

The evolution of lncRNA repertoires and expression patterns in tetrapods

Anamaria Necsulea^{1,2†}, Magali Soumillon^{1,2†}, Maria Warnefors^{1,2}, Angélica Liechti^{1,2}, Tasman Daish³, Ulrich Zeller⁴, Julie C. Baker⁵, Frank Grützner³ & Henrik Kaessmann^{1,2}

Only a very small fraction of long noncoding RNAs (lncRNAs) are well characterized. The evolutionary history of lncRNAs can provide insights into their functionality, but the absence of lncRNA annotations in non-model organisms has precluded comparative analyses. Here we present a large-scale evolutionary study of lncRNA repertoires and expression patterns, in 11 tetrapod species. We identify approximately 11,000 primate-specific lncRNAs and 2,500 highly conserved lncRNAs, including approximately 400 genes that are likely to have originated more than 300 million years ago. We find that lncRNAs, in particular ancient ones, are in general actively regulated and may function predominantly in embryonic development. Most lncRNAs evolve rapidly in terms of sequence and expression levels, but tissue specificities are often conserved. We compared expression patterns of homologous lncRNA and protein-coding families across tetrapods to reconstruct an evolutionarily conserved co-expression network. This network suggests potential functions for lncRNAs in fundamental processes such as spermatogenesis and synaptic transmission, but also in more specific mechanisms such as placenta development through microRNA production.

Evolutionary analyses of protein-coding gene sequences¹ and expression patterns² have provided important insights into the genetic basis of lineage-specific phenotypes and into individual gene functions. For lncRNAs, such analyses remain scarce, despite growing interest in these genes. Recent studies have identified thousands of lncRNAs in human^{3–5}, mouse^{6–9}, fruitfly¹⁰, nematode¹¹ and zebrafish¹². Although most lncRNAs have unknown functions, some are involved in fundamental processes like X-chromosome dosage compensation¹³, genomic imprinting¹⁴, cellular pluripotency and differentiation¹⁵. As a class, lncRNAs seem to be versatile expression regulators that recruit chromatin-modifying complexes to specific locations¹⁶, enhance transcription in *cis*¹⁷ or provide decoy targets for microRNAs (miRNAs)¹⁸. Thus, lncRNA evolutionary studies can also be informative in the wider scope of regulatory networks evolution.

Although several highly conserved lncRNAs are known¹⁹, lncRNAs generally have modest sequence conservation^{6,20,21}. Furthermore, in mouse liver, lncRNA transcription undergoes rapid evolutionary turnover²². These observations suggest that many lncRNAs may have no biological relevance. Detailed evolutionary analyses can clarify lncRNA functionality, but such analyses have been hampered by lack of annotations in non-model organisms.

The evolutionary history of lncRNAs in 11 tetrapods

We used RNA sequencing (RNA-seq) to determine lncRNA repertoires of 11 tetrapod species. We analysed 185 samples and approximately 6 billion RNA-seq reads (Supplementary Table 1), representing the polyadenylated transcriptomes of 8 organs (cortex or whole brain, cerebellum, heart, kidney, liver, placenta, ovary and testes) and 11 species (human, chimpanzee, bonobo, gorilla, orangutan, macaque, mouse, opossum, platypus, chicken and frog), which diverged approximately 370 million years (Myr) ago²³. We included 47 strand-specific samples (approximately 2 billion reads), which allowed us to confirm gene orientation and to predict antisense transcripts (Methods).

Using this data set, we recovered spliced transcripts for most known genes (Extended Data Table 1a and Supplementary Discussion). We evaluated the protein-coding potential of transcripts using genome-wide codon substitution frequency scores (CSF²⁴) and the presence of sequence similarity with known proteins and protein domains (Methods), obtaining correct classifications for approximately 96% of protein-coding genes and 97% of known noncoding RNAs, on average (Extended Data Table 1b). We thus identified between approximately 3,000 and 15,000 multi-exonic lncRNAs in each species, including known lncRNAs for human^{4,5} and mouse⁶, as well as approximately 10,000 novel human and 9,000 novel murine lncRNAs (Fig. 1a and Extended Data Table 2). Although part of the variability in lncRNA repertoire size may be biologically meaningful, much is likely to be explained by unequal sequencing depth and by variable genome sequence and assembly quality (Supplementary Discussion).

We reconstructed homologous families based on sequence similarity and we inferred a stringent minimum evolutionary age of lncRNAs, requiring transcription evidence as an additional criterion (Methods). We also estimated a ‘maximum’ evolutionary age by explicitly accounting for between-species variations in RNA-seq coverage and annotation quality (Methods and Extended Data Table 3a). We thus identified 13,533 lncRNA families transcribed in at least 3 species. Most (81%) lncRNA families were primate-specific, but 2,508 (19%) families likely originated more than 90 Myr ago and 425 (3%) more than 300 Myr ago (Fig. 1a). Most homologous lncRNAs were found in conserved synteny, even for distantly related species (Extended Data Table 3b).

The large proportion of inferred young lncRNAs may be due to fast lncRNA evolution, which prevents detection of distant homologues. Furthermore, the phylogenetic distribution of the species in our data set may contribute to the skewed distribution of estimated ages. To investigate these possibilities, we evaluated DNA sequence conservation across placental mammals²⁵ and variation within populations²⁶.

¹Center for Integrative Genomics, University of Lausanne, 1015 Lausanne, Switzerland. ²Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland. ³The Robinson Institute, School of Molecular and Biomedical Science, University of Adelaide, Adelaide, South Australia 5005, Australia. ⁴Department of Systematic Zoology, Faculty of Agriculture and Horticulture, Humboldt University Berlin, 10099 Berlin, Germany. ⁵Department of Genetics, Stanford University School of Medicine, Stanford, California 94305, USA. [†]Present addresses: Laboratory of Developmental Genomics, Ecole Polytechnique Fédérale de Lausanne (EPFL), 1015 Lausanne, Switzerland (A.N.); Harvard Department of Stem Cell and Regenerative Biology, Harvard University, Cambridge, Massachusetts 02138, USA, and Broad Institute, Cambridge, Massachusetts 02142, USA (M.S.).

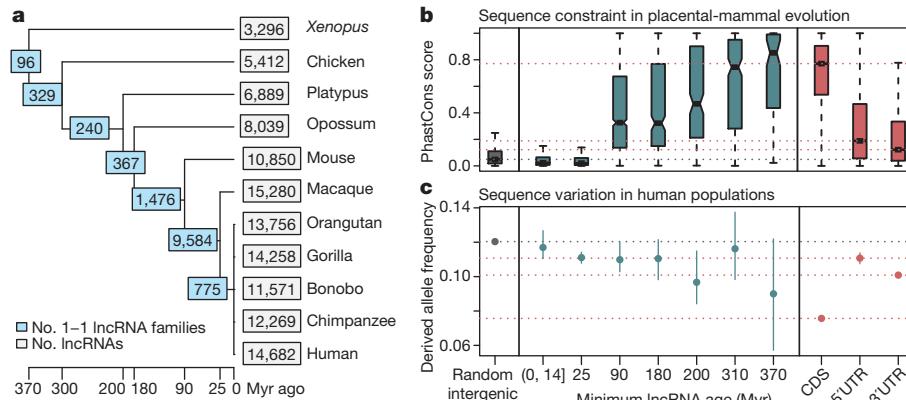


Figure 1 | Evolutionary age and genomic characteristics of lncRNA families.

a, Simplified phylogenetic tree. Internal branches and root, numbers of 1–1 orthologous lncRNA families for each minimum evolutionary age. Tree tips, lncRNA numbers for each species. **b**, Exonic sequence conservation (placental PhastCons score), for random intergenic regions, lncRNA evolutionary age classes, coding (CDS) and untranslated exons of protein-coding genes. **c**, Mean derived allele frequency of autosomal non-CpG single nucleotide

for human lncRNAs (Fig. 1b, c and Methods). We found that young lncRNAs (inferred minimum ages of 25 Myr or younger) have low levels of long-term exonic sequence conservation (median score ~ 0.02), significantly lower than random intergenic regions (median score ~ 0.05 , Wilcoxon test, $P < 10^{-10}$). However, single nucleotide polymorphisms found in primate-specific (minimum evolutionary age 25 Myr) lncRNA exons have significantly lower derived allele frequencies (mean 0.11) than those found in intergenic regions (mean 0.12, randomization test, $P < 0.01$), consistent with recent purifying selection²⁷. The same conclusions were reached using maximum evolutionary age estimates (Extended Data Fig. 1a, b), and when controlling for GC-biased gene conversion²⁸ (Extended Data Fig. 1c) and for linkage to protein-coding genes (Extended Data Fig. 1d). The presence of selective constraint in recent evolution, but not on a broader timescale, is compatible with a recent origination or acquisition of novel functions for a fraction of primate-specific lncRNAs.

Overall, the two measures of selective constraint correlate with evolutionary age estimates (Fig. 1c, d). Remarkably, older lncRNAs (minimum age 90 Myr) have higher levels of long-term exonic sequence conservation than untranslated regions (UTRs), and the oldest age classes are comparable with coding exons (Fig. 1c, Wilcoxon test, $P > 0.05$). Furthermore, lncRNA promoters are as conserved as protein-coding gene promoters even for younger classes (Extended Data Fig. 1e, f), suggesting stronger selective constraints at the transcriptional level, as previously observed⁸.

Active regulation of ancient lncRNAs

We next asked whether lncRNA expression patterns vary with evolutionary age. We found that lncRNAs are lowly transcribed, highly tissue-specific and preferentially expressed in testes (Fig. 2a–c and Extended Data Fig. 2), consistent with previous observations^{4,5}. However, the testes specificity is stronger for young lncRNAs (55%) than for old lncRNAs (46%, Fig. 2a, chi-squared test, $P < 10^{-10}$), in agreement with the hypothesis that the permissive testes chromatin favours new gene origination²⁹. After testes, neural tissues generally express the largest numbers of lncRNAs (Fig. 2a and Extended Data Fig. 2), consistent with a previously reported enrichment of lncRNAs in mouse brain⁹. Surprisingly, for platypus, ovary appears to be the second most favourable tissue for lncRNA expression (Extended Data Fig. 2).

The low expression levels and the testes specificity raise the question of whether lncRNAs are actively regulated, or whether they result from non-specific transcription in open chromatin regions. To test these hypotheses, we analysed the occurrence of transcription-factor-binding

polymorphisms (SNPs) segregating in African populations (1000 Genomes Project²⁶). Intergenic SNPs were randomly drawn in regions matching lncRNA recombination rates (Methods). Error bars, 95% confidence intervals based on 100 bootstrap resampling replicates. Round brackets indicate that the boundary is excluded from the interval; square brackets indicate that the boundary is included in the interval.

sites as an indicator of active regulation. Using a genome-wide set of evolutionarily conserved binding sites predicted *in silico*³⁰ and ChIP-seq transcription-factor-binding data³¹ (Methods), we found that lncRNA promoters were more frequently associated with transcription factors than random intergenic regions (Fig. 2d and Extended Data Fig. 3a, c). Moreover, binding site sequence conservation was stronger in lncRNA promoters than in random intergenic regions and even protein-coding gene promoters, in particular for ancient lncRNAs (Fig. 2e, Wilcoxon test $P < 10^{-10}$). Consistently, the evolutionary turnover of CEBPA and HNF4A binding³² between human and mouse is significantly slower for lncRNAs than expected by chance (Extended Data Fig. 3f, g, Fisher's exact test $P < 10^{-10}$). Taken together, these results suggest that lncRNA

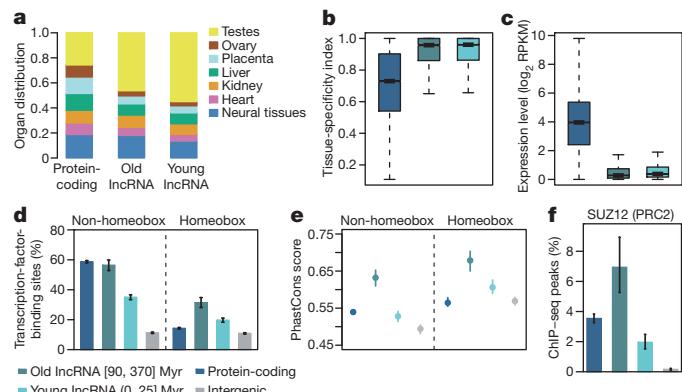


Figure 2 | lncRNA expression patterns and evidence for developmental regulation of old lncRNAs. **a**, Distribution of the organ in which maximum expression is observed, for human protein-coding genes, old lncRNAs (minimum age 90–370 Myr, 2,556 lncRNAs) and young lncRNAs (minimum age 0–25 Myr, 12,126 lncRNAs). **b**, Tissue-specificity index. Values close to 1 represent high tissue specificity. **c**, Distribution of the maximum expression level (\log_2 -transformed RPKM). **d**, Frequency of *in silico*-predicted binding sites for homeobox and non-homeobox transcription factors, in human gene promoters (2 kb upstream) and in random intergenic regions. Error bars, 95% binomial proportion confidence intervals. **e**, Mean sequence conservation (PhastCons score) for transcription-factor-binding sites. Error bars, 95% confidence intervals based on 100 bootstrap replicates. **f**, Frequency of SUZ12 (part of the PRC2 complex) binding (ENCODE ChIP-seq). Error bars, 95% binomial proportion confidence intervals. We analysed 793 ‘old’ lncRNAs, 3,418 ‘young’ lncRNAs and 16,566 protein-coding genes for which the predicted transcription start site was within 100 bp of a cap analysis gene expression (CAGE) tag.

transcription is overall actively regulated, in particular for ancient lncRNAs.

