# Human lncRNAs harbor conserved modules embedded in different sequence contexts

Francesco Ballesio[1], Gerardo Pepe[2], Gabriele Ausiello[2], Andrea Novelletto[2], Manuela Helmer-Citterich[2,*], Pier Federico Gherardini[2,*]


[1] PhD Program in Cellular and Molecular Biology, Department of Biology, University of Rome "Tor Vergata", Rome, Italy.

[2] Department of Biology, University of Rome "Tor Vergata", Rome, Italy.


* To whom correspondence should be addressed. Tel: 00390672594324; Fax:0039062023500; Email: citterich@uniroma2.it; Email: pier.federico.gherardini@uniroma2.it

# Abstract

We analyzed the structure of human long non-coding RNA (lncRNAs) genes to investigate whether the non-coding transcriptome is organized in modular domains, as is the case for protein-coding genes. To this aim, we compared all known human lncRNA exons and identified 340 pairs of exons with high sequence and/or secondary structure similarity but embedded in a dissimilar sequence context. We grouped these pairs in 106 clusters based on their reciprocal similarities. These shared modules are highly conserved between humans and the four great ape species, display evidence of purifying selection and likely arose as a result of recent segmental duplications. Our analysis contributes to the understanding of the mechanisms driving the evolution of the non-coding genome and suggests additional strategies towards deciphering the functional complexity of this class of molecules.

## Author summary

The Human genome includes more than 18,000 genes coding for RNAs that are not translated into proteins, called long non-coding RNAs (lncRNA). Mounting evolutionary and experimental evidence shows that a large amount of these RNAs have a specific function, mainly as regulators of a diverse set of biological processes. Here we set out to investigate whether these genes have a modular organization similar to that of protein-coding genes. Accordingly, we compared the sequence of all the exonic regions of human lncRNAs and identified 106 clusters of non-repetitive exonic modules shared between this class of genes. These modules display evidence of purifying selection, are highly conserved between humans and the four great ape species, and may represent distinct functional units that have been shuffled among multiple lncRNA genes, in a manner similar to the exon-shuffling process that is observed in the coding genome.

# Introduction

41

Many eukaryotic proteins are composed of a discrete number of domains, endowed with autonomous folding capacity and/or characteristic functions. This type of organization is defined as modular, and the process by which this set of modules is recombined into a variety of different protein products is known as "exon-shuffling" [1].

Long noncoding RNAs (lncRNAs) represent a heterogeneous class of RNAs that are not translated into functional protein products but, similar to messenger RNAs, are transcribed from genes that may have an exon/intron structure. These RNAs are generally defined as non-coding RNAs of more than 200 nucleotides in length and can be capped, polyadenylated and spliced [2], much in the same way as the transcripts of protein coding genes. The human genome contains about 18,000 lncRNA genes and 47,000 transcripts [3], most of which are of unknown function. lncRNAs exhibit evidence of purifying selection and experimental evidence shows that at least a portion of them is indeed functional (287 eukaryotic lncRNAs associated with a biological function are collected by [4], 1,273 human lncRNAs by [5]). Some lncRNAs have been characterized in depth and they may function as regulatory molecules both in the nucleus and the cytoplasm, through a variety of mechanisms, including interaction with transcription factors, recruitment of chromatin modifying complexes, modulation of the expression of their neighboring genes, control of mRNA stability and translation and competition for the binding of specific miRNAs [6–8]. Individual lncRNAs have been found to have a role in promotion of metastasis [9], neuronal differentiation [10], regulation of the accumulation of beta amyloid peptide in Alzheimer's disease [11], and many other processes in a diverse array of pathological and physiological contexts. However the identification of the function of lncRNAs on a global scale remains elusive [12], also because their definition likely encompasses an extremely heterogeneous set of genes, whose main, and possibly only, common characteristic is the fact that they do not produce a functional protein product [13].

In general, lncRNAs are significantly less conserved than protein-coding sequences [14], which also suggests that the relationship between sequence and function is particularly complex in this class of molecules. Examples of lncRNA such as *Xist, Megamind, Cyrano* and *Miat* have been described, which have conserved functions throughout multiple organisms, and yet display a level of sequence divergence that challenges sequence homology search tools [13,15]. A corollary of this observation is

70   that similarity amongst lncRNA within a given organism is also limited, and, unlike coding sequences,

71   most lncRNAs appear in single copies in vertebrate genomes [13].

72   However, lncRNAs are significantly more likely to contain repetitive sequences, particularly

73   transposable elements (TEs) [15,16]. On one hand, this could simply indicate that lncRNAs are more

74   prone to transposon insertion, because of their aforementioned looser association between sequence

75   and function [13]. On the other hand, this observation implies the existence of stretches of homologous

76   sequences that are shared among different lncRNAs, even when the lncRNAs themselves are not

77   related by descent.

78   Because TEs are often enriched in sequences with regulatory function, and may contribute to

79   their "spread" within a genome [17], Johnson and Guigò [18] hypothesized that the presence of TEs

80   may result in the sharing of functional cassettes among evolutionarily unrelated lncRNA, possibly

81   implying a modularization of function for this class of molecules [6,12], reminiscent of the notion of

82   domains in the protein-coding world. In support of this hypothesis, it has been reported that TE-derived

83   sequences within lncRNAs are more conserved compared with non-TE sequences [19].

84   Here we set out to expand the identification of modules in lncRNAs that could have contributed

85   to increasing the diversity of the non-coding genome, similar to the exon-shuffling phenomenon that is

86   well known for protein sequences. Our work extends previous observations in three ways, namely by i)

87   focusing on the sharing of individual exons among unrelated lncRNAs within the human genome, ii)

88   specifically excluding exons that contain repetitive sequences, and iii) including secondary structure as

89   an additional criterion to define similarity, as lncRNAs with similar functions often lack linear sequence

90   homology [20], and many examples of ncRNA are known whose function is tied to their secondary

91   structure [21–24].

92

# Results

94

## Exon sequence and secondary structure comparison

96

97   In order to search for similarities among lncRNAs, we performed a pairwise comparison of both

98   the sequence and the predicted secondary structure of 12,097 non-overlapping human lncRNA exons

99    that do not contain repetitive sequences, performing a total of more than 73 million sequence alignments

100   and an equal number of structure alignments. The distributions of the corresponding scores are shown

101   in Figure 1A, B (1A-B Fig.).

102   To identify pairs or groups of exons representing shared sequence elements, hereafter referred

103   to as "modules", it was necessary to select a threshold above which their sequence or structure similarity

104   would be considered significant.

105   We thus investigated the conservation of lncRNA exons in four non-human primates (see

106   Materials and Methods), with the goal of identifying shared sequence elements in the human genome

107   that are also conserved in other primate genomes.
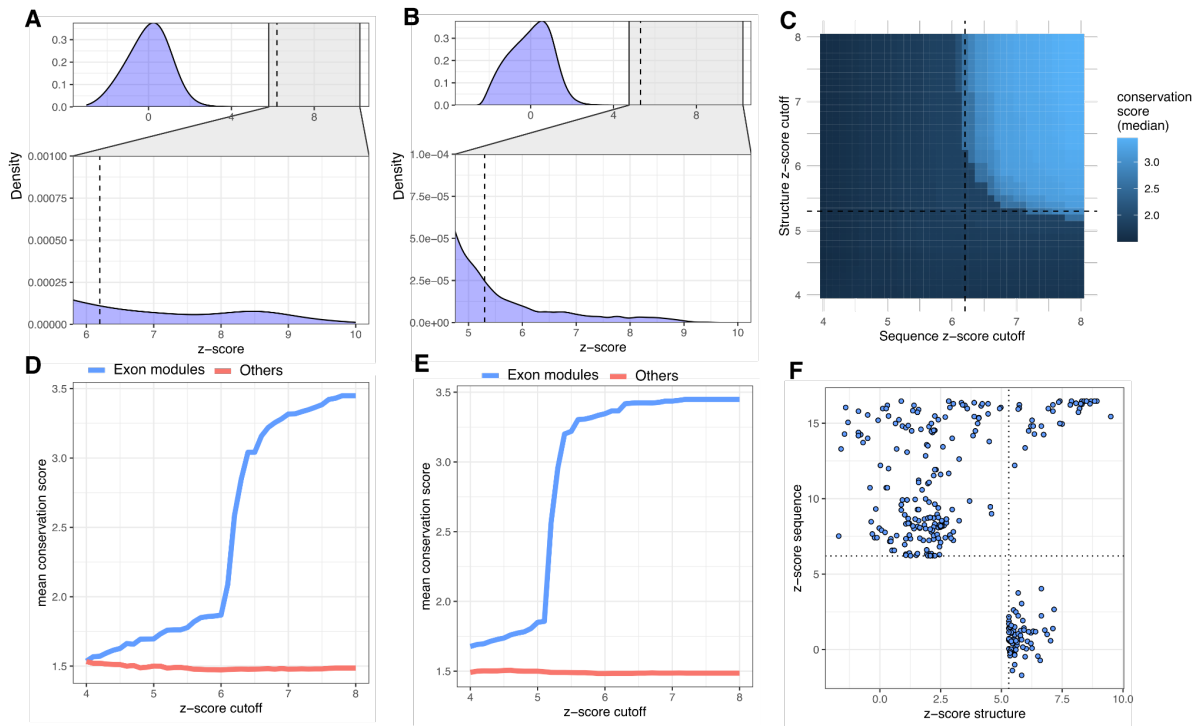
108   Accordingly, we calculated the mean conservation scores of sequence modules across these

109   species, as a function of the similarity score threshold used to define the modules themselves. Using

110   this procedure, we observed a sharp transition in conservation at Z-score similarity thresholds of 6.2

111   and 5.3 for sequence and structure alignments, respectively (1C-D-E Fig.). We consider this increase in

112   conservation, coupled with the high Z-score similarity threshold, as a strong indication that the shared

113   sequence elements we identified represent significant similarities. As a further benchmark, we repeated

114   the entire procedure by aligning exons against random sequences with the same length and base

115   composition. None of the alignments produced z-scores above the 6.2 threshold.

