchell Guttman., et al. 2016. Evolutionary Analysis across Mammals Reveals

698     Distinct Classes of Long Non-Coding RNAs. *Genome Biology* 17 (1). BioMed

699     Central: 19. doi:10.1186/s13059-016-0880-9.

700  Chen, Jiunn-Liang., Maria A Blasco., and Carol W Greider. 2000. Secondary Structure

701     of Vertebrate Telomerase RNA. *Cell* 100 (5): 503–14. doi:10.1016/S0092-

702     8674(00)80687-X.

703  Derrien, Thomas., Rory Johnson., Giovanni Bussotti., Andrea Tanzer., Sarah Djebali.,

704     Hagen Tilgner., et al. 2012. The GENCODE v7 Catalog of Human Long Noncoding

705     RNAs: Analysis of Their Gene Structure, Evolution, and Expression. *Genome*

706     *Research* 22 (9): 1775–89. doi:10.1101/gr.132159.111.

707  Durand, Dannie., Bjarni V Halldórsson., and Benjamin Vernot. 2006. A Hybrid Micro-

708     Macroevolutionary Approach to Gene Tree Reconstruction. *Journal of*

709     *Computational Biology : A Journal of Computational Molecular Cell Biology* 13 (2).

710     Mary Ann Liebert, Inc.  2 Madison Avenue Larchmont, NY 10538 USA: 320–35.

711     doi:10.1089/cmb.2006.13.320.

712  Glazko, G. V. 2003. Estimation of Divergence Times for Major Lineages of Primate

713     Species. *Molecular Biology and Evolution* 20 (3): 424–34.

714     doi:10.1093/molbev/msg050.

715  Guttman, Mitchell., Ido Amit., Manuel Garber., Courtney French., Michael F. Lin., David

716     Feldser., et al. 2009. Chromatin Signature Reveals over a Thousand Highly

717     Conserved Large Non-Coding RNAs in Mammals. *Nature* 458 (7235). Nature

718     Publishing Group: 223–27. doi:10.1038/nature07672.

719  Hangauer, Matthew J., Ian W. Vaughn., and Michael T. McManus. 2013. Pervasive

720      Transcription of the Human Genome Produces Thousands of Previously

721      Unidentified Long Intergenic Noncoding RNAs. *PLoS Genetics* 9 (6).

722      doi:10.1371/journal.pgen.1003569.

723  Hawkes, Emily J., Scott P. Hennelly., Irina V. Novikova., Judith A. Irwin., Caroline

724      Dean., Karissa Y. Sanbonmatsu., et al. 2016. COOLAIR Antisense RNAs Form

725      Evolutionarily Conserved Elaborate Secondary Structures. *Cell Reports* 16 (12).

726      Elsevier: 3087–96. doi:10.1016/j.celrep.2016.08.045.

727  Johnson, Rory., and Roderic Guigó. 2014. The RIDL Hypothesis: Transposable

728      Elements as Functional Domains of Long Noncoding RNAs. *RNA (New York, N.Y.)*

729      20 (7): 959–76. doi:10.1261/rna.044560.114.

730  Kapusta, Aurélie., Zev Kronenberg., Vincent J Lynch., Xiaoyu Zhuo., LeeAnn Ramsay.,

731      Guillaume Bourque., et al. 2013. Transposable Elements Are Major Contributors to

732      the Origin, Diversification, and Regulation of Vertebrate Long Noncoding RNAs.

733      *PLoS Genetics* 9 (4). Public Library of Science: e1003470.

734      doi:10.1371/journal.pgen.1003470.

735  Katoh, Kazutaka., and Daron M Standley. 2013. MAFFT Multiple Sequence Alignment

736      Software Version 7: Improvements in Performance and Usability. *Molecular Biology*

737      *and Evolution* 30 (4): 772–80. doi:10.1093/molbev/mst010.

738  Kearse, Matthew., Richard Moir., Amy Wilson., Steven Stones-Havas., Matthew

739      Cheung., Shane Sturrock., et al. 2012. Geneious Basic: An Integrated and

740      Extendable Desktop Software Platform for the Organization and Analysis of

741      Sequence Data. *Bioinformatics (Oxford, England)* 28 (12): 1647–49.

742     doi:10.1093/bioinformatics/bts199.