Using *in silico* binding-site predictions, we also uncovered a remarkable difference between two transcription-factor classes: homeobox transcription factors, which function in embryonic development, bind preferentially in lncRNA promoters, whereas non-homeobox transcription factors bind more frequently in protein-coding promoters (Fig. 2d and Extended Data Fig. 3b). Notably, 31% of old lncRNA promoters have homeobox transcription-factor-binding sites, more than twice the frequency observed for protein-coding genes (14%, Fisher's exact test, $P < 10^{-10}$). The ChIP-seq data set consisted largely (95%) of non-homeobox transcription factors, 117 (98%) of which were associated significantly more often with protein-coding than with lncRNA promoters (Extended Data Fig. 3d). However, two factors bound more frequently in old lncRNA than in protein-coding promoters: SUZ12, a member of the polycomb repressive complex 2 (PRC2) that functions in pluripotency and differentiation³³ (Fig. 2f) and OCT4 (also known as POU5F1), a homeobox transcription factor that controls pluripotency³⁴ (Extended Data Fig. 3e). The association with homeobox transcription factors and PRC2 suggests that lncRNAs (especially ancient ones) may be important for embryonic development, pluripotency and differentiation¹⁵.

Rapid evolution of lncRNA expression patterns

We next assessed the evolutionary conservation of lncRNA expression patterns. We first estimated the presence of shared transcription across species. To reduce the impact of weak lncRNA sequence conservation, we compared intergenic lncRNAs across closely related primate species (Fig. 3a) and we analysed lncRNAs transcribed in antisense of

protein-coding exons (Extended Data Fig. 4a). We found that lncRNA transcription evolves rapidly: only approximately 92% of human intergenic lncRNAs were also detected as expressed in chimpanzee or bonobo and only approximately 72% were expressed in macaque, whereas more than 98% of conservation was observed for protein-coding genes, for all primates (Fig. 3a). Likewise, the evolutionary turnover of anti-sense lncRNAs is rapid compared to protein-coding genes (Extended Data Fig. 4a). The discrepancy between lncRNAs and protein-coding genes remained considerable when controlling for low lncRNA expression with a read resampling procedure (Fig. 3a and Extended Data Fig. 4a), indicating that rapid transcription evolution is a genuine feature of lncRNAs²².

We also measured correlations of lncRNA expression levels between pairs of species (Fig. 3b). The difference between lncRNAs and protein-coding genes is marked (Fig. 3c): Spearman's correlation coefficient for lncRNA brain expression between human and chimpanzee (which diverged 6 Myr ago) is approximately 0.55, lower than the correlation (0.66) observed for protein-coding genes between human and *Xenopus* (which diverged ~370 Myr ago). However, low lncRNA expression levels explain much of this discrepancy, as differences between correlation coefficients for the two classes of genes were much lower after resampling controls (Fig. 3c). For both protein-coding genes and lncRNAs, the testes have the fastest rates of evolution (Extended Data Fig. 4b).

We also observed that lncRNA tissue specificity is well conserved among primates, but not beyond. Indeed, a hierarchical clustering of samples based on pairwise correlations for eutherian lncRNA families revealed preferential grouping among related organs for primates, though all mouse samples clustered together (Fig. 3c and Extended Data Fig. 4f, g). Moreover, 47% of human tissue-specific lncRNAs had

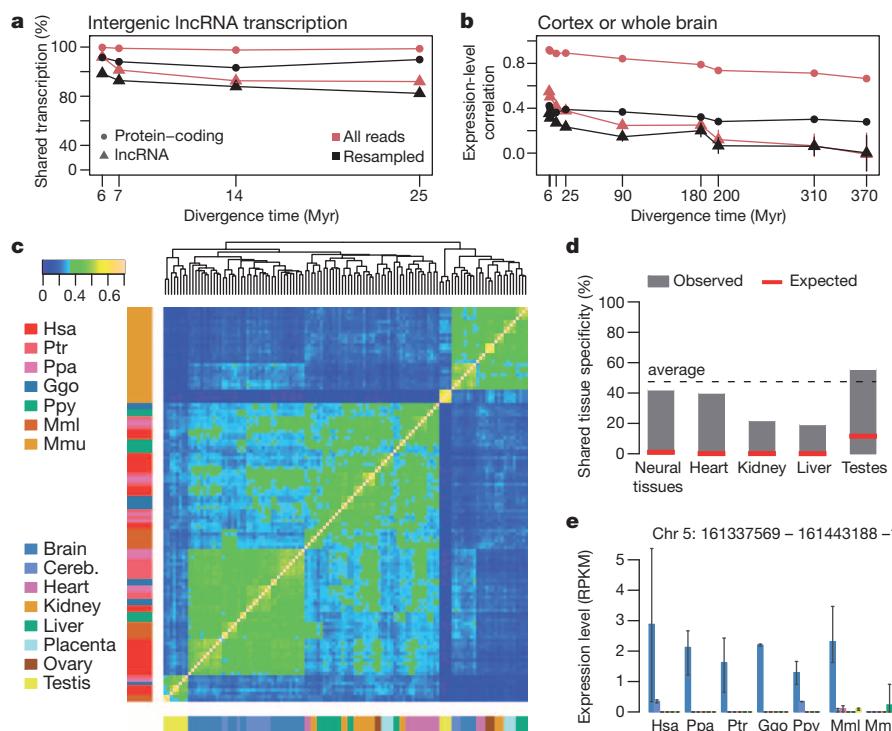


Figure 3 | Evolution of lncRNA expression patterns in tetrapods.

a, Percentage of human lncRNAs (4,430 intergenic primate lncRNA families) transcribed in other primates, in pool of 5 somatic tissues (Methods). **b**, Pairwise Spearman correlations between human and other species, for cortex or whole brain. In **a** and **b**, 'all reads' represents estimates obtained with all reads, and 're-sampled' represents estimates obtained after resampling identical numbers of mapped reads per species and tissue; error bars, 95% confidence intervals obtained with 100 bootstrap resampling replicates. **c**, Hierarchical clustering of pairwise Spearman correlations, for 1,716 lncRNA families with 1–1 orthologues in all eutherians. Samples are colour-coded according to the

organ (horizontal) and species (vertical). Hsa, *Homo sapiens*; Ppa, *Pan troglodytes* (chimpanzee); Ppy, *Pan paniscus* (bonobo); Ggo, *Gorilla gorilla*; Ppy, *Pongo pygmaeus* (orangutan); Mml, *Macaca mulatta* (macaque); Mmu, *Mus musculus* (mouse). **d**, Proportion of human organ-specific lncRNAs (771 lncRNAs with minimum evolutionary age >90 Myr, tissue-specificity index >0.9, RPKM >0.1) for which organ specificity is shared across primates. Red lines, random expectation; dashed line, average conserved specificity across organs. **e**, A lncRNA with conserved neural tissue specificity across primates. Error bars, range observed in biological replicates. Chromosomal coordinates are given in the plot title.

conserved specificity in all primates, while only 28% had conservation across all eutherians (Fig. 3d and Extended Data Fig. 4c–e). These proportions are significantly lower than for protein-coding genes, for which 81% are conserved across all primates and 72% across all eutherians (Fisher's exact test, $P < 10^{-10}$), but higher than randomly expected (randomization test, $P < 0.01$). The extent of conservation varies among tissues (Fig. 3d and Extended Data Fig. 4c–e), but is always significantly higher than expected by chance (randomization test, $P < 0.01$). These observations are illustrated by a lncRNA identified within a cluster of GABA (γ -aminobutyric acid) receptors on human chromosome 5, expressed in neural tissues for primates, but detected only in liver in mouse (Fig. 3e).

Evolutionarily conserved co-expression network

Finally, we evaluated the co-expression of lncRNAs and protein-coding genes, which can indicate functional relatedness³⁵ or regulatory relationships³⁶. As co-expression may also arise spuriously, we used evolutionary conservation as a criterion for significance³⁵. We analysed a set of 16,076 protein-coding gene families and 1,770 lncRNA families expressed in at least 3 species (Methods). We evaluated expression correlations for all gene pairs and tested if the combination of correlation coefficients across species was significantly higher (for positive associations) or lower (for negative associations) than expected by chance³⁵ (Methods). The conserved co-expression relationships formed a network with 9,388 nodes (8,971 protein-coding and 417 lncRNAs) and 97,556 edges (Supplementary Table 2). The same criteria applied on randomized gene families identified only approximately 160 co-expression relationships, proving the reconstruction specificity (Supplementary Discussion).

The co-expression network can predict functional relatedness, as illustrated by the high frequency of connections within gene ontology

(GO) categories: out of 115 GO categories with at least 100 members, 101 (88%) had within-category connections more often than randomly expected (Fig. 4a). To verify if the direction of network connections may also predict regulatory associations, we analysed 710 connections annotated as expression activation/inhibition relationships in the String³⁷ database. We found that approximately 70% of positive connections are annotated as activation relationships, significantly more than negative connections (30%, Fisher's exact test, $P = 0.01$; Extended Data Fig. 4a). Consistent with this, we found an overwhelming majority of negative connections for the REST and *HBP1* transcriptional repressors (Fig. 4b). Positive co-expression also often arises for genes that participate in complexes, such as the sodium channel subunit *SCNN1B* (Fig. 4b). Most (72%) network connections are positive co-expression cases. However, whereas lncRNAs have more negative connections (Fig. 4b). Interestingly, the imprinted lncRNA *H19*, which functions as a miRNA precursor³⁸, has a majority of negative connections (Fig. 4b).