116   By using these thresholds, we identified a total of 340 exon pairs (219 identified by sequence,

117   75 by structure and 46 by both), involving 338 different exons and 218 different genes (1F Fig.). Starting

118   from these pairwise similarities, we identified 106 clusters (exon modules) defined by homologous

119   lncRNA exons represented in at least two copies in the same or different genes (2 Fig., S1 Fig. and S1

120   Table).

121   To rule out the possibility that similarity between exons in a pair of genes is simply due to

122   paralogy, we aligned the entire genes using BLAST and excluded pairs with alignment coverage on the

123   smallest gene of the pair greater than 80%. Measuring the alignment coverage of the entire genes,

124   including introns, allowed us to identify and exclude cases of complete paralogy even in the presence

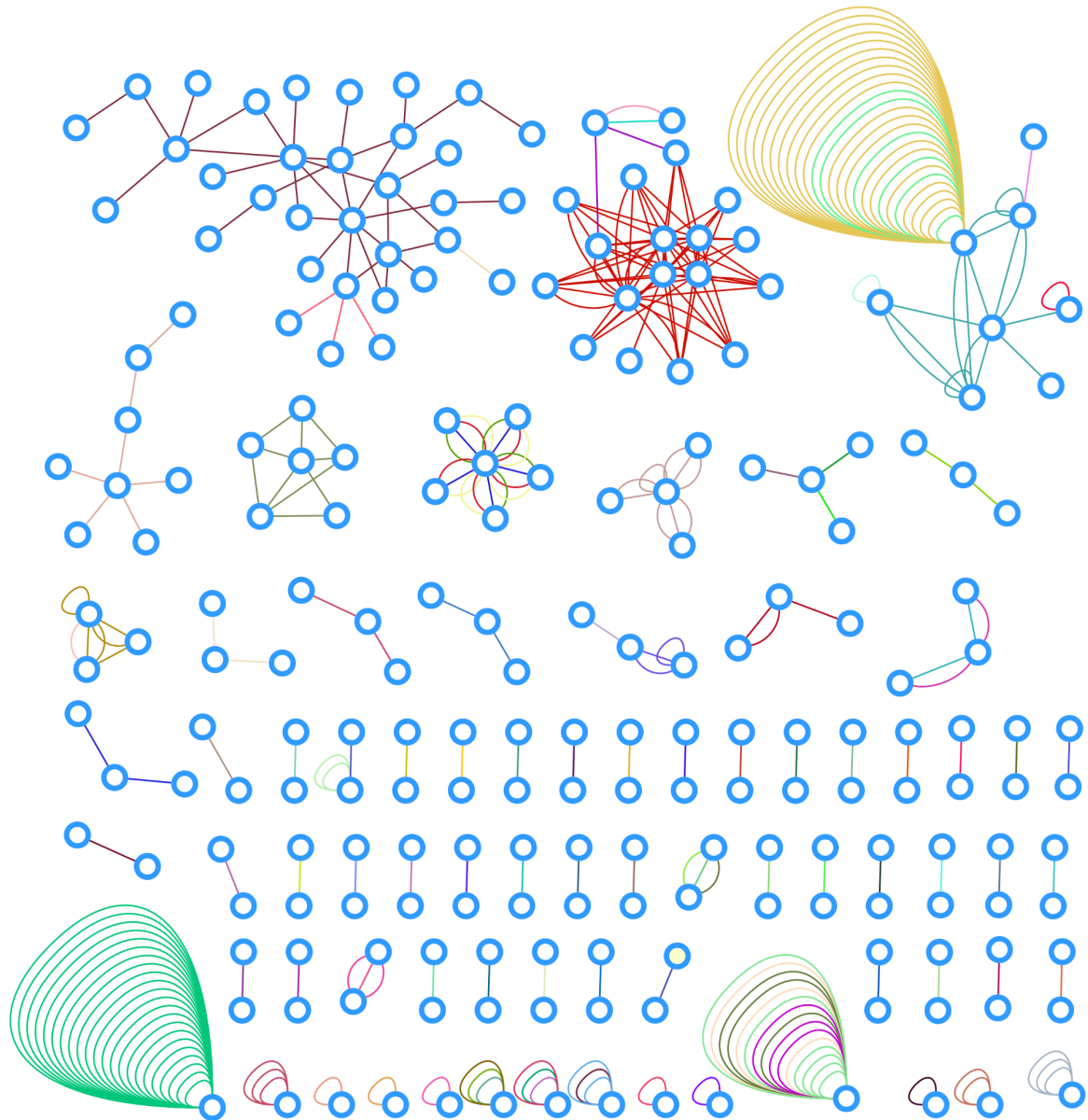125   of intronization or imprecise exon annotation.

126   We note that, in general, our analysis is dependent on the reliability of the reconstruction of the

127   whole transcript structure, which is used to define the exons themselves. This is summarized by the

128   Transcript Support Level (TSL, S1 Table).

129



130

**Fig.1 Sequence and structure alignments results.** A) Distribution of z-transformed pairwise alignment scores for sequence; B) Distribution of z-transformed pairwise alignment scores for structures, for these distributions, a close-up around the proposed cutoff thresholds is also shown; C) heatmap representing the conservation scores in the four non-human primates of all pairs selected at the different z-score thresholds of sequence and structure alignments; D, E) Mean conservation scores (within four non-human primates) of members of clusters defined by different z-score thresholds of pairwise similarity for sequence (D) and structure (E). Note the steep increase in evolutionary conservation for the z-score cutoff of 6.2 (sequence) and 5.3 (structure), respectively; F) Scatter plot of sequence and structure similarity z-scores of the exon pairs (for the sake of clarity, the more than 73 million pairs below the thresholds are not shown).

141

**Fig.2 Network representation of the exon-sharing gene clusters and the corresponding exon modules.** Each node represents a lncRNA gene and each edge an exonic module shared between two genes. Same color edges within a gene cluster represent a module. Self-loops represent instances where the same module occurs multiple times in a single gene. The network representation was generated using Cytoscape [25].
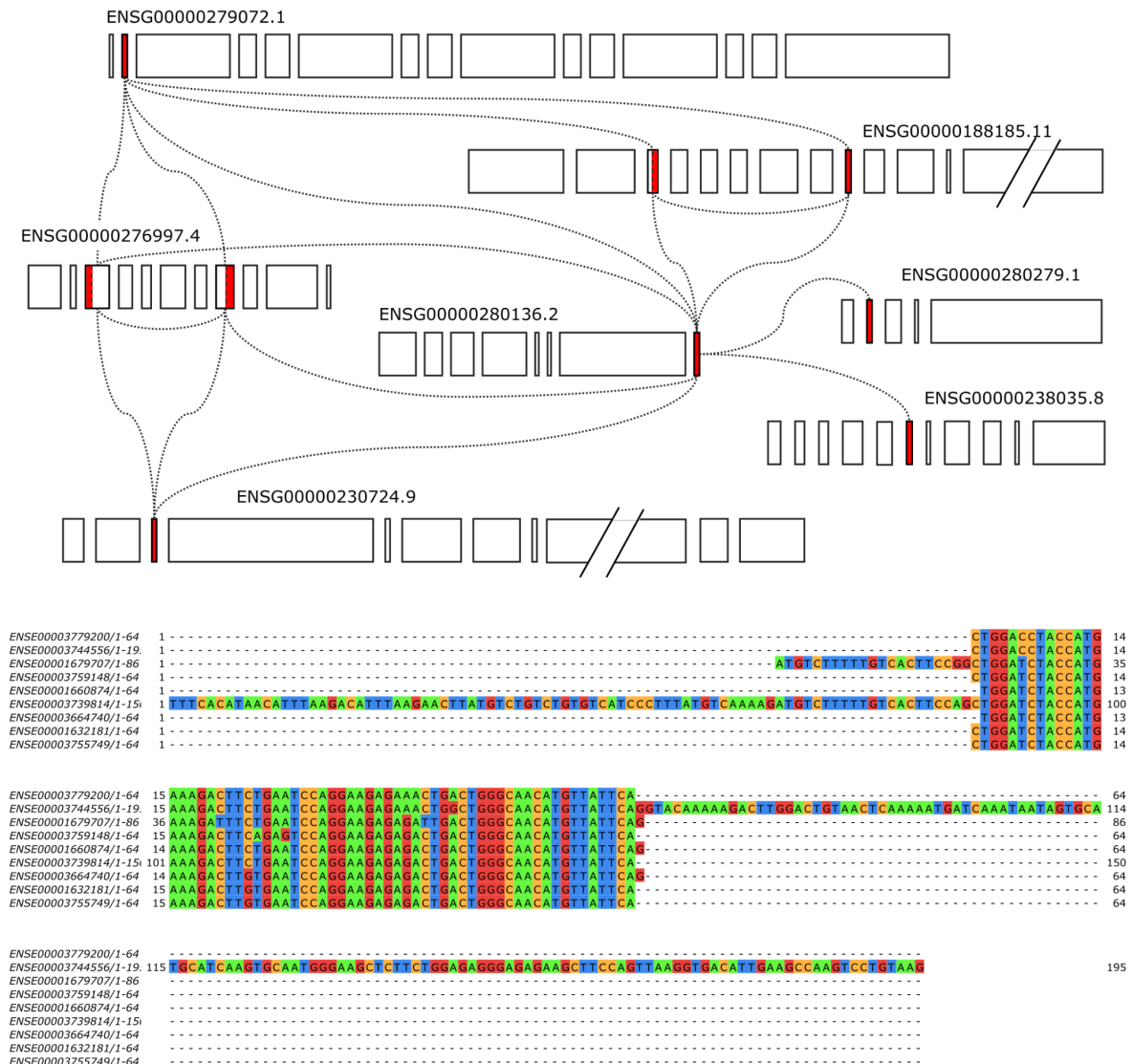
147

Figure 3A (3A Fig.) is an example of one of the identified exon modules shared by a group of 7 lncRNA genes: ENSG00000279072.1, ENSG00000188185.11, ENSG00000276997.4, ENSG00000280136.2, ENSG00000280279.1, ENSG00000230724.9, ENSG00000238035.8. This cluster consists of 9 exons that contain a region of ~65 nucleotides with high sequence similarity

152    (external gap trimmed sequence identity 92-98%, 3B Fig.) embedded in different genes. It is worth noting

153    that, in some cases, the module constitutes an exon on its own, whereas in other cases it is part of a

154    larger exon.

155



156

157    **Fig.3 An example of the identified exon modules.** A) Schematic representation of 7 genes

158    containing representatives (in red) of exons contributing to a module cluster. Each box represents an

159    exon, with width proportional to its length (intron length not to scale); B) multiple alignment of the 9 exons

160    contributing to the cluster.

