



HAL
open science

A comprehensive catalog of LncRNAs expressed in T-cell acute lymphoblastic leukemia

Yasmina Kermezli, Wiam Saadi, Mohamed Belhocine, E. L. Mathieu,
Marc-Antoine Garibal, Vahid Asnafi, Mourad Aribi, Salvatore Spicuglia,
Denis Puthier

► **To cite this version:**

Yasmina Kermezli, Wiam Saadi, Mohamed Belhocine, E. L. Mathieu, Marc-Antoine Garibal, et al..
A comprehensive catalog of LncRNAs expressed in T-cell acute lymphoblastic leukemia. *Leukemia &
lymphoma*, 2019, pp.1-13. 10.1080/10428194.2018.1551534 . hal-02078367

HAL Id: hal-02078367

<https://amu.hal.science/hal-02078367>

Submitted on 19 Jun 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Title: A comprehensive catalog of LncRNAs expressed in T-Cell Acute**
2 **Lymphoblastic Leukemia.**

3
4 **Authors:** Yasmina Kermezli¹⁻³, Wiam Saadi¹⁻³, Mohamed Belhocine^{1,2,5}, Eve-Lyne
5 Mathieu^{1,2}, Marc-Antoine Garibal^{1,6}, Vahid Asnafi⁴, Mourad Aribi^{3,#}, Salvatore Spicuglia^{1,2*}
6 and Denis Puthier^{1,2*}
7
8

9 **AUTHOR AFFILIATIONS**

10 ¹Aix-Marseille University, Inserm, TAGC, UMR1090, Marseille, France.

11 ²Equipe Labéllisée Ligue Nationale Contre le Cancer

12 ³Tlemcen University, The Laboratory of Applied Molecular Biology and Immunology, Algeria.

13 ⁴Université Paris Descartes Sorbonne Cité, Institut Necker-Enfants Malades (INEM), Institut National
14 de la Santé et de la Recherche Médicale (Inserm) U1151, and Laboratory of Onco-Haematology,
15 Assistance Publique-Hôpitaux de Paris (AP-HP), Hôpital Necker Enfants-Malades, Paris, France

16 ⁵Present address: Molecular Biology and Genetics Laboratory, Dubai, United Arab Emirates.

17 ⁶Present address: Aix Marseille University, Inserm, MMG, Marseille France

18 #Senior author

19
20 *Correspondence: salvatore.spicuglia@inserm.fr, denis.puthier@univ-amu.fr
21

22 **Abstract**
23

24 Several studies have demonstrated that LncRNAs can play major roles in cancer
25 development. The creation of a catalogue of LncRNAs expressed in T cell acute
26 lymphoblastic leukemia (T-ALL) is thus of particular importance. However, this task is
27 challenging as LncRNA expression is highly restricted in a time and space manner and may
28 thus greatly differ between samples. We performed a systematic transcript discovery in RNA-
29 Seq data obtained from T-ALL primary cells and cell lines. This led to the identification of
30 2560 novel LncRNAs. After the integration of these transcripts into a large compendium of
31 LncRNAs (n=30478) containing both known LncRNAs and those previously described in T-
32 ALLs, we then performed a systematic genomic and epigenetic characterization of these
33 transcript models demonstrating that these novel LncRNAs share properties with known
34 LncRNAs. Finally, we provide evidences that these novel transcripts could be enriched in
35 LncRNAs with potential oncogenic effects and identified a subset of LncRNAs coregulated
36 with T-ALL oncogenes. Overall, our study represents a comprehensive resource of LncRNAs
37 expressed in T-ALL and might provide new cues on the role of lncRNAs in this type of
38 leukemia.
39

40 **Keywords:** Large non-coding RNA, LncRNA, T-cell acute leukemia, T-ALL, oncogenes
41
42
43
44
45
46
47

48 **Introduction**

49

50 Long non-coding RNAs (LncRNAs) are a novel class of untranslated RNA species
51 defined as transcripts with poor coding potential and size above 200 nucleotides[1,2]. They
52 can lie in both sense and antisense direction of exonic or intronic elements, or in intergenic
53 regions (a subclass termed ‘long intergenic non coding RNAs’, LincRNA), or even in the
54 promoter regions of coding [3]. LncRNAs are transcribed by RNA polymerase II and mirror
55 the features of protein-coding genes, such as polyadenylation and splicing, without
56 containing a functional open reading frame. They are often transcribed at lower abundance
57 than coding genes and in a more tissue-specific manner. In this regard, the function of these
58 transcripts is suggested to be particularly important to shape cell identity. Several studies
59 have demonstrated that lncRNAs are functional and regulate both the expression of
60 neighboring genes and distant genomic sequences by a variety of mechanisms [4]. A growing
61 number of examples also demonstrated that LncRNAs play a major role in cancer
62 development by acting on different levels of regulation to disrupt cellular regulatory
63 networks including proliferation, immortality and motility [5].

64

65 T-cell acute lymphoblastic leukemia (T-ALL) is an aggressive hematological cancer
66 arising from the transformation of T cell [6,7]. Cytogenetic and global transcriptomic
67 analyses led to the classification of T-ALL into molecular groups characterized by the
68 abnormal expression of specific transcription factors (TAL; LMO1/2; TLX1/3; LYL; HOXA;
69 MEF2c, respectively) and their block of differentiation at specific stages [8,9]. Although, the
70 outcome of T-ALLs has globally improved by modern poly-chemotherapy, T-ALL remains
71 of poor prognosis notably in relapsing cases. A major obstacle to understanding the
72 mechanisms of T-ALL oncogenesis is the heterogeneous cellular and molecular nature of the
73 disease, which is driven by a complex interplay of multiple oncogenic events. In this context,
74 LncRNA signatures have been shown to define oncogenic subtypes [10] and several
75 LncRNAs regulated by key T-ALL oncogenes have been identified [11–14].

76

77 The time and space restricted expression of LncRNAs makes it challenging to
78 envision the creation of a complete catalog of LncRNAs. Yet such a catalog appears as a
79 prerequisite to better characterize LncRNAs involved in pathological processes. We thus
80 performed systematic transcripts discovery in a collection of T-ALL samples [15] and
81 integrate previously created catalogs into a non-redundant set. The subsequent list of
82 LncRNAs was thoroughly characterized regarding genomic structure and epigenetic features.
83 Finally, we set up a strategy to prioritize LncRNAs having potential oncogenic effects. This
84 approach allowed us to point out LncRNA candidates potentially relevant in the leukemia
85 pathogenesis.

86

87

88

89

90

91

92 **Material and methods**

93

94 ***De novo* LncRNAs discovery**

95 RNA-Seq experiment from Atak *et al.* [15], were retrieved in BAM format from European
96 Genome-phenome Archive under accession number EGAS00001000536. The alignments
97 (BAM files) were provided to cufflinks (v2.2.1) which aims at assembling reads
98 into transcript models. Cufflinks was used with default settings, except for arguments “-j/
99 pre-mrna-fraction” (set to 0.6) and “-a/--junc-alpha” (set to 0.00001) in order to reduce the
100 number of intronic transcripts (that may correspond to fragments of immature transcripts) and
101 to include well-supported exons into transcript models [16]. The transcript models obtained
102 from the 50 samples (31 primary T-ALL patients, 18 T-ALL cell lines and 1 pool of 5
103 thymuses) were subjected to a cleaning procedure using bedtools (v2.17.0) in order to remove
104 all transcripts described in hg19 RefSeq annotation (Illumina iGenomes web site) [17].
105 Cuffcompare v2.1.1 was then used to merge all files and remove transcript model redundancy
106 [18]. The subsequent gtf file was then filtered to eliminate any transcript model defined in
107 RefSeq, Gencode V19 and new lincRNAs discovered by Trimarchi and his coworkers [11].
108 Transcripts expression levels were estimated using cufflink ('-G' option) and only those with
109 FPKM greater than 1 in at least one sample were kept. Filters on transcripts size (at least 200
110 nucleotides as defined for LncRNAs), number of exons (at least 2 exons) and poor coding
111 potential (CPAT score lower than 0.2) as expected from LncRNAs [19] were subsequently
112 applied.

113

114 **Genomic annotation of transcripts**

115 The subsequent gtf, including all transcripts categories, ([Dataset 1](#)) was then used to
116 perform genomic annotation. All analyses were done using R software or Python scripts.

117

118 **Assessment of LncRNA Tissue Specificity**

119 We processed a set of fastq files corresponding to 20 human tissues (SRA accession
120 number SRP056969). After read mapping (tophat2), the genes expression levels were
121 quantified using Cuffdiff [18]. The gene expression specificity of each gene was computed
122 across all tissues using the tau score [20]:

123

$$124 \quad \tau = \sum_{i=0}^n \frac{(1-\hat{x}_i)^n}{n-1} ; \hat{x}_i = \frac{x_i}{\max_{1 \leq i \leq n}(x_i)}$$

125

126 Where (i) n corresponds to the number of samples, (ii) x_i corresponds to the expression level
127 (log2-transformed FPKM values) in condition i (iii) and $\max(x)$ corresponds to the
128 maximum expression level through all tissues.

129

129 **Epigenetic characterization of LncRNA**

130 The ChIP-Seq datasets obtained from thymus and 3 T-ALL cell lines (DND41, Jurkat
131 and RPMI-8402) and corresponding to H3K4me3 and H3K27ac were obtained from
132 ENCODE and GEO databases. The H3K4me3 and H3K27ac ChIP-Seq in RPMI-8402 cell
133 line were sequenced in our laboratory (SRA accession numbers SRX3437292 and
134 SRX3437293, see Table S1). To measure the ChIP-Seq signal around the TSS ([-3000,

135 +3000] pb) we focused on genes with FPKM above 1. Coverage analyses were performed
136 using a Python script making calls to the pyBigWig python library.

137

138 **Search for potential oncogene**

139 We computed the variance of Log2-transformed FPKM values as a score to find
140 genes displaying high dispersion of expression levels across samples. Pearson's correlation
141 coefficients were computed, using R software, between the top 10% variants LncRNAs and
142 known coding T-ALL oncogenes to bring out coding-noncoding pairs.

