Original Research Article

# Human lncRNAs harbor conserved modules embedded in different sequence contexts

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

[a] *PhD Program in Cellular and Molecular Biology, Department of Biology, University of Rome "Tor Vergata", Rome, Italy*
[b] *Department of Biology, University of Rome "Tor Vergata", Rome, Italy*

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.

## 1. Introduction

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 non-coding 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 (ncRNAs) 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 in lncRNAdb v2.0 [4], 1273 human lncRNAs in Lnc2Cancer 3.0 [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 that similarity amongst lncRNA within a given organism

---

is also limited, and, unlike coding sequences, most lncRNAs appear in single copies in vertebrate genomes [13].

However, lncRNAs are significantly more likely to contain repetitive sequences, particularly transposable elements (TEs) [15,16]. On one hand, this could simply indicate that lncRNAs are more prone to transposon insertion, because of their aforementioned looser association between sequence and function [13]. On the other hand, this observation implies the existence of stretches of homologous sequences that are shared among different lncRNAs, even when the lncRNAs themselves are not related by descent.

Because TEs are often enriched in sequences with regulatory function, and may contribute to their "spread" within a genome [17], Johnson and Guigò [18] hypothesized that the presence of TEs may result in the sharing of functional cassettes among evolutionarily unrelated lncRNA, possibly implying a modularization of function for this class of molecules [6,12], reminiscent of the notion of domains in the protein-coding world. In support of this hypothesis, it has been reported that TE-derived sequences within lncRNAs are more conserved compared with non-TE sequences [19].

Here we set out to expand the identification of modules in lncRNAs that could have contributed to increasing the diversity of the non-coding genome, similar to the exon-shuffling phenomenon that is well known for protein sequences. Our work extends previous observations in three ways, namely by i) focusing on the sharing of individual exons among unrelated lncRNAs within the human genome, ii) specifically excluding exons that contain repetitive sequences, and iii) including secondary structure as an additional criterion to define similarity, as lncRNAs with similar functions often lack linear sequence homology [20], and many examples of ncRNAs are known whose function is tied to their secondary structure [21–24].

## 2. Materials and methods

### 2.1. Dataset

We used gencode version 29 [3], to select 34,509 exons annotated as long intergenic non-coding RNA, which do not have overlaps with protein-coding genes, and downloaded their chromosomal coordinates as a gtf file. We then used these coordinates to obtain the corresponding sequences from the hg38 version of the human genome (UCSC genome browser), converting the gtf to bed file and using the getfasta tool from the bedtools suite [25], with repetitive sequences masked by RepeatMasker (Smit et al., unpublished data, www.repeatmasker.org) and Tandem Repeats Finder [26]. We removed 18,703 exons containing repetitive sequences and retained 15,806 exons. 3709 of these were shared by different isoforms of the same lncRNA gene. In such cases we only considered the longest isoform, thus obtaining a final set of 12,097 non-overlapping exons that do not contain repetitive sequences. These exons belong to 5423 different lncRNA genes.

### 2.2. Sequence alignments

All exon sequences were compared to each other using the Needleman and Wunsch global alignment algorithm [27], using the same default gap penalties scores as the EMBOSS Needle tool for global alignments of nucleic acids sequences [28] (−10 for gap insertions, −0.5 for gap extensions) and the EDNAFULL substitution matrix.

### 2.3. Structure alignments

The secondary structure of each exon was calculated using RNAfold [28,29], as the minimum free energy (MFE) structure, and represented by its dot-bracket notation. These representations were converted into the BEAR alphabet for RNA secondary structure notation [30]. The BEAR alphabet is an encoding method for RNA secondary structure, whose characters encode for a specific secondary structure element

(loop, stem, bulge and internal loop) with specific length (e.g. a nucleotide that is part of a stem of length 5 is represented by one character and a different character is used to represent a stem of a different length). The global structure alignments were performed using the BEAGLE algorithm [31], with default parameters (−2 for gap insertions, −0.7 for gap extensions, +0.6 for the sequence match bonus) and the substitution matrix for RNA structural elements (MBR, Matrix of Bear-encoded RNAs) described by Mattei et al. [30]. To avoid favoring alignments between unstructured regions we modified the original MBR, assigning a score of 0 to matches in these regions. BEAGLE is an algorithm for pairwise RNA secondary structure global comparison similar to the Needleman and Wunsch algorithm for sequence alignments.

For both sequence and structure alignments we considered the scores of the aligned sequences after trimming external gaps. The score of each alignment was normalized by its length, to avoid biases towards longer sequences. We selected only alignments of a length of at least 50 nucleotides after the external gap trimming. The final distributions consisted of approx. 73 million values, with z-scores ranging from ~-36 to ~16 and from ~-3 to ~9, respectively.

### 2.4. Repetitive elements and cis-regulatory elements

Repetitive sequences were mapped using the rmsk table from the UCSC genome browser, which is derived from RepeatMasker (Smit et al., unpublished data, http://www.repeatmasker.org).

Cis-regulatory elements coordinates were derived from the ENCODE Registry of candidate cis-Regulatory Elements (cCREs) combined from all human cell types [32]. The enrichments were calculated using a Fisher's exact test between modules containing a particular CRE and the other lncRNA exons of the dataset with a Benjamini-Hochberg correction.

### 2.5. Evolutionary conservation score

The evolutionary conservation score for each exon was calculated using an approach similar to Ref. [33], using the BLAST + suite of command-line tools [34]. More specifically, the BLASTn algorithm was used to perform an alignment of all the lncRNA exons of our dataset (12, 097). In view of the pattern of the evolutionary conservation of lncRNA sequences [14], we used the genomes of four primate species closely related to *H. sapiens*: *Pan troglodytes* (Chimpanzee, taxid:9598), *Pan paniscus* (Bonobo, taxid:9597), *Pongo pygmaeus* (Orangutan, taxid:9601) and *Gorilla gorilla* (Gorilla, taxid:9592). For each lncRNA exon we then calculated a comprehensive conservation score as the sum of the best match bit-score over the four species, divided by the length of the query sequence. Though the four organisms are phyletically related, we used this procedure to buffer lineage-specific effects and potential genome annotation errors.

For both sequence and structure similarity scores, the resulting distributions were compared with the inter-specific degree of sequence conservation, under the hypothesis that constraints on exon variation acted both intra- and inter-specifically. These comparisons were used to explore the relationship between intra- and inter-specific conservation scores around the z-score value of 6.0 proposed by Mitrophanov and Borodovsky [35] as the threshold to distinguish homologous sequences (Fig. 1).

We excluded from this comparison exon pairs located in genes that are globally similar as the similarity of the exons would simply reflect gene paralogy. To do so we performed a pairwise alignment of the genes containing the exon pairs using BLASTn. The genomic coordinates of the whole genes, including the introns, were retrieved from the gencode version 29 gtf file [3], and we used the same procedure described above for the exons to obtain their sequences. Local alignments were performed considering the smallest gene of the pair as the query and the longest as the subject, and excluding pairs presenting a total query
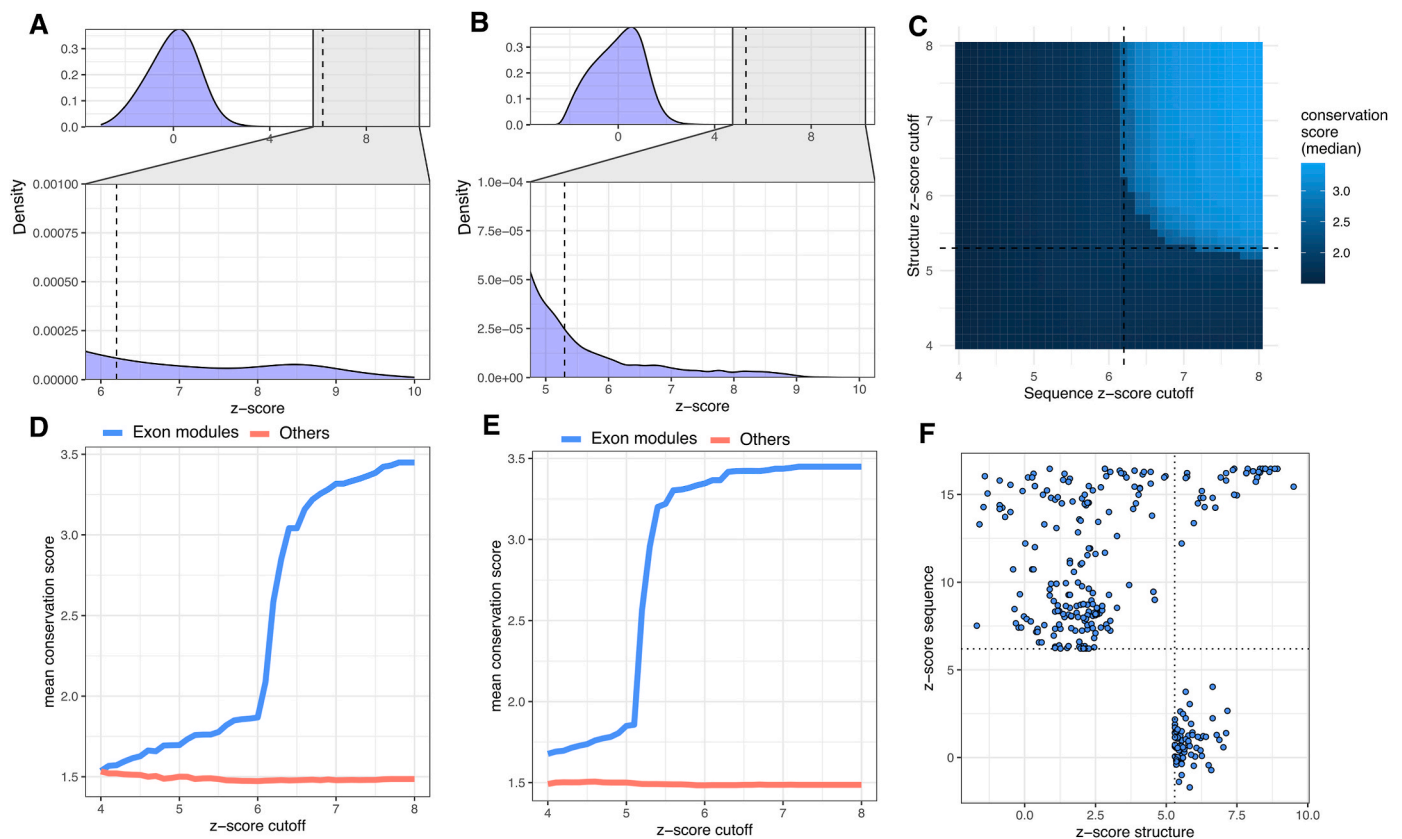
**Fig. 1. Sequence and structure alignments results.** A) Distribution of z-transformed pairwise alignment scores for sequence; B) Distribution of z-transformed pairwise alignment scores for structures, for these distributions, a close-up around the proposed cutoff thresholds is also shown; C) heatmap 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).

coverage greater than or equal to 80 %. For each exon pair, we also checked the coordinates from the bed file, excluding overlapping pairs.