161

162    We then analyzed in more detail the sequence context of exon modules. More specifically, we

163    looked at the sequence similarity of additional exons flanking the modules, to rule out the possibility that
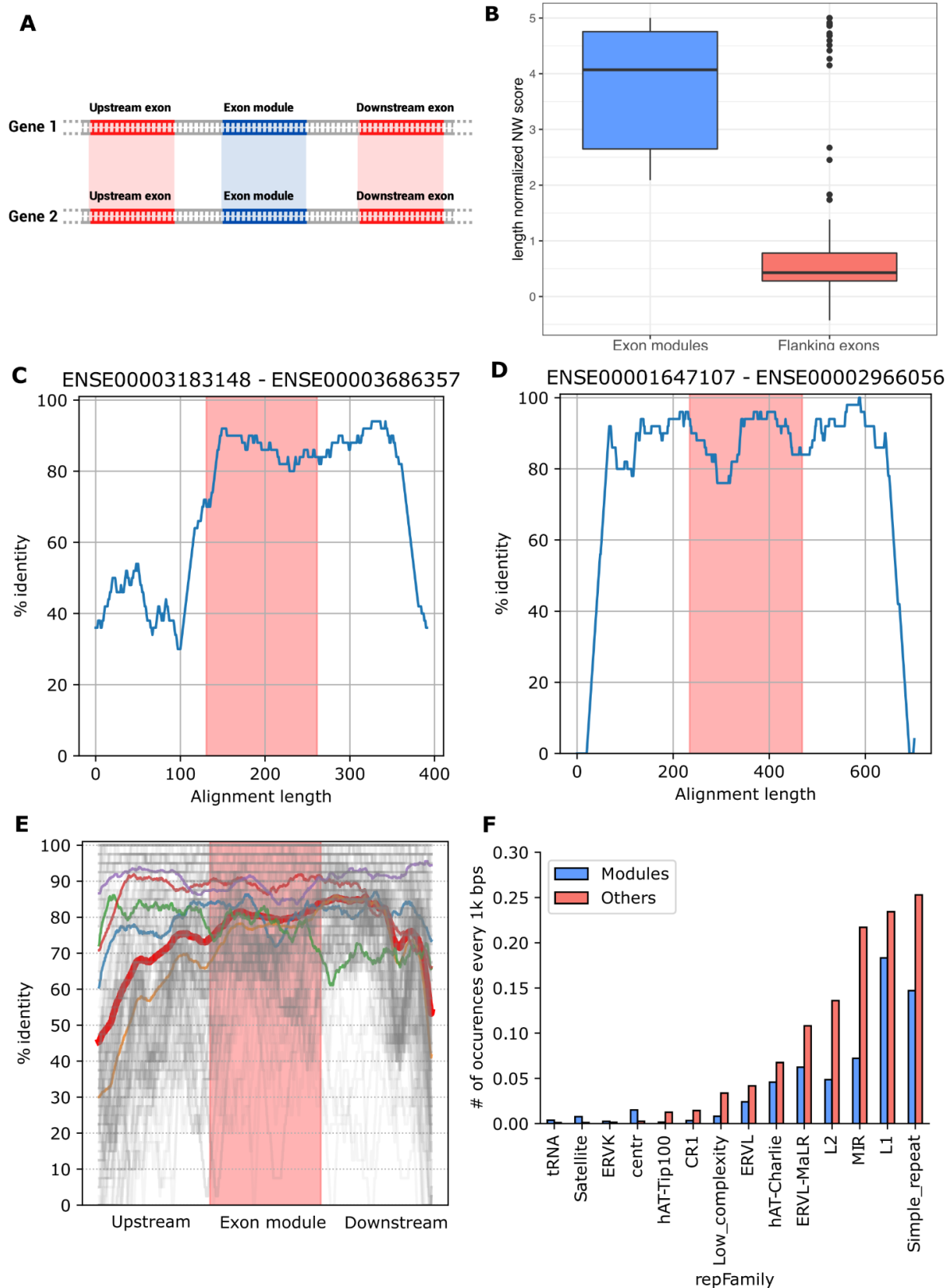
164 the similarity between modules in different genes simply reflects global sequence similarity between the

165 exonic components of genes (see Materials and Methods and 4A Fig.). The alignment scores for exons

166 flanking the putative module in the same gene, upstream and downstream (4B Fig.), showed that the

167 similarity between exon modules is significantly higher than that of the sequence context in which they

168 are embedded. We also observed a small proportion of cases in which the flanking exons are also similar

169 (outliers in 4B Fig.). These cases fall outside the criteria used to define exon modules: in 17 cases

170 because they are less than 50 nucleotides in length, and in another 17 cases because they contain

171 repetitive sequences.

172 We then analyzed the sequence similarity of the intronic sequences flanking the exon modules.

173 To this end, we defined genomic regions of interest by extending upstream and downstream the

174 sequence of each candidate exon pair, until we obtained two sequences with a length equal to three

175 times that of the longest exon of the pair (4C-E Fig.). We limited the analysis to pairs with sequence

176 similarity above the z-score threshold of 6.2 and excluded modules repeating within the same gene. For

177 each pair of genomic regions of interest, we performed a global alignment using the same parameters

178 used to identify the exon modules, and calculated the percentage identity of the pairs using overlapping

179 windows of 50 nucleotides with a single nucleotide shift, to generate graphs depicting the extent of the

180 similarity. We found that, in the majority of instances, sequence similarity extends into the flanking

181 intronic regions. More specifically, in approximately one third of the cases, the similarity encompassed

182 both the upstream and downstream intron, in another third of the cases the similarity extended to a

183 single intron, while the remainder of cases lacked a clear pattern. We did not observe any cases where

184 the similarity was confined to the boundaries of the candidate exon modules.

185 The extension of the similarity through the flanking introns suggests that the most common

186 mechanism responsible for the origin of exon modules is segmental duplication of a genomic DNA

187 stretch encompassing the parental copy of an exon. This is the same mechanism suggested as a driver

188 of exon shuffling in protein coding genes [26]. To further confirm these findings, we compared our results

189 with the data present in the UCSC Segmental Dups track (genomicSuperDups) which contains regions

190 detected as putative genomic duplications within the human genome. These regions represent large

191 recent duplications (>= 1 kb and >= 90% identity) that originated over the last ~40 million years of

192 human evolution, based on neutral expectation of divergence [26]. For 84 of the 340 lncRNA exon pairs

193    identified here, we found a match in the segmental duplications identified by Bailey et al. In 81 of these

194    cases the duplicated stretch includes the entire exons of the pair, while in 3 cases the duplication is

195    interrupted within the exon. We also observed a higher frequency of pairs located on the same

196    chromosome (~20.5%) compared with what is observed when the same exons are randomly paired

197    (~3.6%). Moreover, pairs of exon modules that are on the same chromosome are closer together when

198    compared to the same random pairing control (Mann-Whitney p-value=9.86e-05). A higher rate of

199    occurrence on the same chromosome has been described for segmental duplications [27]. To further

200    extend the analysis of flanking regions, we compared the rate of occurrence of multiple families of

201    repetitive elements in the introns flanking candidate exonic modules vs other lncRNA exons (for exons

202    located at the ends of a gene, we included a region of 10k bps in the genome). We calculated the

203    number of occurrences per 1,000 base pairs of each family of repetitive elements on the set of regions

204    flanking the exon modules vs the other lncRNA exons (4F Fig., S2 Table) thus obtaining a distribution

205    of occurrences where the observations correspond to the individual sequence regions. We then

206    compared these distributions using a Mann-Whitney U test, with Bonferroni correction for multiple

207    hypothesis testing. We observed significant differences for 15 of 46 families (padj<0.05). Interestingly,

208    centromere and satellite repeats are among the few classes of repeats enriched in regions flanking the

209    exon modules, while most classes of transposon- or endogenous retrovirus-derived repeats are

210    depleted. Since the genomic regions proximal to centromeres and telomeres are enriched with

211    segmental duplications [28],  this observation further points at segmental duplication as the main driver

212    of  the appearance of these exon modules, as opposed to, for instance, transposition. The enrichment

213    of this type of repetitive sequences can be explained by the localization near the centromeres or

214    telomeres of a portion of the modules (S2 Fig.). Moreover, searching for transposase domains using a

215    procedure similar to the one described in [29] did not reveal significant differences in their occurrence

216    among genes containing exon modules (data not shown), further highlighting that transposition is not
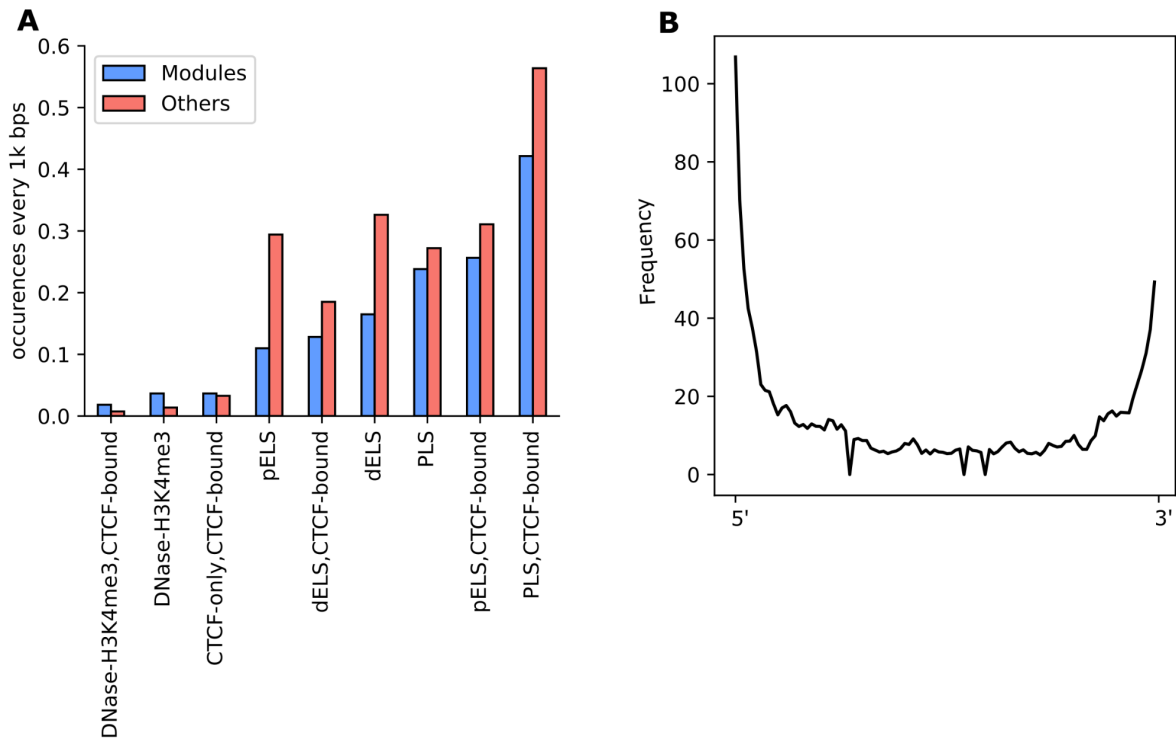
217    the main driver of this process.

218

**Fig.4 Analysis of the sequence regions flanking exon modules.** A) For each pair of genes

containing a shared exon module we compared the similarities of the upstream and downstream flanking

221    exons (when present); B) Distributions of the length-normalized Needleman and Wunsch scores of

222    exonic modules (in blue) and of their upstream and downstream flanking exons (in red); C) A pair of

223    exons in which the similarity only extends to  the downstream flanking intron; D) A pair of exons in which

224    the similarity extends upstream and downstream into both flanking introns; E) Overall representation of

225    all the length-scaled similarities between all the exon pairs and their flanking introns (in grey), the median

226    identity percentage is represented in red. The other colored lines represent five clusters of similarity

227    patterns as defined by grouping individual lines; F) Number of occurrences per thousand base pairs of

228    families of repetitive sequences in flanking introns with significant differences (padj<0.05) between the

229    exonic modules and the other lncRNA exons.

230

231

232         To further investigate the characteristics of these modules we looked at the distribution of cis-

233    regulatory elements (CRE) within their sequences (5A Fig.). This research highlighted a depletion in

234    exon modules of the most frequent CREs (Fisher's exact test padj=6.2e10-3, 1.2e10-2, 1.2e10-2 for

235    pELS, dELS and PLS, CTCF-bound respectively). One of the few elements that are not depleted are

236    H3K4me3 marks, which are characteristic of transcriptionally active regions (Howe et al. 2017).

237    Interestingly this histone modification is usually found in the region corresponding to the beginning of

238    the transcript [30]. Accordingly, when we investigated the position of the exonic modules within their

239    transcripts (5B Fig.), we detected a higher frequency of the modules at the 5' end. This finding is

240    consistent with what is observed in protein coding genes, which in vertebrates tend to increase their

241    length over time by gaining recently evolved domains, primarily through the addition of sequences at

242    the 5' end of genes [31]. The insertion of these modules at the extremities of the transcript  presumably

243    allows  the addition of genetic material with minimal disruption to the existing sequence.