143

144 **LncRNA expression analyses**

145 Total RNA was extracted using TRIzol (Invitrogen) according to the manufacturer's
146 instruction. 1 µg of RNA was treated with 1 U of DNase I (Ambion) and incubated at 37°C
147 for 30 min. DNase I was then inactivated (15mM of EDTA and incubation at 75°C for 10
148 minutes). DNase-treated RNAs were reversed transcribed using SuperScript II (Invitrogen)
149 and oligo (dT) or random primers according to the manufacturer's instruction. Control
150 genomic DNA was purified from RPMI-8402 cells using the DNeasy Kit (Qiagen) according
151 to the manufacturer's specifications. Sequences of primers used for PCR of LncRNA
152 *XLOC_00017544_Atak*, *XLOC_00009269_Atak* and *XLOC_00012823_Atak* are provided in
153 **Table S2**. PCR using 1 µL of cDNA was performed with Herculase II Fusion kit (Agilent,
154 Waldbronn, Germany) following manufacturer instructions. Amplifications were carried out
155 with 40 cycles (95°C for 1 minute, denaturation at 95°C for 20 seconds, annealing for 20
156 seconds, extension at 68°C for 1 minutes), followed by a final extension step (68°C for 4
157 minutes).

158

159 **Quantitative reverse transcriptase polymerase chain reaction (qRT-PCR)**

160 The qPCR with Power SYBR green mix (Thermo Fisher) was performed on a
161 Mx3000P real-time PCR system. Each reaction was performed with 2µl of cDNA. *GAPDH*
162 was used as reference for normalization.

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178 Results

179

180 Building a catalog of LncRNAs expressed in T-ALLs

181

182 In order to get an exhaustive catalog of LncRNAs expressed in T-ALLs we first
183 performed a systematic transcript discovery on 50 RNA-Seq samples (a pool of 5 normal
184 thymuses, 31 T-ALLs primary blasts and 18 T-ALLs cell lines) previously described by Atak
185 *et al* [15] (**Figure 1A**). *De novo* transcripts were filtered and only multi-exonic transcripts
186 with size greater than 200 bp, FPKM greater than 1 and coding potential lower than 0.2 were
187 kept. This transcript models were then merged with known LncRNAs obtained from
188 GENCODE version 19 [21] and T-ALL LncRNAs described by Trimarchi *et al* [11]. The
189 final catalog contains a non-redundant list of 30478 LncRNAs. This encompasses 26092
190 from GENCODE (LncRNA_Known), 1826 LncRNAs from the Trimarchi dataset
191 (LncRNA_Trimarchi) and 2560 new LncRNAs from the Atak dataset (LncRNA_Atak). The
192 corresponding GTF file is provided as supplementary ([Dataset 1](#)).

193

194 Genomic characterization

195

196 We next performed a thorough genomic comparison of the three different sets of
197 LncRNA transcripts obtained from GENCODE, Trimarchi *et al* and our own analysis of Atak
198 *et al* RNA-Seq data. They will be denoted hereafter as LncRNA_Known,
199 LncRNA_Trimarchi and LncRNA_Atak respectively. Throughout the analysis, these three
200 sets of LncRNAs were compared to coding transcripts (mRNA) in order to underlie their
201 specific properties. **Figure 1B** shows that the number of exons differs between mRNA
202 transcripts (which mainly contain 5 exons or more) and the three sets of LncRNAs. This
203 underscores the unusual exonic structure of LncRNAs that tends to be limited to two exons as
204 already reported by others [21]. Note that the *de novo* LncRNA_Atak dataset lack mono-
205 exonic transcripts as they were discarded during the filtering process. Regarding transcript
206 size, the reference LncRNAs were found to be shorter than mRNAs (average size of 0.8 KB
207 compared to 2.5 KB) as observed by Derrien *et al* [21]. The mean size of *de novo*
208 LncRNA_Atak was close to that of the reference (LncRNA_Known) validating our
209 procedure of transcript reconstruction (**Figure 1C**). In contrast, the mean size of LncRNAs
210 defined by Trimarchi *et al* were greater (4 kb) than the mean size of mRNAs (2.5 kb), which
211 may point out an intrinsic difference in the procedure used for reconstruction of underlying
212 transcript models. Concerning chromosomal distribution, the LncRNAs tend to be similarly
213 spread throughout the chromosomes while several differences were observed (**Figure 1D**).
214 LncRNA_Trimarchi dataset is enriched in transcripts from chromosome 13 and Y while
215 depleted of transcripts located on chromosome 19. Such a result may probably highlight the
216 representation of some particular tumor karyotypes and gender distribution in the samples
217 used by Trimarchi *et al*. In contrast, the proportion of LncRNAs from LncRNAs_Atak
218 dataset is very similar throughout the chromosomes although their representation is slightly
219 increased on chromosome 21 and 22. LncRNAs are generally classified based on their
220 location with respect to protein-coding genes. We defined five types of LncRNAs (**Figure**
221 **1E**): (i) ‘Intergenic’, transcribed outside of any known coding gene; (ii) ‘Divergent’,

222 produced in promoter regions of coding genes on opposite strand; (iii) ‘Convergent’ whose
223 transcription ends in 3' regions of coding genes on opposite strand (iv) “Sense” and (v)
224 “Antisense” whose transcription takes place inside the gene body of a coding gene in sense
225 or antisense direction, respectively. Regarding these five classes, the composition of
226 LncRNA_Atak dataset was rather close to the reference with 55.38% and 55.56% of
227 intergenic transcripts respectively although less transcripts were classified as divergent
228 (6.57% versus 14.44% for GENCODE) and more transcripts were labelled as antisense
229 (32.30% versus 21.34%) (**Figure 1F**). In contrast, the Trimarchi dataset was found to be
230 mainly composed of intergenic transcripts (90.47%) and divergent transcripts (8.12%) since
231 the other classes were discarded during the building steps of the catalog [11].

232

233 **LncRNAs expression in T-ALL and thymus**

234

235 We next intended to assess the expression of these LncRNAs in several sample
236 groups including normal thymus, T-ALL cell lines and patient samples. We computed, for
237 each transcript, its median expression level across each sample group (**Figure 2A**). In
238 agreement with the weak expression level of LncRNAs reported earlier [21], the mRNAs
239 were more highly expressed than any of the three LncRNA sets. Transcripts from
240 LncRNA_Atak were found to be more highly expressed than LncRNA_Known in Thymus
241 and T-ALL patients while slightly less expressed in cell lines. Of note, however, weaker
242 expression was observed in the LncRNA_Trimarchi dataset suggesting that the lack of
243 accuracy in transcript reconstruction step may also impair proper quantification of these
244 transcripts. LncRNAs are known to be highly tissue specific compared to mRNAs. To verify
245 this, we computed the *tau* tissue-specificity score [20] using a public RNA-Seq dataset
246 encompassing 20 human tissues [22]. This score ranges from 0 for housekeeping genes to 1
247 for highly tissue-specific genes. As expected, mRNAs displayed a bimodal signal with a
248 major fraction of genes behaving as ubiquitous genes and a minor fraction having high. In
249 contrast, a clear shift toward high tissue-specificity scores were observed for LncRNAs
250 regardless of the underlying groups (**Figure 2B**). Moreover, assessment of expression in
251 individual tissues demonstrated a strong bias for thymus-specific LncRNAs in the
252 LncRNA_Trimarchi and LncRNA_Atak dataset (**Supplementary Figure S1**). This also
253 underlines that while LncRNA_Atak were selected against LncRNA_Trimarchi, numerous
254 LncRNAs with strong expression bias in the thymus remained to be discovered. Altogether,
255 these results underscore the enrichment that exists in our catalog for LncRNA biases toward
256 tissue-specificity.

257

258 **Functional annotation**

259 Many LncRNAs have been shown to regulate the expression of neighbor genes in cis
260 [11,23,24]. Therefore, to characterize the functional relevance of the LncRNA sets we
261 performed a functional annotation of their closest neighbor genes using GREAT [25]. For the
262 LncRNA_Known class, significant enrichment was observed for annotation terms related to
263 ubiquitous processes including ‘genes expression’ and ‘metabolism’ (**Figure 3A**). In contrast,
264 for the LncRNA_Trimarchi set, annotation terms related to immune system were

265 significantly enriched, including: ‘regulation of interleukin 4 production’, ‘leukocyte
266 activation’ and ‘T cell receptor V(D)J recombination’ (**Figure 3B**). In the same way, closest
267 genes for LncRNA_Atak dataset were related to ‘negative regulation of Notch signaling
268 pathway’ or ‘regulation of leukocyte degranulation’ for instance (**Figure 3C**). This indicates
269 that both LncRNA_Atak and LncRNA_Trimarchi datasets are enriched for LncRNAs located
270 close to coding genes having major role in normal immune processes and leukemia
271 development. As some LncRNA have been shown to regulate protein-coding genes in cis,
272 this would suggest that some of our newly discovered transcripts could potentially act on key
273 genes regulating immune response and oncogenic processes.

274

275 **Epigenetic features of LncRNAs**

276

277 LncRNAs are known to share epigenetic features with coding genes. Both H3K4me3
278 and H3K27ac have been described as epigenetic marks strongly associated to the promoter
279 region of expressed genes. In order to compare epigenetic features across all transcript sets,
280 we used H3K4me3 and H3K27ac ChIP-Seq obtained from normal thymus and three T-ALL
281 cell lines (DND41, RPMI-8402 and Jurkat). We filtered LncRNAs and mRNAs based on
282 their expression in the corresponding samples by selecting transcripts with FPKM above 1.
283 Using ChIP-seq datasets for H3K4me3 and H3K27ac, we then computed the number of reads
284 falling in binned regions around the promoter (defined as [-3000, 3000] pb around the TSS)
285 for each gene. The mean number of reads for each bin across all genes of a class was used to
286 compute the meta profile shown in Figure 4. Although the results slightly differ between the
287 samples, the LncRNAs and mRNAs sets displayed consistent epigenetic profiles. A striking
288 difference is observed for the LncRNA_Trimarchi set, which display high levels of H3K27ac
289 likely indicating a location bias toward enhancer regions [26].