### 2.6. Syntenies

Synteny data were collected from SyntDB [36], which takes into account positional conservation and sequence similarity to identify syntenic regions of human lncRNAs across primates. This database comprises synteny information for 55,632 transcripts. From this dataset we selected conservation data in Chimpanzee, Bonobo, Orangutan and Gorilla for the 8390 lncRNA transcripts containing the 12,097 exons in our dataset.

### 2.7. 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) [37] and population frequencies were obtained from the ALFA allele frequency aggregator project [38]. 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 [39] to select SNPs located within each of the 12,097 exons in our dataset, obtaining ~764,000 variants with associated allele frequency information.

### 2.8. Transition/transversion ratio

The transition 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).

### 2.9. 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.

### 2.10. 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 [40]. This tool enables the identification of RNA sequence and structure motifs involved in protein binding in one or more input RNA molecules, by measuring, through a Fisher's exact test, their enrichment compared to a background of RNAs from Rfam with similar length and structure content, defined as the fraction of paired nucleotides in the RNA secondary structure. The database of motifs that is included in BRIO is derived from high throughput protein-RNA binding experiments (PAR-CLIP, eCLIP and HITS) analyzed by Adinolfi et al. [41]. For this analysis, we considered the default enrichment significance threshold of p-value<0.05 to evaluate the enrichment of a motif in a group of exon modules. We chose to use this algorithm

because in addition to identifying common motifs on some particular modules, it allows us to associate them with motifs enriched in RNA that interact with specific proteins from experimental data.

## 3. Results

### 3.1. Exon sequence and secondary structure comparison

In order to search for similarities among lncRNAs, we performed a pairwise comparison of both the sequence and the predicted secondary structure of 12,097 non-overlapping human lncRNA exons that do not contain repetitive sequences, performing a total of more than 73 million sequence alignments and an equal number of structure alignments. The distributions of the corresponding scores are shown in (Fig. 1A and B).

To identify pairs or groups of exons representing shared sequence elements, hereafter referred to as "modules", it was necessary to select a threshold above which their sequence or structure similarity would be considered significant.

We thus investigated the conservation of lncRNA exons in four non-human primates (see Materials and Methods), with the goal of identifying shared sequence elements in the human genome that are also conserved in other primate genomes.

Accordingly, we calculated the mean conservation scores of sequence modules across these species, as a function of the similarity score threshold used to define the modules themselves. Using this procedure, we observed a sharp transition in conservation at Z-score similarity thresholds of 6.2 and 5.3 for sequence and structure alignments, respectively (Fig. 1C-E). We consider this increase in conservation,

coupled with the high Z-score similarity threshold, as a strong indication that the shared sequence elements we identified represent significant similarities. As a further benchmark, we repeated the entire procedure by aligning exons against random sequences with the same length and base composition. None of the alignments produced z-scores above the 6.2 threshold.

By using these thresholds, we identified a total of 340 exon pairs (219 identified by sequence, 75 by structure and 46 by both), involving 338 different exons and 218 different genes (Fig. 1F). Starting from these pairwise similarities, we identified 106 clusters (exon modules) defined by homologous lncRNA exons represented in at least two copies in the same or different genes (Fig. 2, S1 and Table S1).

To rule out the possibility that similarity between exons in a pair of genes is simply due to paralogy, we aligned the entire genes using BLAST and excluded pairs with alignment coverage on the smallest gene of the pair greater than 80 %. Measuring the alignment coverage of the entire genes, including introns, allowed us to identify and exclude cases of complete paralogy even in the presence of intronization or imprecise exon annotation.

We note that, in general, our analysis is dependent on the reliability of the reconstruction of the whole transcript structure, which is used to define the exons themselves. This is summarized by the Transcript Support Level (TSL, Table S1).

Fig. 3A 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
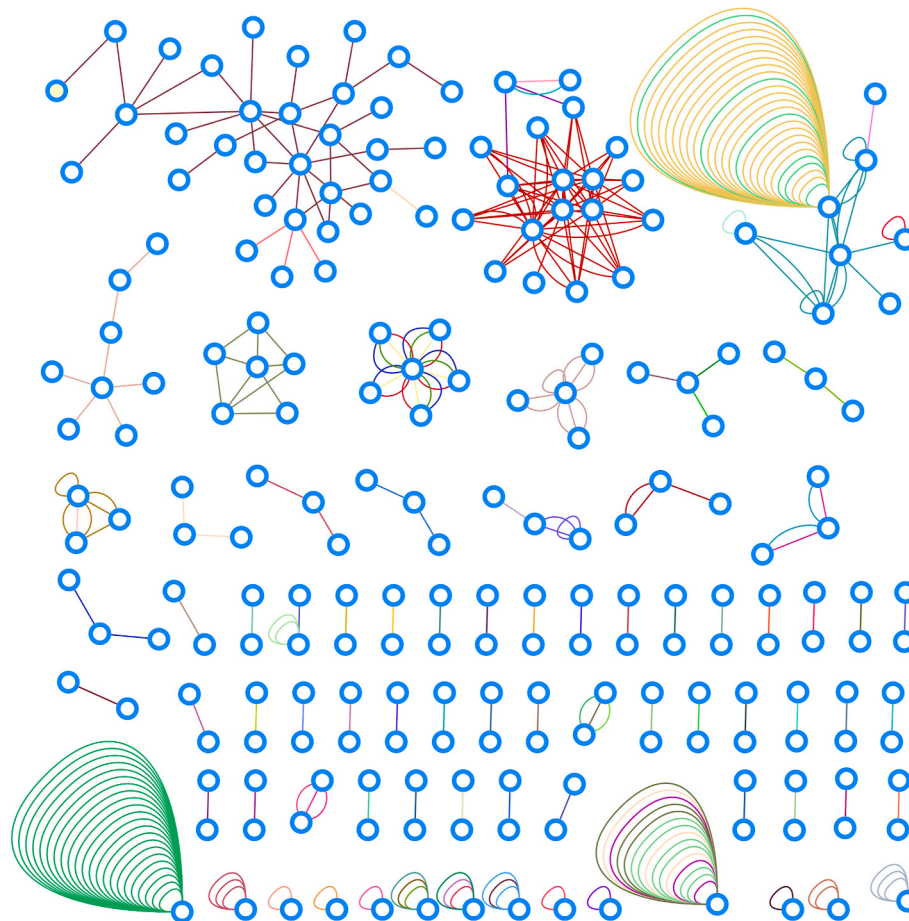


**Fig. 2. Network representation of the exon-sharing gene clusters and the corresponding exon modules.** Each node represents a lncRNA gene and each edge an exonic module shared between two genes. Same color edges within a gene cluster represent a module. Self-loops represent instances where the same module occurs multiple times in a single gene. The network representation was generated using Cytoscape [42].
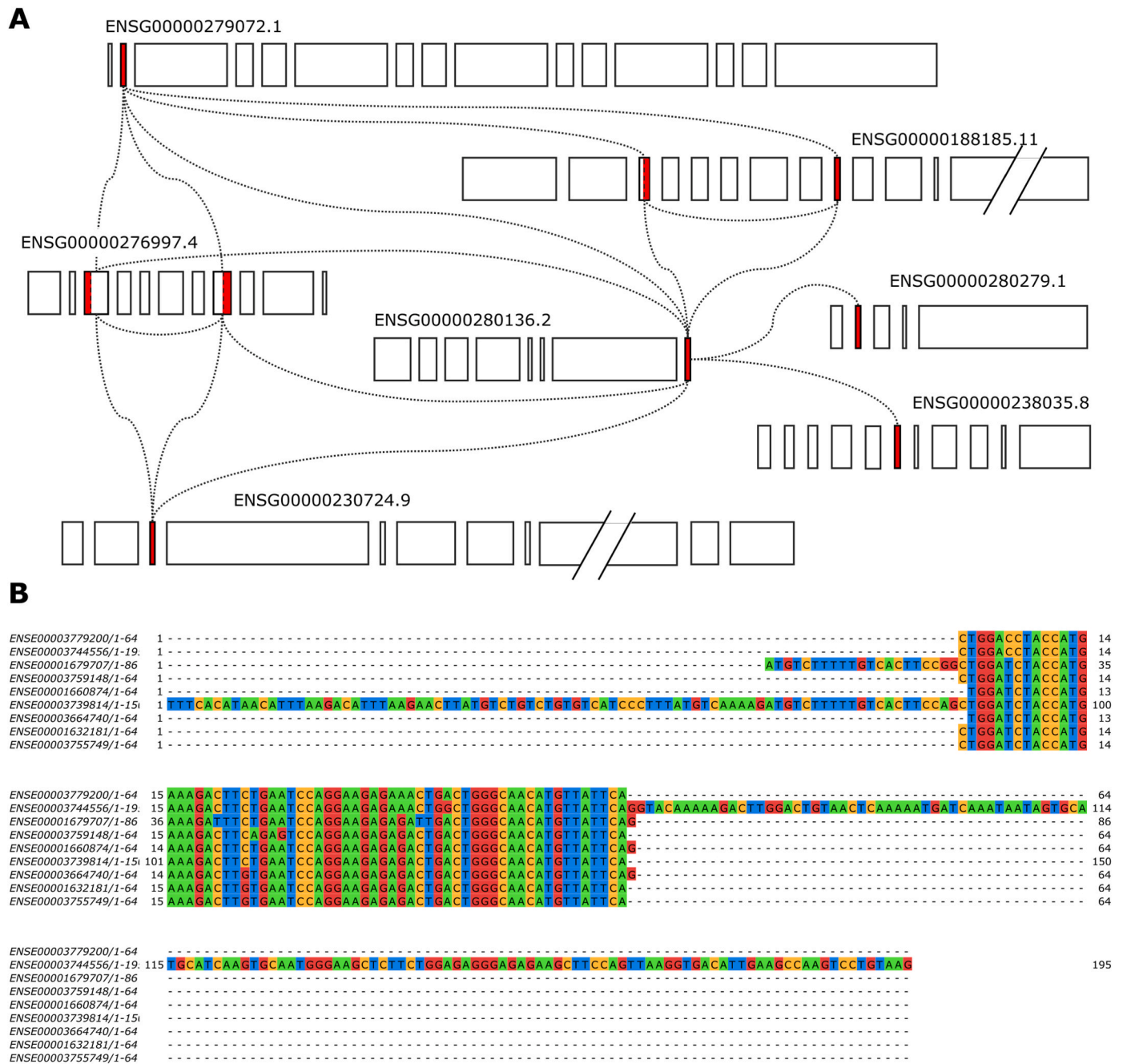
**Fig. 3. An example of the identified exon modules.** A) Schematic representation of 7 genes containing representatives (in red) of exons contributing to a module cluster. Each box represents an exon, with width proportional to its length (intron length not to scale); B) multiple alignment of the 9 exons contributing to the cluster.