**Fig 5 Cis-regulatory elements (CRE) and position of the modules**. A) number of occurrences of the different CREs from the annotation present in ENCODE every thousand nucleotides in the modules (in blue) and in the other lncRNA exons of the dataset (in red); B) the y axis indicates the frequency of regions containing modules relative to their position on their transcript (which is indicated on the Y axis, see Methods), as the sum of modules present in that region. The higher y value therefore indicates that there is a greater number of modules at the ends of the transcripts, particularly at the level of the 5' end.

## Evolutionary conservation of exon modules

To analyze in detail the inter-specific conservation of exon modules, we compared their conservation scores (see Materials and Methods) with the conservation scores of functionally annotated lncRNA exons, using the conservation scores of other lncRNA exons as control. Functionally annotated lncRNA genes were collected from the lnc2Cancer database [5], which contains experimentally supported annotations of lncRNA associated with a biological function, as derived from the literature (see Materials and Methods). The comparison of these three categories revealed that the conservation score of exon modules was higher than that of exons belonging to functionally annotated lncRNA genes
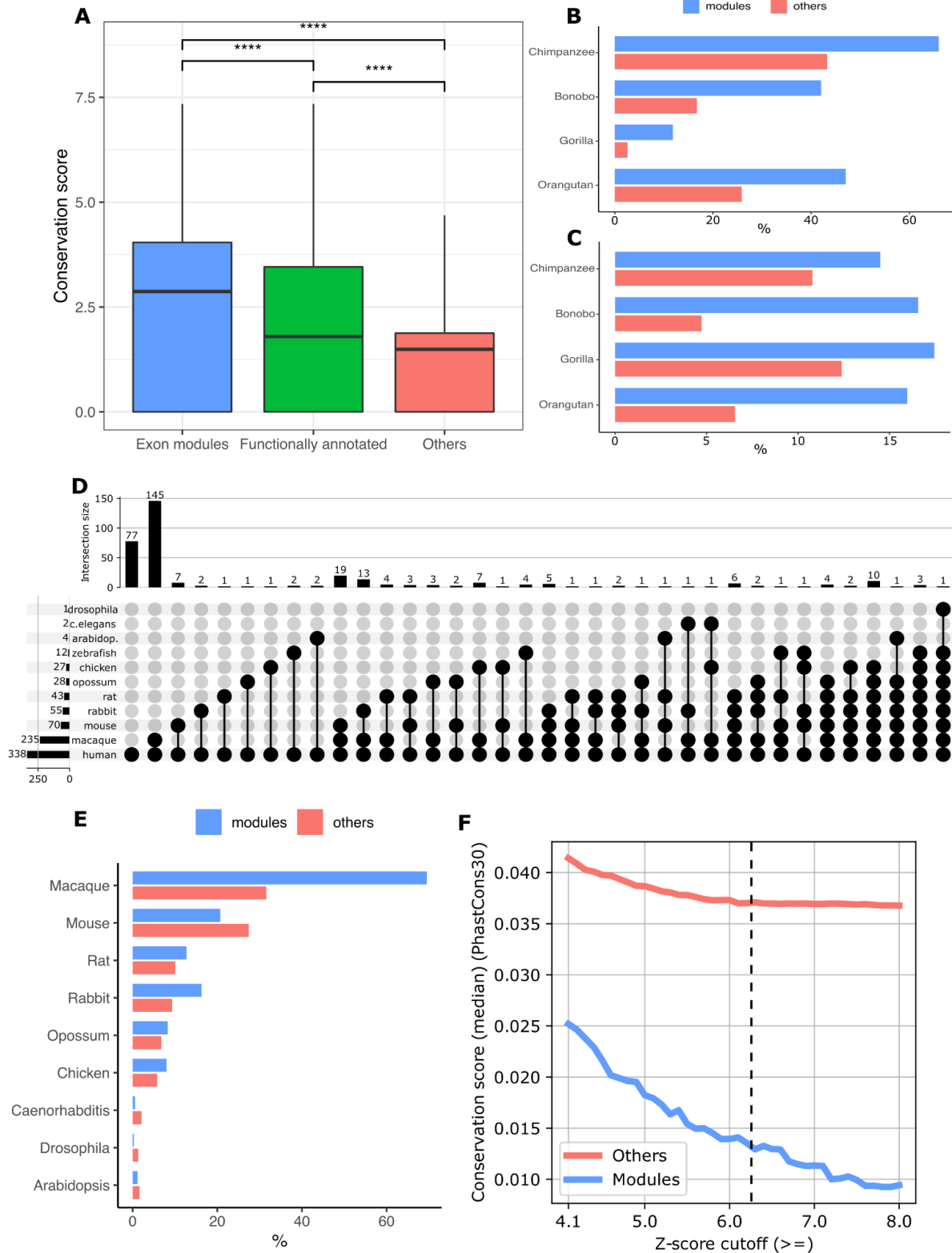
262    (Mann-Whitney p-value=6.3e-5), and both the conservation score of exon modules and of exons

263    belonging to functionally annotated lncRNA were significantly higher than the conservation score of the

264    remaining lncRNA exons (Mann-Whitney p-value=7.4e-27 and 3.5e-26 respectively, 6A Fig.). When

265    looking at the conservation of exon modules in four higher primate species, we also observed a greater

266    proportion of exons with a BLAST hit among exon modules vs the remaining exons. More specifically,

267    65.97% of the exon modules have a BLAST hit in Chimpanzee, 42.01% in Bonobo, 11.83% in Gorilla

268    and 47.04% in Orangutan. Conversely, only 43.23%, 16.72%, 2.61%, 25.86% of the control exons (i.e

269    the portion of the 12097 lncRNA exons that have no repetitive and non-overlapping sequences and that

270    are not modules) have BLAST hits on the same species respectively (6B Fig.) To evaluate the

271    significance of these results we performed a Fisher's exact test on the aggregated data from the different

272    species, which confirmed that these results are significant (p-value=4.10e-15).

273        Since the BLAST similarity score with non-human primates does not take into account the

274    genomic position of exons in different organisms, i.e. it cannot distinguish between the similarity of true

275    orthologs vs in- and out- paralogs, we investigated whether exon modules are located in regions of

276    synteny between non-human primates more often than other exons. To this end we leveraged the

277    SynthDB [32] database, which provides data on orthology relationships between humans and other

278    primates. We observed that the percentage of genes located in a syntenic region is higher for genes

279    that contain at least one exon module, compared with those which do not. Accordingly, 14.50% of the

280    exon modules are located in genes that have an ortholog in Chimpanzee, 16.57% in Bonobo, 17.46%

281    in Gorilla and 15.98% in Orangutan. While for the other lncRNA exons we observed percentages of

282    10.79, 4.74, 12.39, 6.55 in the same species respectively. We then performed a Fisher's exact test

283    comparing exons modules that belong to genes with an ortholog in at least one of the species mentioned

284    above to the other exons which confirmed the significance of our results (p-value=5.63e-05) (6C Fig.).

285        To strengthen the evolutionary conservation analysis, and to compare our results with the

286    analysis by Sarropoulos et al. 2019 [33], we extended it by including additional species. To this end, we

287    aligned all lncRNA exons using blastn against the genomes of the organisms used in Sarropoulos et al.

288    2019 (Macaque, taxid: 9544; Rabbit, taxid: 9986; Chicken, taxid: 9031; Opossum, taxid: 13616; Rat,

289    taxid: 10116; Mouse, taxid: 10090), and other model organisms (Danio rerio, taxid: 7955; Drosophila

290    melanogaster, taxid: 7227; Caenorhabditis elegans, taxid: 6239; Arabidopsis thaliana, taxid: 3702),

291    using an e-value threshold of 0.01 to identify hits (5D-E Fig., S3 Fig.). Figure 6E (6E Fig.) displays the

292     percentage of exonic modules vs other lncRNA exons that have at least one hit in the species indicated

293     above. This analysis shows  a rapid decay in the number of similar exons as the evolutionary distance

294     from humans increases. Figure 5F (5F Fig.) shows the 30 mammal PhastCons scores of the exon

295     modules, as a function of the z-score similarity threshold used to define the modules themselves (i.e.

296     the threshold described in 1A Fig.). This analysis demonstrates that the exon modules identified in this

297     work, which are highly similar as they were selected on the basis of having a Z-score of at least 6.2 and

298     5.3 in the sequence and structure alignment respectively, represent duplications that are recent (as

299     implied by the high levels of sequence similarity) and that are exclusively found in humans and higher

300     primates, and thus have lower PhastCons scores on the entire set of 30 mammals (6F Fig.) .

301     Overall, the above results reveal that roughly 4% of lncRNA genes (218 lncRNA genes/5,423

302     total lncRNAs genes which contain at least one exon without repetitive sequences, see Materials and

303     Methods) include one or more exons having significant similarity with exonic portions of other lncRNAs.

304     To our knowledge, this represents the first draft of a genome-wide catalog of shared lncRNA exons.

**Fig.6 Evolutionary conservation of exon modules.** A) Box-plot of the conservation scores in four non-human primates for exon modules, functionally annotated exons from the lnc2Cancer database, and controls; B) Percentage of exon modules (in blue) and other exons (in red) that showed a BLAST hit (e-value <0.001) in the primate species considered; C) Percentage of genes showing a

310  conserved syntenic region (as defined in SynthDB) among those containing exon modules (in blue) vs

311  genes not containing an exon module (in red); D) Upset plot representing the exons that have a BLAST

312  hit  in the species analyzed in Sarropoulos et al. and in other model organisms; E) Percentages of

313  modules (in blue) and other exons (in red) showing a BLAST hit in the indicated species  F) PhastCons

314  30 mammals scores of members of clusters defined by different z-score thresholds of pairwise similarity

315  from sequence alignments (in blue) and the other lncRNA exons of the dataset (in red).

316

## Nucleotide variation in modules

318

319  To further investigate whether exon modules may represent conserved functional units, we

320  analyzed the occurrence and frequency of single nucleotide polymorphisms (SNPs) in these regions, as

321  a lower incidence of variants may indicate the existence of constraints associated with  functional

322  sequences, due to the effects of purifying selection [34].  Accordingly, we collected SNP data from the

323  1000 Genome project from dbSNP 153 [35] and we observed 12.87 variants per thousand bases in

324  control exons (which are not modules) and 11.83 in modules. We then obtained from the ALFA allele

325  frequencies aggregator [36] a total of 764,005 SNPs  located in lncRNA exons, [36] and their associated

326  frequencies. For each exon, we calculated the index of nucleotide diversity $\theta\pi$ [37] as

327
$$\theta\pi = \frac{\sum_{i=1}^{l} 2f_i(1 - f_i)}{l}$$

328  where $f_i$ represents the frequency of variants in the $i$ th position of the exon sequence in the population,

329  and $l$ represents the length of the exon.