The network connectivity depends on expression levels, as more connections were detected for highly expressed genes (Extended Data Fig. 5b, c). Expectedly, lncRNAs generally had lower connectivity (median degree 2) than protein-coding genes (median degree 5, Wilcoxon test $P < 10^{-10}$; Extended Data Fig. 5d), and transcription factors were less well connected (median degree 4) than non-transcription-factor protein-coding genes. However, when resampling genes with similar expression levels, lncRNAs had higher degrees (median 3) than protein-coding genes (median 2, randomization test $P < 0.01$), and transcription factors had higher connectivity than other protein-coding genes (median 3, randomization test $P = 0.02$; Extended Data Fig. 5d), consistent with their central roles in regulatory networks. The highly connected lncRNAs may represent interesting candidates for further studies of gene expression regulation. Notably, lncRNAs had connections in

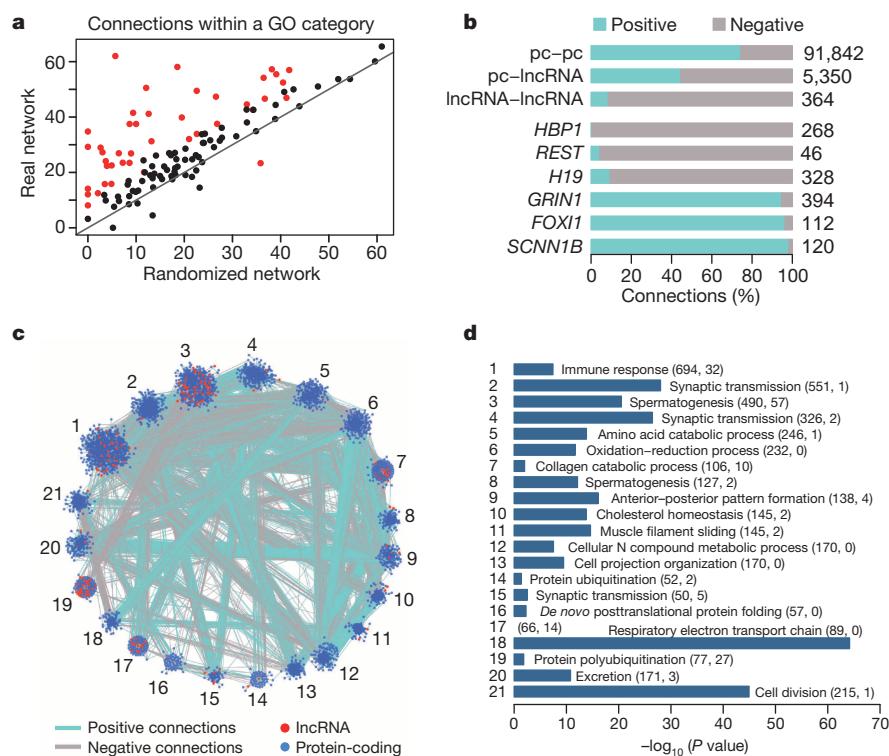


Figure 4 | Evolutionarily conserved co-expression network of protein-coding genes and lncRNAs. **a**, Percentage of genes with connections within the same GO category, in real and randomized co-expression networks, for 115 biological process categories. Red, significant difference between real and randomized data ($P < 0.05$). **b**, Percentage of positive connections, for the entire network and for six genes with extreme positive:negative ratios. pc-*pc*, connections between two protein-coding genes; pc-lncRNA, connections

between a protein-coding and a lncRNA gene; lncRNA-lncRNA, connections between two lncRNAs. Total numbers of connections per category are indicated on the right side of the plot. **c**, Cytoscape⁴⁷ representation of the 21 largest MCL clusters in the co-expression network. **d**, GO enrichment for the 21 largest MCL clusters; only the most significant GO category is displayed. Numbers of coding and lncRNA genes in each cluster are shown in brackets.

cis more often than protein-coding genes (Extended Data Fig. 5e). An excess of connections in *cis* was also found for protein-coding genes acting in body plan development, in particular for HOX genes (Extended Data Fig. 5e and Supplementary Table 3).

Finally, we used the co-expression network to infer potential functions for lncRNAs. Using the Markov clustering algorithm (MCL³⁹), we identified 1,326 groups of highly inter-connected genes, including 21 clusters with at least 50 genes (Fig. 4c and Supplementary Table 4). The proportion of lncRNAs in these clusters varied between 0 and 26% (Fig. 4d). The clusters were enriched for organ-specific functions, such as spermatogenesis (testes), synaptic transmission (neural tissues), catabolic processes (liver), muscle functions (heart) (Methods, Fig. 4c and Supplementary Table 4). We also recovered specific processes, such as anterior-posterior pattern formation in a cluster that includes HOX genes (Fig. 4c). The clusters with highest lncRNA proportions were enriched in spermatogenesis functions (Fig. 4c), in agreement with the predominant lncRNA testes specificity. GO enrichment analyses for individual nodes suggested potential lncRNA involvement in, for example, nervous system development, cell adhesion, transcription (Supplementary Table 5).

miRNA precursors in the *H19* co-expression network

The only MCL cluster without significant GO enrichments (Fig. 4d) contains a high proportion (17.5%) of lncRNAs, including *H19*. As *H19* is a precursor for *miR-675*, which targets *IGF1R* and thus stalls placenta growth during late gestation³⁸, we scanned the network for other potential miRNA precursors (Methods). Unexpectedly, genes positively connected with *H19* had the highest average density of embedded miRNAs (Extended Data Fig. 6a). These include one exceptional case: a lncRNA that could potentially promote the transcription of between 2 and 7 miRNAs in different species (Fig. 5a, Supplementary Table 6 and Supplementary Discussion). This lncRNA (that we name *H19X*, for *H19* X-linked co-expressed lncRNA) is transcribed in all studied species and thus likely originated at least 370 Myr ago, in the tetrapod ancestor. Notably, its expression pattern appears to have dramatically shifted during evolution, from an ancestral testes-predominant pattern to preferential expression in the chorioallantoic placenta of eutherians (Fig. 5a).

The miRNAs associated with *H19X* comprise two conserved tetrapod families, four placental-mammal-specific families and one rodent-specific miRNA (Supplementary Discussion). Interestingly, the two oldest families (with representative members *miR-503*, *miR-16c*, and *miR-424*, *miR322*, *mir-15c*, respectively) seem to have undergone accelerated sequence evolution in the eutherian ancestor (Extended Data Fig. 6b). In human and mouse, these miRNAs are in general highly expressed in the placenta (Extended Data Fig. 6c, d).

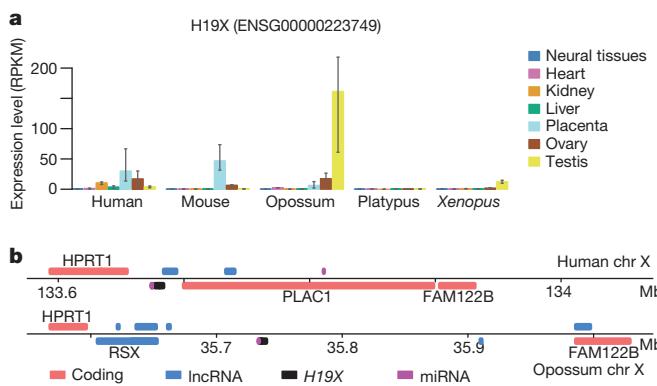


Figure 5 | *H19* co-expression network and miRNA precursors. **a**, Expression pattern for an X-linked *H19* co-expressed lncRNA (*H19X*, identified as ENSG00000223749 in the Ensembl database), in five tetrapod species. The error bars represent the range observed in biological replicates. **b**, Genomic neighbourhood of *H19X* in human and opossum.

Finally, *H19X* is a neighbour of *Rsx*, the lncRNA that drives imprinted X-inactivation in marsupials⁴⁰ (Fig. 5b), suggesting that *H19X* may itself be imprinted. These results suggest that *H19X* may function like *H19*, by promoting miRNA transcription, preferentially in the placenta and in an imprinted manner. Although validation is needed, this illustrates how the reconstruction of a conserved co-expression network, enabled by the broad evolutionary perspective of our study, can predict lncRNA functions and stimulate further investigations.

METHODS SUMMARY

We sequenced polyadenylated transcriptomes of 11 species and 8 tissues with Illumina GAI and HiSeq2000 technologies. We detected multi-exonic transcripts based on transcribed island and splice junction coordinates, using TopHat⁴¹ and Cufflinks⁴². Protein-coding potential was inferred using codon substitution frequency scores (CSF²⁴) and sequence similarity with known proteins⁴³ and protein domains⁴⁴. We included published lncRNA annotations for human and mouse⁴⁵ and projected annotations across species. We reconstructed homologous families based on DNA sequence similarity, with single-link clustering. We inferred lncRNA evolutionary ages based on the phylogenetic distribution of species with transcription evidence, or for which its absence was due to low coverage or incomplete annotation. We computed RPKM (reads per kilobase per million mapped reads) levels using non-overlapping exonic regions and unambiguously mapped reads, and we normalized them through median-scaling². We computed tissue-specificity indexes as previously described⁴⁶. To control for unequal coverage, we simulated read distributions by resampling identical numbers of reads per species and tissue, keeping proportions among genes unchanged. We reconstructed an evolutionarily conserved co-expression network by computing expression correlations between gene pairs and identifying cross-species combinations that are significantly higher or lower than randomly expected³⁵. Network analysis was done with MCL³⁹ and Cytoscape⁴⁷. For all statistics and graphics we used R⁴⁸.

Online Content Any additional Methods, Extended Data display items and Source Data are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 31 December 2012; accepted 5 December 2013.

Published online 19 January 2014.

- Kosiol, C. et al. Patterns of positive selection in six mammalian genomes. *PLoS Genet.* **4**, e1000144 (2008).
- Brawand, D. et al. The evolution of gene expression levels in mammalian organs. *Nature* **478**, 343–348 (2011).
- Khalil, A. M. et al. Many human large intergenic noncoding RNAs associate with chromatin-modifying complexes and affect gene expression. *Proc. Natl. Acad. Sci. USA* **106**, 11667–11672 (2009).
- Cabili, M. N. et al. Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses. *Genes Dev.* **25**, 1915–1927 (2011).
- Derrien, T. et al. The GENCODE v7 catalog of human long noncoding RNAs: Analysis of their gene structure, evolution, and expression. *Genome Res.* **22**, 1775–1789 (2012).
- Guttman, M. et al. Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in mammals. *Nature* **458**, 223–227 (2009).
- Guttman, M. et al. *Ab initio* reconstruction of cell type-specific transcriptomes in mouse reveals the conserved multi-exonic structure of lncRNAs. *Nature Biotechnol.* **28**, 503–510 (2010).
- Carninci, P. et al. The transcriptional landscape of the mammalian genome. *Science* **309**, 1559–1563 (2005).
- Mercer, T. R., Dinger, M. E., Sunkin, S. M., Mehler, M. F. & Mattick, J. S. Specific expression of long noncoding RNAs in the mouse brain. *Proc. Natl. Acad. Sci. USA* **105**, 716–721 (2008).
- Young, R. S. et al. Identification and properties of 1,119 candidate lncRNA loci in the *Drosophila melanogaster* genome. *Genome Biol. Evol.* **4**, 427–442 (2012).
- Nam, J. W. & Bartel, D. Long non-coding RNAs in *C. elegans*. *Genome Res.* (2012).
- Ulitsky, I., Shkumatava, A., Jan, C. H., Sive, H. & Bartel, D. P. Conserved function of lncRNAs in vertebrate embryonic development despite rapid sequence evolution. *Cell* **147**, 1537–1550 (2011).
- Chow, J. C., Yen, Z., Ziesche, S. M. & Brown, C. J. Silencing of the mammalian X chromosome. *Annu. Rev. Genomics Hum. Genet.* **6**, 69–92 (2005).
- Sleutels, F., Zwart, R. & Barlow, D. P. The non-coding Air RNA is required for silencing autosomal imprinted genes. *Nature* **415**, 810–813 (2002).
- Dinger, M. E. et al. Long noncoding RNAs in mouse embryonic stem cell pluripotency and differentiation. *Genome Res.* **18**, 1433–1445 (2008).
- Rinn, J. L. et al. Functional demarcation of active and silent chromatin domains in human HOX loci by noncoding RNAs. *Cell* **129**, 1311–1323 (2007).
- Ørom, U. A. et al. Long noncoding RNAs with enhancer-like function in human cells. *Cell* **143**, 46–58 (2010).
- Cesana, M. et al. A long noncoding RNA controls muscle differentiation by functioning as a competing endogenous RNA. *Cell* **147**, 358–369 (2011).