290

291 **Experimental validation expression of three LncRNAs in RPMI-8402, and Jurkat cell 292 lines**

293

294 We next aimed at validating the expression of *de novo* identified LncRNAs from the
295 LncRNA_Atak dataset. We selected three LncRNAs (*XLOC_00017544_Atak*,
296 *XLOC_00009269_Atak* and *XLOC_00012823_Atak*) located on chromosome 9, 2 and 3 and
297 containing 5, 8 and 2 exons respectively (see [Dataset 1](#) for coordinates). RNA-Seq and ChIP-
298 Seq signals indicated that all three LncRNAs were expressed in RPMI-8402 and Jurkat cell
299 lines (**Figure 5A**). This result was confirmed for the 3 candidates by RT-PCR performed on
300 RNA isolated from RPMI-8402, and Jurkat cell lines. PCR products corresponding to DNA
301 fragments of expected sizes were observed in all three cases (**Figure 5B**).

302

303 **Variability of gene expression among T-ALL samples predict potential oncogenic 304 LncRNAs**

305

306 T cell transformation is related to many genomic and chromosomal abnormalities,
307 which can lead to aberrant gene transcription [6,27]. Many of the described oncogenes are
308 not expressed in normal T cell development [28–30] and only restricted to a subset of T-ALL

309 samples. Therefore, it is expected that the expression of these oncogenes should be associated
310 with a high variance across leukemic samples. Based on this hypothesis, we aimed at mining
311 our LncRNA dataset for potentially new oncogenes using the variance as a proxy.

312

313 We first computed the variance of coding genes across T-ALL cell lines and patient
314 samples and check our ability to recover known T-ALL oncogenes [28]. As depicted in
315 **Figure 6A**, typical leukemia oncogenes (TAL1, TLX1, TLX3, HOXA9, NKX3-1, LMO2)
316 were ranked within the top 10 % of genes with the highest variance in the cell lines and
317 patients. A statistical analysis demonstrated that both in cell lines and patients, leukemia
318 oncogenes have significantly higher variance compared to non-oncogenic genes or to a
319 random list of genes matched for expression distribution (**Figure 6B**). The same prioritizing
320 strategy was applied to LncRNAs in order to identify potential oncogene candidates.
321 Strikingly, numerous LncRNAs known for their implication in cancer (*e.g.* H19, XIST,
322 LUNAR1, MIAT and NEAT1) were ranked within the top 10% of LncRNAs displaying the
323 highest variance in both cell lines and patients. Interestingly, several *de novo* LncRNAs from
324 the Atak dataset, including the 3 LncRNAs validated in **Figure 5**, were found among the
325 highest variable transcripts (**Figure 6C**). Moreover, the variance of LncRNA_Atak list was
326 found to be significantly higher when compared to the two other sets, suggesting that it may
327 be enriched for potential oncogenic LncRNAs (**Figure 6D**). As an example, we validated the
328 variable expression of *XLOC_00000871_Atak*, one of the LncRNAs with the highest variance
329 in both T-ALL patients and cell lines (**Figure 6B and S2**), by RT-qPCR across a panel of T-
330 ALL cell lines (**Figure 6E**).

331

332 **Correlation between T-ALL oncogenes and LncRNAs**

333

334 To address the possibility that some of the highly variable LncRNAs might be
335 associated with the regulation of key oncogenes, we computed the expression correlation
336 between the 10% of LncRNAs with highest variance and the set of known T-ALL oncogenes
337 and retrieved the correlated gene-lncRNA pairs. 60.5 % (1146) and 59% (1116) of these
338 LncRNAs were correlated ($r > 0,5$) with at least one oncogene in T-ALL cell lines and
339 patients, respectively (**Supplementary dataset 2-3**). For instances, the expression of
340 LUNAR-1 (*XLOC_LNC_TALL01_Trimarchi*), a Notch1-regulated LncRNA in T-ALL (11),
341 was highly correlated with *NOTCH1* ($r=0.75$; **Supplementary dataset 2-3**). About 14%
342 (patients) to 17% (cell lines) of the correlated gene-lncRNA pairs were located within the
343 same chromosome, while 10% were separated by less than 1 Mb (**Supplementary Figure**
344 **S3**), suggesting potential cis-regulation for a substantial number of oncogenes. Thirty percent
345 of correlated LncRNAs come from the LncRNA_Atak dataset, with some being highly
346 correlated with T-ALL oncogenes (**Figure 7A**). Two examples are shown in Figure 7B-C.
347 Interestingly, some T-ALL oncogenes, such as EML1 and OLIG2 (Figure 7B-C), correlated
348 with several LncRNAs and form complex regulatory networks (**Supplementary Figure S4**).
349 Altogether, these results suggest that some LncRNA from the LncRNA_Atak set might be
350 potential regulators of oncogenic genes in T-ALL and pave the way for more detailed studies.

351

352

Discussion

LncRNAs are transcribed weakly in a large fraction of the genome and display remarkably restricted expression in a space and time-dependent manner [2], making challenging to draw up an exhaustive list of transcripts expressed in all cellular conditions. This is especially true in the case of tumor cells where genomic alterations are expected to alter transcriptional programs, generally leading to large heterogeneous cancer subtypes [6]. In order to establish a comprehensive catalogue of LncRNAs expressed in T-ALL, we analyzed a set of 50 RNA-Seq samples produced by Atak *et al.* [15] (31 primary T-ALL patients, 18 T-ALL cell lines and 1 pool of 5 thymuses) and performed *de novo* transcript discovery in order to systematically identify transcript models. This approach led to the discovery of 2560 novel LncRNAs. Subsequently, we perform a deep characterization of the genomic and epigenetic properties of these transcripts and showed they are comparable to previously identified LncRNAs.

Several approaches have been suggested to identify functionally relevant LncRNAs, including guilty-by-association or correlation-based approaches [31]. Master oncogenes in T-ALL are generally ectopically expressed in a restricted number of patients resulting in highly variable expression among tumor samples. Indeed, genes displaying high variance throughout the T-ALL samples were demonstrated to be significantly enriched in known T-ALL oncogenes. We thus used the expression variance as a proxy to estimate oncogenic potential of the LncRNA expressed in T-ALL. We observed that LncRNAs with known implication in cancer (e.g. LUNAR1) were ranked among those with the highest variance. Interestingly, many newly identified LncRNAs were found to have highly variable expression. Combined variance and correlation analysis also suggest that a fraction of these LncRNAs could have oncogenic properties by functionally interacting with known oncogenes.

One of the key features of LncRNAs is that their expression pattern is highly tissue and cell type specific [2]. This is consistent with our finding that *de novo* LncRNAs discovered in T-ALL demonstrated high tissue-specificity (**Figure 2**) and that many LncRNAs found in T-ALL are expressed in few leukemic samples and (**Figure 6**). Consequently, molecules targeting either their expression or their interactions with chromatin or protein complexes would represent therapeutic targets able to kill cancer cells while sparing normal cells [5]. Additionally, correlation of expression patterns with leukemia progression and outcome could lead to novel prognosis markers and help classification and stratification of the patients.

T-ALL comprises several molecular subgroups characterized by the aberrant expression of distinct oncogenic transcription factors, unique gene expression signatures, and different prognoses [6]. While the existence of specific molecular subtypes of T-ALL has long been established, therapeutic strategies are applied uniformly across subtypes, leading to variable responses between patients coupled with high toxicity. Our comprehensive resource of LncRNAs expressed in T-ALL should allow further exploration of LncRNAs potentially involved in leukemia and provide new rationales for patients/risk stratification.

397

Acknowledgements

398

399

400

401

402

403

404

405

406

407

408

409

Legends to figures

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

Figure 1: Genomic characterization of LncRNA transcripts. (A) Schematic illustration of the procedure used to create our LncRNA catalog. (B) Bar plots displaying the number of exons per transcript set. (C) Distribution of transcript sizes. Each vertical line indicates the mean transcript sizes of the corresponding set. (D) Chromosomal distribution of transcript sets. (E) Schematic illustration of LncRNAs categories. LncRNA exons appear as red and coding genes as blue. (F) Pie chart representing the fraction of LncRNA categories across the transcripts sets.

Figure 2: Expression levels and tissue-specificity of LncRNAs classes. (A) Violin plot showing expression levels of transcripts. For each transcript, expression level was computed by calculating the median of its expression values among 31 T-ALL blasts, 18 T-ALL cell lines and a pool of 5 normal thymus. Wilcoxon test was used to assess differences. (B) Representative human tissues [22] were used to compute gene expression and assess tissue-specificity of each transcript. The density plots shows the distributions of the tissue-specificity score (see material and method section). Each vertical line indicates the mean tissue specificity score of the corresponding class.

Figure 3: Functional annotation of neighboring genes for the LncRNAs. (A) LncRNA_Known. (B) LncRNA_Trimarchi. (C) LncRNA_Atak. The x -axis is corresponding to $-\log_{10}(\text{corrected } p\text{-value})$ and the y -axis shows the corresponding biological processes.

Figure 4: TSS coverage plot of ChIP-Seq signal for H3K4me3 and. Signals are shown for the four transcript in one tissue (total thymus) and three human cell lines (DND41, RPMI-8402 and Jurkat).

Figure 5: Experimental validation of expression for three lncRNAs. (A) Integrated genomics viewer (IGV) screenshots displaying H3K4me3 ChIP-Seq signals as well as RNA-Seq signals for genomic regions corresponding to *XLOC_00017544_Atak*, *XLOC_00009269_Atak* and *XLOC_00012823_Atak* in RPMI-8402, and Jurkat cell lines. (B) PCR validation of *XLOC_00017544_Atak*, *XLOC_00009269_Atak* and *XLOC_00012823_Atak* in RPMI-8402 and Jurkat cell lines. *MALAT1* was used as positive control.

Figure 6: New candidate oncogenes identification by variance. (A) Variance for coding genes in T-ALL cell lines and patients. (B) Box plots showing the distribution of variance of

446 coding genes in T-ALL cell lines and patients. Wilcoxon test was used to assess differences.
447 (C) Variance of non-coding genes in T-ALL cell lines and patients. Arrows highlight
448 leukemic oncogenes. (D) Box plots showing the distribution of variance of non-coding in T-
449 ALL cell lines and patients. Wilcoxon test was used to assess differences. (E) Variability of
450 XLOC_00000871_Atak expression normalized against the GAPDH gene (n = 20) in thymus
451 and cell lines.