330  After comparing the distributions of $\theta\pi$ scores with the Mann-Whitney U test, we obtained a p-

331  value of 2.14e-02 in the comparison between modules and exons from functionally annotated genes, a

332  p-value of 2.77e-02 from the comparison between exon modules and other lncRNA exons and a non-

333  significant p-value (7.45e-01) from the comparison between functionally annotated and others,

334  confirming a significant lower propensity to harbor variation in exon modules as compared to the other

335  two groups. These findings indicate the existence of evolutionary constraints which limit the occurrence

336  of variants with polymorphic frequencies in exon modules, which in turn may reduce the rate of

337  evolutionary change in the long-term. We also looked at the frequency of polymorphic complete exon

338  deletions, but the results were not statistically significant (data not shown).

339

## Search for characteristics shared with protein coding genes

341

342    To confirm that exon modules do not simply represent mis-annotated protein domains, we
343    compared their sequence characteristics with those of known coding genes.

344    França et al. [38] observed that symmetric shuffling units (exons whose length is an exact
345    multiple of three) are strongly over-represented in human protein coding genes, due to their lower impact
346    on the reading frame when transposed. We found an opposite trend in lncRNA exon modules, with only
347    25% having a length that is a multiple of three, which confirms the lack of relevance of the reading frame.
348    By contrast, in the remainder of the exons, this proportion is 33%, i.e. what would be expected under a
349    random model.

350    The transition/transversion ratio (Ti/tv) among polymorphic variants should be 0.5 under a purely
351    random model, resulting from four possible transitions/eight possible transversions. However, real data
352    depart remarkably from this expectation, with functional regions and protein coding regions presenting
353    values higher than 0.5, since transitions are more likely to result in non-synonymous substitutions (e.g.
354    when they occur in the third base of a codon) [39]. Exon modules displayed values of 1.9, in line with
355    previous results for lncRNAs [40]. As a reference, these values contrast sharply with those for protein
356    coding genes, which range between 2.8-2.9 +/- 0.1 [40].

357

358

## Functional hypothesis and organization of putative modules

## in clusters of lncRNA genes

361

362    To further describe exon modules, here we show some examples of their organization within
363    the structure of their lncRNA genes. Only 12 of 218 genes containing exon modules are associated with
364    a known biological function in the lnc2Cancer database [5]. For most of them, the specific region of the
365    lncRNA molecule responsible for that function is unknown. In the next two paragraphs we will provide a
366    more detailed description for two of the identified modules, in an attempt to capture their putative

367    functions. The first example refers to an exon module recognized by virtue of sequence similarity, and

368    the second one refers to an exon module recognized by virtue of structure similarity.

369

## Identification of a putative YBX1 binding module

371

372    Figure 7A (7A Fig.) shows an example of a putative module represented in a pair of exons as a

373    sequence of ~200 nucleotides sharing a high sequence similarity (>87%). The exons involved are

374    ENSE00003710224.1 and ENSE00003838358.1 which belong to genes ENSG00000182165.17 (also

375    known as TP53TG1) and ENSG00000285540.1, respectively. TP53TG1 is a lncRNA involved in the p53

376    network response to DNA damage [41], which has a role as tumor suppressor by blocking the

377    tumorigenic activity of the RNA binding protein (RBP) YBX1 [42]. More in detail, the expression of

378    TP53TG1 is induced by p53 under cellular stress conditions that involve the induction of double-strand

379    breaks [41], while the interaction in the cytoplasm between TP53TG1 and YBX1 prevents the migration

380    of the latter inside the nucleus where it might promote the transcription of a series of oncogenes [43].

381    *Diaz-Lagares et al.* [42] demonstrated that a central region of TP53TG1, which includes the putative

382    module in the exon ENSE00003710224.1, is responsible for YBX1 binding. Moreover, they proved that

383    YBX1 binding motifs CACC are necessary to ensure the tumor-suppressor function of TP53TG1. We

384    identified two occurrences of the CACC motif in ENSE00003710224.1 and one in ENSE00003838358.1,

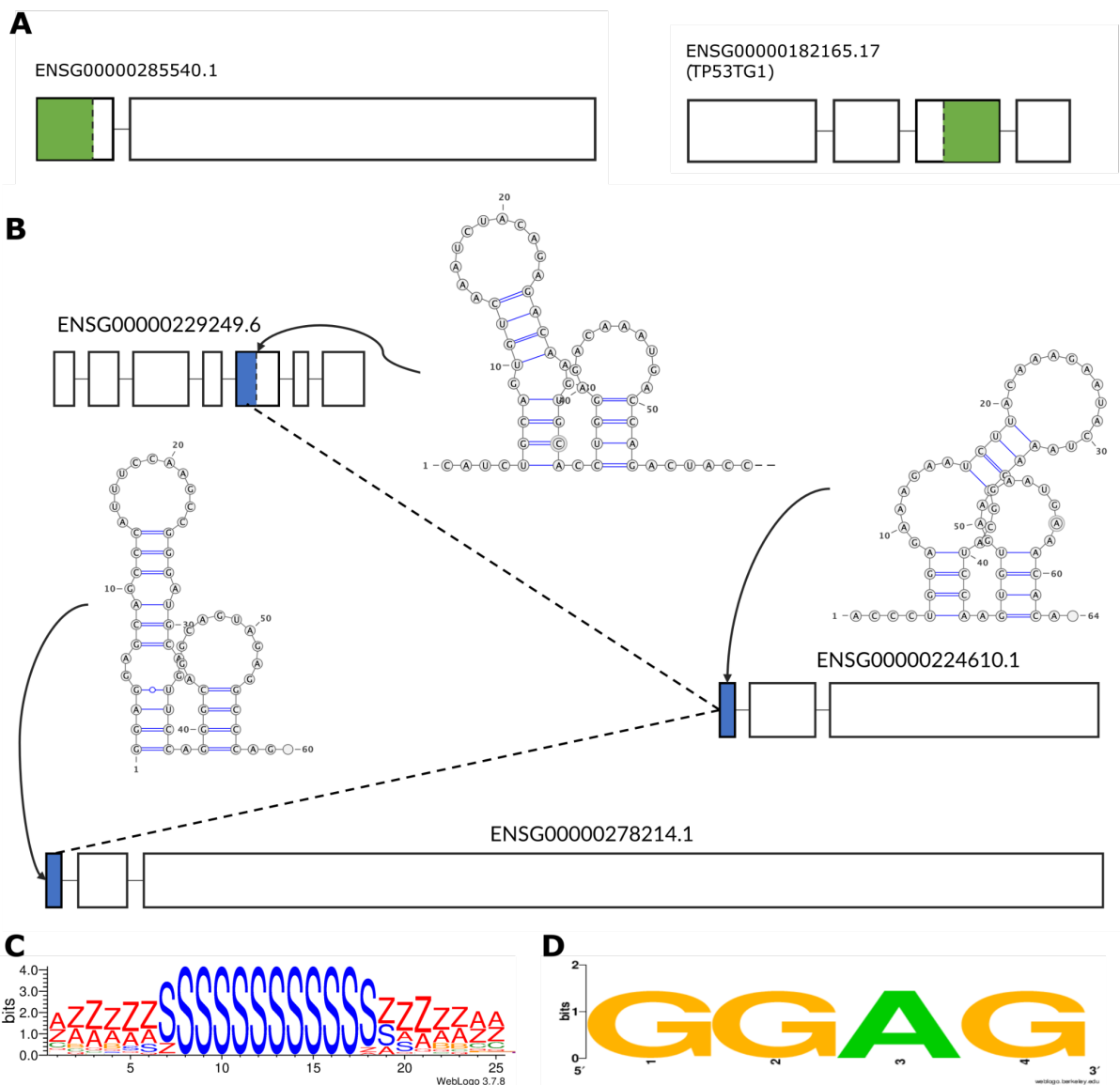385    suggesting a common role for this module.

386

## Identification of a putative LIN28B binding module

388

389    Figure 7B-D (7B-D Fig.) shows an example of a module with high structure similarity, embedded

390    in dissimilar sequence contexts. The exons involved are ENSE00003741285.1, ENSE00001800736.1

391    and ENSE00001782399.1 which belong to ENSG00000278214.1, ENSG00000224610.1 and

392    ENSG00000229249.6, respectively (7B Fig.). These three exons fold into a similar secondary structure,

393    composed of two stems ending with a hairpin loop, with one of the two stems having one or two internal

394    loops.

395    To detect a possible function, common to the three representatives of this exon module, we

396    searched for the presence of enriched structure and sequence motifs using the BRIO web server (see

397     Materials and Methods). BRIO identified a significantly enriched (Fisher's exact test padj<0.05) structure

398     motif shared between all the exons of the group (7C Fig.). This particular motif was associated by Adinolfi

399     *et al.* [44] with a series of different RNAs capable of binding some RBPs including LIN28B. This is an

400     evolutionary conserved RBP involved in several cellular processes, which acts as a critical oncogene

401     activated in cancer [45]. LIN28B is known to be able to bind different mRNAs, including a set of mRNAs

402     for splicing factors [46], miRNAs [47] and lncRNAs such as NEAT1 [48]. Furthermore, LIN28B C-terminal

403     zinc knuckle (ZnK) mediates specific binding to a conserved GGAG motif [49] which is also a sequence

404     motif present in all the three representatives of this module (7D Fig.). These observations suggest a

405     possible role of this module in binding LIN28B.



406

407

408       **Fig.7 Organization of a sequence and a structure module and identified motifs.** A)

409   Schematic representation of the lncRNA genes containing the putative YBX1 binding module (in green);

410   B) Representation of the lncRNA genes containing the exons with the putative LIN28B binding module

411   and their secondary structures. The blue boxes represent the exons with high structural similarity that

412   form the module; C) secondary structure motif revealed by BRIO represented with the BEAR alphabet

413   [50]; D) sequence motif recognized by ZnK in the three modules. The RNA secondary structure

414   representations were generated using VARNA (Darty et al. 2009); Sequence and structure logos were

415   generated using WebLogo [51].

416

# Discussion

418

419       This work identified a set of lncRNA exons with high sequence and/or structure similarity that

420   are embedded within globally dissimilar genes, confirming the hypothesis of exon sharing between this

421   class of molecules similar to protein-coding genes. This set contains a total of 340 pairs of exons that

422   can be grouped, on the basis of their reciprocal connections, in 106 clusters. In contrast to previous

423   work [18], our analysis focused on exons that do not contain repetitive sequences. The resulting dataset

424   of exon modules likely represents the result of recent segmental duplications that are almost exclusively

425   found in humans and higher primates. These findings support the hypothesis that the non-coding

426   transcriptome is structured into modular domains, similar to the organization observed in protein-coding

427   genes.

428       Approximately 4% (218 out of 5,423) of all the lncRNA genes in our dataset contain an exon

429   module. Even though we cannot assign a specific function to each of these modules, as it has been

430   done for the majority of protein coding domains, it is tempting to infer that sharing of functional modules

431   between different lncRNAs may contribute to expanding the functional repertoire of the non-coding

432   genome, similar to the shuffling of functional exons in coding sequences [12].