19. Chodroff, R. A. et al. Long noncoding RNA genes: conservation of sequence and brain expression among diverse amniotes. *Genome Biol.* **11**, R72 (2010).
20. Marques, A. C. & Ponting, C. P. Catalogues of mammalian long noncoding RNAs: modest conservation and incompleteness. *Genome Biol.* **10**, R124 (2009).
21. Cabili, M. N. et al. Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses. *Genes Dev.* **25**, 1915–1927 (2011).
22. Kutter, C. et al. Rapid turnover of long noncoding RNAs and the evolution of gene expression. *PLoS Genet.* **8**, e1002841 (2012).
23. Hedges, S. B., Dudley, J. & Kumar, S. TimeTree: a public knowledge-base of divergence times among organisms. *Bioinformatics* **22**, 2971–2972 (2006).
24. Lin, M. F. et al. Revisiting the protein-coding gene catalog of *Drosophila melanogaster* using 12 fly genomes. *Genome Res.* **17**, 1823–1836 (2007).
25. Siepel, A. et al. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* **15**, 1034–1050 (2005).
26. The 1000 Genomes Project Consortium. An integrated map of genetic variation from 1,092 human genomes. *Nature* **135**, 56–65 (2012).
27. Ward, L. D. & Kellis, M. Evidence of abundant purifying selection in humans for recently-acquired regulatory functions. *Science* **337**, 1675–1678 (2012).
28. Galtier, N. & Duret, L. Adaptation or biased gene conversion? Extending the null hypothesis of molecular evolution. *Trends Genet.* **23**, 273–277 (2007).
29. Soumilon, M. et al. Cellular source and mechanisms of high transcriptome complexity in the mammalian testis. *Cell Rep.* **3**, 2179–2190 (2013).
30. Lindblad-Toh, K. et al. A high-resolution map of human evolutionary constraint using 29 mammals. *Nature* **478**, 476–482 (2011).
31. The Encode Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
32. Schmidt, D. et al. Five-vertebrate ChIP-seq reveals the evolutionary dynamics of transcription factor binding. *Science* **328**, 1036–1040 (2010).
33. Walker, E., Manias, J. L., Chang, W. Y. & Stanford, W. L. PCL2 modulates gene regulatory networks controlling self-renewal and commitment in embryonic stem cells. *Cell Cycle* **10**, 45–51 (2011).
34. Chambers, I. & Tomlinson, S. R. The transcriptional foundation of pluripotency. *Development* **136**, 2311–2322 (2009).
35. Stuart, J. M., Segal, E., Koller, D., Kim, S. K. & A. Gene-coexpression network for global discovery of conserved genetic modules. *Science* **302**, 249–255 (2003).
36. Shkumatava, A., Stark, A., Sive, H. & Bartel, D. P. Coherent but overlapping expression of microRNAs and their targets during vertebrate development. *Genes Dev.* **23**, 466–481 (2009).
37. Franceschini, A. et al. STRING v9.1: protein–protein interaction networks, with increased coverage and integration. *Nucleic Acids Res.* **41**, D808–D815 (2013).
38. Keniry, A. et al. The *H19* lincRNA is a developmental reservoir of miR-675 that suppresses growth and *Igf1r*. *Nature Cell Biol.* **14**, 659–665 (2012).
39. Van Dongen, S. & Abreu-Goodger, C. Using MCL to extract clusters from networks. *Methods Mol. Biol.* **804**, 281–295 (2012).
40. Grant, J. et al. *Rsx* is a metatherian RNA with *Xist*-like properties in X-chromosome inactivation. *Nature* **487**, 254–258 (2012).
41. Trapnell, C., Pachter, L. & Salzberg, S. L. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* **25**, 1105–1111 (2009).
42. Trapnell, C. et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature Biotechnol.* **28**, 511–515 (2010).
43. UniProt. Reorganizing the protein space at the Universal Protein Resource (UniProt). *Nucleic Acids Res.* **40**, D71–D75 (2012).
44. Punta, M. et al. The Pfam protein families database. *Nucleic Acids Res.* **40**, D290–D301 (2012).
45. Flicek, P. et al. Ensembl 2012. *Nucleic Acids Res.* **40**, D84–D90 (2012).
46. Yanai, I. et al. Genome-wide midrange transcription profiles reveal expression level relationships in human tissue specification. *Bioinformatics* **21**, 650–659 (2005).
47. Smoot, M. E., Ono, K., Ruscheinski, J., Wang, P.-L. & Ideker, T. Cytoscape 2.8: new features for data integration and network visualization. *Bioinformatics* **27**, 431–432 (2011).
48. R Development Core Team. *R: A language and environment for statistical computing* <http://www.r-project.org> (R Foundation for Statistical Computing, 2011).

Supplementary Information is available in the online version of the paper.

Acknowledgements We thank L. Froidevaux and D. Cortéz for help with genome sequencing, J. Meunier for help with preliminary miRNA analyses, K. Harshman and the Lausanne Genomics Technology Facility for high-throughput sequencing support, I. Xenarios for computational support, S. Bergmann and Z. Kutalik for advice on co-expression analyses. Human embryonic and fetal material was provided by the Joint MRC/Wellcome Trust (grant 099175/Z/12/Z) Human Developmental Biology Resource (<http://www.hdbi.org>). The computations were performed at the Vital-IT (<http://www.vital-it.ch>) Center for high-performance computing of the SIB Swiss Institute of Bioinformatics. This research was supported by grants from the European Research Council (Starting Independent Researcher Grant 242597, SexGenTransEvolution) and the Swiss National Science Foundation (grant 31003A_130287) to H.K.A.N. was supported by a FEBS long-term postdoctoral fellowship.

Author Contributions A.N. conceived and performed all biological analyses and wrote the manuscript, with input from all authors. A.N. and M.W. processed RNA-seq data. M.S. and A.L. generated RNA-seq data. T.D. and F.G. collected platypus samples. U.Z. collected opossum samples. J.C.B. provided mouse placenta samples and contributed to *H19X* analyses. The project was supervised and originally designed by H.K.

Author Information The sequencing data have been deposited in the Gene Expression Omnibus (accession GSE43520) and SRA (PRJNA186438 and PRJNA202404). Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to A.N. (anamaria.necsulea@epfl.ch) or H.K. (henrik.kaessmann@unil.ch).

METHODS

RNA sequencing and initial analysis. Our main data set consists of 185 RNA-seq (135 previously published² and 50 new) samples, amounting to approximately 6 billion raw reads (Supplementary Table 1). The libraries were prepared using standard Illumina protocols and sequenced with Illumina GAI or HiSeq2000, as single-end reads, after poly(A) selection. After ensuring data comparability (Supplementary Discussion), we included 47 samples that we generated with a strand-specific RNA-seq protocol, for six species (human, mouse, opossum, platypus, chicken and *Xenopus*). To gain statistical power for co-expression network reconstruction, we incorporated 44 Illumina and 4 Applied Biosystems (ABI) Solid RNA-seq samples published by other groups (Supplementary Table 1 and Supplementary Discussion). We aligned the reads and detected splice junctions *de novo* using TopHat⁴¹ v1.4.0 and Bowtie⁴⁹ v0.12.5. The genome sequences were retrieved from Ensembl⁴⁵ v62. Given the genetic similarity between chimpanzee and bonobo and the unavailability of the bonobo genome sequence when we started our project, we used the chimpanzee genome as a reference for all bonobo analyses.

lncRNA detection. To detect genes *de novo* with RNA-seq, we developed an algorithm that predicts multi-exonic transcribed loci based on transcribed island and splice junction coordinates and we used Cufflinks⁴² to assemble transcripts from genomic read alignments (Supplementary Discussion). We combined multi-exonic transcripts detected with the two methods and Ensembl 62 annotations (including GENCODE lncRNAs⁵) into non-redundant data sets for each species. For human, we included approximately 8,000 lncRNAs predicted with RNA-seq²¹. To assess the evolution of sense-antisense transcripts, we repeated the detection procedure using only strand-specific samples. After the initial detection procedure, which used mainly in-house generated samples, we added to our analyses several previously published RNA-seq samples, mainly from the human ENCODE⁵ and Illumina Human Body Map⁴ projects, as well as several strand-specific samples that we generated at a later stage to increase coverage for the placenta, ovary and testes for several species (Supplementary Table 1). We did not repeat the entire detection procedure with these new samples, but we used the additional splice junction information to join fragmented lncRNA loci. We also discarded *de novo* detected loci which thus appeared to be unannotated UTRs, as they were joined with protein-coding genes. We determined the coding potential of genes based on the codon substitution frequency (CSF²⁴) score and on the presence of sequence similarity with known proteins (SwissProt⁴³ database) or protein domains (Pfam-A⁴⁴ database). As *de novo* gene predictions can be incomplete or fragmented, we chose to assess the coding potential genome-wide rather than only for predicted exonic regions. We used the CSF score to define potential coding regions on a genome-wide scale, by scanning multiple species alignments (available through the UCSC Genome Browser⁵⁰). Genes were said to be potentially noncoding if they were sufficiently distant (>2 kb away) from a CSF-predicted coding region. Several distance thresholds were tested (Supplementary Discussion). We evaluated two additional methods (reading frame conservation⁵¹ and presence of open reading frames), but these performed less well and were not used in our final analyses (Supplementary Discussion). After estimating the coding potential independently for each species, we verified that the classifications of the members of homologous families agreed, thus further reducing the possibility of misclassifications.

Cross-species annotation projection. To reduce the inequalities in annotation depth among species, we projected the annotations across species and included the projected gene models in each species' data set. To do this, we searched for sequence similarity (blastn⁵²) between the complementary DNAs of a reference species and the repeat-masked genomes of the target species. We accepted projections without rearrangements or internal repeats and with inferred intron sizes below 100 kb. To avoid redundancy, the projections were added recursively, and only if they did not overlap with already annotated genes (Supplementary Methods).

We reduced the occurrence of fragmented gene predictions (a single gene is annotated as multiple neighbour loci), using a homology-directed defragmentation procedure that takes advantage of the availability of multiple species. We searched for sequence similarity (blastn⁵²) between the cDNA sequences of each species and classified as potentially 'fragmented' those neighbouring loci that could be reliably aligned with different regions of a single locus in another species (Supplementary Methods). For our final lncRNA data set, we excluded candidates that clustered with protein-coding sequences (thus reducing the possibility of misclassifying UTRs as lncRNAs) and we used 'de-fragmented' lncRNA annotations as controls for our analyses.

lncRNA filtering. We applied several filters to ensure reliability of the lncRNA data set. For species-specific lncRNAs we required: minimum exonic length 200 bp, at least 75% or 500 bp of non-overlapping exonic sequence, minimum 5 kb distance between lncRNA exons and Ensembl-annotated protein-coding gene exons, support by at least 5 non-strand-specific and 5 strand-specific reads (including splice junction reads), Ensembl gene biotypes (when available) 'lncRNA' or 'processed_transcript', no clustering (fragmentation) with protein-coding genes. For families

of lncRNAs with n species, we required noncoding classification with both CSF and sequence similarity in at least $n - 1$ species and with at least one of the two criteria in all species, minimum exonic length 200 bp (50 bp for projected genes) in all species, support by at least two reads in at least two species, minimum distance 5 kb to protein-coding gene exons for all species. For families that included Ensembl-annotated lncRNAs, we required the above criteria to be satisfied in at least $n - 1$ out of n species. For genes that overlapped on the antisense strand with other genes, we required support with strand-specific reads. We note that the list of lncRNAs provided for each species includes projected genes for which transcription evidence could not be found in the corresponding species, if these genes belonged to homologous families in which at least two species had transcription evidence.

Reconstruction of homologous lncRNA families and lncRNA evolutionary age. We reconstructed homologous lncRNA families based on DNA sequence similarity. We searched for similarity between the cDNA sequences of each species, using blastn⁵². As in Ensembl Compara⁵³, we extracted reciprocal best hits for each pair of species and significant self-hits for each species and we clustered genes with single-linkage. As lncRNAs can overlap with protein-coding genes or with transposable elements, we repeated the procedure after masking these regions, with no significant change. For improved sensitivity, we searched for alignments of wider regions, including 5 kb of flanking sequences, in whole genome alignments generated with blastz and multiz⁵⁴ (available through the UCSC Genome Browser). Potential homologues were called for alignments that mapped to a single target species gene. This homology inference was used as a control for our analyses. We inferred the minimum lncRNA evolutionary age with parsimony, based on the phylogenetic distribution of the species with transcription evidence in the homologous gene families. We note that this estimate represents a strict lower boundary, since transcription may be undetectable for lowly expressed genes, in particular for the species with lower overall read coverage.