743    Kim, Daehwan., Geo Pertea., Cole Trapnell., Harold Pimentel., Ryan Kelley., Steven L

744        Salzberg., et al. 2013. TopHat2: Accurate Alignment of Transcriptomes in the

745        Presence of Insertions, Deletions and Gene Fusions. *Genome Biology* 14 (4).

746        BioMed Central: R36. doi:10.1186/gb-2013-14-4-r36.

747    Lamesch, Philippe., Tanya Z Berardini., Donghui Li., David Swarbreck., Christopher

748        Wilks., Rajkumar Sasidharan., et al. 2012. The Arabidopsis Information Resource

749        (TAIR): Improved Gene Annotation and New Tools. *Nucleic Acids Research* 40

750        (Database issue). Oxford University Press: D1202-10. doi:10.1093/nar/gkr1090.

751    Li, Lin., Steven R Eichten., Rena Shimizu., Katherine Petsch., Cheng-Ting Yeh., Wei

752        Wu., et al. 2014. Genome-Wide Discovery and Characterization of Maize Long

753        Non-Coding RNAs. *Genome Biology* 15 (2): R40. doi:10.1186/gb-2014-15-2-r40.

754    Liu, Jun., Choonkyun Jung., Jun Xu., Huan Wang., Shulin Deng., Lucia Bernad., et al.

755        2012. Genome-Wide Analysis Uncovers Regulation of Long Intergenic Noncoding

756        RNAs in Arabidopsis. *The Plant Cell* 24 (11): 4333–45.

757        doi:10.1105/tpc.112.102855.

758    Lorenz, Ronny., Stephan H Bernhart., Christian Höner Zu Siederdissen., Hakim Tafer.,

759        Christoph Flamm., Peter F Stadler., et al. 2011. ViennaRNA Package 2.0.

760        *Algorithms for Molecular Biology : AMB* 6 (1): 26. doi:10.1186/1748-7188-6-26.

761    Lyons, Eric., Matthew Bomhoff., Fan Li., and Brian D. Gregory. 2014. EPIC-CoGe:

762        Functional and Diversity Comparative Genomics.

763        https://pag.confex.com/pag/xxii/webprogram/Paper10615.html.

764  Merchant, Nirav., Eric Lyons., Stephen Goff., Matthew Vaughn., Doreen Ware., David

765      Micklos., et al. 2016. The iPlant Collaborative: Cyberinfrastructure for Enabling

766      Data to Discovery for the Life Sciences. *PLOS Biology* 14 (1). Public Library of

767      Science: e1002342. doi:10.1371/journal.pbio.1002342.

768  Migeon, Barbara R., Ethan Kazi., Camille Haisley-Royster., Jie Hu., Roger Reeves.,

769      Linda Call., et al. 1999. Human X Inactivation Center Induces Random X

770      Chromosome Inactivation in Male Transgenic Mice. *Genomics* 59 (2): 113–21.

771      doi:10.1006/geno.1999.5861.

772  Mohammadin, Setareh., Patrick P Edger., J Chris Pires., and Michael Eric Schranz.

773      2015. Positionally-Conserved but Sequence-Diverged: Identification of Long Non-

774      Coding RNAs in the Brassicaceae and Cleomaceae. *BMC Plant Biology* 15 (1):

775      217. doi:10.1186/s12870-015-0603-5.

776  Necsulea, Anamaria., Magali Soumillon., Maria Warnefors., Angélica Liechti., Tasman

777      Daish., Ulrich Zeller., et al. 2014. The Evolution of lncRNA Repertoires and

778      Expression Patterns in Tetrapods. *Nature* 505 (7485). Nature Publishing Group, a

779      division of Macmillan Publishers Limited. All Rights Reserved.: 635–40.

780      doi:10.1038/nature12943.

781  Nelson, Andrew D. L., and Dorothy E. Shippen. 2015. Evolution of TERT-Interacting

782      lncRNAs: Expanding the Regulatory Landscape of Telomerase. *Frontiers in*

783      *Genetics* 6 (September): 1–6. doi:10.3389/fgene.2015.00277.