433       LncRNA exon modules identified in this work display a higher degree of sequence conservation

434   and synteny in four primate great ape species than the remainder of lncRNA exons. A high level of

435   conservation between related species is suggestive of purifying selection and is a landmark

436   characteristic of functional genetic elements [52]. Exon modules also harbor a lower frequency of SNPs

437   compared with control sequences, which suggests that purifying selection also persists intra-specifically

438 in human populations. Our set included 46 exon pairs highly similar in both sequence and structure (1F

439 Fig.), which are associated with the highest conservation scores. Even though we cannot infer the age

440 of the duplication/shuffling event based on our analysis, our results show that the exons involved are

441 subjected to extreme purifying selection, which preserved both sequence and structure. Taken together,

442 this evidence suggests that these modules play an important role within their respective lncRNA genes,

443 even though their exact function is yet to be characterized.

444 Some of the modules may also have undergone an accelerated divergence. Our set includes

445 219 and 75 exon pairs similar only in sequence or structure, respectively, and since our inclusion criteria

446 considered both similarity and evolutionary conservation, examples of accelerated evolution may have

447 escaped our search. This mechanism is equally relevant, especially when searching for evolutionary

448 innovations specific to the human lineage. However, different methods than those used here are

449 required to identify such cases. Finally, it was reported that homologous lncRNAs can, in some cases,

450 conserve their function over long evolutionary times, despite having diverged in both their nucleotide

451 sequences and their secondary structures [53]. The above considerations suggest that our analysis may

452 underestimate the extent of module sharing in lncRNAs. Other limitations include the fact that the correct

453 identification of exons within lncRNAs is strongly dependent on the reliability of the reconstruction of the

454 whole transcript structure. This is usually summarized by the TSL parameter (Transcript Support Level)

455 which we included, for every exon, in the Supporting information (S1 Table).

456 In the few cases for which functional information on a lncRNA is available, it may be possible to

457 infer the function of the shared module. We report two examples of modules conserved in either

458 sequence or structure. In both cases, the ability to bind specific targets is the inferred associated

459 function.

460 Overall, our results highlight the presence of groups of exons sharing high sequence or structure

461 similarity within dissimilar lncRNA genes. These exons are highly conserved across primate species

462 and depleted of inter-individual variation among humans (SNPs), and we suggest that they may

463 represent functional modules.

464 The identification of these modules could constitute a tool for decoding the function of the many

465 lncRNAs that are currently uncharacterized. Membership in a shared exon cluster represents a feature

466 that deserves annotation, even though conclusive proof of shared function will require experimental

467 evidence.

468

# Materials and Methods

470

## Dataset

472

473     We used gencode version 29 [3], to select 34,509 exons annotated as long intergenic non-

474     coding RNA, which do not have overlaps with protein coding genes, and downloaded their chromosomal

475     coordinates as a gtf file. We then used these coordinates to obtain the corresponding sequences from

476     the hg38 version of the human genome (UCSC genome browser), converting the gtf to bed file and

477     using the getfasta tool from the bedtools suite [54], with repetitive sequences masked by RepeatMasker

478     (Smit et al., unpublished data, www.repeatmasker.org) and Tandem Repeats Finder [55]. We removed

479     18,703 exons containing repetitive sequences and retained 15,806 exons. 3,709 of these were shared

480     by different isoforms of the same lncRNA gene. In such cases we only considered the longest isoform,

481     thus obtaining a final set of 12,097 non-overlapping exons that do not contain repetitive sequences.

482     These exons belong to 5,423 different lncRNA genes.

483

## Sequence alignments

485

486     All exon sequences were compared to each other using the Needleman and Wunsch global

487     alignment algorithm [56], using the same default gap penalties scores as the EMBOSS Needle tool for

488     global alignments of nucleic acids sequences [57] (-10 for gap insertions, -0.5 for gap extensions) and

489     the EDNAFULL substitution matrix.

490

## Structure alignments

492

493     The secondary structure of each exon was calculated using RNAfold [57,58], as the minimum

494     free energy (MFE) structure, and represented by its dot-bracket notation. These representations were

495     converted into the BEAR alphabet for RNA secondary structure notation (Mattei et al. 2014). The BEAR

496  alphabet is an encoding method for RNA secondary structure, whose characters encode for a specific

497  secondary structure element (loop, stem, bulge and internal loop) with specific length (e.g. a nucleotide

498  that is part of a stem of length 5 is represented by one character and a different character is used to

499  represent a stem of a different length). The global structure alignments were performed using the

500  BEAGLE algorithm [59] , with default parameters (-2 for gap insertions, -0.7 for gap extensions, +0.6 for

501  the sequence match bonus) and the substitution matrix for RNA structural elements (MBR, Matrix of

502  Bear-encoded RNAs) described in [50]. To avoid favoring alignments between unstructured regions we

503  modified the original MBR, assigning a score of 0 to matches in these regions. BEAGLE is an algorithm

504  for pairwise RNA secondary structure global comparison similar to the Needleman and Wunsch

505  algorithm for sequence alignments.

506  For both sequence and structure alignments we considered the scores of the aligned sequences

507  after trimming external gaps. The score of each alignment was normalized by its length, to avoid biases

508  towards longer sequences. We selected only alignments of a length of at least 50 nucleotides after the

509  external gap trimming. The final distributions consisted of approx. 73 million values, with z-scores

510  ranging from ~-36 to ~16 and from ~-3 to ~9, respectively.

511

512  # Repetitive elements and cis-regulatory elements

513

514  Repetitive sequences were mapped using the rmsk table from the UCSC genome browser,

515  which is derived from RepeatMasker (Smit et al., unpublished data, www.repeatmasker.org).

516  Cis-regulatory elements coordinates are derived from the ENCODE Registry of candidate cis-

517  Regulatory Elements (cCREs) combined from all human cell types [60]. The enrichments are calculated

518  using a Fisher's exact test between modules containing a particular CRE and the other lncRNA exons

519  of the dataset with a Benjamini-Hochberg correction.

520

521  # Evolutionary conservation score

522

523  The evolutionary conservation score for each exon was calculated using an approach similar to

524  [61], using the BLAST+ suite of command-line tools [62]. More specifically, the BLASTn algorithm was

525    used to perform an alignment of all the lncRNA exons of our dataset (12,097). In view of the pattern of

526    the evolutionary conservation of lncRNA sequences [14], we used the genomes of four primate species

527    closely related to *H. sapiens*: *Pan troglodytes* (Chimpanzee, taxid:9598), *Pan paniscus* (Bonobo,

528    taxid:9597), *Pongo pygmaeus* (Orangutan, taxid:9601) and *Gorilla gorilla* (Gorilla, taxid:9592). For each

529    lncRNA exon we then calculated a comprehensive conservation score as the sum of the best match bit-

530    score over the four species, divided by the length of the query sequence. Though the four organisms

531    are phyletically correlated, we used this procedure to buffer lineage-specific effects and potential

532    genome annotation errors.

533          For both sequence and structure similarity scores, the resulting distributions were compared

534    with the inter-specific degree of sequence conservation, under the hypothesis that constraints on exon

535    variation acted both intra- and inter-specifically. These comparisons were used to explore the

536    relationship between intra- and inter-specific conservation scores around the z-score value of 6.0

537    proposed by [63] as the threshold to distinguish homologous sequences (1 Fig.).

538          We excluded from this comparison exon pairs located in genes that are globally similar as the

539    similarity of the exons would simply reflect gene paralogy. To do so we performed a pairwise alignment

540    of the genes containing the exon pairs using BLASTn. The genomic coordinates of the whole genes,

541    including the introns, were retrieved from the gencode version 29 gtf file [3], and we used the same

542    procedure described above for the exons to obtain their sequences. Local alignments were performed

543    considering the smallest gene of the pair as the query and the longest as the subject, and excluding

544    pairs presenting a total query coverage greater than or equal to 80%. For each exon pair, we also

545    checked the coordinates from the bed file, excluding overlapping pairs.

546

## Syntenies

548

549          Synteny data were collected from SyntDB [32], which takes into account positional conservation

550    and sequence similarity to identify syntenic regions of human lncRNAs across primates. This database

551    comprises synteny information for 55632 transcripts. From this dataset we selected conservation data

552    in Chimpanzee, Bonobo, Orangutan and Gorilla for the 8,390 lncRNA transcripts containing the 12,097

553    exons in our dataset.

554

## Single nucleotide polymorphisms (SNPs)

SNPs locations were retrieved from common dbSNP 153 (variants with a minor allele frequency (MAF) of at least 1% (0.01) in the 1000 Genomes Phase 3 dataset) [35] and population frequencies were obtained from the ALFA allele frequency aggregator project [36]. The release 2 vcf format file contains variant frequency data aggregated from 79 different studies on more than 900 million SNPs. We used the tabix tool from the SAMtools suite of programs [64] to select SNPs located within one of the 12,097 exons in our dataset, obtaining ~764,000 variants with associated allele frequency information.

## Transition/transversion ratio

The transitions to transversion ratio (Ti/Tv) was calculated by using the variant data present in the common dbSNP 153 (see above) for all the 12,097 lncRNA exons in our dataset, as the number of pyrimidine-pyrimidine or purine-purine substitutions (transitions), divided by the number of purine-pyrimidine or pyrimidine-purine substitutions (transversions).

## Protein coding exons

The protein coding exon coordinates were obtained from the gencode version 29 annotation and mapped on the hg38 version of the human genome using the same procedure described for the lncRNA exons.

## Motifs scan

The search for sequence and structure motifs in the putative LIN28B binding module was performed using the BRIO (BEAM RNA Interaction mOtifs) web server [65]. This tool enables the identification of RNA sequence and structure motifs involved in protein binding in one or more input RNA

583 molecules, by measuring, through a Fisher's exact test, their enrichment compared to a background of

584 RNAs from Rfam with similar length and structure content, defined as the fraction of paired nucleotides

585 in the RNA secondary structure. The database of motifs that is included in BRIO is derived from high

586 throughput protein-RNA binding experiments (PAR-CLIP, eCLIP and HITS) analyzed in [44]. For this

587 analysis, we considered the default enrichment significance threshold of p-value<0.05 to evaluate the

588 enrichment of a motif in a group of exon modules. We chose to use this algorithm because in addition

589 to identifying common motifs on some particular modules, it allows us to associate them with motifs

590 enriched in RNA that interact with specific proteins from experimental data.

591

# Funding

593

597

# Author Contributions

599

600 Conceptualization: Francesco Ballesio, Gerardo Pepe, Gabriele Ausiello, Andrea Novelletto, Manuela

601 Helmer-Citterich, Pier Federico Gherardini.

602 Data curation: Francesco Ballesio.

603 Investigation: Francesco Ballesio.

604 Funding acquisition: Manuela Helmer-Citterich, Pier Federico Gherardini.

605 Writing – original draft: Francesco Ballesio.

606 Writing – review & editing: Gerardo Pepe, Gabriele Ausiello, Manuela Helmer-Citterich, Pier Federico

607 Gherardini, Andrea Novelletto.

608

609 Gabriele Ausiello, Manuela Helmer-Citterich, Pier Federico Gherardini, Andrea Novelletto are senior

610 authors.

611

# Data Availability

613

614    The annotation was obtained from GENCODE v29

615    (https://www.gencodegenes.org/human/release_29.html).