In addition, we tested whether the absence of transcription in some species can be simply attributed to differences in RNA-seq read coverage, and we provide an additional estimate of the potential evolutionary age of lncRNAs. We estimated the proportion of mapped reads assigned to a given lncRNA, separately for each species and tissue. For each lncRNA family and for each tissue, we then estimated the minimum such proportion (p_{\min}), over all species in which the lncRNA was detected as transcribed. Given that for projected genes we often recover only a limited fraction of the original exonic length, the p_{\min} probability was further adjusted to reflect the difference in exonic length between the species with no transcription evidence and the species in which p_{\min} was observed (p_{\min} was multiplied by the ratio of the two exonic lengths). We then assessed the probability of observing 0 reads out of the total n mapped reads, given a theoretical detection probability of p_{\min} and assuming a binomial distribution, in the species for which transcription could not be detected in that tissue. If the tissue was not sampled for a given species (such as orangutan testes or non-human great ape placenta), the probability was set to 1. Finally, these probabilities were multiplied over all available tissues, to obtain a combined estimate of the likelihood that the absence of transcription in that species is simply due to differences in read coverage and/or annotated exonic length. We then re-estimated the evolutionary age of the lncRNA family, taking into account the phylogenetic distribution of the species in which transcription was either detected, or for which the absence of transcription could be attributed to read coverage and/or exonic length issues. This additional age estimate is termed the 'maximum' evolutionary age.

Selective constraint on DNA sequences. We computed average PhastCons²⁵ scores for exons and promoter regions, using genome-wide nucleotide resolution scores from the UCSC Genome Browser⁵⁰. We downloaded SNP data from the 1000 Genomes Project²⁶, we filtered the SNPs to exclude potential CpG sites and we computed the average derived allele frequency (DAF) for the African population. For DAF comparisons, we derived 95% confidence intervals from 100 bootstrap resampling replicates (parametric statistics cannot be applied due to non-normal distributions). We analysed only autosomal SNPs, residing in regions of moderate recombination (<2 cM per Mb), as measured using the DECODE⁵⁵ sex-averaged recombination maps in 20 kb windows centred on the SNP. As a neutral control, we resampled intergenic SNPs (>5 kb away from coding or noncoding genes) found in regions of similar recombination rates as lncRNAs (Supplementary Discussion). For overlapping genes (for example, sense-antisense transcripts), both measures of selective constraint were estimated using non-overlapping exonic regions.

Expression-level estimation and normalization. We estimated RPKM values from unambiguous read alignments obtained with TopHat⁴¹. To ensure an unbiased measurement, we considered only exonic regions that could be unambiguously assigned to a single gene. We also measured expression levels with Cufflinks v2.0.0, using all mapped reads, with the embedded multi-read and fragment bias correction methods (Supplementary Discussion). For projected genes, for which exon annotations are often incomplete, we included 1-kb flanking sequences on each side in the expression computation, if this extended region did not overlap with

other transcribed loci. We normalized expression levels among samples with a median scaling, using the 1000 least-varying genes as a reference, as described previously².

Transcription-factor-binding analysis. We used a genome-wide set of human transcription-factor-binding sites (~2.7 million sites, for 375 transcription factors), predicted *in silico*³⁰, as well as ChIP-seq peaks for 127 transcription factors (excluding those directly associated with PolII or PolIII) from the human ENCODE project³⁶. We analysed the occurrence of transcription-factor-binding sites or peaks in promoter regions, exclusively for genes for which the predicted transcription start site was found within 100 bp of a CAGE tag cluster (data from the FANTOM project⁵⁷). Two promoter region sizes were tested (2 kb and 5 kb), reaching similar conclusions. We also used ChIP-seq data for HNF4A and CEBPA for human and mouse³². We aligned promoter regions for the two species and considered that transcription-factor binding was conserved if peaks were found in both species within 10 kb of the aligned transcription start site. As a control, we analysed transcription-factor binding and binding conservation for 20,000 randomly drawn intergenic regions.

Expression evolution analyses. For the qualitative assessment of transcription conservation, we analysed 4,430 intergenic lncRNAs (>5 kb away from protein-coding genes) that had 1–1 orthologues in all primate species and which had at least 2 mapped reads in human in a pool of brain, cerebellum, heart, kidney and liver samples, as well as 2,492 human lncRNAs that overlapped on the antisense strand with exons of protein-coding genes, which had orthologues in at least one of the other species with strand-specific data (mouse, opossum, platypus, chicken and *Xenopus*). These antisense lncRNAs were further filtered to extract genes that were expressed in human brain and testes. We evaluated Spearman's correlation coefficients between pairs of samples, on lncRNA or protein-coding gene RPKM values. All available 1–1 orthologues were used. As a control for our expression evolution analyses, we resampled the same average number of reads per gene for each species and tissue, keeping the proportions among genes identical to the original distribution.

Tissue-specific expression. We evaluated the tissue specificity of the expression pattern with a previously proposed index⁴⁶, which varies between 0 for housekeeping genes and 1 for tissue-restricted genes:

$$\frac{\sum_{i=1}^n \left(1 - \frac{\exp_i}{\exp_{\max}}\right)}{n-1}$$

where n is the number of tissues, \exp_i is the expression value in tissue i , and \exp_{\max} the maximum expression level over all tissues. We used RPKM and log₂-transformed RPKM for expression values, reaching the same conclusions. The randomly expected proportion of conserved specificity across species was computed as the product of the observed proportions of tissue-specific genes in each species, for each tissue.

Reconstruction and analysis of the co-expression network. We reconstructed the evolutionarily conserved co-expression network for lncRNAs and protein-coding genes following a previously proposed method³⁵ (Supplementary Discussion). For each species and for each pair of genes (lncRNA or protein-coding), we computed the Pearson correlation coefficients of expression patterns. Given two homologous families, we examined whether the combination of correlation coefficients measured in each species was significantly higher or lower than expected by chance. The statistical tests were carried out by comparing the observed ranks of the correlation coefficients with a random n -dimensional order statistics³⁵. We computed correlations only for genes expressed in at least three samples for each species, and we computed P values only if correlations were evaluated in at least three species. We allow negative connections, which have lower than expected rank combinations. We considered only lncRNAs estimated to have originated in the Eutherian ancestor or earlier, but without requiring representatives in all descendant

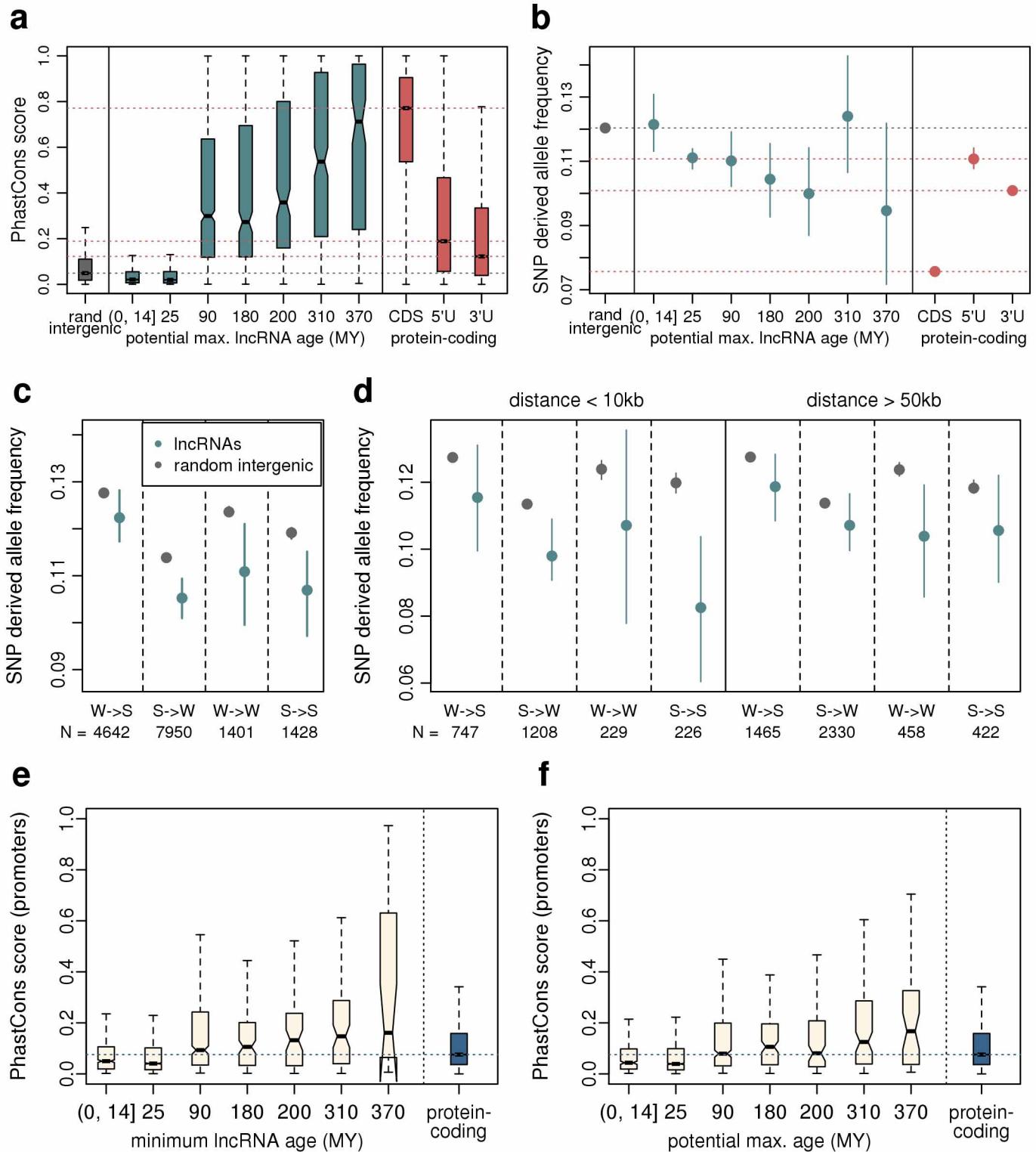
species. As P -value computations are highly time-consuming with a large number of species, analyses were carried out using a representative subset of seven species: human, macaque, mouse, opossum, platypus, chicken and *Xenopus*. For greater accuracy of the reconstruction we extended our in-house generated data set to include previously published, comparable RNA-seq samples (Supplementary Table 1). We visualized the network with Cytoscape⁴⁷ and we detected clusters of highly inter-connected genes with the Markov Cluster (MCL) algorithm³⁹.

Defining potential miRNA precursors. To search for lncRNAs that may promote transcription of miRNAs or are potentially processed into miRNAs, we extracted all miRNA hairpin sequences from miRBase⁵⁸ 18 and searched for sequence similarity (blastn⁵²) against all annotated gene regions, including 10 kb of flanking sequences. Genes with at least one miRNA hairpin alignment (95% identity, aligned on the entire length) on the same strand were considered potential miRNA precursors.

Statistical analyses. All statistical analyses and graphical representations (including gene expression clustering, principal component analysis, randomization tests for statistical significance) were done in R⁴⁸. For statistical tests involving the co-expression network, we generated a set of 100 randomized networks by permuting the gene identifiers of the nodes for each edge. The randomized networks had the same distribution of edges types (positive, negative, coding-coding, coding-noncoding), and the node degree was preserved. To test the significance of the network properties (for example, *cis* connections), we derived a P value by comparing the values observed in real and randomized networks. To compare the degrees of connectivity among gene types by controlling for unequal expression levels, we extracted lncRNAs with maximum expression levels (\log_2 RPKM) between 3 and 6, and divided them into 6 discrete expression classes ([3, 3.5], [3.5, 4], ..., [5.5, 6] \log_2 RPKM) (round brackets represent open (excluded) boundaries of intervals, square brackets represent closed (included) boundaries). We then drew transcription-factor and non-transcription-factor protein-coding genes matching the relative proportions of lncRNAs in each expression class. The resampling was repeated 100 times.