with high sequence similarity (external gap trimmed sequence identity 92–98 %, Fig. 3B) embedded in different genes. It is worth noting that, in some cases, the module constitutes an exon on its own, whereas in other cases it is part of a larger exon.

We then analyzed in more detail the sequence context of exon modules. More specifically, we looked at the sequence similarity of additional exons flanking the modules, to rule out the possibility that the similarity between modules in different genes simply reflects global sequence similarity between the exonic components of genes (see Materials and Methods and Fig. 4A). The alignment scores for exons flanking the putative module in the same gene, upstream and downstream (Fig. 4B), showed that the similarity between exon modules is significantly higher than that of the sequence context in which they are

embedded. We also observed a small proportion of cases in which the flanking exons are also similar (outliers in Fig. 4B). These cases fall outside the criteria used to define exon modules: in 17 cases because they are less than 50 nucleotides in length, and in another 17 cases because they contain repetitive sequences.

We then analyzed the sequence similarity of the intronic sequences flanking the exon modules. To this end, we defined genomic regions of interest by extending upstream and downstream the sequence of each candidate exon pair, until we obtained two sequences with a length equal to three times that of the longest exon of the pair (Fig. 4C-E). We limited the analysis to pairs with sequence similarity above the z-score threshold of 6.2 and excluded modules repeating within the same gene. For each pair of genomic regions of interest, we performed a global
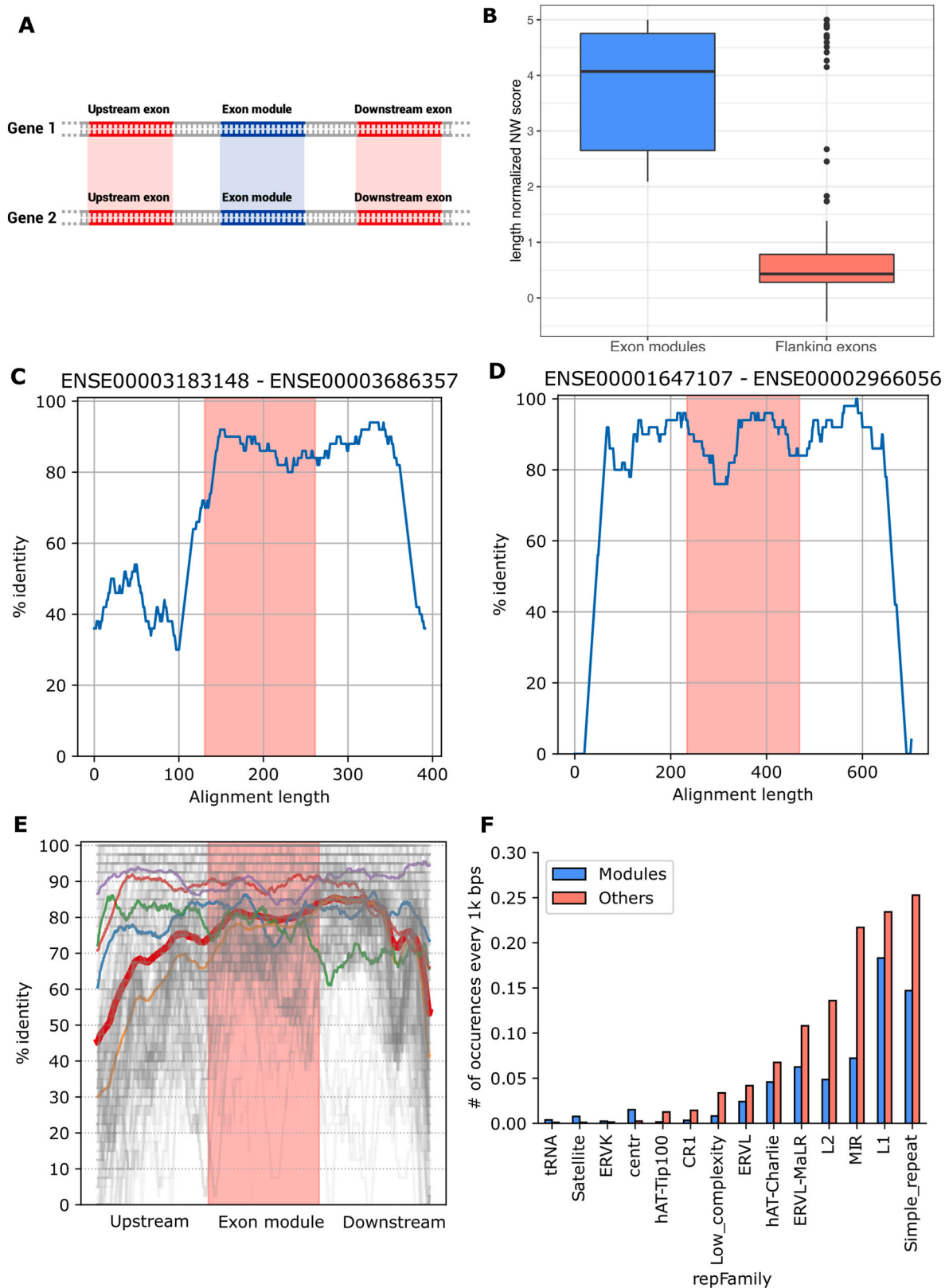
**Fig. 4. Analysis of the sequence regions flanking exon modules.** A) For each pair of genes containing a shared exon module we compared the similarities of the upstream and downstream flanking exons (when present); B) Distributions of the length-normalized Needleman and Wunsch scores of exonic modules (in blue) and of their upstream and downstream flanking exons (in red); C) A pair of exons in which the similarity only extends to the downstream flanking intron; D) A pair of exons in which the similarity extends upstream and downstream into both flanking introns; E) Overall representation of all the length-scaled similarities between all the exon pairs and their flanking introns (in gray), the median identity percentage is represented in red. The other colored lines represent five clusters of similarity patterns as defined by grouping individual lines; F) Number of occurrences per thousand base pairs of families of repetitive sequences in flanking introns with significant differences (padj<0.05) between the exonic modules and the other lncRNA exons.

alignment using the same parameters used to identify the exon modules, and calculated the percentage identity of the pairs using overlapping windows of 50 nucleotides with a single nucleotide shift, to generate graphs depicting the extent of the similarity. We found that, in the majority of instances, sequence similarity extends into the flanking intronic regions. More specifically, in approximately one third of the cases, the similarity encompassed both the upstream and downstream intron, in another third of the cases the similarity extended to a single intron, while the remainder of cases lacked a clear pattern. We did not observe any cases where the similarity was confined to the boundaries of the candidate exon modules.

The extension of the similarity through the flanking introns suggests that the most common mechanism responsible for the origin of exon modules is segmental duplication of a genomic DNA stretch encompassing the parental copy of an exon. This is the same mechanism suggested as a driver of exon shuffling in protein-coding genes [43]. To further confirm these findings, we compared our results with the data present in the UCSC Segmental Dups track (genomicSuperDups) which contains regions detected as putative genomic duplications within the human genome. These regions represent large recent duplications (≥1 kb and ≥90 % identity) that originated over the last ~40 million years along the human lineage, based on neutral expectation of divergence [43]. For 84 of the 340 lncRNA exon pairs identified here, we found a match in the segmental duplications identified by Bailey et al. [43], in 81 of these cases the duplicated stretch includes the entire exons of the pair, while in 3 cases the duplication is interrupted within the exon. We also observed a higher frequency of pairs located on the same chromosome (~20.5 %) compared with what is observed when the same exons are randomly paired (~3.6 %). Moreover, pairs of exon modules that are on the same chromosome are closer together when compared to the same random pairing control (Mann-Whitney p-value = 9.86e-05). A higher rate of occurrence on the same chromosome has been described for segmental duplications [44]. To further extend the analysis of

flanking regions, we compared the rate of occurrence of multiple families of repetitive elements in the introns flanking candidate exonic modules vs other lncRNA exons (for exons located at the ends of a gene, we included a region of 10k bps in the genome). We calculated the number of occurrences per 1000 base pairs of each family of repetitive elements on the set of regions flanking the exon modules vs the other lncRNA exons (Fig. 4F and Table S2) thus obtaining a distribution of occurrences where the observations correspond to the individual sequence regions. We then compared these distributions using a Mann-Whitney *U* test, with Bonferroni correction for multiple hypothesis testing. We observed significant differences for 15 of 46 families (padj<0.05). Interestingly, centromere and satellite repeats are among the few classes of repeats enriched in regions flanking the exon modules, while most classes of transposon- or endogenous retrovirus-derived repeats are depleted. Since the genomic regions proximal to centromeres and telomeres are enriched with segmental duplications [45], this observation further points at segmental duplication as the main driver of the appearance of these exon modules, as opposed to, for instance, transposition. The enrichment of this type of repetitive sequences can be explained by the localization near the centromeres or telomeres of a portion of the modules (Fig. S2). Moreover, searching for transposase domains using a procedure similar to the one described by Koch [46] did not reveal significant differences in their occurrence among genes containing exon modules (data not shown), further highlighting that transposition is not the main driver of this process.