452 **Figure 7: Co-expression between oncogenes and LncRNAs from Atak dataset.** (A):
453 Heatmaps showing the correlation between oncogenes (columns) and the most correlated
454 transcript from the LncRNA_Atak dataset (rows). Correlation are shown both for cell lines
455 (left panel) and patients (right panel). (B) Screenshots obtained from IGV displaying two
456 examples of strong co-expression between an oncogene and a LncRNA from the Atak
457 dataset. Tracks corresponding to cell lines are shown in red while patients are shown in
458 green. (C) Scatter plots showing the expression of pairs oncogene-LncRNA_Atak in cell
459 lines (red) and patients (green).

460

461 References

- 462 [1] Schadt EE, Edwards SW, GuhaThakurta D, et al. A comprehensive transcript index of
463 the human genome generated using microarrays and computational approaches.
464 *Genome Biol.* 2004;5:R73.
- 465 [2] Guttman M, Rinn JL. Modular regulatory principles of large non-coding RNAs. *Nature.*
466 2012;482:339–346.
- 467 [3] Harrow J, Frankish A, Gonzalez JM, et al. GENCODE: the reference human genome
468 annotation for The ENCODE Project. *Genome Res.* 2012;22:1760–1774.
- 469 [4] Bonasio R, Shiekhattar R. Regulation of transcription by long noncoding RNAs. *Annu.*
470 *Rev. Genet.* 2014;48:433–455.
- 471 [5] Schmitt AM, Chang HY. Long Noncoding RNAs in Cancer Pathways. *Cancer Cell.*
472 2016;29:452–463.
- 473 [6] Ferrando AA, Neuberg DS, Staunton J, et al. Gene expression signatures define novel
474 oncogenic pathways in T cell acute lymphoblastic leukemia. *Cancer Cell.* 2002;1:75–
475 87.
- 476 [7] Soulier J, Clappier E, Cayuela J-M, et al. HOXA genes are included in genetic and
477 biologic networks defining human acute T-cell leukemia (T-ALL). *Blood.*
478 2005;106:274–286.
- 479 [8] Asnafi V. HiJAKing T-ALL. *Blood.* 2014;124:3038–3040.
- 480 [9] Ntziachristos P, Abdel-Wahab O, Aifantis I. Emerging concepts of epigenetic
481 dysregulation in hematological malignancies. *Nat. Immunol.* 2016;17:1016–1024.

- 482 [10] Wallaert A, Durinck K, Van Loocke W, et al. Long noncoding RNA signatures define
483 oncogenic subtypes in T-cell acute lymphoblastic leukemia. *Leukemia*. 2016;30:1927–
484 1930.
- 485 [11] Trimarchi T, Bilal E, Ntziachristos P, et al. Genome-wide mapping and characterization
486 of Notch-regulated long noncoding RNAs in acute leukemia. *Cell*. 2014;158:593–606.
- 487 [12] Wang Y, Wu P, Lin R, et al. LncRNA NALT interaction with NOTCH1 promoted cell
488 proliferation in pediatric T cell acute lymphoblastic leukemia. *Sci Rep*. 2015;5:13749.
- 489 [13] Ngoc PCT, Tan SH, Tan TK, et al. Identification of novel lncRNAs regulated by the TAL1
490 complex in T-cell acute lymphoblastic leukemia. *Leukemia*. 2018;
- 491 [14] Wallaert A, Durinck K, Taghon T, et al. T-ALL and thymocytes: a message of noncoding
492 RNAs. *J Hematol Oncol*. 2017;10:66.
- 493 [15] Atak ZK, Gianfelici V, Hulselmans G, et al. Comprehensive analysis of transcriptome
494 variation uncovers known and novel driver events in T-cell acute lymphoblastic
495 leukemia. *PLoS Genet*. 2013;9:e1003997.
- 496 [16] Roberts A, Pimentel H, Trapnell C, et al. Identification of novel transcripts in
497 annotated genomes using RNA-Seq. *Bioinformatics*. 2011;27:2325–2329.
- 498 [17] Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic
499 features. *Bioinformatics*. 2010;26:841–842.
- 500 [18] Trapnell C, Roberts A, Goff L, et al. Differential gene and transcript expression analysis
501 of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc*. 2012;7:562–578.
- 502 [19] Wang L, Park HJ, Dasari S, et al. CPAT: Coding-Potential Assessment Tool using an
503 alignment-free logistic regression model. *Nucleic Acids Res*. 2013;41:e74.
- 504 [20] Kryuchkova-Mostacci N, Robinson-Rechavi M. A benchmark of gene expression tissue-
505 specificity metrics. *Brief. Bioinformatics*. 2017;18:205–214.
- 506 [21] Derrien T, Johnson R, Bussotti G, et al. The GENCODE v7 catalog of human long
507 noncoding RNAs: analysis of their gene structure, evolution, and expression. *Genome*
508 *Res*. 2012;22:1775–1789.
- 509 [22] Duff MO, Olson S, Wei X, et al. Genome-wide identification of zero nucleotide
510 recursive splicing in *Drosophila*. *Nature*. 2015;521:376–379.
- 511 [23] Pandey RR, Mondal T, Mohammad F, et al. Kcnq1ot1 antisense noncoding RNA
512 mediates lineage-specific transcriptional silencing through chromatin-level regulation.
513 *Mol. Cell*. 2008;32:232–246.
- 514 [24] Zhao J, Sun BK, Erwin JA, et al. Polycomb proteins targeted by a short repeat RNA to
515 the mouse X chromosome. *Science*. 2008;322:750–756.

- 516 [25] McLean CY, Bristor D, Hiller M, et al. GREAT improves functional interpretation of cis-
517 regulatory regions. *Nat. Biotechnol.* 2010;28:495–501.
- 518 [26] Ørom UA, Shiekhattar R. Long noncoding RNAs usher in a new era in the biology of
519 enhancers. *Cell.* 2013;154:1190–1193.
- 520 [27] Van Vlierberghe P, Ferrando A. The molecular basis of T cell acute lymphoblastic
521 leukemia. *J. Clin. Invest.* 2012;122:3398–3406.
- 522 [28] Xia Y, Brown L, Yang CY, et al. TAL2, a helix-loop-helix gene activated by the
523 (7;9)(q34;q32) translocation in human T-cell leukemia. *Proc. Natl. Acad. Sci. U.S.A.*
524 1991;88:11416–11420.
- 525 [29] Brown L, Cheng JT, Chen Q, et al. Site-specific recombination of the tal-1 gene is a
526 common occurrence in human T cell leukemia. *EMBO J.* 1990;9:3343–3351.
- 527 [30] Armstrong SA, Staunton JE, Silverman LB, et al. MLL translocations specify a distinct
528 gene expression profile that distinguishes a unique leukemia. *Nat. Genet.*
529 2002;30:41–47.
- 530 [31] Liao Q, Liu C, Yuan X, et al. Large-scale prediction of long non-coding RNA functions in
531 a coding-non-coding gene co-expression network. *Nucleic Acids Res.* 2011;39:3864–
532 3878.