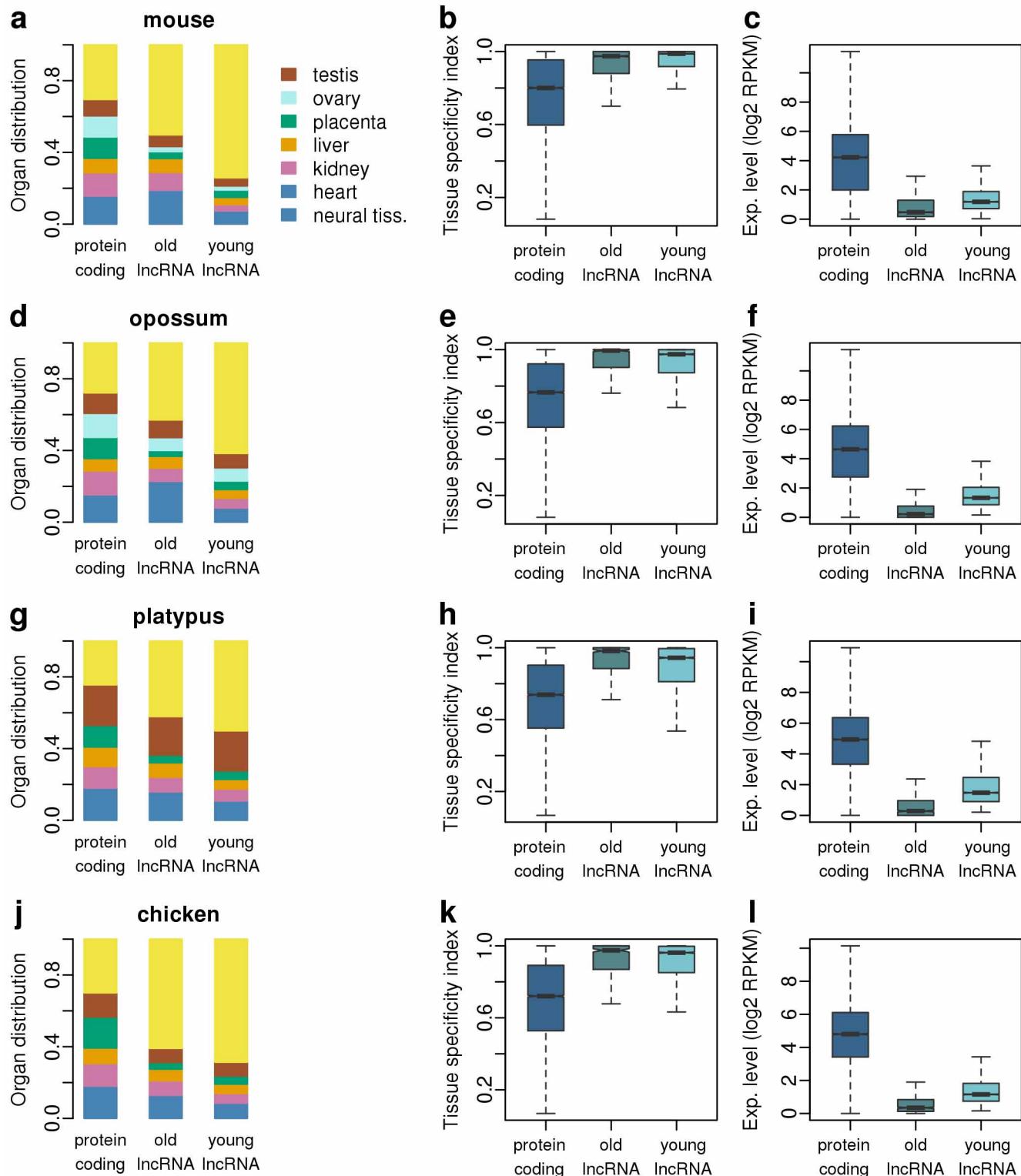
Data availability. The sequencing data have been submitted to GEO (accession GSE43520) and SRA (PRJNA186438 and PRJNA202404). The lncRNA annotations and homologous families have been made available on the publisher's website (Supplementary Data 1 and 2), as well as gene expression levels for lncRNAs and protein-coding genes (Supplementary Data 3) and miRNAs (Supplementary Data 4).

49. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.* **10**, R25 (2009).
50. Rhead, B. *et al.* The UCSC Genome Browser database: update 2010. *Nucleic Acids Res.* **38**, D613–D619 (2010).
51. Kellis, M., Patterson, N., Birren, B., Berger, B. & Lander, E. S. Methods in comparative genomics: genome correspondence, gene identification and regulatory motif discovery. *J. Comput. Biol.* **11**, 319–355 (2004).
52. Altschul, S. F., Gish, W., Miller, W., Myers, E. & Lipman, D. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
53. Vilella, A. J. *et al.* EnsemblCompara GeneTrees: complete, duplication-aware phylogenetic trees in vertebrates. *Genome Res.* **19**, 327–335 (2009).
54. Blanchette, M. *et al.* Aligning multiple genomic sequences with the threaded blockset aligner. *Genome Res.* **14**, 708–715 (2004).
55. Kong, A. *et al.* Fine-scale recombination rate differences between sexes, populations and individuals. *Nature* **467**, 1099–1103 (2010).
56. Gerstein, M. B. *et al.* Architecture of the human regulatory network derived from ENCODE data. *Nature* **489**, 91–100 (2012).
57. Carninci, P. *et al.* Genome-wide analysis of mammalian promoter architecture and evolution. *Nature Genet.* **38**, 626–635 (2006).
58. Kozomara, A. & Griffiths-Jones, S. miRBase: integrating microRNA annotation and deep-sequencing data. *Nucleic Acids Res.* **39**, D152–D157 (2011).



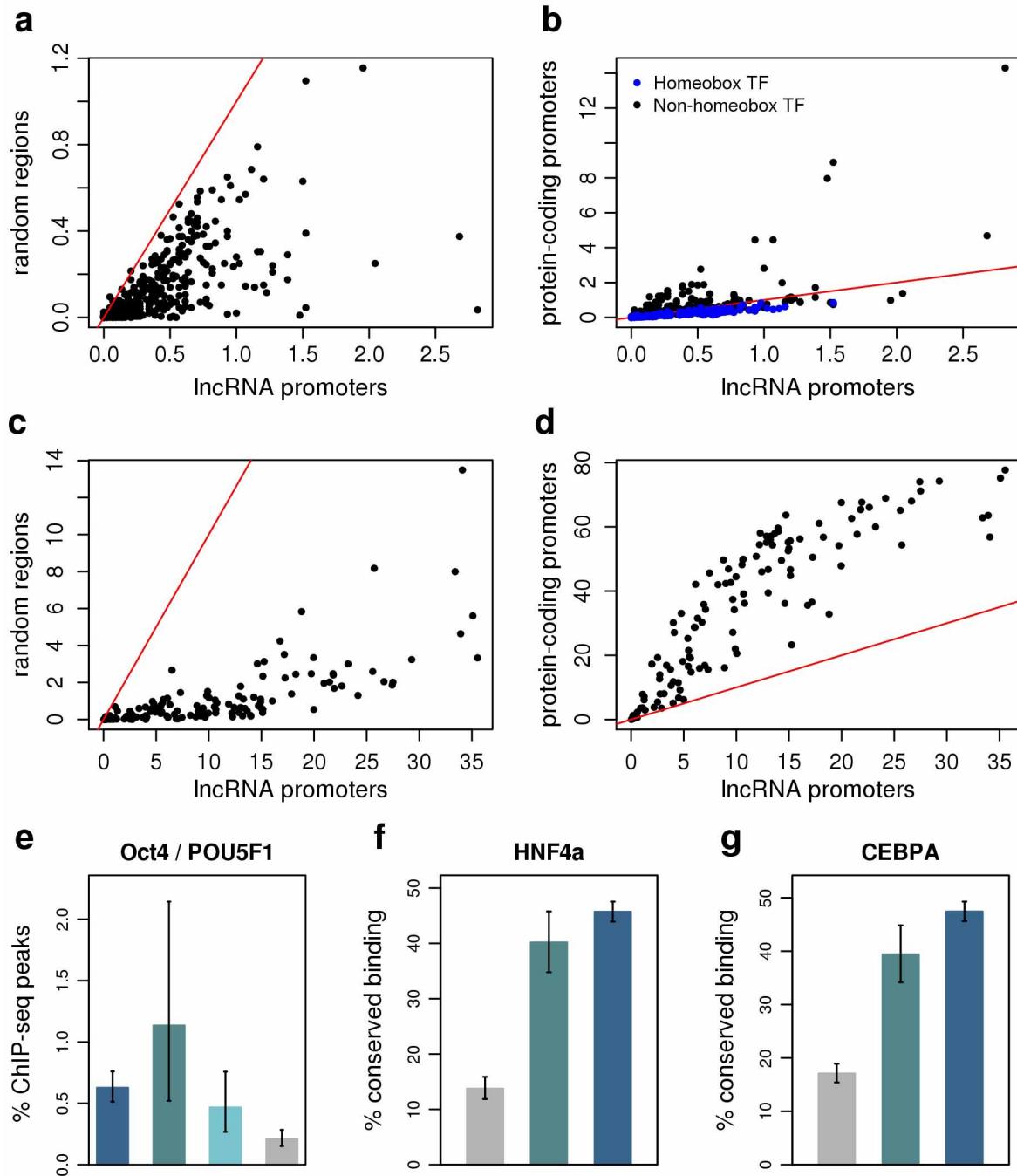
Extended Data Figure 1 | IncRNA evolutionary age and sequence conservation patterns. **a**, Exonic sequence conservation (mean placental PhastCons score), for random intergenic regions, IncRNA maximum evolutionary age classes, coding and untranslated exons of protein-coding genes. **b**, Mean DAF of autosomal non-CpG SNPs segregating in African populations (1000 Genomes project²⁰). Intergenic SNPs were randomly drawn in regions matching IncRNA recombination rates (Methods). **c**, Mean DAF for the four classes of mutation orientation (W to S (W→S) or AT to GC; S to W (S→W) or GC to AT; W to W (W→W), or AT to AT; and S to S (S→S), or GC to GC) for autosomal non-CpG SNPs found in primate-specific (age 25 Myr) IncRNA exonic regions (blue) or in intergenic regions with matching

recombination rates (grey). The W→S and S→W mutation classes are known to be affected by GC-biased gene conversion. **d**, Same as c but for IncRNAs that are found close to (left panel, maximum distance 10 kb) or far from (right panel, minimum distance 50 kb) Ensembl-annotated coding or noncoding genes. **e**, Mean placental PhastCons score for promoter regions (1 kb upstream) of IncRNA minimum evolutionary age classes (beige) and protein-coding genes (blue). **f**, Mean placental PhastCons score for promoter regions (1 kb upstream) of IncRNA maximum evolutionary age classes (beige) and protein-coding genes (blue). Error bars, 95% confidence intervals based on 100 bootstrap resampling replicates.