To further investigate the characteristics of these modules we looked at the distribution of cis-regulatory elements (CRE) within their sequences (Fig. 5A). This analysis highlighted a depletion in exon modules of the most frequent CREs (Fisher's exact test padj = 6.2e10-3, 1.2e10-2, 1.2e10-2 for pELS, dELS and PLS, CTCF-bound respectively). Non-depleted elements include H3K4Me1, which represent markers for enhancer lncRNAs [47], and H3K4me3 marks, which are characteristic of transcriptionally active regions (Howe et al., 2017). Interestingly this
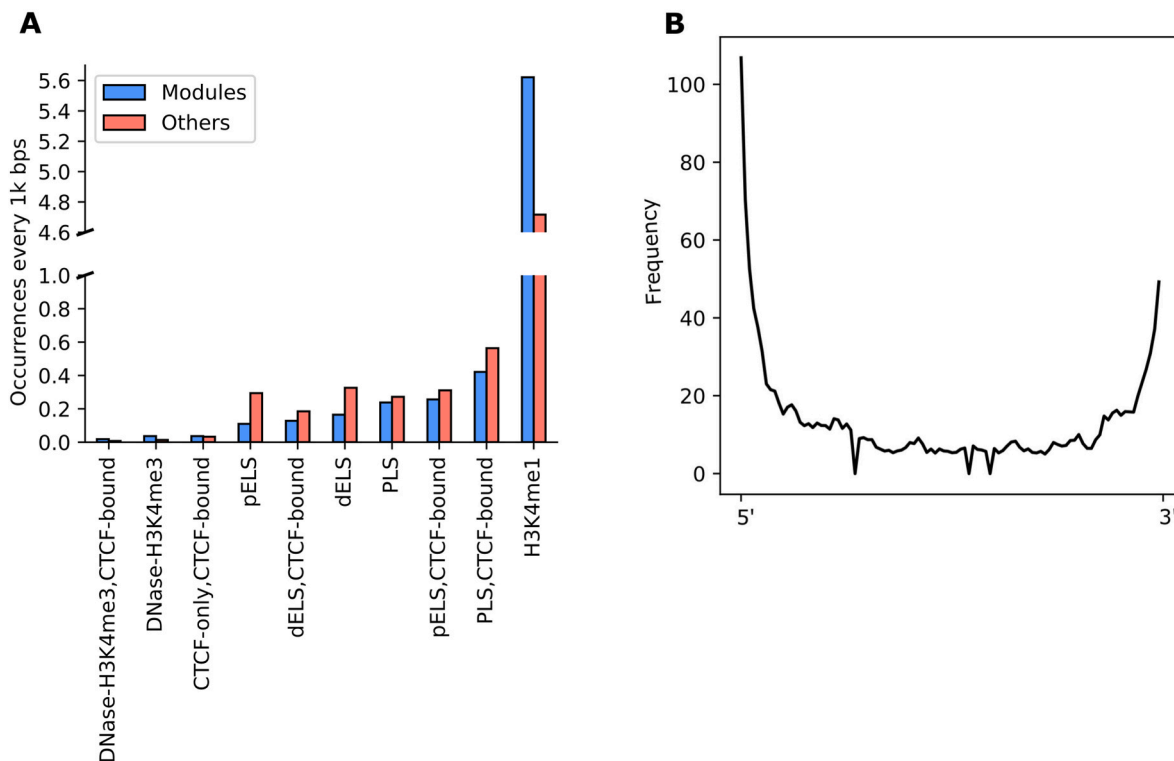


**Fig. 5. Cis-regulatory elements (CRE) and position of the modules.** A) number of occurrences of the different CREs from the annotation present in ENCODE every thousand nucleotides in the modules (in blue) and in the other lncRNA exons of the dataset (in red); B) the y axis indicates the frequency 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.

histone modification is usually found in the region corresponding to the beginning of the transcript [48]. Accordingly, when we investigated the position of the exonic modules within their transcripts (Fig. 5B), we detected a higher frequency of the modules at the 5′ end. This finding is consistent with what is observed in protein-coding genes, which in vertebrates tend to increase their length over time by gaining recently evolved domains, primarily through the addition of sequences at the 5′ end of genes [49]. The insertion of these modules at the extremities of the transcript presumably allows the addition of genetic material with minimal disruption to the existing sequence.

### 3.2. 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 (Mann-Whitney p-value = 6.3e-5), and both the conservation score of exon modules and of exons belonging to functionally annotated lncRNA were significantly higher than the conservation score of the remaining lncRNA exons (Mann-Whitney p-value = 7.4e-27 and 3.5e-26 respectively, Fig. 6A). When looking at the conservation of exon modules in four higher primate species, we also observed a greater proportion of exons with a BLAST hit among exon modules vs the remaining exons. More specifically, 65.97 % of the exon modules have a BLAST hit in Chimpanzee, 42.01 % in Bonobo, 11.83 % in Gorilla and 47.04 % in Orangutan. Conversely, only 43.23 %, 16.72 %, 2.61 %, 25.86 % of the control exons (i.e the portion of the 12,097 lncRNA exons that have no repetitive and non-overlapping sequences and that are not modules) have BLAST hits on the same species, respectively (Fig. 6B) To evaluate the significance of these results we performed a Fisher's exact test on the aggregated data from the different species, which confirmed that these results are significant (p-value = 4.10e-15).

Since the BLAST similarity score with non-human primates does not take into account the genomic position of exons in different organisms, i. e. it cannot distinguish between the similarity of true orthologs vs in- and out-paralogs, we investigated whether exon modules are located in regions of synteny between non-human primates more often than other exons. To this end we leveraged the SynthDB [36] database, which provides data on orthology relationships between humans and other primates. We observed that the percentage of genes located in a syntenic region is higher for genes that contain at least one exon module, compared with those which do not. Accordingly, 14.50 % of the exon modules are located in genes that have an ortholog in Chimpanzee, 16.57 % in Bonobo, 17.46 % in Gorilla and 15.98 % in Orangutan. While for the other lncRNA exons we observed percentages of 10.79, 4.74, 12.39, 6.55 in the same species respectively. We then performed a Fisher's exact test comparing exons modules that belong to genes with an ortholog in at least one of the species mentioned above to the other exons which confirmed the significance of our results (p-value = 5.63e-05) (Fig. 6C).

To strengthen the evolutionary conservation analysis, and to compare our results with the analysis by Sarropoulos et al., 2019 [50], we extended it by including additional species. To this end, we aligned all lncRNA exons using blastn against the genomes of the organisms used in that work (Macaque, taxid: 9544; Rabbit, taxid: 9986; Chicken, taxid: 9031; Opossum, taxid: 13,616; Rat, taxid: 10,116; Mouse, taxid: 10, 090), and other model organisms (Danio rerio, taxid: 7955; Drosophila melanogaster, taxid: 7227; Caenorhabditis elegans, taxid: 6239; Arabidopsis thaliana, taxid: 3702), using an e-value threshold of 0.01 to identify hits (Fig. 6D–E and Fig. S3). Figure 6E (Fig. 6E) displays the percentage of exonic modules vs other lncRNA exons that have at least one hit in the species indicated above. This analysis shows a rapid decay in the number of similar exons as the evolutionary distance from humans increases. Figure 6F (Fig. 6F) shows the 30 mammal PhastCons scores of the exon modules, as a function of the z-score similarity threshold used to define the modules themselves (i.e. the threshold described in Fig. 1A). This analysis demonstrates that the exon modules identified in this work, which are highly similar as they were selected on the basis of having a Z-score of at least 6.2 and 5.3 in the sequence and structure alignment respectively, represent duplications that are recent (as implied by the high levels of sequence similarity) and that are exclusively found in humans and higher primates, and thus have lower PhastCons scores on the entire set of 30 mammals (Fig. 6F).

Overall, the above results reveal that roughly 4 % of lncRNA genes (218 lncRNA genes/5423 total lncRNAs genes which contain at least one exon without repetitive sequences, see Materials and Methods) include one or more exons having significant similarity with exonic portions of other lncRNAs. To our knowledge, this represents the first draft of a genome-wide catalog of shared lncRNA exons.

### 3.3. Nucleotide variation in modules

To further investigate whether exon modules may represent conserved functional units, we analyzed the occurrence and frequency of single nucleotide polymorphisms (SNPs) in these regions, as a lower incidence of variants may indicate the existence of constraints associated with functional sequences, due to the effects of purifying selection [51]. Accordingly, we collected SNP data from the 1000 Genome project from dbSNP 153 [37] and we observed 12.87 variants per thousand bases in control exons (which are not modules) and 11.83 in modules. We then obtained from the ALFA allele frequencies aggregator [38] a total of 764,005 SNPs located in lncRNA exons [38], and their associated frequencies. For each exon, we calculated the index of nucleotide diversity $\theta\pi$ [52] as

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

where $f_i$ represents the frequency of variants in the $i$ th position of the exon sequence in the population, and $l$ represents the length of the exon.

After comparing the distributions of $\theta\pi$ scores with the Mann-Whitney $U$ test, we obtained a p-value of 2.14e-02 in the comparison between modules and exons from functionally annotated genes, a p-value of 2.77e-02 from the comparison between exon modules and other lncRNA exons and a non-significant p-value (7.45e-01) from the comparison between functionally annotated and others, confirming a significant lower propensity to harbor variation in exon modules as compared to the other two groups. These findings indicate the existence of evolutionary constraints which limit the occurrence of variants with polymorphic frequencies in exon modules, which in turn may reduce the rate of evolutionary change in the long-term. We also looked at the frequency of polymorphic complete exon deletions, but the results were not statistically significant (data not shown).