533

Figure 1

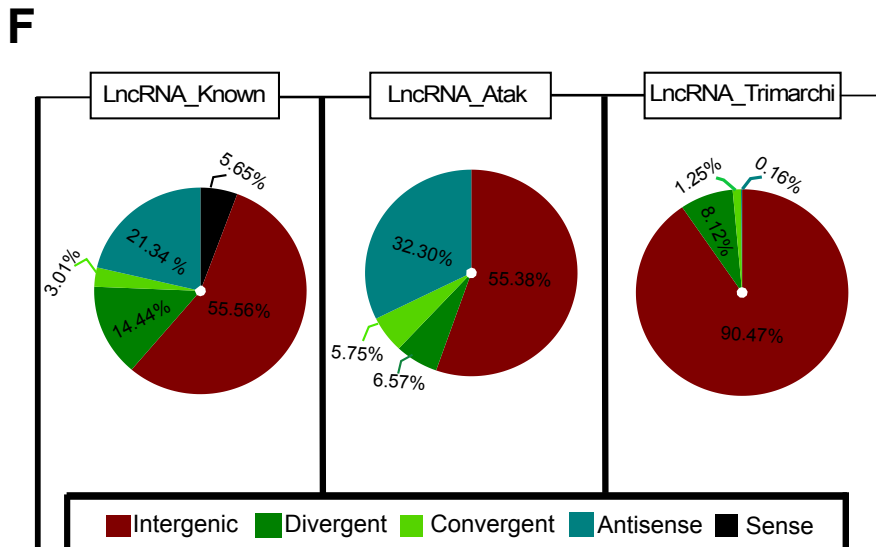
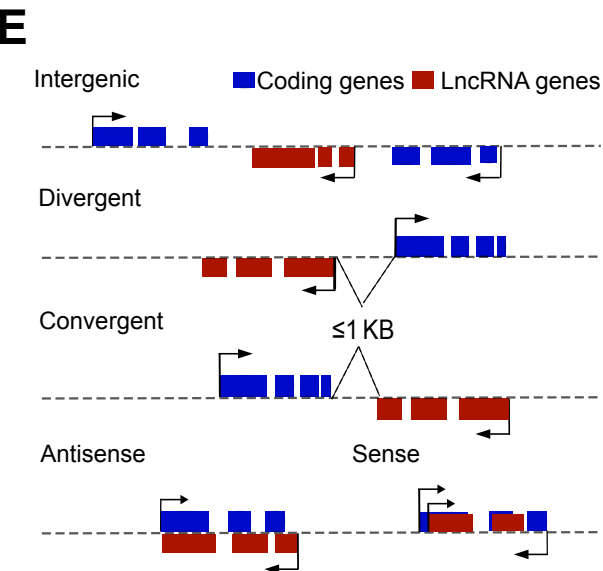
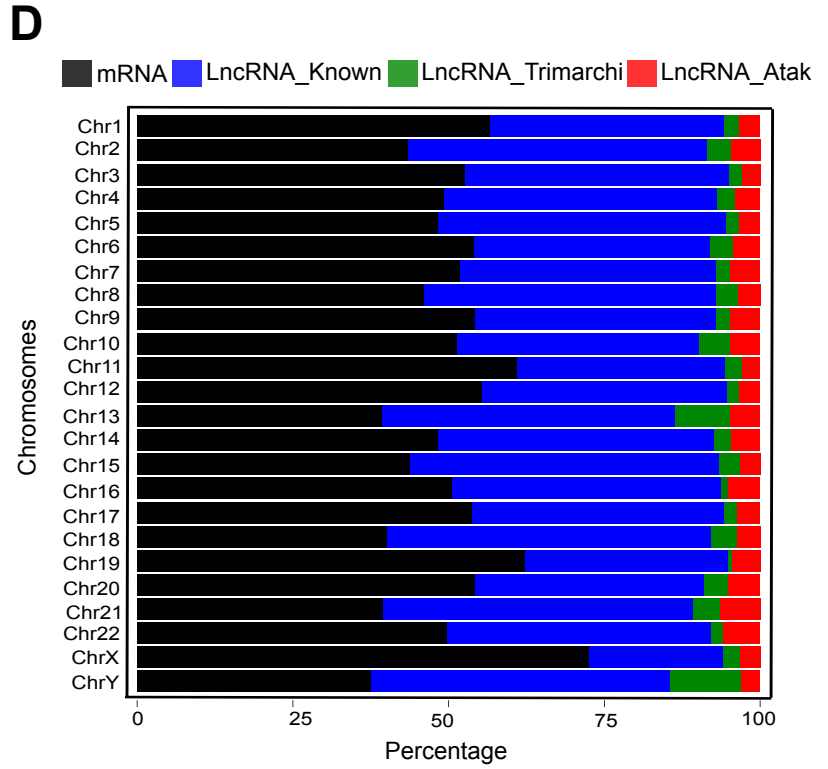
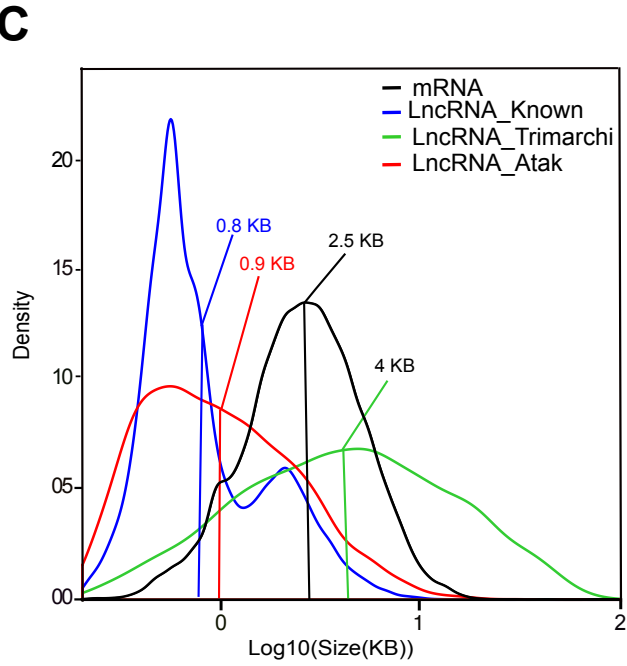
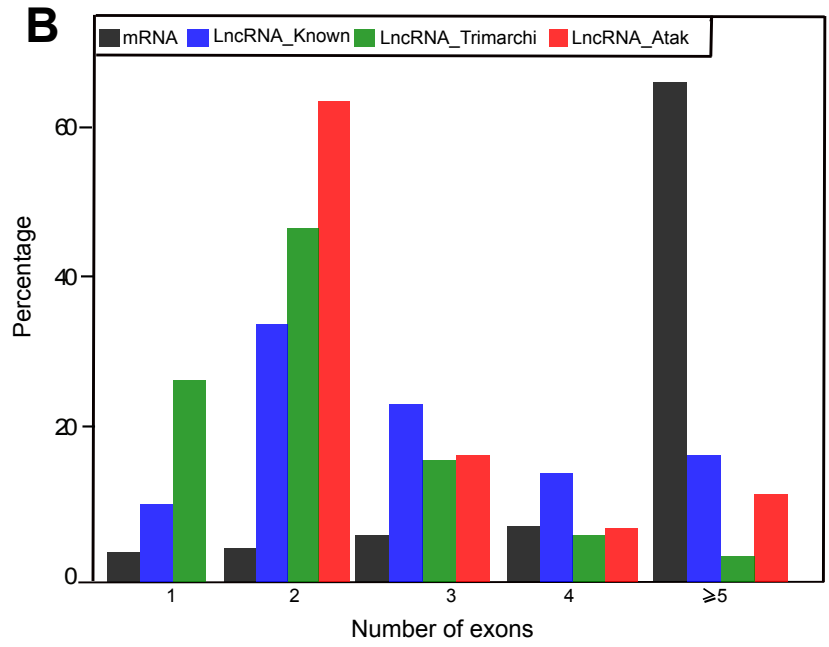
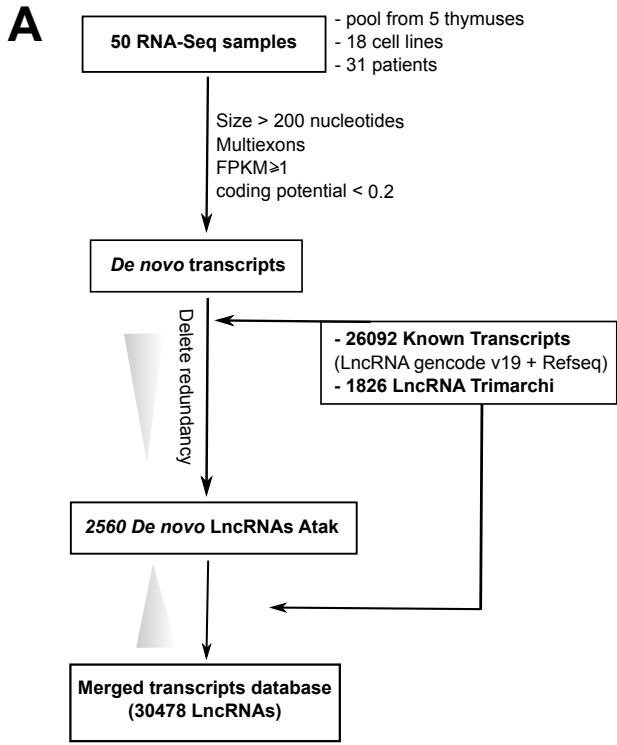