784  Nelson, Andrew D L A.D.L., Evan S. E.S. Forsythe., Xiangchao Gan., Miltos Tsiantis.,

785      and M.A. Mark a. Beilstein. 2014. Extending the Model of Arabidopsis Telomere

Length and Composition across Brassicaceae. *Chromosome Research* 22 (2):

153–66. doi:10.1007/s10577-014-9423-y.

Nelson, Andrew D L., Evan S Forsythe., Upendra K Devisetty., David S Clausen., Asher

K Haug-Batzell., Ari M R Meldrum., et al. 2016. A Genomic Analysis of Factors

Driving lincRNA Diversification: Lessons from Plants. *G3&amp;#58;*

*Genes|Genomes|Genetics* 6 (September). Genetics Society of America: 2881–91.

doi:10.1534/g3.116.030338.

Pertea, Mihaela., Daehwan Kim., Geo M Pertea., Jeffrey T Leek., and Steven L

Salzberg. 2016. Transcript-Level Expression Analysis of RNA-Seq Experiments

with HISAT, StringTie and Ballgown. *Nat Protocols* 11 (9): 1650–67.

doi:10.1038/nprot.2016.095.

Quinlan, Aaron R., and Ira M Hall. 2010. BEDTools: A Flexible Suite of Utilities for

Comparing Genomic Features. *Bioinformatics (Oxford, England)* 26 (6). Oxford

University Press: 841–42. doi:10.1093/bioinformatics/btq033.

Quinn, Jeffrey J., Qiangfeng C Zhang., Plamen Georgiev., Ibrahim A Ilik., Asifa Akhtar.,

and Howard Y Chang. 2016. Rapid Evolutionary Turnover Underlies Conserved

lncRNA-Genome Interactions. *Genes & Development* 30 (2). Cold Spring Harbor

Laboratory Press: 191–207. doi:10.1101/gad.272187.115.

Robinson, Mark D., and Alicia Oshlack. 2010. A Scaling Normalization Method for

Differential Expression Analysis of RNA-Seq Data. *Genome Biology* 11 (3). BioMed

Central: R25. doi:10.1186/gb-2010-11-3-r25.

Spannagl, Manuel., Thomas Nussbaumer., Kai C. Bader., Mihaela M. Martis., Michael

808         Seidel., Karl G. Kugler., et al. 2016. PGSB plantsDB: Updates to the Database

809         Framework for Comparative Plant Genome Research. *Nucleic Acids Research* 44

810         (D1): D1141–47. doi:10.1093/nar/gkv1130.

811 Stamatakis, Alexandros. 2014. RAxML Version 8: A Tool for Phylogenetic Analysis and

812         Post-Analysis of Large Phylogenies. *Bioinformatics (Oxford, England)* 30 (9): 1312–

813         13. doi:10.1093/bioinformatics/btu033.

814 Tang, Haibao., and Eric Lyons. 2012. Unleashing the Genome of Brassica Rapa.

815         *Frontiers in Plant Science* 3 (January). Frontiers: 172. doi:10.3389/fpls.2012.00172.

816 Theimer, Carla A., and Juli Feigon. 2006. Structure and Function of Telomerase RNA.

817         *Current Opinion in Structural Biology* 16 (3): 307–18. doi:10.1016/j.sbi.2006.05.005.

818 Trapnell, Cole., Brian A Williams., Geo Pertea., Ali Mortazavi., Gordon Kwan., Marijke J

819         van Baren., et al. 2010. Transcript Assembly and Quantification by RNA-Seq

820         Reveals Unannotated Transcripts and Isoform Switching during Cell Differentiation.

821         *Nature Biotechnology* 28 (5). Nature Publishing Group: 511–15.

822         doi:10.1038/nbt.1621.

823 Ulitsky, Igor. 2016. Evolution to the Rescue: Using Comparative Genomics to

824         Understand Long Non-Coding RNAs. *Nature Reviews Genetics* 17 (10). Nature

825         Research: 601–14. doi:10.1038/nrg.2016.85.

826 Ulitsky, Igor., and David P Bartel. 2013. Ulitsky, Igor, and David P Bartel. 2013.

827         "lincRNAs: Genomics, Evolution, and Mechanisms." Cell 154 (1): 26–46.