### 3.4. Search for characteristics shared with protein-coding genes

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

França et al. [53] observed that symmetric shuffling units (exons whose length is an exact multiple of three) are strongly over-represented in human protein-coding genes, due to their lower impact on the reading frame when transposed. We found an opposite trend in lncRNA exon modules, with only 25 % having a length that is a multiple of three,
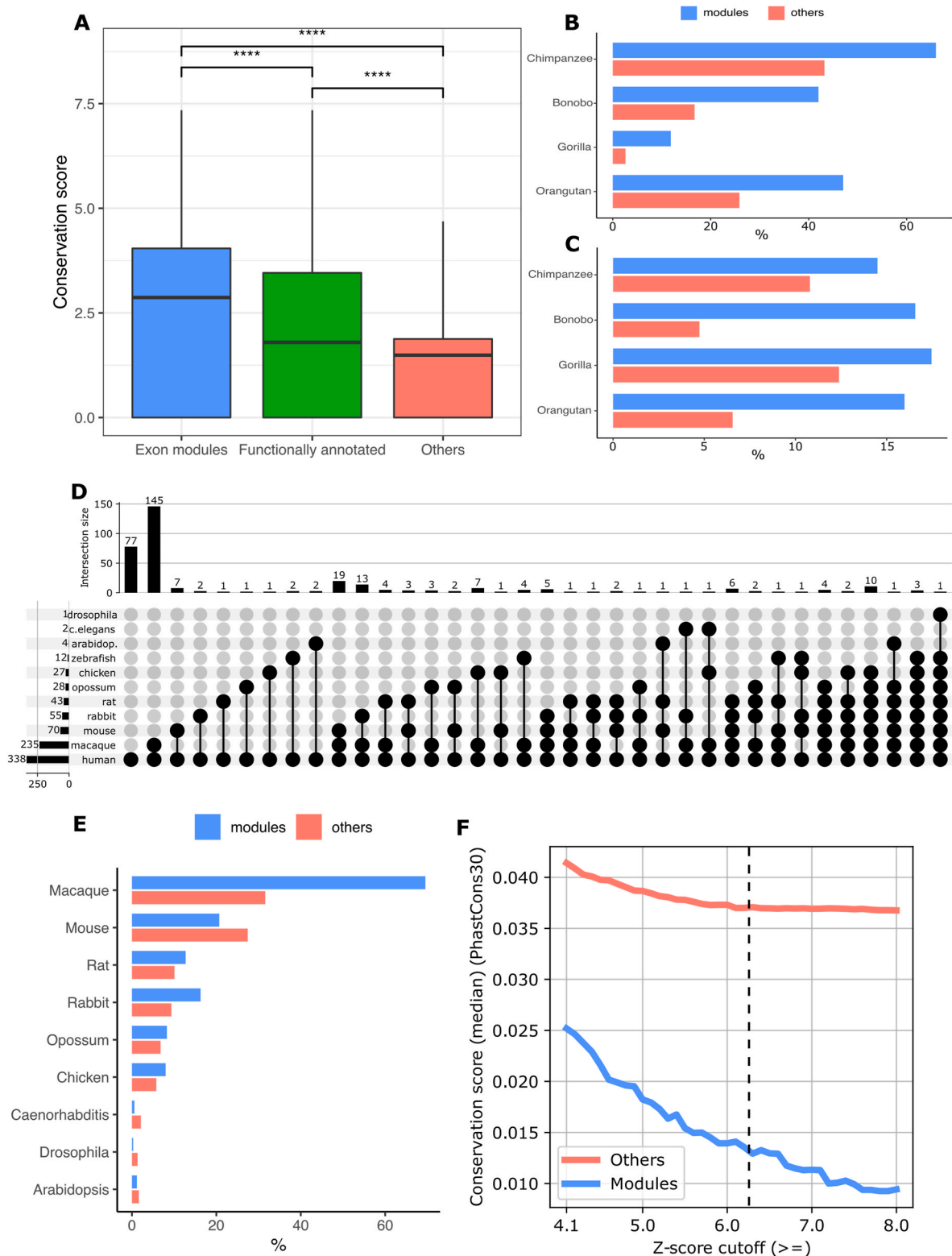
**Fig. 6. Evolutionary conservation of exon modules.** A) Box-plot of the conservation scores in four non-human primates for exon modules, functionally annotated exons from the lnc2Cancer database, and controls; B) Percentage of exon modules (in blue) and other exons (in red) that showed a BLAST hit (e-value <0.001) in the primate species considered; C) Percentage of genes showing a conserved syntenic region (as defined in SynthDB) among those containing exon modules (in blue) vs genes not containing an exon module (in red); D) Upset plot representing the exons that have a BLAST hit in the species analyzed in Sarropoulos et al. [50] and in other model organisms; E) Percentages of modules (in blue) and other exons (in red) showing a BLAST hit in the indicated species F) PhastCons 30 mammals scores of members of clusters defined by different z-score thresholds of pairwise similarity from sequence alignments (in blue) and the other lncRNA exons of the dataset (in red).

which confirms the lack of relevance of the reading frame. By contrast, in the remainder of the exons, this proportion is 33 %, i.e. what would be expected under a random model.

The transition/transversion ratio (Ti/Tv) among polymorphic variants should be 0.5 under a purely random model, resulting from four possible transitions/eight possible transversions. However, real data depart remarkably from this expectation, with functional regions and protein-coding regions presenting values higher than 0.5, since transitions are more likely to result in non-synonymous substitutions (e.g. when they occur in the third base of a codon) [54]. Exon modules displayed values of 1.9, in line with previous results for lncRNAs [55]. As a reference, these values contrast sharply with those for protein-coding genes, which range between 2.8 and $2.9 \pm 0.1$ [55].

To further explore how these exons differ from protein-coding sequences, we analyzed their coding potential using the Coding Potential Calculator 2.0 (CPC2) tool [56]. Among the module-containing exons, only 2 out of 338 were identified as having coding potential. We conducted a parallel analysis on all the remaining 11,759 lncRNA exons without modules. This comparison yielded a similar result, with only 49 exons being designated as "coding". A comparison of these frequencies using Fisher's exact test was not significant (p-value = 0.08).

We also compared the coding probability score distributions obtained with the same software, using a Mann-Whitney test, and did not identify significant differences (p-value >0.56).

### 3.5. Analysis of the characteristics of exon modules that are not conserved in primates

We broadened our study to compare exon modules that are not conserved in primates against other exon modules, aiming to highlight their differences. After evaluating the conservation scores for all exon modules in primates (detailed in the Materials and Methods section), we identified 86 exon modules with a conservation score of zero, indicating a lack of conservation. Proximal enhancer-like sequences (pELS) where the only cis-regulatory elements (CREs) with a different abundance between conserved and non-conserved exons (more frequent in non-conserved exons, p-value = 0.038, Fisher's exact test, Fig. S4).

We also assessed the coding potential of these non-conserved exon modules compared to their conserved counterparts using the Coding Potential Calculator 2.0 (CPC2) tool [56], The results showed slightly significantly higher coding probability scores in non-conserved exons, as indicated by a Mann-Whitney p-value of 0.001. However, only two out of 250 conserved exons were classified as "coding," and none of the 86 non-conserved exons were, suggesting that both groups generally have low coding potential.

Furthermore, the ratio of transitions to transversions, calculated in section 3.4, was similar for both groups (1.9 for conserved exon modules and 2 for non-conserved exon modules), which contrasts with the typical range of $2.8–2.9 \pm 0.1$ observed in protein-coding genes. We also noted a slight increase in exons whose length is a multiple of 3 in non-conserved modules (29 %) compared to conserved ones (24 %). However, these percentages are still lower than those seen in other long non-coding RNA exons that do not contain modules (33 %).

### 3.6. Functional hypothesis and organization of putative modules in clusters of lncRNA genes

To further describe exon modules, here we show some examples of their organization within the structure of their lncRNA genes. Only 12 of 218 genes containing exon modules are associated with a known biological function in the lnc2Cancer database [5]. For most of them, the specific region of the lncRNA molecule responsible for that function is unknown. In the next two paragraphs we will provide a more detailed description for two of the identified modules, in an attempt to capture their putative functions. The first example refers to an exon module recognized by virtue of sequence similarity, and the second one refers to

an exon module recognized by virtue of structure similarity.

### 3.7. Identification of a putative YBX1 binding module

Fig. 7A shows an example of a putative module represented in a pair of exons as a sequence of ~200 nucleotides sharing a high sequence similarity (>87 %). The exons involved are ENSE00003710224.1 and ENSE00003838358.1 which belong to genes ENSG00000182165.17 (also known as *TP53TG1*) and ENSG00000285540.1, respectively. *TP53TG1* is a lncRNA involved in the p53 network response to DNA damage [57], which has a role as tumor suppressor by blocking the tumorigenic activity of the RNA binding protein (RBP) YBX1 [58]. More in detail, the expression of *TP53TG1* is induced by p53 under cellular stress conditions that involve the induction of double-strand breaks [57], while the interaction in the cytoplasm between *TP53TG1* and YBX1 prevents the migration of the latter inside the nucleus where it might promote the transcription of a series of oncogenes [59]. *Diaz-Lagares* et al. [58] demonstrated that a central region of *TP53TG1*, which includes the putative module in the exon ENSE00003710224.1, is responsible for YBX1 binding. Moreover, they proved that YBX1 binding motifs CACC are necessary to ensure the tumor-suppressor function of *TP53TG1*. We identified two occurrences of the CACC motif in ENSE00003710224.1 and one in ENSE00003838358.1, suggesting a common role for this module.

### 3.8. Identification of a putative LIN28B binding module

Fig. 7B–D shows an example of a module with high structure similarity, embedded in dissimilar sequence contexts. The exons involved are ENSE00003741285.1, ENSE00001800736.1 and ENSE00001782399.1 which belong to ENSG00000278214.1, ENSG00000224610.1 and ENSG00000229249.6, respectively (Fig. 7B). These three exons fold into a similar secondary structure, composed of two stems ending with a hairpin loop, with one of the two stems having one or two internal loops.

To detect a possible function, common to the three representatives of this exon module, we searched for the presence of enriched structure and sequence motifs using the BRIO web server (see Materials and Methods). BRIO identified a significantly enriched (Fisher's exact test padj<0.05) structure motif shared between all the exons of the group (Fig. 7C). This particular motif was associated by Adinolfi et al. [41] with a series of different RNAs capable of binding some RBPs including LIN28B. This is an evolutionary conserved RBP involved in several cellular processes, which acts as a critical oncogene activated in cancer [62]. LIN28B is known to be able to bind different mRNAs, including a set of mRNAs for splicing factors [63], miRNAs [64] and lncRNAs such as *NEAT1* [65]. Furthermore, LIN28B C-terminal zinc knuckle (ZnK) mediates specific binding to a conserved GGAG motif [66] which is also a sequence motif present in all the three representatives of this module (Fig. 7D). These observations suggest a possible role of this module in binding LIN28B. To validate the interaction between these RNA binding proteins and the identified exons, we also searched whether these pairs were present in the data obtained from three public CLIP-seq datasets: CLIPdb [67], eCLIP [68], Starbase [69]. Accordingly, we found experimental evidence for the interaction between the lncRNA encoded by ENSG00000224610 and LIN28B in the PAR-CLIP experiment GSM1087851 [70].