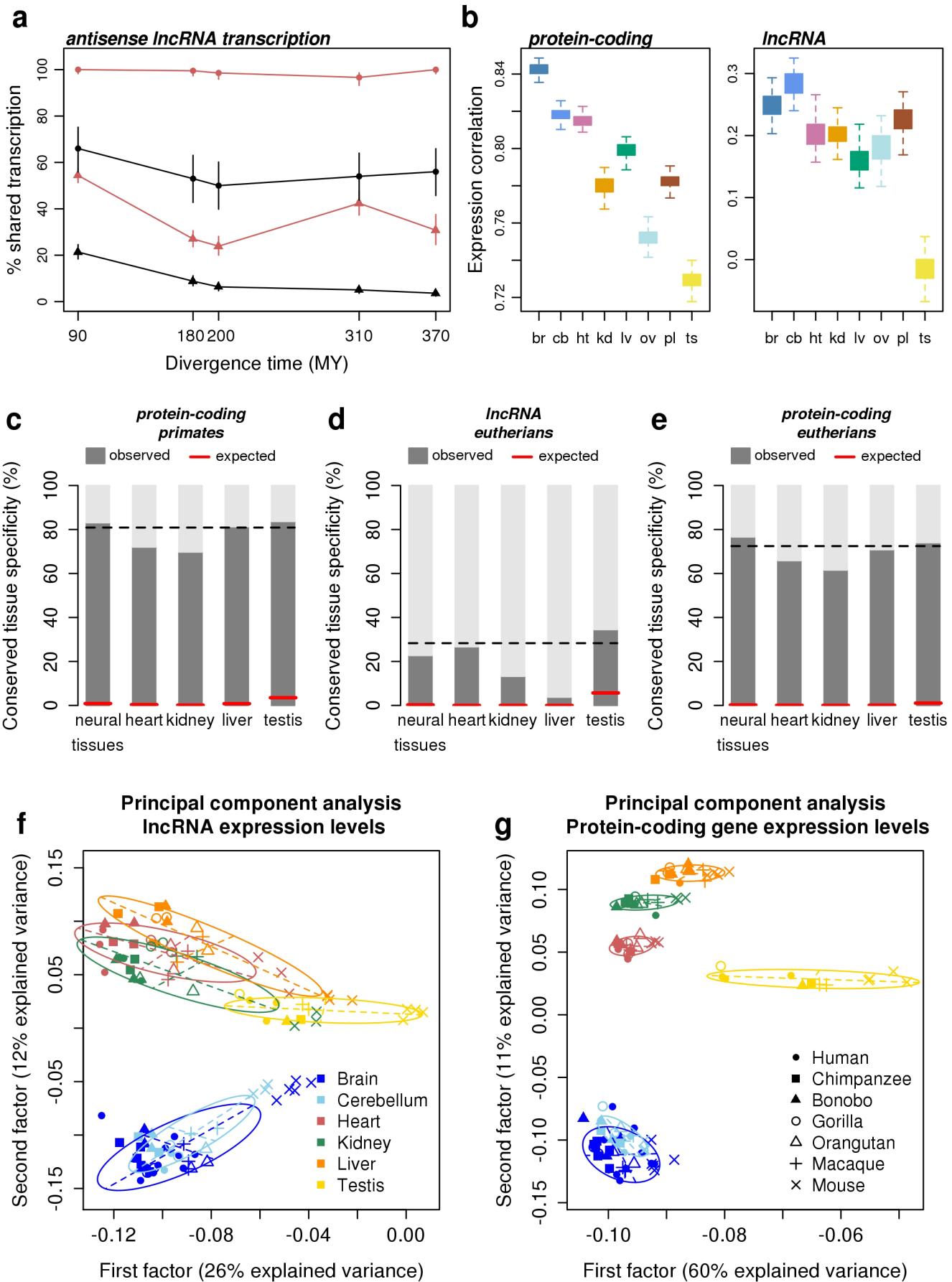
Extended Data Figure 2 | IncRNA expression patterns in four tetrapod species. **a**, Proportions of genes with observed maximum expression in different organs for mouse protein-coding genes, old lncRNAs (shared across at least two species) and young lncRNAs (species-specific). **b**, Tissue-specificity

index, for the same classes of mouse genes. Values close to 1 represent high tissue specificity. **c**, Distribution of the maximum expression level (\log_2 -transformed RPKM). **d–f**, Same as **a–c** but for the opossum. **g–i**, Same as **a–c** but for the platypus. **j–l**, Same as **a–c** but for the chicken.



Extended Data Figure 3 | Transcription-factor binding at lncRNA promoters. **a**, Comparison between the frequencies of *in silico*-predicted transcription-factor (TF)-binding sites in lncRNA promoters (2 kb upstream) and in random intergenic regions. **b**, Comparison between the frequencies of *in silico*-predicted TF-binding sites in lncRNA and protein-coding gene promoters (2 kb upstream). Homeobox TFs are shown in blue. **c**, Comparison between the frequencies of experimentally determined (ChIP-seq ENCODE) TF-binding sites in lncRNA promoters (2 kb upstream) and in random

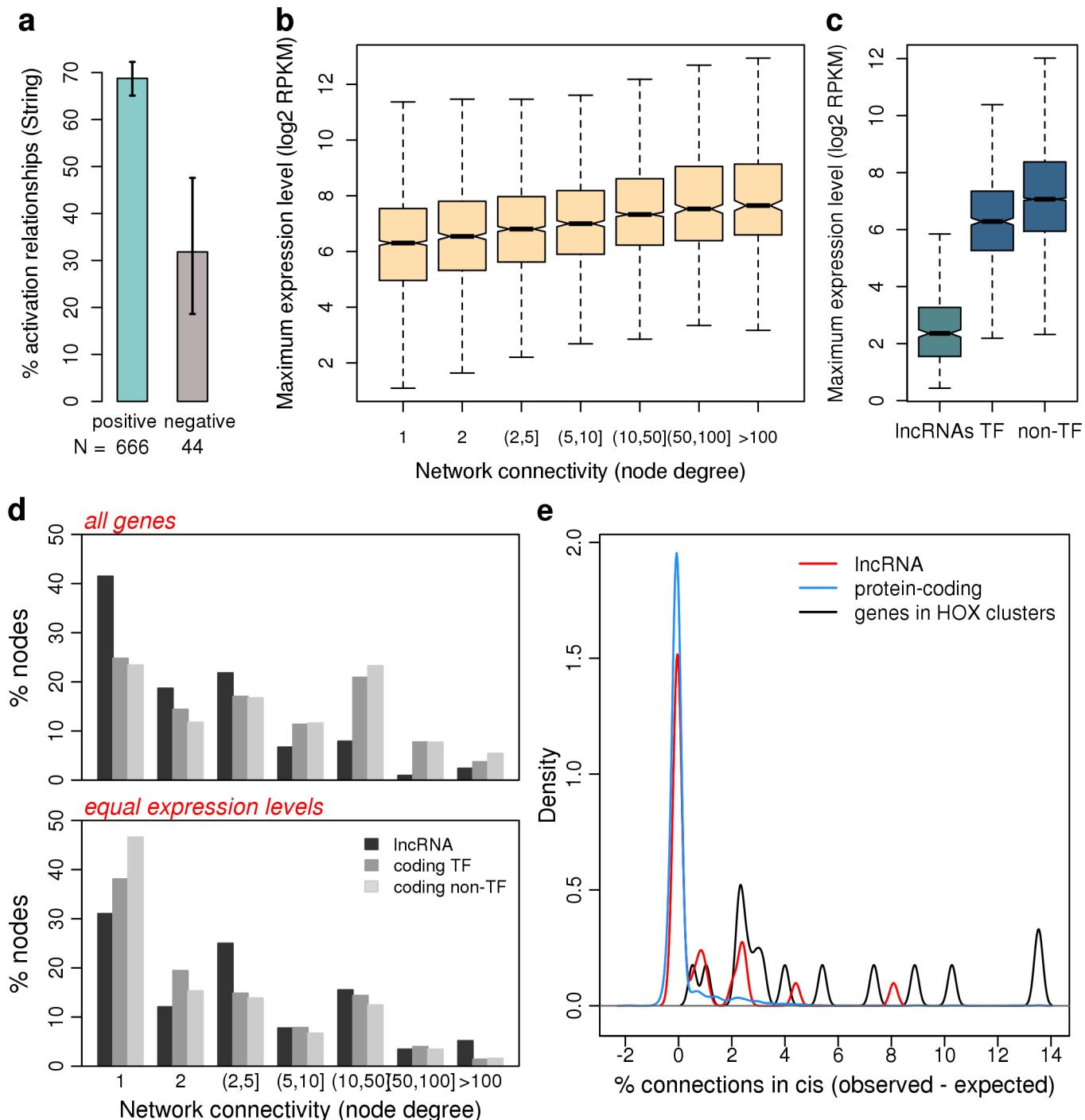
intergenic regions. **d**, Comparison between the frequencies of experimentally determined (ChIP-seq ENCODE) predicted TF-binding sites in lncRNA and protein-coding gene promoters (2 kb upstream). **e**, Frequency of binding (Encode ChIP-seq data) for OCT4 (also known as POU5F1). **f**, **g**, Proportion of HNF4A- CEBPA-binding events shared between human and mouse, for random intergenic regions, lncRNA (321 lncRNAs with binding events and liver expression, supported by CAGE data) and protein-coding gene promoters (5 kb upstream).



Extended Data Figure 4 | Evolution of lncRNA expression patterns.

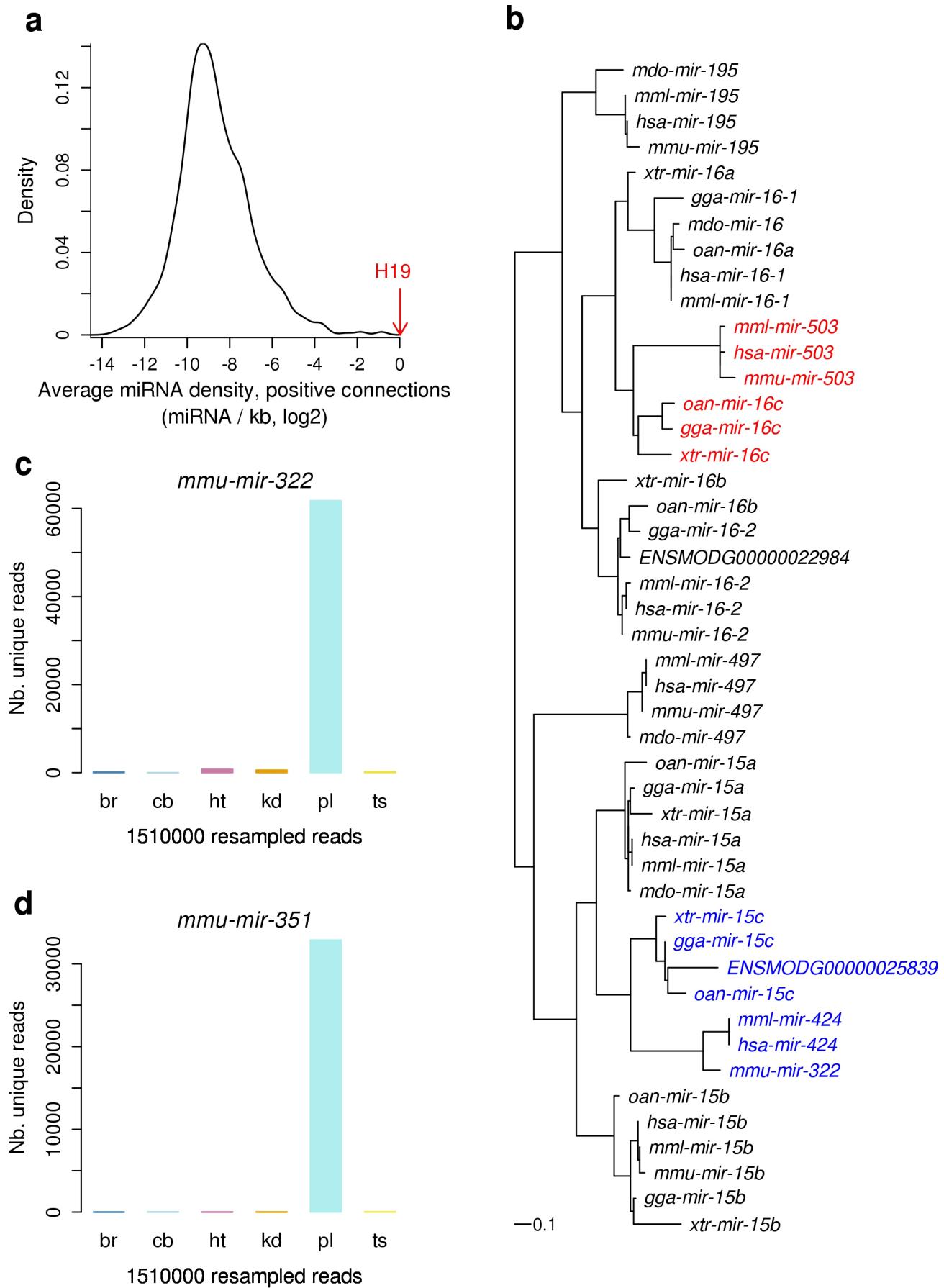
a, Percentage of human lncRNAs (found in antisense of protein-coding genes) that have transcription evidence in other species, as a function of the divergence time. Transcription evidence was assessed in a pool of brain and testes strand-specific RNA-seq data, for 2,535 human antisense lncRNAs that had 1–1 orthologues in at least one other species and transcription evidence in human (Methods). **b**, Spearman correlation of human and mouse expression levels, in different tissues. The boxplots represent the variation observed in 100 bootstrap replicates. **c**, Proportion of human organ-specific protein-coding genes (tissue-specificity index >0.9 , RPKM >0.1) for which the organ

specificity is shared across primates. Red lines, random expectation of shared organ specificity; horizontal black line, average conserved specificity for all organs. **d**, Proportion of human organ-specific lncRNAs (minimum evolutionary age >90 Myr, tissue-specificity index >0.9 , RPKM >0.1) for which the organ specificity is shared across eutherians. Red lines, random expectation of shared organ specificity; horizontal black line, average conserved specificity for all organs. **e**, Same as **c**, conservation across eutherian species. **f**, Principal component analysis of lncRNA expression levels for families of eutherian 1–1 orthologues. **g**, Principal component analysis of protein-coding gene expression levels for families of eutherian 1–1 orthologues.