616    The hg38 version of the human genome was downloaded from UCSC genome browser

617    (http://hgdownload.cse.ucsc.edu/goldenPath/hg38/bigZips/).

618    The BEAGLE webserver for RNA structure alignments is available at:

619    http://beagle.bio.uniroma2.it.

620    Functionally annotated lncRNA was downloaded from: http://bio-

621    bigdata.hrbmu.edu.cn/lnc2cancer.

622    Variant frequencies in human populations are available in the ncbi website

623    (https://www.ncbi.nlm.nih.gov/snp/docs/gsr/alfa/#ftp-download).

624    The BRIO webserver for RNA interaction motif search is available at:

625    http://brio.bio.uniroma2.it.

626    SynthDB is available at: http://syntdb.amu.edu.pl.

627    For a list of the 340 exon pairs identified see S1 Table.

628

# Supporting information

630

631    S1 Fig. Numerosity of lncRNA exons per exon cluster.

632    S2 Fig. Positions of the exon modules on the human chromosomes

633    S3 Fig. Comparison of BLAST hit frequencies at different evolutionary divergence ages of

634    exonic modules (A) and lncRNAs genes analyzed by Sarropoulos et al. (B).

635    S1 Table A list of identified exonic modules and their properties.

636    S2 Table Repetitive elements in the regions flanking the exonic modules.

637    S3 Table Percentage of transposase domains identified in genes containing exon modules

638    and in other lncRNA genes.

639

640 # References

641    1.    Gilbert W. Why genes in pieces? In: Nature Publishing Group UK [Internet]. 1 Feb 1978

642          [cited 23 Oct 2023]. doi:10.1038/271501a0

643    2.    Engreitz JM, Haines JE, Perez EM, Munson G, Chen J, Kane M, et al. Local regulation

644          of gene expression by lncRNA promoters, transcription and splicing. Nature. 2016;539:

645          452–455.

646    3.    Frankish A, Diekhans M, Jungreis I, Lagarde J, Loveland JE, Mudge JM, et al.

647          GENCODE 2021. Nucleic Acids Res. 2021;49: D916–D923.

648    4.    Quek XC, Thomson DW, Maag JLV, Bartonicek N, Signal B, Clark MB, et al. lncRNAdb

649          v2.0: expanding the reference database for functional long noncoding RNAs. Nucleic

650          Acids Res. 2015;43: D168–73.

651    5.    Gao Y, Shang S, Guo S, Li X, Zhou H, Liu H, et al. Lnc2Cancer 3.0: an updated

652          resource for experimentally supported lncRNA/circRNA cancer associations and web

653          tools based on RNA-seq and scRNA-seq data. Nucleic Acids Res. 2021;49: D1251–

654          D1258.

655    6.    Fort V, Khelifi G, Hussein SMI. Long non-coding RNAs and transposable elements: A

656          functional relationship. Biochim Biophys Acta Mol Cell Res. 2021;1868: 118837.

657    7.    Ponting CP, Oliver PL, Reik W. Evolution and functions of long noncoding RNAs. Cell.

658          2009;136: 629–641.

659    8.    Statello L, Guo C-J, Chen L-L, Huarte M. Gene regulation by long non-coding RNAs and

its biological functions. Nat Rev Mol Cell Biol. 2021;22: 96–118.

9. Gupta RA, Shah N, Wang KC, Kim J, Horlings HM, Wong DJ, et al. Long non-coding RNA HOTAIR reprograms chromatin state to promote cancer metastasis. Nature. 2010;464: 1071–1076.

10. Ng S-Y, Johnson R, Stanton LW. Human long non-coding RNAs promote pluripotency and neuronal differentiation by association with chromatin modifiers and transcription factors. EMBO J. 2012;31: 522–533.

11. Faghihi MA, Modarresi F, Khalil AM, Wood DE, Sahagan BG, Morgan TE, et al. Expression of a noncoding RNA is elevated in Alzheimer's disease and drives rapid feed-forward regulation of beta-secretase. Nat Med. 2008;14: 723–730.

12. Guttman M, Rinn JL. Modular regulatory principles of large non-coding RNAs. Nature. 2012;482: 339–346.

13. Ulitsky I, Bartel DP. lincRNAs: genomics, evolution, and mechanisms. Cell. 2013;154: 26–46.

14. Johnsson P, Lipovich L, Grandér D, Morris KV. Evolutionary conservation of long non-coding RNAs; sequence, structure, function. Biochim Biophys Acta. 2014;1840: 1063–1071.

15. Ulitsky I, Shkumatava A, Jan CH, Sive H, Bartel DP. Conserved function of lincRNAs in vertebrate embryonic development despite rapid sequence evolution. Cell. 2011;147: 1537–1550.

16. Kelley D, Rinn J. Transposable elements reveal a stem cell-specific class of long noncoding RNAs. Genome Biol. 2012;13: R107.

17. Fueyo R, Judd J, Feschotte C, Wysocka J. Roles of transposable elements in the regulation of mammalian transcription. Nat Rev Mol Cell Biol. 2022;23: 481–497.

684  18. Johnson R, Guigó R. The RIDL hypothesis: transposable elements as functional

685      domains of long noncoding RNAs. RNA. 2014;20: 959–976.

686  19. Kapusta A, Kronenberg Z, Lynch VJ, Zhuo X, Ramsay L, Bourque G, et al. Transposable

687      elements are major contributors to the origin, diversification, and regulation of vertebrate

688      long noncoding RNAs. PLoS Genet. 2013;9: e1003470.

689  20. Kirk JM, Kim SO, Inoue K, Smola MJ, Lee DM, Schertzer MD, et al. Functional

690      classification of long non-coding RNAs by k-mer content. Nat Genet. 2018;50: 1474–

691      1482.

692  21. Martin L, Meier M, Lyons SM, Sit RV, Marzluff WF, Quake SR, et al. Systematic

693      reconstruction of RNA functional motifs with high-throughput microfluidics. Nat Methods.

694      2012;9: 1192–1194.

695  22. Muckenthaler MU, Galy B, Hentze MW. Systemic iron homeostasis and the iron-

696      responsive element/iron-regulatory protein (IRE/IRP) regulatory network. Annu Rev Nutr.

697      2008;28: 197–213.

698  23. Zhang C, Lee K-Y, Swanson MS, Darnell RB. Prediction of clustered RNA-binding

699      protein motif sites in the mammalian genome. Nucleic Acids Res. 2013;41: 6793–6807.

700  24. Oberstrass FC, Lee A, Stefl R, Janis M, Chanfreau G, Allain FH-T. Shape-specific

701      recognition in the structure of the Vts1p SAM domain with RNA. Nat Struct Mol Biol.

702      2006;13: 160–167.

703  25. Gustavsen JA, Pai S, Isserlin R, Demchak B, Pico AR. RCy3: Network biology using

704      Cytoscape from within R. F1000Res. 2019;8: 1774.

705  26. Bailey JA, Gu Z, Clark RA, Reinert K, Samonte RV, Schwartz S, et al. Recent segmental

706      duplications in the human genome. Science. 2002;297: 1003–1007.

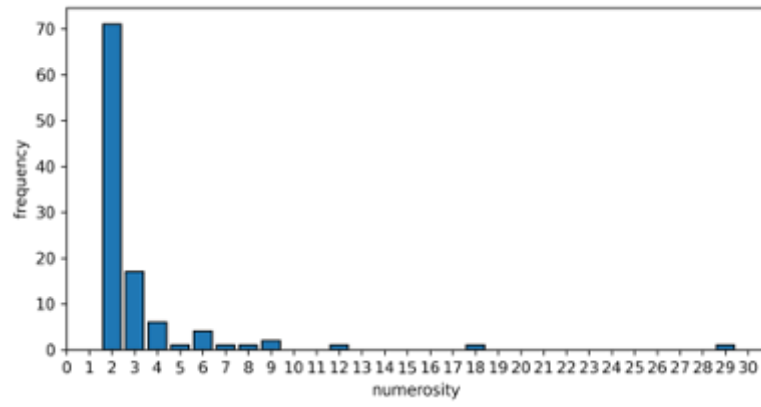707  27. Antonarakis SE. Content and variation of the human genome. Medical and Health

708  Genomics. Elsevier; 2016. pp. 161–177.

709 28. Abdullaev ET, Umarova IR, Arndt PF. Modelling segmental duplications in the human

710   genome. BMC Genomics. 2021;22: 496.

711 29. Koch L. Capturing transposases for new proteins. Nature reviews. Genetics. 2021. pp.

712   266–267.

713 30. Li B, Carey M, Workman JL. The role of chromatin during transcription. Cell. 2007;128:

714   707–719.

715 31. Toll-Riera M, Albà MM. Emergence of novel domains in proteins. BMC Evol Biol.

716   2013;13: 1–10.

717 32. Bryzghalov O, Szcześniak MW, Makałowska I. SyntDB: defining orthologues of human

718   long noncoding RNAs across primates. Nucleic Acids Res. 2020;48: D238–D245.

719 33. Sarropoulos I, Marin R, Cardoso-Moreira M, Kaessmann H. Developmental dynamics of

720   lncRNAs across mammalian organs and species. Nature. 2019;571: 510–514.

721 34. Cvijović I, Good BH, Desai MM. The Effect of Strong Purifying Selection on Genetic

722   Diversity. Genetics. 2018;209: 1235–1278.

723 35. Sherry ST, Ward M-H, Kholodov M, Baker J, Phan L, Smigielski EM, et al. dbSNP: the

724   NCBI database of genetic variation. Nucleic Acids Res. 2001;29: 308–311.

725 36. Phan, Jin, Zhang, Qiang, Shekhtman, Shao, et al. ALFA: allele frequency aggregator.

726   National Center for Biotechnology Information, US National Library of Medicine.

727 37. Nei M. Molecular Evolutionary Genetics. New York Chichester, West Sussex: Columbia

728   University Press; 1987.

729 38. França GS, Cancherini DV, de Souza SJ. Evolutionary history of exon shuffling.

730   Genetica. 2012;140: 249–257.

731   39. Yang Z, Bielawski JP. Statistical methods for detecting molecular adaptation. Trends

732         Ecol Evol. 2000;15: 496–503.

733   40. Wang J, Raskin L, Samuels DC, Shyr Y, Guo Y. Genome measures used for quality

734         control are dependent on gene function and ancestry. Bioinformatics. 2015;31: 318–323.

735   41. Takei Y, Ishikawa S, Tokino T, Muto T, Nakamura Y. Isolation of a novel TP53 target

736         gene from a colon cancer cell line carrying a highly regulated wild-type TP53 expression

737         system. Genes Chromosomes Cancer. 1998;23: 1–9.

738   42. Diaz-Lagares A, Crujeiras AB, Lopez-Serra P, Soler M, Setien F, Goyal A, et al.

739         Epigenetic inactivation of the p53-induced long noncoding RNA TP53 target 1 in human

740         cancer. Proc Natl Acad Sci U S A. 2016;113: E7535–E7544.

741   43. Finkbeiner MR, Astanehe A, To K, Fotovati A, Davies AH, Zhao Y, et al. Profiling YB-1

742         target genes uncovers a new mechanism for MET receptor regulation in normal and

743         malignant human mammary cells. Oncogene. 2009;28: 1421–1431.