## 4. Discussion and conclusions

This work identified a set of lncRNA exons with high sequence and/ or structure similarity that are embedded within globally dissimilar genes, confirming the hypothesis of exon sharing between this class of molecules, similarly to protein-coding genes. This set contains a total of 340 pairs of exons that can be grouped, on the basis of their reciprocal connections, in 106 clusters. In contrast to previous work [18], our
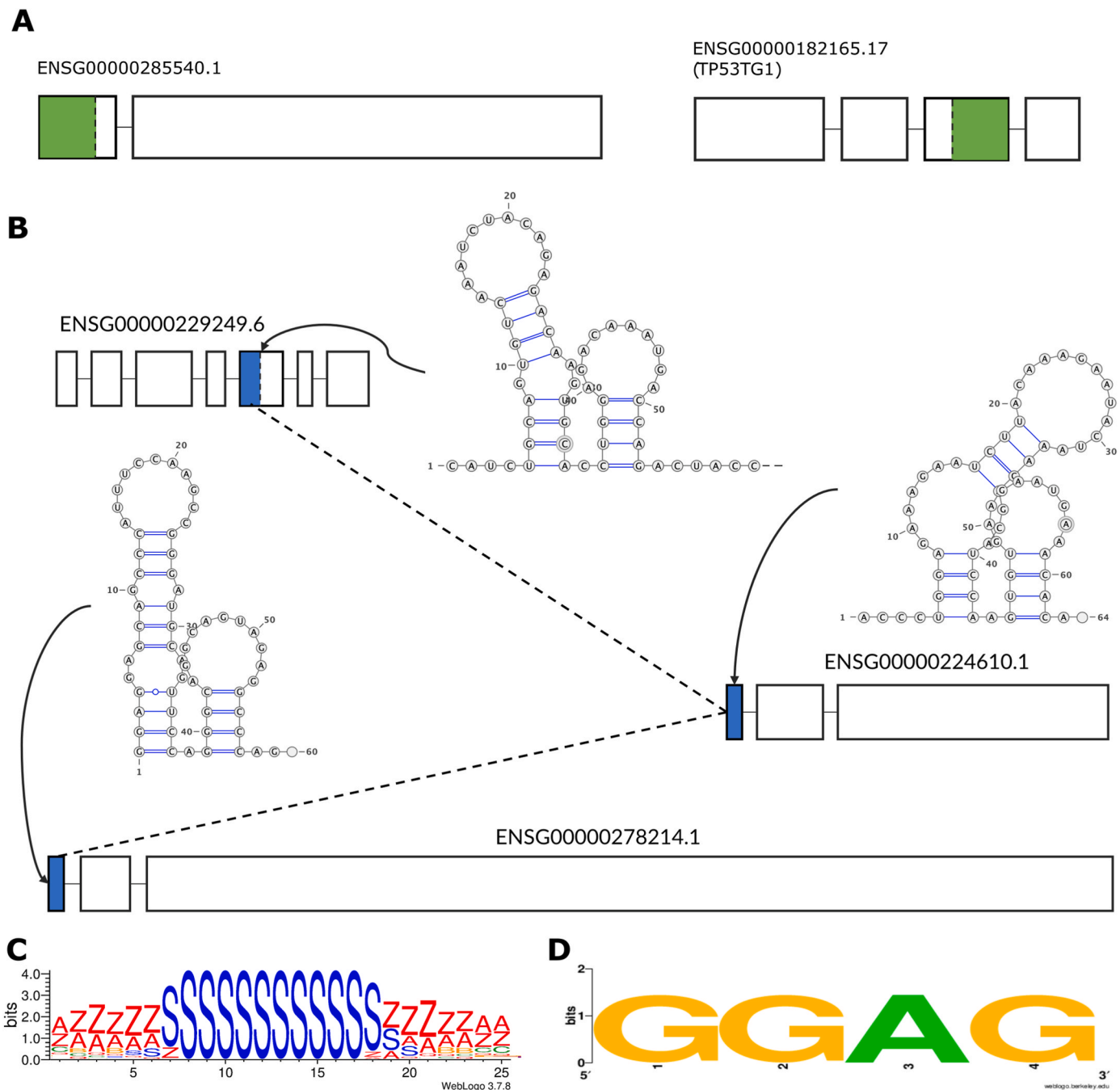
**Fig. 7. Organization of a sequence and a structure module and identified motifs.** A) Schematic representation of the lncRNA genes containing the putative YBX1 binding module (in green); B) Representation of the lncRNA genes containing the exons with the putative LIN28B binding module and their secondary structures. The blue boxes represent the exons with high structural similarity that form the module; C) secondary structure motif revealed by BRIO represented with the BEAR alphabet [30]; D) sequence motif recognized by ZnK in the three modules. The RNA secondary structure representations were generated using VARNA [60]; Sequence and structure logos were generated using WebLogo [61].

analysis focused on exons that do not contain repetitive sequences. The resulting dataset of exon modules likely represents the result of recent segmental duplications that are almost exclusively found in humans and higher primates. These findings support the hypothesis that the non-coding transcriptome is structured into modular domains, similar to the organization observed in protein-coding genes.

Approximately 4 % (218 out of 5423) of all the lncRNA genes in our dataset contain an exon module. Despite the different features preventing a direct comparison, it is interesting to note that this figure is comparable to the 6.4 % of protein coding genes with evidence of exon shuffling events in *H. sapiens* according to França et al. [53].

Even though we cannot assign a specific function to each of these modules, as it has been done for the majority of protein-coding domains, it is tempting to infer that sharing of functional modules between different lncRNAs may contribute to expanding the functional repertoire of the non-coding genome, similar to the shuffling of functional exons in coding sequences [12].

LncRNA exon modules identified in this work display a higher degree of sequence conservation and synteny in four primate great ape species than the remainder of lncRNA exons. A high level of conservation between related species is suggestive of purifying selection and is a landmark characteristic of functional genetic elements [71]. Exon modules

also harbor a lower frequency of SNPs compared with control sequences, which suggests that purifying selection also persists intra-specifically in human populations. Our set included 46 exon pairs highly similar in both sequence and structure (Fig. 1F), which are associated with the highest conservation scores. Even though we cannot infer the age of the duplication/shuffling event based on our analysis, our results show that the exons involved are subjected to extreme purifying selection, which preserved both sequence and structure. Taken together, this evidence suggests that these modules play an important role within their respective lncRNA genes, even though their exact function is yet to be characterized.

Because the Z-score thresholds used for module identification were based on the identification of an inflection point in the conservation scores of the sequence regions classified as modules (see Fig. 1D), our analysis is biased towards the identification of modules that are both similar within the human genome and also conserved across other genomes. As such, our analysis may underestimate examples of accelerated divergence, similar to the ones discussed in the section about exon modules that are not conserved in primates.

Finally, it was reported that homologous lncRNAs can, in some cases, conserve their function over long evolutionary times, despite having diverged in both their nucleotide sequences and their secondary structures [72]. The above considerations suggest that our analysis may underestimate the extent of module sharing in lncRNAs. Other limitations include the fact that the correct identification of exons within lncRNAs is strongly dependent on the reliability of the reconstruction of the whole transcript structure. This is usually summarized by the TSL parameter (Transcript Support Level) which we included, for every exon, in the Supporting information (Table S1). To evaluate the accuracy of the transcript annotations, we have additionally included in the table details regarding the confidence level provided by GENCODE, the potential association of the exon with a transcript marked as Ensembl Canonical, and whether these transcripts are part of the GENCODE basic set.

In the few cases for which functional information on a lncRNA is available, it may be possible to infer the function of the shared module. We report two examples of modules conserved in either sequence or structure. In both cases, the ability to bind specific targets is the inferred associated function.

Overall, our results highlight the presence of groups of exons sharing high sequence or structure similarity within dissimilar lncRNA genes. These exons are highly conserved across primate species and depleted of inter-individual variation among humans (SNPs), and we suggest that they may represent functional modules.

The identification of these modules could constitute a tool for decoding the function of the many lncRNAs that are currently uncharacterized. Membership in a shared exon cluster represents a feature that deserves annotation, even though conclusive proof of shared function will require experimental evidence.

## Funding

## Data availability

The annotation was obtained from GENCODE v29 (https://www.gencodegenes.org/human/release_29.html).

The hg38 version of the human genome was downloaded from UCSC genome browser (http://hgdownload.cse.ucsc.edu/goldenPath/hg38/bigZips/).

The BEAGLE web-server for RNA structure alignments is available at: http://beagle.bio.uniroma2.it.

Functionally annotated lncRNA was downloaded from: http://bio-bi gdata.hrbmu.edu.cn/lnc2cancer.

Variant frequencies in human populations are available in the ncbi website (https://www.ncbi.nlm.nih.gov/snp/docs/gsr/alfa/#ftp-download).

The BRIO web-server for RNA interaction motif search is available at: http://brio.bio.uniroma2.it.

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

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

## CRediT authorship contribution statement

**Francesco Ballesio:** Writing – original draft, Software, Investigation, Data curation, Conceptualization. **Gerardo Pepe:** Writing – review & editing, Conceptualization. **Gabriele Ausiello:** Writing – review & editing, Supervision, Conceptualization. **Andrea Novelletto:** Writing – review & editing, Supervision, Conceptualization. **Manuela Helmer-Citterich:** Writing – review & editing, Supervision, Funding acquisition, Conceptualization. **Pier Federico Gherardini:** Writing – review & editing, Supervision, Funding acquisition, Conceptualization.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Appendix A. Supplementary data

Supplementary data to this article can be found online at https://doi.org/10.1016/j.ncrna.2024.06.013.