Figure 2

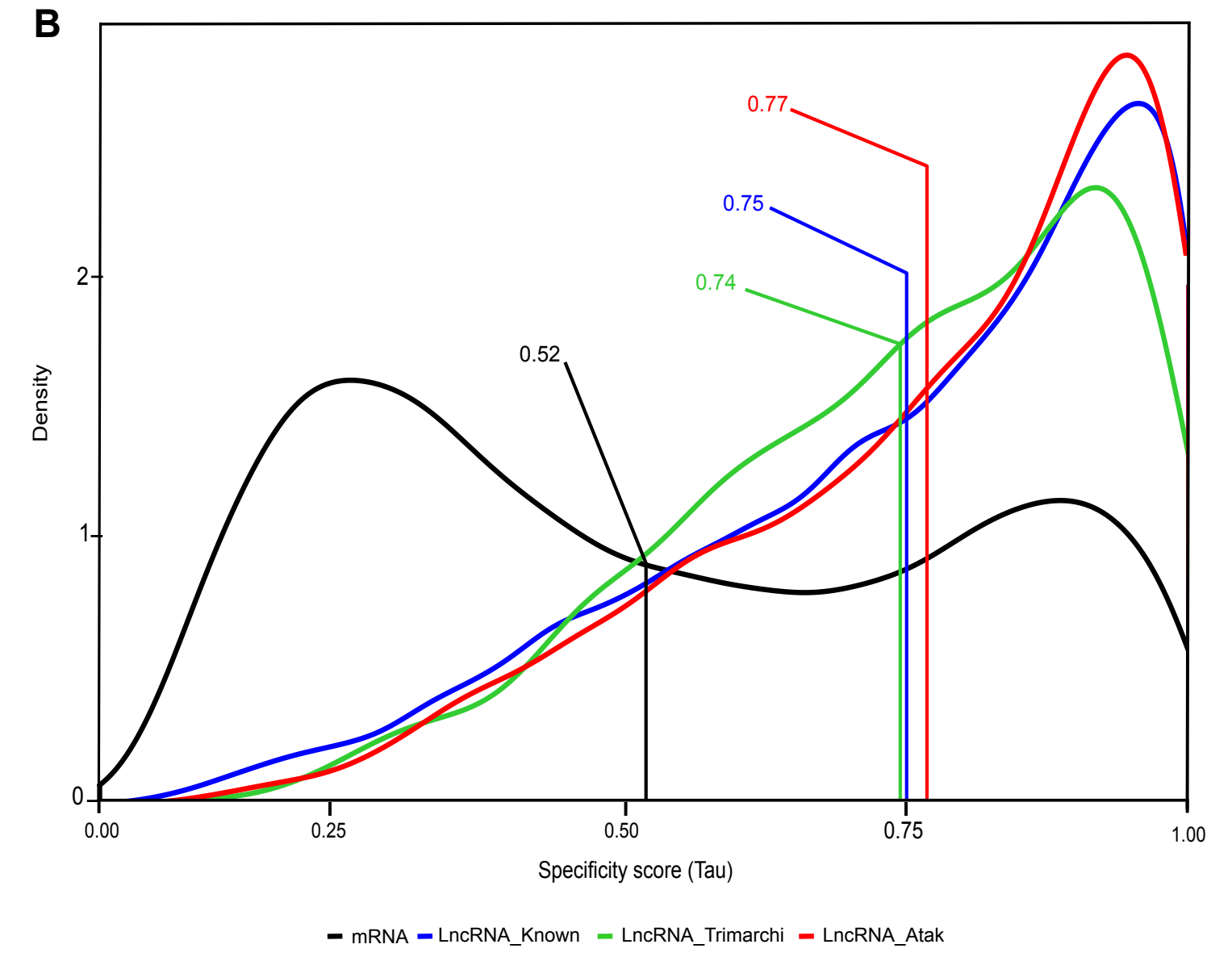
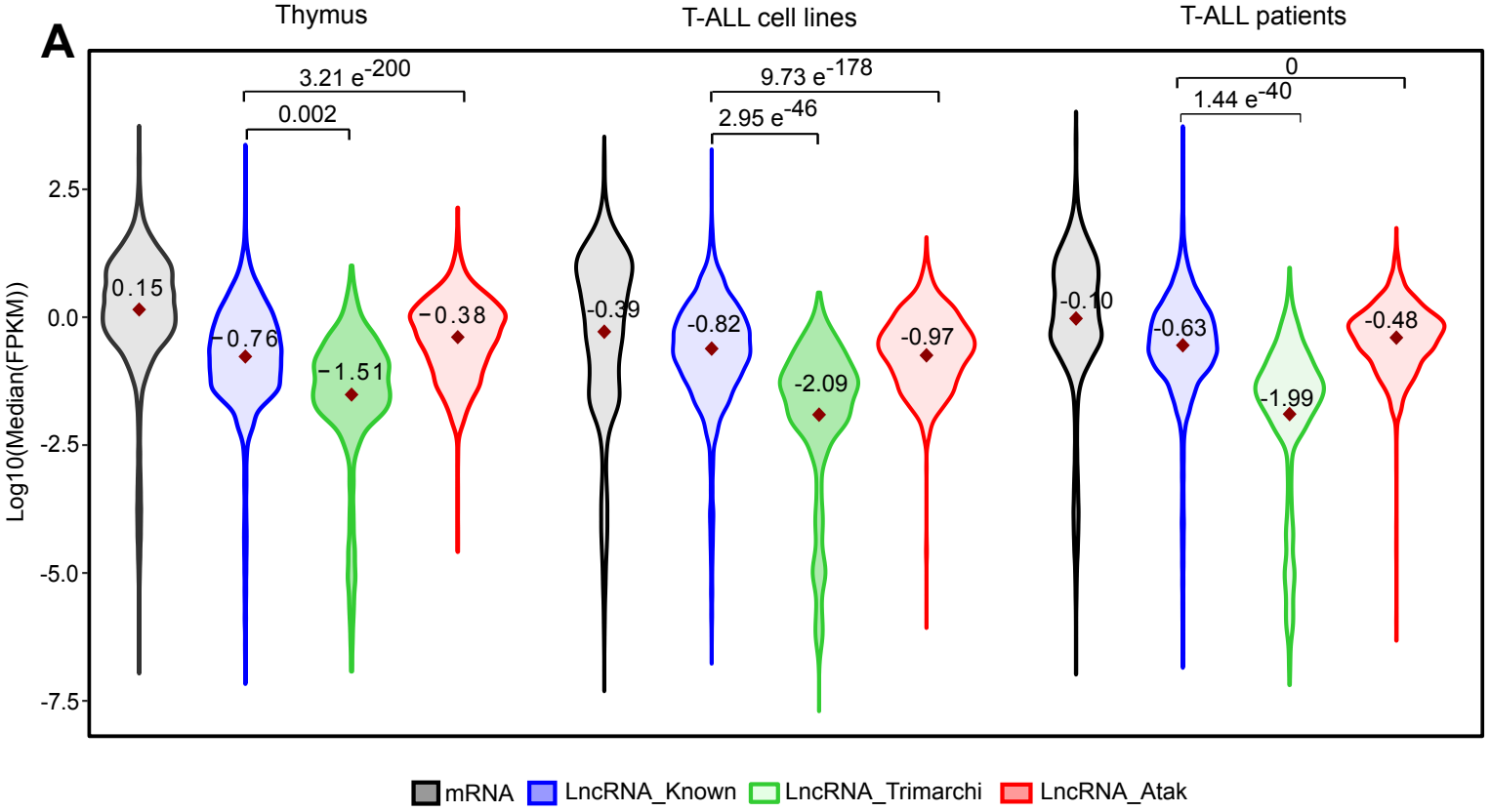
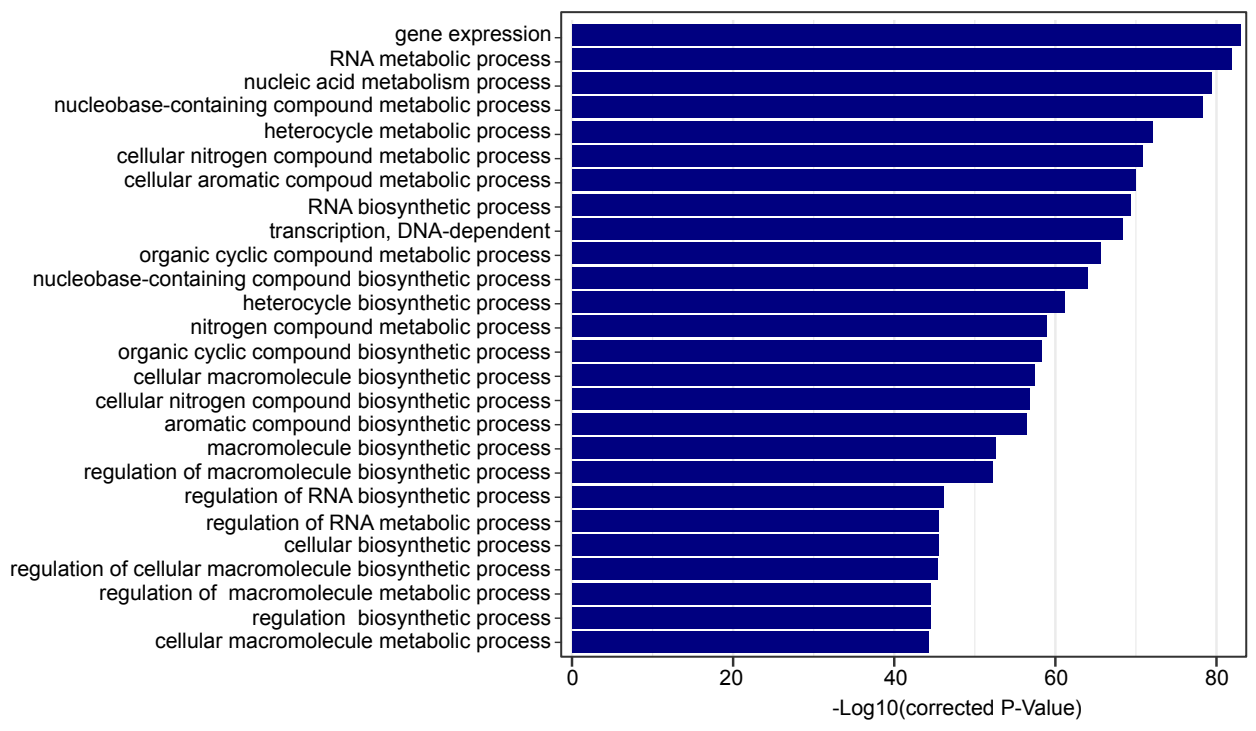
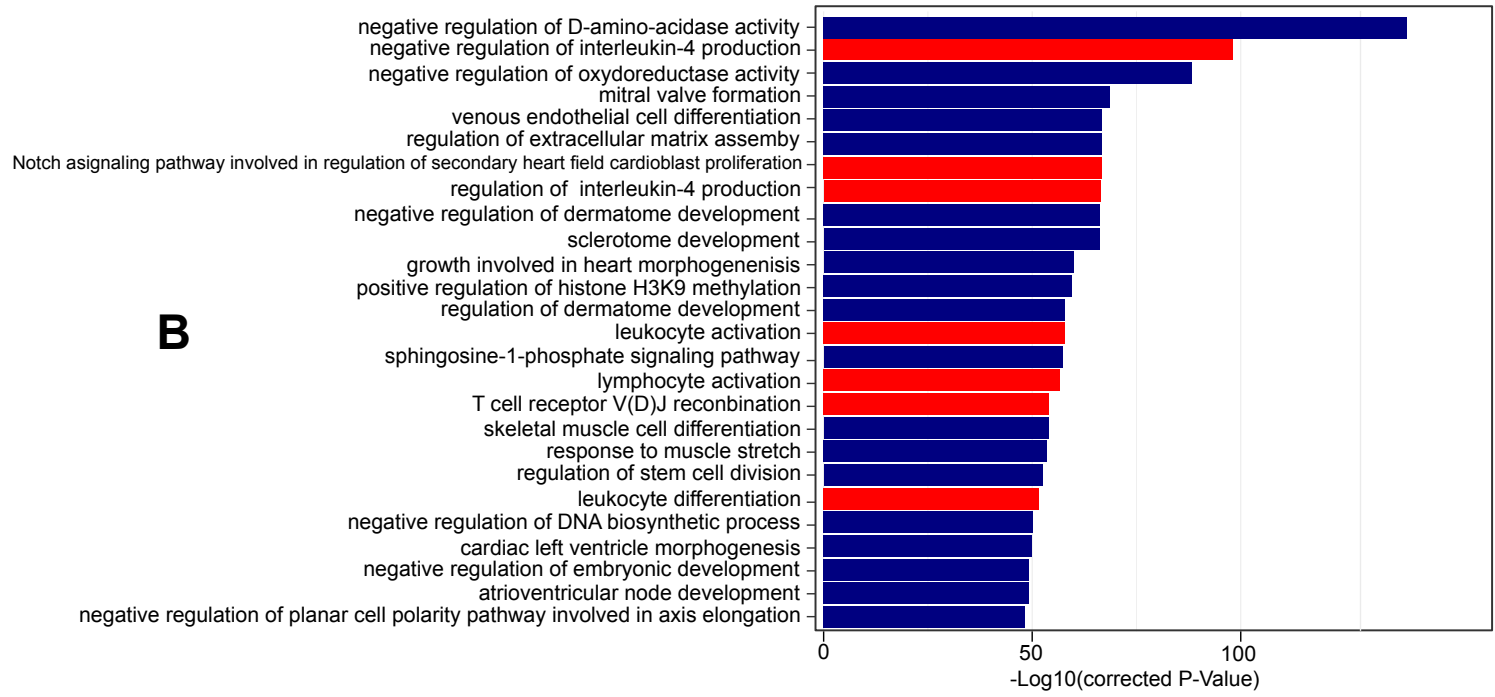