828         doi:10.1016/j.cell.2013.06.020.lincRNAs: Genomics, Evolution, and Mechanisms.

829         *Cell* 154 (1): 26–46. doi:10.1016/j.cell.2013.06.020.

830    Ulitsky, Igor., Alena Shkumatava., Calvin H Jan., Hazel Sive., and David P Bartel. 2011.

831        Conserved Function of lincRNAs in Vertebrate Embryonic Development despite

832        Rapid Sequence Evolution. *Cell* 147 (7): 1537–50. doi:10.1016/j.cell.2011.11.055.

833    Usdin, Karen. 2008. The Biological Effects of Simple Tandem Repeats: Lessons from

834        the Repeat Expansion Diseases. *Genome Research* 18 (7). Cold Spring Harbor

835        Laboratory Press: 1011–19. doi:10.1101/gr.070409.107.

836    Wang, Dong., Zhipeng Qu., Lan Yang., Qingzhu Zhang., Zhi-Hong Liu., Trung Do., et al.

837        2017. Transposable Elements (TEs) Contribute to Stress-Related Long Intergenic

838        Noncoding RNAs in Plants. *The Plant Journal*, January. doi:10.1111/tpj.13481.

839    Wang, Huan., Qi-Wen Niu., Hui-Wen Wu., Jun Liu., Jian Ye., Niu Yu., et al. 2015.

840        Analysis of Non-Coding Transcriptome in Rice and Maize Uncovers Roles of

841        Conserved lncRNAs Associated with Agriculture Traits. *The Plant Journal* 84 (2):

842        404–16. doi:10.1111/tpj.13018.

843    Wang, Kevin C., and Howard Y Chang. 2011. Molecular Mechanisms of Long

844        Noncoding RNAs. *Molecular Cell* 43 (6): 904–14. doi:10.1016/j.molcel.2011.08.018.

845    Washietl, Stefan., Sven Findeiss., Stephan A Müller., Stefan Kalkhof., Martin von

846        Bergen., Ivo L Hofacker., et al. 2011. RNAcode: Robust Discrimination of Coding

847        and Noncoding Regions in Comparative Sequence Data. *RNA (New York, N.Y.)* 17

848        (4): 578–94. doi:10.1261/rna.2536111.

849    Wyatt, Haley D M., Stephen C West., and Tara L Beattie. 2010. InTERTpreting

850        Telomerase Structure and Function. *Nucleic Acids Research* 38 (17). Oxford

851        University Press: 5609–22. doi:10.1093/nar/gkq370.

852    Xu, Hengyi;, Andrew D L; A.D.L. Nelson., and Dorothy E; D.E. Shippen. 2015. A

853        Transposable Element within the Non-Canonical Telomerase RNA of Arabidopsis

854        Thaliana Modulates Telomerase Activity in Response to DNA Damage. *PLoS*

855        *Genetics* In press (6). Public Library of Science: e1005281.

856        doi:10.1371/journal.pgen.1005281.

857    Zhang, Qi., Nak-Kyoon Kim., and Juli Feigon. 2011. Architecture of Human Telomerase

858        RNA. *Proceedings of the National Academy of Sciences of the United States of*

859        *America* 108 (51): 20325–32. doi:10.1073/pnas.1100279108.

860    Zhang, Yu-Chan., Jian-You Liao., Ze-Yuan Li., Yang Yu., Jin-Ping Zhang., Quan-Feng

861        Li., et al. 2014. Genome-Wide Screening and Functional Analysis Identify a Large

862        Number of Long Noncoding RNAs Involved in the Sexual Reproduction of Rice.

863        *Genome Biology* 15 (12). BioMed Central: 512. doi:10.1186/s13059-014-0512-1.

864

865

866
867
868

Figure 01.TIF



Nelson et al, Figure 1

Figure 02.TIF



Nelson et al, Figure 2

Figure 03.TIF

To regenerate the micro-synteny analyses
https://genomevolution.org/r/lxvp
https://genomevolution.org/r/lxvo
https://genomevolution.org/r/lxvn
https://genomevolution.org/r/lxz6