## References

[1] W. Gilbert, Why Genes in Pieces? Nature Publishing Group UK, 1978 https://doi.org/10.1038/271501a0.
[2] J.M. Engreitz, J.E. Haines, E.M. Perez, G. Munson, J. Chen, M. Kane, P.E. McDonel, M. Guttman, E.S. Lander, Local regulation of gene expression by lncRNA promoters, transcription and splicing, Nature 539 (2016) 452–455.
[3] A. Frankish, M. Diekhans, I. Jungreis, J. Lagarde, J.E. Loveland, J.M. Mudge, C. Sisu, J.C. Wright, J. Armstrong, I. Barnes, A. Berry, A. Bignell, C. Boix, S. Carbonell Sala, F. Cunningham, T. Di Domenico, S. Donaldson, I.T. Fiddes, C. García Girón, J.M. Gonzalez, T. Grego, M. Hardy, T. Hourlier, K.L. Howe, T. Hunt, O.G. Izuogu, R. Johnson, F.J. Martin, L. Martínez, S. Mohanan, P. Muir, F. C.P. Navarro, A. Parker, B. Pei, F. Pozo, F.C. Riera, M. Ruffier, B.M. Schmitt, E. Stapleton, M.-M. Suner, I. Sycheva, B. Uszczynska-Ratajczak, M.Y. Wolf, J. Xu, Y. T. Yang, A. Yates, D. Zerbino, Y. Zhang, J.S. Choudhary, M. Gerstein, R. Guigó, T.J. P. Hubbard, M. Kellis, B. Paten, M.L. Tress, P. Flicek, Gencode 2021, Nucleic Acids Res. 49 (2021) D916–D923.
[4] X.C. Quek, D.W. Thomson, J.L.V. Maag, N. Bartonicek, B. Signal, M.B. Clark, B. S. Gloss, M.E. Dinger, lncRNAdb v2.0: expanding the reference database for functional long noncoding RNAs, Nucleic Acids Res. 43 (2015) D168–D173.
[5] Y. Gao, S. Shang, S. Guo, X. Li, H. Zhou, H. Liu, Y. Sun, J. Wang, P. Wang, H. Zhi, X. Li, S. Ning, Y. Zhang, Lnc2Cancer 3.0: an updated resource for experimentally supported lncRNA/circRNA cancer associations and web tools based on RNA-seq and scRNA-seq data, Nucleic Acids Res. 49 (2021) D1251–D1258.
[6] V. Fort, G. Khelifi, S.M.I. Hussein, Long non-coding RNAs and transposable elements: a functional relationship, Biochim. Biophys. Acta Mol. Cell Res. 1868 (2021) 118837.
[7] C.P. Ponting, P.L. Oliver, W. Reik, Evolution and functions of long noncoding RNAs, Cell 136 (2009) 629–641.
[8] L. Statello, C.-J. Guo, L.-L. Chen, M. Huarte, Gene regulation by long non-coding RNAs and its biological functions, Nat. Rev. Mol. Cell Biol. 22 (2021) 96–118.
[9] R.A. Gupta, N. Shah, K.C. Wang, J. Kim, H.M. Horlings, D.J. Wong, M.-C. Tsai, T. Hung, P. Argani, J.L. Rinn, Y. Wang, P. Brzoska, B. Kong, R. Li, R.B. West, M. J. van de Vijver, S. Sukumar, H.Y. Chang, Long non-coding RNA HOTAIR reprograms chromatin state to promote cancer metastasis, Nature 464 (2010) 1071–1076.
[10] S.-Y. Ng, R. Johnson, L.W. Stanton, Human long non-coding RNAs promote pluripotency and neuronal differentiation by association with chromatin modifiers and transcription factors, EMBO J. 31 (2012) 522–533.
[11] M.A. Faghihi, F. Modarresi, A.M. Khalil, D.E. Wood, B.G. Sahagan, T.E. Morgan, C. E. Finch, G. St Laurent, P.J. Kenny, C. Wahlestedt, Expression of a noncoding RNA is elevated in Alzheimer's disease and drives rapid feed-forward regulation of beta-secretase, Nat. Med. 14 (2008) 723–730.
[12] M. Guttman, J.L. Rinn, Modular regulatory principles of large non-coding RNAs, Nature 482 (2012) 339–346.

[13] I. Ulitsky, D.P. Bartel, lincRNAs: genomics, evolution, and mechanisms, Cell 154 (2013) 26–46.

[14] P. Johnsson, L. Lipovich, D. Grandér, K.V. Morris, Evolutionary conservation of long non-coding RNAs; sequence, structure, function, Biochim. Biophys. Acta 1840 (2014) 1063–1071.

[15] I. Ulitsky, A. Shkumatava, C.H. Jan, H. Sive, D.P. Bartel, Conserved function of lincRNAs in vertebrate embryonic development despite rapid sequence evolution, Cell 147 (2011) 1537–1550.

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

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

[18] R. Johnson, R. Guigó, The RIDL hypothesis: transposable elements as functional domains of long noncoding RNAs, RNA 20 (2014) 959–976.

[19] A. Kapusta, Z. Kronenberg, V.J. Lynch, X. Zhuo, L. Ramsay, G. Bourque, M. Yandell, C. Feschotte, Transposable elements are major contributors to the origin, diversification, and regulation of vertebrate long noncoding RNAs, PLoS Genet. 9 (2013) e1003470.

[20] J.M. Kirk, S.O. Kim, K. Inoue, M.J. Smola, D.M. Lee, M.D. Schertzer, J.S. Wooten, A.R. Baker, D. Sprague, D.W. Collins, C.R. Horning, S. Wang, Q. Chen, K.M. Weeks, P.J. Mucha, J.M. Calabrese, Functional classification of long non-coding RNAs by k-mer content, Nat. Genet. 50 (2018) 1474–1482.

[21] L. Martin, M. Meier, S.M. Lyons, R.V. Sit, W.F. Marzluff, S.R. Quake, H.Y. Chang, Systematic reconstruction of RNA functional motifs with high-throughput microfluidics, Nat. Methods 9 (2012) 1192–1194.

[22] M.U. Muckenthaler, B. Galy, M.W. Hentze, Systemic iron homeostasis and the iron-responsive element/iron-regulatory protein (IRE/IRP) regulatory network, Annu. Rev. Nutr. 28 (2008) 197–213.

[23] C. Zhang, K.-Y. Lee, M.S. Swanson, R.B. Darnell, Prediction of clustered RNA-binding protein motif sites in the mammalian genome, Nucleic Acids Res. 41 (2013) 6793–6807.

[24] F.C. Oberstrass, A. Lee, R. Stefl, M. Janis, G. Chanfreau, F.H.-T. Allain, Shape-specific recognition in the structure of the Vts1p SAM domain with RNA, Nat. Struct. Mol. Biol. 13 (2006) 160–167.

[25] A.R. Quinlan, I.M. Hall, BEDTools: a flexible suite of utilities for comparing genomic features, Bioinformatics 26 (2010) 841–842.

[26] G. Benson, Tandem repeats finder: a program to analyze DNA sequences, Nucleic Acids Res. 27 (1999) 573–580.

[27] S.B. Needleman, C.D. Wunsch, A general method applicable to the search for similarities in the amino acid sequence of two proteins, J. Mol. Biol. 48 (1970) 443–453.

[28] P. Rice, I. Longden, A. Bleasby, EMBOSS: the European molecular biology open software suite, Trends Genet. 16 (2000) 276–277.

[29] R. Lorenz, S.H. Bernhart, C. Höner Zu Siederdissen, H. Tafer, C. Flamm, P. F. Stadler, I.L. Hofacker, ViennaRNA package 2.0, Algorithm Mol. Biol. 6 (2011) 26.

[30] E. Mattei, G. Ausiello, F. Ferrè, M. Helmer-Citterich, A novel approach to represent and compare RNA secondary structures, Nucleic Acids Res. 42 (2014) 6146–6157.

[31] E. Mattei, M. Pietrosanto, F. Ferrè, M. Helmer-Citterich, Web-Beagle: a web server for the alignment of RNA secondary structures, Nucleic Acids Res. 43 (2015) W493–W497.

[32] J.E. Moore, M.J. Purcaro, H.E. Pratt, C.B. Epstein, N. Shoresh, J. Adrian, T. Kawli, C.A. Davis, A. Dobin, R. Kaul, J. Halow, E.L. Van Nostrand, P. Freese, D.U. Gorkin, Y. Shen, Y. He, M. Mackiewicz, F. Pauli-Behn, B.A. Williams, A. Mortazavi, C. A. Keller, X.-O. Zhang, S.I. Elhajjajy, J. Huey, D.E. Dickel, V. Snetkova, X. Wei, X. Wang, J.C. Rivera-Mulia, J. Rozowsky, J. Zhang, S.B. Chhetri, J. Zhang, A. Victorsen, K.P. White, A. Visel, G.W. Yeo, C.B. Burge, E. Lécuyer, D.M. Gilbert, J. Dekker, J. Rinn, E.M. Mendenhall, J.R. Ecker, M. Kellis, R.J. Klein, W.S. Noble, A. Kundaje, R. Guigó, P.J. Farnham, J.M. Cherry, R.M. Myers, B. Ren, B. R. Graveley, M.B. Gerstein, L.A. Pennacchio, M.P. Snyder, B.E. Bernstein, B. Wold, R.C. Hardison, T.R. Gingeras, J.A. Stamatoyannopoulos, Z. Weng, Expanded encyclopaedias of DNA elements in the human and mouse genomes, Nature 583 (2020) 699–710.

[33] A. Jha, M. Quesnel-Vallières, D. Wang, A. Thomas-Tikhonenko, K.W. Lynch, Y. Barash, Identifying common transcriptome signatures of cancer by interpreting deep learning models, Genome Biol. 23 (2022) 117.

[34] C. Camacho, G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer, T. L. Madden, BLAST+: architecture and applications, BMC Bioinf. 10 (2009) 421.

[35] A.Y. Mitrophanov, M. Borodovsky, Statistical significance in biological sequence analysis, Briefings Bioinf. 7 (2006) 2–24.

[36] O. Bryzghalov, M.W. Szcześniak, I. Makałowska, SyntDB: defining orthologues of human long noncoding RNAs across primates, Nucleic Acids Res. 48 (2020) D238–D245.