744   44. Adinolfi M, Pietrosanto M, Parca L, Ausiello G, Ferrè F, Helmer-Citterich M. Discovering

745         sequence and structure landscapes in RNA interaction motifs. Nucleic Acids Res.

746         2019;47: 4958–4969.

747   45. Lin X, Shen J, Dan Peng, He X, Xu C, Chen X, et al. RNA-binding protein LIN28B

748         inhibits apoptosis through regulation of the AKT2/FOXO3A/BIM axis in ovarian cancer

749         cells. Signal Transduct Target Ther. 2018;3: 23.

750   46. Wilbert ML, Huelga SC, Kapeli K, Stark TJ, Liang TY, Chen SX, et al. LIN28 binds

751         messenger RNAs at GGAGA motifs and regulates splicing factor abundance. Mol Cell.

752         2012;48: 195–206.

753   47. Piskounova E, Polytarchou C, Thornton JE, LaPierre RJ, Pothoulakis C, Hagan JP, et al.

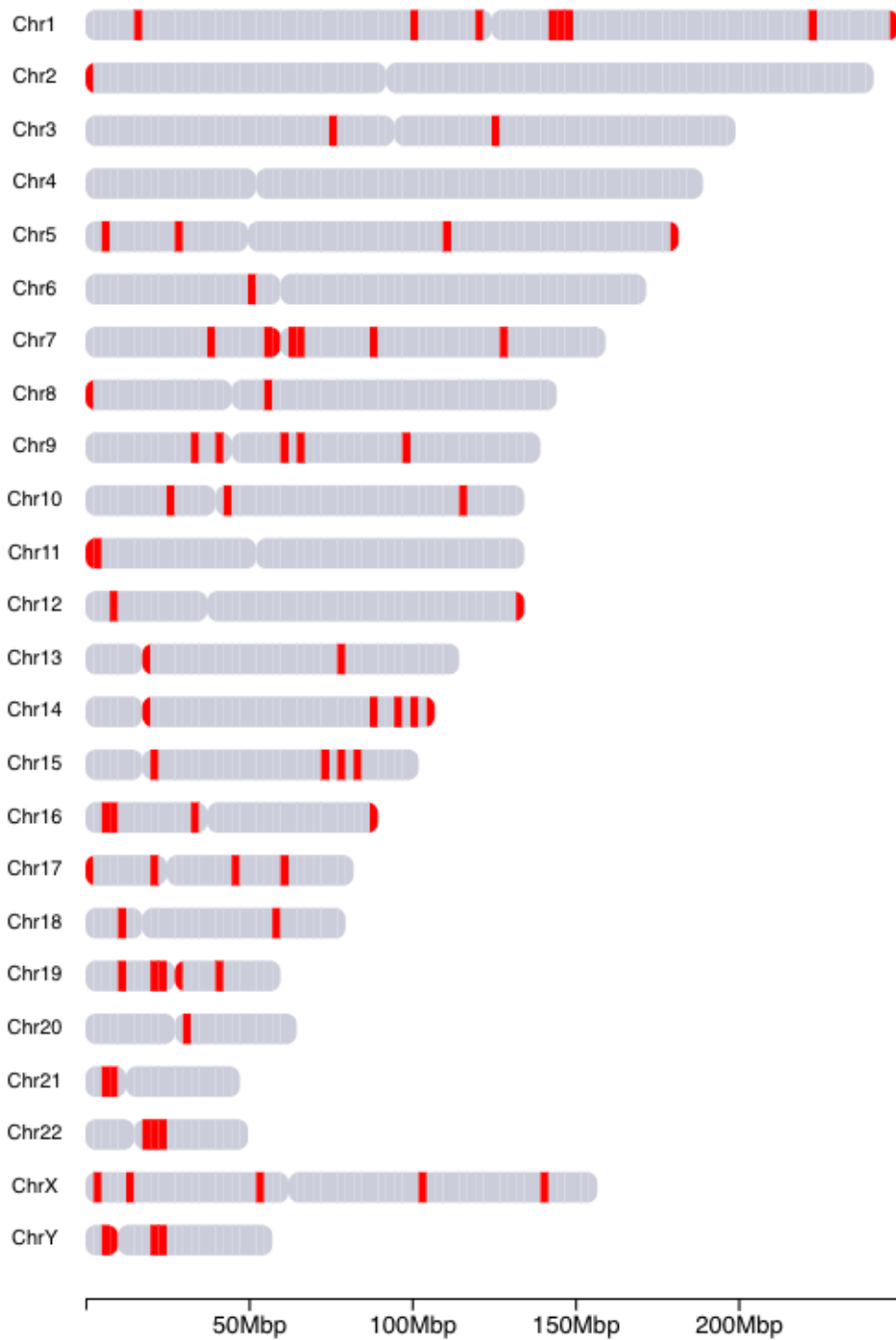754         Lin28A and Lin28B inhibit let-7 microRNA biogenesis by distinct mechanisms. Cell.

755      2011;147: 1066–1079.

756    48.  Yong W, Yu D, Jun Z, Yachen D, Weiwei W, Midie X, et al. Long noncoding RNA

757         NEAT1, regulated by LIN28B, promotes cell proliferation and migration through sponging

758         miR-506 in high-grade serous ovarian cancer. Cell Death Dis. 2018;9: 861.

759    49.  Peters DT, Fung HKH, Levdikov VM, Irmscher T, Warrander FC, Greive SJ, et al.

760         Human Lin28 Forms a High-Affinity 1:1 Complex with the 106~363 Cluster miRNA miR-

761         363. Biochemistry. 2016;55: 5021–5027.

762    50.  Mattei E, Ausiello G, Ferrè F, Helmer-Citterich M. A novel approach to represent and

763         compare RNA secondary structures. Nucleic Acids Res. 2014;42: 6146–6157.

764    51.  Crooks GE, Hon G, Chandonia J-M, Brenner SE. WebLogo: a sequence logo generator.

765         Genome Res. 2004;14: 1188–1190.

766    52.  Kellis M, Wold B, Snyder MP, Bernstein BE, Kundaje A, Marinov GK, et al. Defining

767         functional DNA elements in the human genome. Proc Natl Acad Sci U S A. 2014;111:

768         6131–6138.

769    53.  Karner H, Webb C-H, Carmona S, Liu Y, Lin B, Erhard M, et al. Functional Conservation

770         of LncRNA JPX Despite Sequence and Structural Divergence. J Mol Biol. 2020;432:

771         283–300.

772    54.  Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic

773         features. Bioinformatics. 2010;26: 841–842.

774    55.  Benson G. Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids

775         Res. 1999;27: 573–580.

776    56.  Needleman SB, Wunsch CD. A general method applicable to the search for similarities

777         in the amino acid sequence of two proteins. J Mol Biol. 1970;48: 443–453.

778   57.   Rice P, Longden I, Bleasby A. EMBOSS: the European Molecular Biology Open

779         Software Suite. Trends Genet. 2000;16: 276–277.

780   58.   Lorenz R, Bernhart SH, Höner Zu Siederdissen C, Tafer H, Flamm C, Stadler PF, et al.

781         ViennaRNA Package 2.0. Algorithms Mol Biol. 2011;6: 26.

782   59.   Mattei E, Pietrosanto M, Ferrè F, Helmer-Citterich M. Web-Beagle: a web server for the

783         alignment of RNA secondary structures. Nucleic Acids Res. 2015;43: W493–7.

784   60.   Moore JE, Purcaro MJ, Pratt HE, Epstein CB, Shoresh N, Adrian J, et al. Expanded

785         encyclopaedias of DNA elements in the human and mouse genomes. Nature. 2020;583:

786         699–710.

787   61.   Jha A, Quesnel-Vallières M, Wang D, Thomas-Tikhonenko A, Lynch KW, Barash Y.

788         Identifying common transcriptome signatures of cancer by interpreting deep learning

789         models. Genome Biol. 2022;23: 117.

790   62.   Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:

791         architecture and applications. BMC Bioinformatics. 2009;10: 421.

792   63.   Mitrophanov AY, Borodovsky M. Statistical significance in biological sequence analysis.

793         Brief Bioinform. 2006;7: 2–24.

794   64.   Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve years of

795         SAMtools and BCFtools. Gigascience. 2021;10. doi:10.1093/gigascience/giab008

796   65.   Guarracino A, Pepe G, Ballesio F, Adinolfi M, Pietrosanto M, Sangiovanni E, et al. BRIO:

797         a web server for RNA sequence and structure motif scan. Nucleic Acids Res. 2021;49:

798         W67–W71.
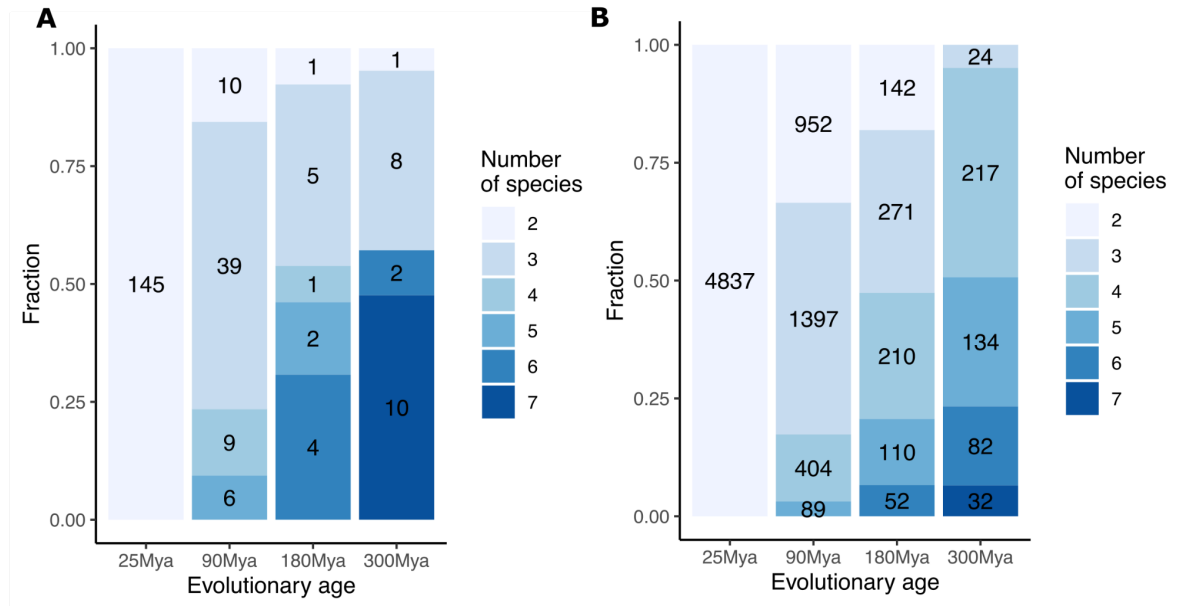
799

800

801



802     Fig. S1 - Numerosity of lncRNA exons per exon cluster.

803

804    Fig. S2 - Positions of the exon modules on the human chromosomes

Fig. S3 - Comparison of BLAST hit frequencies at different evolutionary divergence ages of exonic modules (A) and lncRNAs genes analyzed by Sarropoulos et al. (B).