Extended Data Figure 5 | Characteristics of the evolutionarily conserved co-expression network. **a**, Proportion of activation/inhibition relationships annotated in the String database, for positive and negative co-expression network connections. **b**, Gene expression levels (maximum over all available sample and species for each co-expression network node) for different network connectivity classes. **c**, Gene expression levels (maximum over all available sample and species for each co-expression network node) for connected IncRNAs, transcription factors (TFs) and non-TF protein-coding genes.

d, Network connectivity (node degree) for IncRNAs (black), transcription factors (medium grey) and for non-transcription factors protein-coding genes (light grey). Top, raw data; bottom, after correcting for expression level differences. **e**, Difference between observed and expected proportions of connections in *cis*, for IncRNAs (red), protein-coding genes (blue) and for genes found in HOX clusters (black). The expected proportions were computed through randomizations (Methods).



Extended Data Figure 6 | Expression patterns and sequence evolution of *H19X*-associated miRNAs. **a**, Distribution of the average embedded miRNA density (miRNA hairpins per kb, in the gene body or 10 kb downstream), for genes that are positively connected with each network node. Red arrow, average miRNA density for genes that are positively connected with *H19*. **b**, Maximum likelihood reconstruction of the phylogeny of the ancient *H19X*-associated miRNA family (representative members *miR-503*, *miR-322*, *miR-424*, *miR-15c*, *miR-16c*). miRNAs associated with *H19X* are displayed in red (subfamily containing *miR-503* and *miR-16c*) and blue (subfamily containing *miR-424*, *miR-322* and *miR-15c*). miRNA names are derived from

miRBase where available, including three-letter species abbreviations. Hsa, *Homo sapiens*; Mdo, *Monodelphis domestica* (opossum); Mml, *Macaca mulatta* (macaque); Mmu, *Mus musculus* (mouse); Oan, *Ornithorhynchus anatinus* (platypus); Gga, *Gallus gallus* (chicken), Xtr, *Xenopus tropicalis*. Ensembl identifiers are given for two opossum miRNAs. **c**, Expression pattern of the mouse miRNA *mmu-miR-322*, associated with *H19X*. The expression level was computed as the number of uniquely mapping reads per miRNA, after resampling the same number of reads per tissue. **d**, Same as **c** but for the mouse miRNA *mmu-miR-351*.

Extended Data Table 1 | Validation of the *de novo* detection and classification methods

(a)

species	protein-coding		lincRNA		processed transcript	
	partial	complete	partial	complete	partial	complete
human	17006 (88%)	9510 (49%)	942 (77%)	442 (36%)	4708 (54%)	2331 (27%)
chimp / bonobo	15457 (93%)	10292 (62%)	-	-	NA	NA
gorilla	14623 (92%)	9693 (61%)	-	-	NA	NA
orangutan	13222 (87%)	8676 (57%)	-	-	NA	NA
macaque	14617 (93%)	9235 (59%)	-	-	NA	NA
mouse	16824 (90%)	12072 (64%)	1000 (78%)	647 (51%)	1568 (64%)	1021 (41%)
opossum	12204 (95%)	7936 (62%)	-	-	NA	NA
platypus	9221 (90%)	3782 (37%)	-	-	NA	NA
chicken	13611 (89%)	9598 (63%)	-	-	NA	NA
Xenopus	13373 (89%)	7052 (47%)	-	-	NA	NA
average	90 %	57 %	78%	44%	56%	30%

(b)

species	protein-coding	lincRNA	processed transcript	tRNA, rRNA
human	19247 (92%)	721 (58%)	6710 (78%)	511 (96%)
chimp / bonobo	17537 (99%)	-	-	505 (97%)
gorilla	16519 (99%)	-	-	508 (97%)
orangutan	16514 (99%)	-	-	516 (97%)
macaque	16807 (100%)	-	-	696 (97%)
mouse	20651 (95%)	1043 (81%)	2062 (84%)	306 (96%)
opossum	13607 (99%)	-	-	170 (97%)
platypus	10521 (99%)	-	-	203 (99%)
chicken	14105 (86%)	-	-	201 (96%)
Xenopus	15492 (96%)	-	-	268 (99%)
average	96 %	70 %	79 %	97 %

a. Proportion of Ensembl-annotated (release 62) multi-exonic protein-coding genes, lncRNAs and processed transcripts recovered with our *de novo* detection methods. Partial overlap: number (percentage) of Ensembl-annotated multi-exonic genes for which at least half of the exons were recovered *de novo*. Complete: number (percentage) of multi-exonic genes for which all exons were recovered *de novo*. Protein-coding genes were filtered to retain those with 'known' or 'known by projection' gene status. **b.** Proportion of Ensembl-annotated protein-coding genes, lncRNAs, processed transcripts and other noncoding RNA genes (transfer RNA (tRNA), ribosomal RNA (rRNA)) that were correctly classified as coding or noncoding with our approach.

Extended Data Table 2 | LncRNA repertoires in 11 tetrapod species

(a)

species	total	orphan	in 1-1 fam.	intergenic	intragenic	de novo	known	projected
Hsa	14682	481	14201 (92%)	12286	2396	2030	4619 (3263)	8032
Ptr/Ppa	14654	347	14307 (90%)	12695	1959	3450	0 (0)	11203
Ggo	14258	501	13757 (91%)	12546	1712	4530	0 (0)	9726
Ppy	13756	229	13527 (61%)	12189	1566	1099	0 (0)	12655
Mml	15280	1060	14220 (88%)	13463	1817	5931	0 (0)	9348
Mmu	10850	7895	2955 (79%)	9045	1805	7485	1580 (1580)	1784
Mdo	8039	6579	1460 (56%)	7171	868	6815	0 (0)	1223
Oan	6889	5890	999 (59%)	6576	313	6097	0 (0)	790
Gga	5412	4730	682 (71%)	4951	461	4857	0 (0)	554
Xtr	3296	3059	237 (49%)	3133	163	3091	0 (0)	204

(b)

species	total	orphan	in 1-1 fam.	intergenic	intragenic	de novo	known	projected
Hsa	12677	8080	4597 (61%)	6823	5854	4161	5646 (4485)	2869
Mmu	15934	11604	4330 (70%)	9138	6796	11448	2482 (2482)	2003
Mdo	9635	6898	2737 (42%)	5847	3788	7416	0 (0)	2218
Oan	8037	5941	2096 (42%)	5704	2333	6346	0 (0)	1689
Gga	8358	6733	1625 (47%)	5095	3263	7099	0 (0)	1258
Xtr	5314	4557	757 (31%)	3748	1566	4683	0 (0)	630

a. IncRNA repertoires determined using all RNA-seq samples available for each species, including both strand-specific and non-strand-specific data. Gga, *Gallus gallus* (chicken); Ggo, *Gorilla gorilla*; Hsa, *Homo sapiens*; Mdo, *Monodelphis domestica* (opossum); Mml, *Macaca mulatta* (macaque); Mmu, *Mus musculus* (mouse); Oan: *Ornithorhynchus anatinus* (platypus); Ppa, *Pan paniscus* (bonobo); Ppy, *Pongo pygmaeus* (orangutan); Ptr, *Pan troglodytes* (chimpanzee); Xtr: *Xenopus tropicalis*. Orphan, IncRNAs for which no orthologues could be detected; 1–1 fam, IncRNAs found in 1–1 orthologous families; Intergenic, IncRNAs found >5 kb away from Ensembl-annotated protein-coding genes; Intragenic, IncRNAs that overlap with Ensembl-annotated protein-coding genes on the opposite strand, but are found at least 5 kb away from their exons; De novo, previously unknown IncRNAs detected with RNA-seq; Known, IncRNAs that confirm previously known loci (including GENCODE/Ensembl human and mouse annotations (numbers in parentheses) and a set of 8,264 human IncRNAs previously detected with RNA-Seq⁴). Projected, IncRNAs derived from cross-species annotation projections. **b.** IncRNA repertoires determined with strand-specific data.

Extended Data Table 3 | LncRNA evolutionary age estimates and synteny conservation

(a)

	hominins	african apes	great apes	primates	eutherians	therians	mammals	amniotes	tetrapods
hominins	47.1%	1%	30.8%	20.7%	0.5%	0%	0%	0%	0%
african apes	0%	32.5%	40.4%	25.4%	0.4%	0.8%	0%	0.2%	0.4%
great apes	0%	0%	80.8%	16.6%	1.2%	0.8%	0.2%	0.2%	0.2%
primates	0%	0%	0%	97.1%	1.3%	0.7%	0.4%	0.3%	0.1%
eutherians	0%	0%	0%	0%	88.6%	5.7%	2.9%	2%	0.9%
therians	0%	0%	0%	0%	0%	81.2%	12.2%	3.4%	3.2%
mammals	0%	0%	0%	0%	0%	0%	92.5%	4.4%	3.1%
amniotes	0%	0%	0%	0%	0%	0%	0%	91.3%	8.7%
tetrapods	0%	0%	0%	0%	0%	0%	0%	0%	100%

(b)

	Hsa	Ptr / Ppa	Ggo	Ppy	Mml	Mmu	Mdo	Oan	Gga	Xtr
Hsa		94.7%	93.6%	89.4%	90.8%	90.9%	67.2%	51.1%	90.1%	79.4%
Ptr / Ppa	97.5%		94.2%	90.5%	92.1%	91.6%	68.1%	55.6%	92.2%	81.5%
Ggo	95.9%	94.9%		89.8%	90.2%	91.2%	68.6%	51.3%	89.4%	82.8%
Ppy	95.4%	94.2%	92.7%		91%	92.3%	69.2%	53.5%	89.5%	87%
Mml	94.4%	93.4%	91.7%	88.2%		89.6%	67.4%	48.8%	90.9%	78.1%
Mmu	86%	83.6%	81.3%	77.9%	79.8%		60.2%	50.4%	87.2%	82.6%
Mdo	88.5%	87.9%	87%	83.6%	84.2%	89.8%		55.2%	90.1%	82%
Oan	59.6%	58.8%	55.9%	54%	55.8%	60.6%	44.7%		84%	75.5%
Gga	54.8%	51.2%	47.6%	47.5%	50%	58.1%	40.9%	48.2%		72.7%
Xtr	61.3%	59.5%	55.7%	51.8%	56.8%	61.1%	46.2%	46.6%	79.4%	

a. Comparison between the minimum evolutionary age of lncRNA families (requiring transcription evidence in all species), and the maximum potential evolutionary age (Methods). The numbers represent the percentage of cases in which a given 'minimum age' estimate (rows) is associated with a given 'maximum age' estimate (columns). **b.** Synteny conservation for pairs of neighbouring genes that contain at least one lncRNA. The neighbouring gene pairs in the reference species (see Extended Data Table 2 legend) were genes with 1–1 orthologues in the target species, separated by 5–100 kb in the reference genome. The numbers represent the percentage of neighbouring gene pairs in the reference species (rows) for which the 1–1 orthologues in the target species (columns) were found on the same chromosome, separated by at most 100 kb.