[37] S.T. Sherry, M.-H. Ward, M. Kholodov, J. Baker, L. Phan, E.M. Smigielski, K. Sirotkin, dbSNP: the NCBI database of genetic variation, Nucleic Acids Res. 29 (2001) 308–311.

[38] Phan, Jin, Zhang, Qiang, Shekhtman, Shao, Revoe, Villamarin, Ivanchenko, Kimura, Others, ALFA: allele frequency aggregator, National Center for Biotechnology Information, US National Library of Medicine (n.d.).

[39] P. Danecek, J.K. Bonfield, J. Liddle, J. Marshall, V. Ohan, M.O. Pollard, A. Whitwham, T. Keane, S.A. McCarthy, R.M. Davies, H. Li, Twelve years of SAMtools and BCFtools, GigaScience 10 (2021) https://doi.org/10.1093/gigascience/giab008.

[40] A. Guarracino, G. Pepe, F. Ballesio, M. Adinolfi, M. Pietrosanto, E. Sangiovanni, I. Vitale, G. Ausiello, M. Helmer-Citterich, BRIO: a web server for RNA sequence and structure motif scan, Nucleic Acids Res. 49 (2021) W67–W71.

[41] M. Adinolfi, M. Pietrosanto, L. Parca, G. Ausiello, F. Ferrè, M. Helmer-Citterich, Discovering sequence and structure landscapes in RNA interaction motifs, Nucleic Acids Res. 47 (2019) 4958–4969.

[42] J.A. Gustavsen, S. Pai, R. Isserlin, B. Demchak, A.R. Pico, RCy3: network biology using Cytoscape from within R, F1000Res 8 (2019) 1774.

[43] J.A. Bailey, Z. Gu, R.A. Clark, K. Reinert, R.V. Samonte, S. Schwartz, M.D. Adams, E.W. Myers, P.W. Li, E.E. Eichler, Recent segmental duplications in the human genome, Science 297 (2002) 1003–1007.

[44] S.E. Antonarakis, Content and variation of the human genome, in: Medical and Health Genomics, Elsevier, 2016, pp. 161–177.

[45] E.T. Abdullaev, I.R. Umarova, P.F. Arndt, Modelling segmental duplications in the human genome, BMC Genom. 22 (2021) 496.

[46] L. Koch, Capturing transposases for new proteins, Nat. Rev. Genet. 22 (2021) 266–267.

[47] I.H. Chowdhury, H.P. Narra, A. Sahni, K. Khanipov, Y. Fofanov, S.K. Sahni, Enhancer associated long non-coding RNA transcription and gene regulation in experimental models of rickettsial infection, Front. Immunol. 9 (2018) 3014.

[48] B. Li, M. Carey, J.L. Workman, The role of chromatin during transcription, Cell 128 (2007) 707–719.

[49] M. Toll-Riera, M.M. Albà, Emergence of novel domains in proteins, BMC Evol. Biol. 13 (2013) 1–10.

[50] I. Sarropoulos, R. Marin, M. Cardoso-Moreira, H. Kaessmann, Developmental dynamics of lncRNAs across mammalian organs and species, Nature 571 (2019) 510–514.

[51] I. Cvijović, B.H. Good, M.M. Desai, The effect of strong purifying selection on genetic diversity, Genetics 209 (2018) 1235–1278.

[52] M. Nei, Molecular Evolutionary Genetics, Columbia University Press, New York Chichester, West Sussex, 1987.

[53] G.S. França, D.V. Cancherini, S.J. de Souza, Evolutionary history of exon shuffling, Genetica 140 (2012) 249–257.

[54] Z. Yang, J.P. Bielawski, Statistical methods for detecting molecular adaptation, Trends Ecol. Evol. 15 (2000) 496–503.

[55] J. Wang, L. Raskin, D.C. Samuels, Y. Shyr, Y. Guo, Genome measures used for quality control are dependent on gene function and ancestry, Bioinformatics 31 (2015) 318–323.

[56] Y.-J. Kang, D.-C. Yang, L. Kong, M. Hou, Y.-Q. Meng, L. Wei, G. Gao, CPC2: a fast and accurate coding potential calculator based on sequence intrinsic features, Nucleic Acids Res. 45 (2017) W12–W16.

[57] Y. Takei, S. Ishikawa, T. Tokino, T. Muto, Y. Nakamura, Isolation of a novel TP53 target gene from a colon cancer cell line carrying a highly regulated wild-type TP53 expression system, Genes Chromosomes Cancer 23 (1998) 1–9.

[58] A. Diaz-Lagares, A.B. Crujeiras, P. Lopez-Serra, M. Soler, F. Setien, A. Goyal, J. Sandoval, Y. Hashimoto, A. Martinez-Cardús, A. Gomez, H. Heyn, C. Moutinho, J. Espada, A. Vidal, M. Paúles, M. Galán, N. Sala, Y. Akiyama, M. Martínez-Iniesta, L. Farré, A. Villanueva, M. Gross, S. Diederichs, S. Guil, M. Esteller, Epigenetic inactivation of the p53-induced long noncoding RNA TP53 target 1 in human cancer, Proc. Natl. Acad. Sci. U.S.A. 113 (2016) E7535–E7544.

[59] M.R. Finkbeiner, A. Astanehe, K. To, A. Fotovati, A.H. Davies, Y. Zhao, H. Jiang, A. L. Stratford, A. Shadeo, C. Boccaccio, P. Comoglio, P.R. Mertens, P. Eirew, A. Raouf, C.J. Eaves, S.E. Dunn, Profiling YB-1 target genes uncovers a new mechanism for MET receptor regulation in normal and malignant human mammary cells, Oncogene 28 (2009) 1421–1431.

[60] K. Darty, A. Denise, Y. Ponty, VARNA: interactive drawing and editing of the RNA secondary structure, Bioinformatics 25 (2009) 1974–1975.

[61] G.E. Crooks, G. Hon, J.-M. Chandonia, S.E. Brenner, WebLogo: a sequence logo generator, Genome Res. 14 (2004) 1188–1190.

[62] X. Lin, J. Shen, Dan Peng, X. He, C. Xu, X. Chen, J.L. Tanyi, K. Montone, Y. Fan, Q. Huang, L. Zhang, X. Zhong, RNA-binding protein LIN28B inhibits apoptosis through regulation of the AKT2/FOXO3A/BIM axis in ovarian cancer cells, Signal Transduct. Targeted Ther. 3 (2018) 23.

[63] M.L. Wilbert, S.C. Huelga, K. Kapeli, T.J. Stark, T.Y. Liang, S.X. Chen, B.Y. Yan, J. L. Nathanson, K.R. Hutt, M.T. Lovci, H. Kazan, A.Q. Vu, K.B. Massirer, Q. Morris, S. Hoon, G.W. Yeo, LIN28 binds messenger RNAs at GGAGA motifs and regulates splicing factor abundance, Mol. Cell 48 (2012) 195–206.

[64] E. Piskounova, C. Polytarchou, J.E. Thornton, R.J. LaPierre, C. Pothoulakis, J. P. Hagan, D. Iliopoulos, R.I. Gregory, Lin28A and Lin28B inhibit let-7 microRNA biogenesis by distinct mechanisms, Cell 147 (2011) 1066–1079.

[65] W. Yong, D. Yu, Z. Jun, D. Yachen, W. Weiwei, X. Midie, J. Xingzhu, W. Xiaohua, Long noncoding RNA NEAT1, regulated by LIN28B, promotes cell proliferation and migration through sponging miR-506 in high-grade serous ovarian cancer, Cell Death Dis. 9 (2018) 861.

[66] D.T. Peters, H.K.H. Fung, V.M. Levdikov, T. Irmscher, F.C. Warrander, S.J. Greive, O. Kovalevskiy, H.V. Isaacs, M. Coles, A.A. Antson, Human Lin28 forms a high-affinity 1:1 complex with the 106~363 cluster miRNA miR-363, Biochemistry 55 (2016) 5021–5027.

[67] Y.-C.T. Yang, C. Di, B. Hu, M. Zhou, Y. Liu, N. Song, Y. Li, J. Umetsu, Z.J. Lu, CLIPdb: a CLIP-seq database for protein-RNA interactions, BMC Genom. 16 (2015) 1–8.

[68] E.L. Van Nostrand, G.A. Pratt, A.A. Shishkin, C. Gelboin-Burkhart, M.Y. Fang, B. Sundararaman, S.M. Blue, T.B. Nguyen, C. Surka, K. Elkins, R. Stanton, F. Rigo, M. Guttman, G.W. Yeo, Robust transcriptome-wide discovery of RNA-binding protein binding sites with enhanced CLIP (eCLIP), Nat. Methods 13 (2016) 508–514.

[69] J.-H. Li, S. Liu, H. Zhou, L.-H. Qu, J.-H. Yang, starBase v2.0: decoding miRNA-ceRNA, miRNA-ncRNA and protein-RNA interaction networks from large-scale CLIP-Seq data, Nucleic Acids Res. 42 (2014) D92–D97.

[70] M. Hafner, K.E.A. Max, P. Bandaru, P. Morozov, S. Gerstberger, M. Brown, H. Molina, T. Tuschl, Identification of mRNAs bound and regulated by human LIN28 proteins and molecular requirements for RNA recognition, RNA 19 (2013) 613–626.

[71] M. Kellis, B. Wold, M.P. Snyder, B.E. Bernstein, A. Kundaje, G.K. Marinov, L. D. Ward, E. Birney, G.E. Crawford, J. Dekker, I. Dunham, L.L. Elnitski, P.

J. Farnham, E.A. Feingold, M. Gerstein, M.C. Giddings, D.M. Gilbert, T.R. Gingeras, E.D. Green, R. Guigo, T. Hubbard, J. Kent, J.D. Lieb, R.M. Myers, M.J. Pazin, B. Ren, J.A. Stamatoyannopoulos, Z. Weng, K.P. White, R.C. Hardison, Defining functional DNA elements in the human genome, Proc. Natl. Acad. Sci. U.S.A. 111 (2014) 6131–6138.

[72] H. Karner, C.-H. Webb, S. Carmona, Y. Liu, B. Lin, M. Erhard, D. Chan, P. Baldi, R. C. Spitale, S. Sun, Functional conservation of LncRNA JPX despite sequence and structural divergence, J. Mol. Biol. 432 (2020) 283–300.