Figure 3

A



B



C

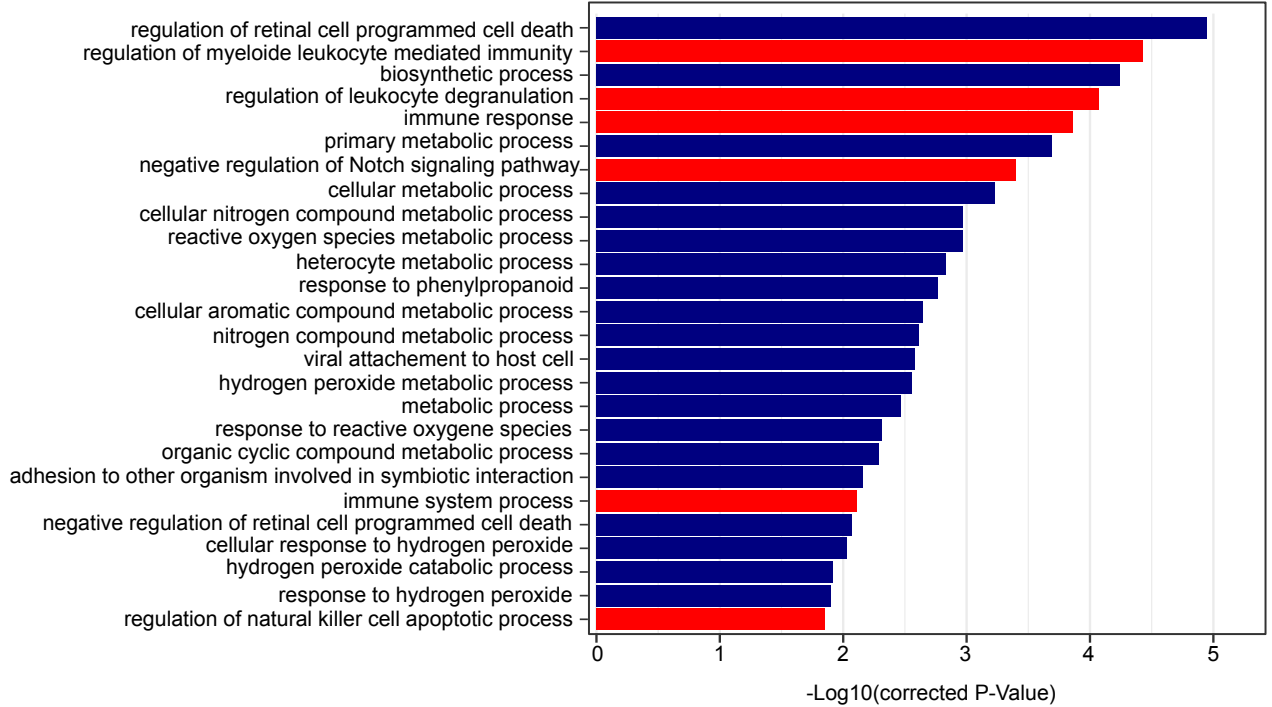


Figure 4

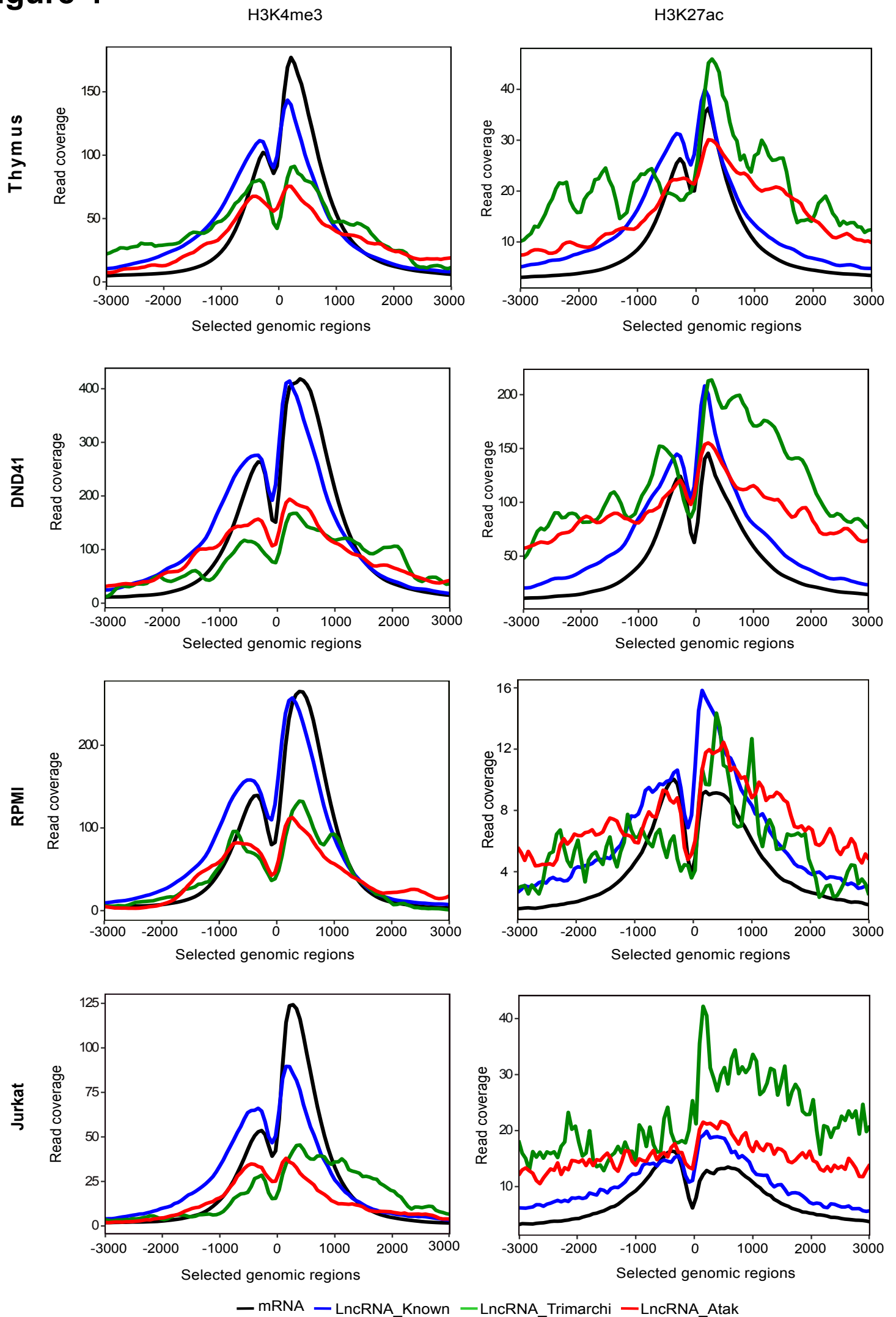


Figure 5

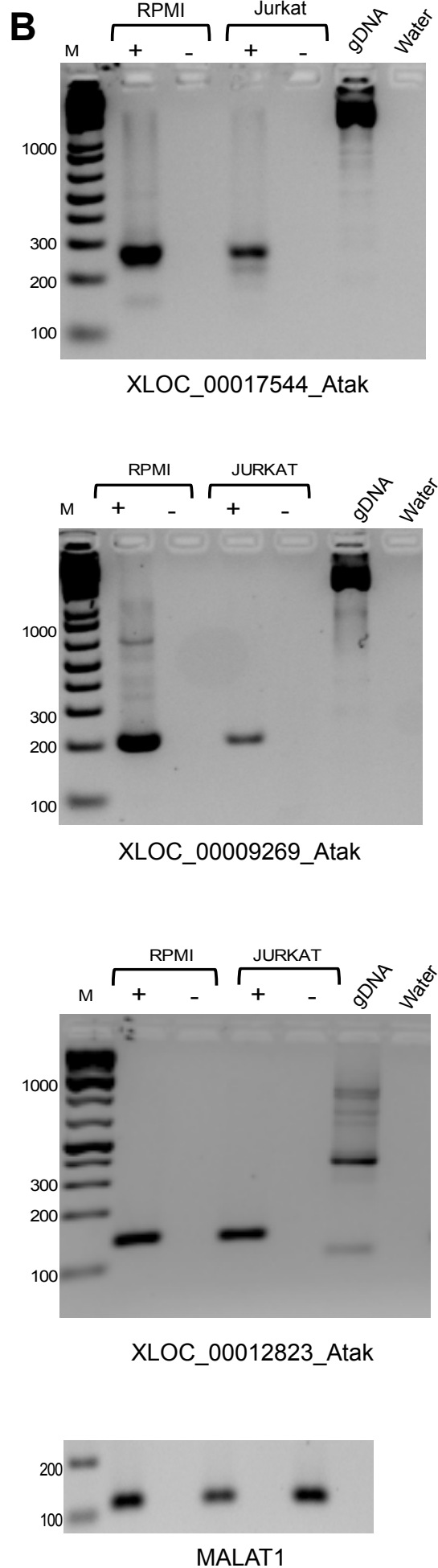
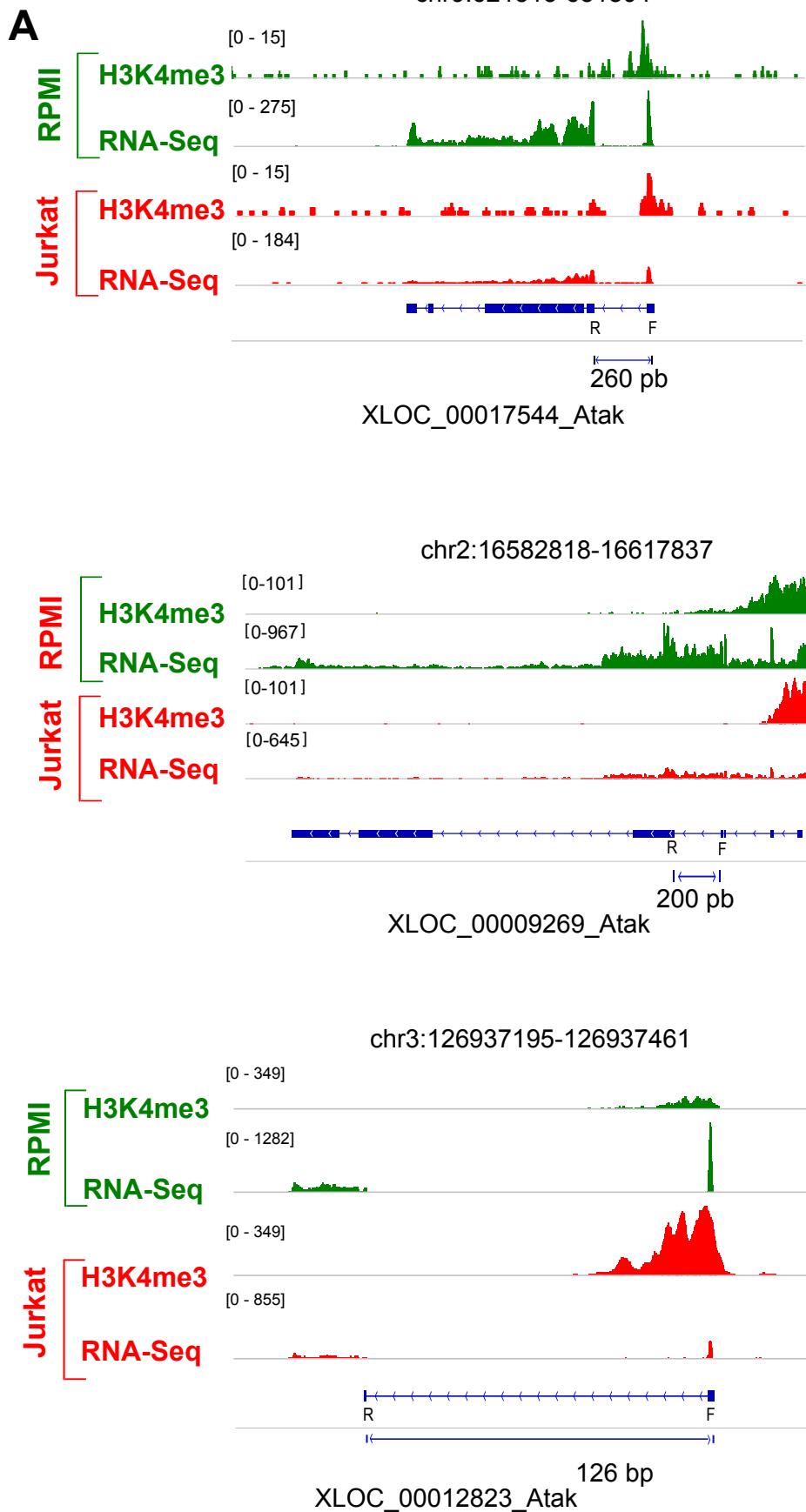
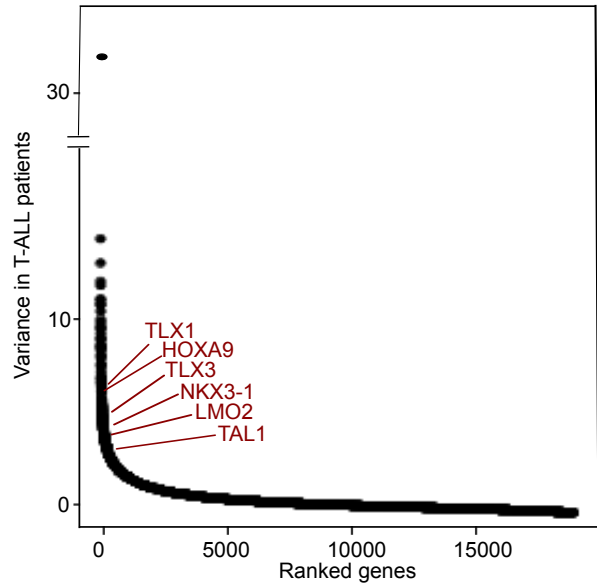
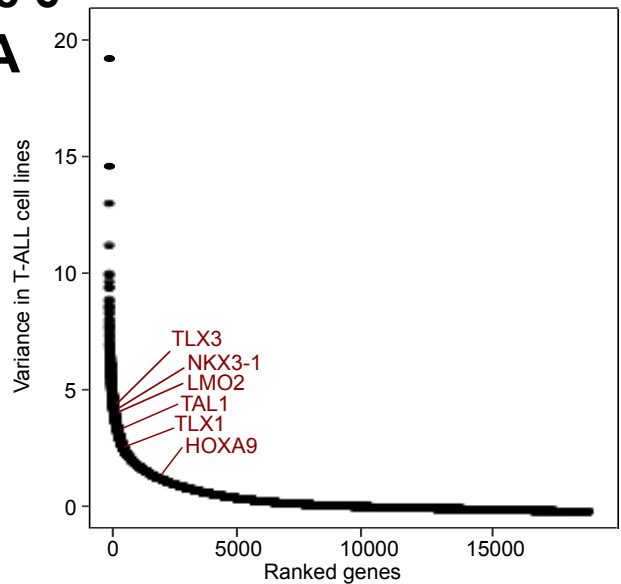
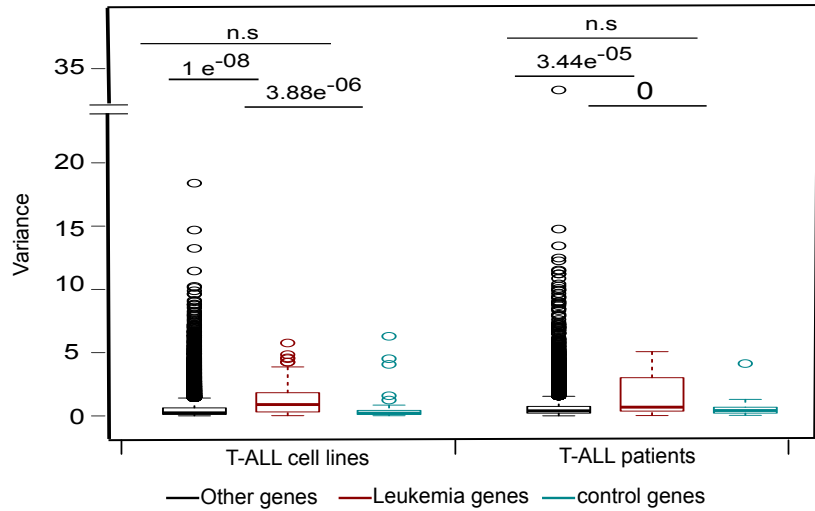


Figure 6

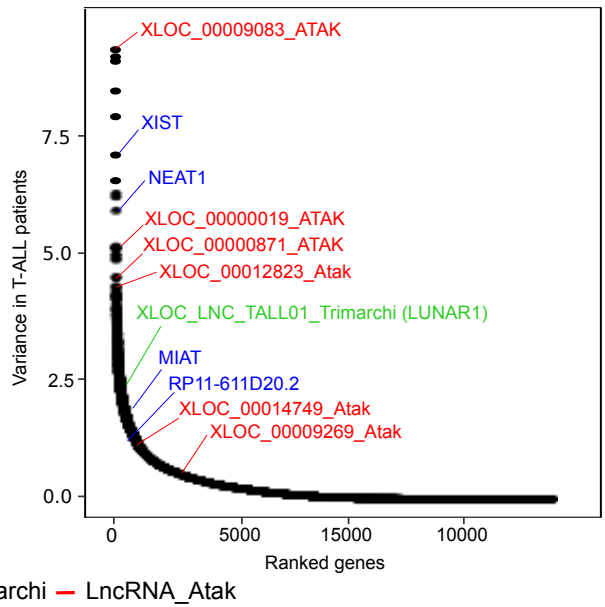
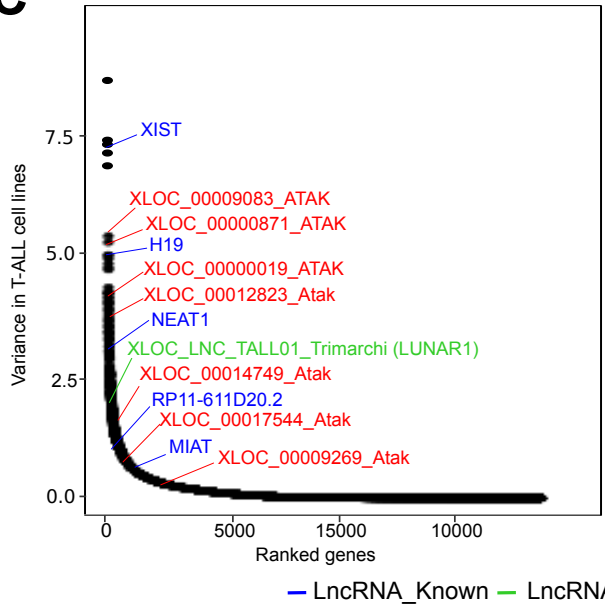
A



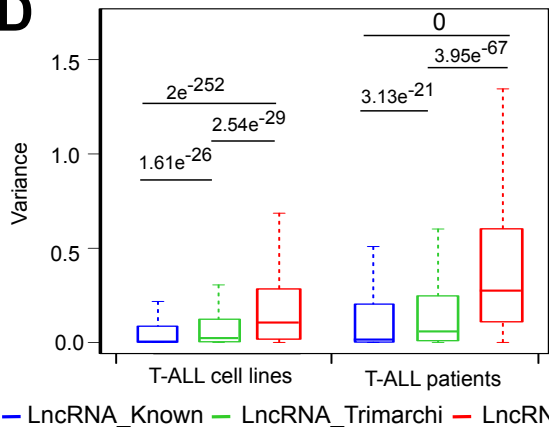
B



C



D



E

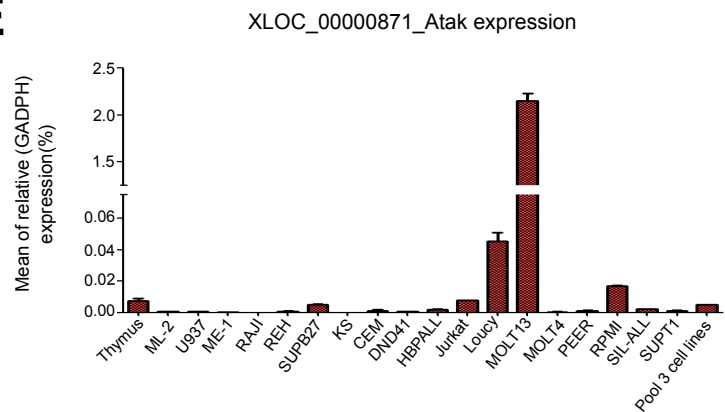


Figure 7

