1  **Comparative transcriptomics analyses across species, organs and developmental stages reveal**

2  **functionally constrained lncRNAs**

3

4  Fabrice Darbellay[1,$,*] and Anamaria Necsulea[1,2,*]

5

6

7  [1]School of Life Sciences, École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

8  [2]Université de Lyon, Université Lyon 1, CNRS, Laboratoire de Biométrie et Biologie Évolutive UMR

9  5558, F-69622 Villeurbanne, France

10

11  [$]Present address: Environmental Genomics and Systems Biology Division, Lawrence Berkeley

12  National Laboratory, 1 Cyclotron Road, Berkeley, California 94720, USA.

13

14

15  [*]Corresponding authors:

16  Fabrice Darbellay (fabrice.darbellay@epfl.ch)

17  Anamaria Necsulea (anamaria.necsulea@univ-lyon1.fr)

18

19  Running title: Functionally constrained lncRNAs in embryonic development

20

21  Keywords: long non-coding RNAs; evolution; development; comparative transcriptomics.

22

23

24

1

**Abstract**

**Background** Transcription of long non-coding RNAs (lncRNAs) is pervasive, but their functionality is disputed. As a class, lncRNAs show little selective constraint and negligible phenotypic effects upon perturbation. However, key biological roles were demonstrated for individual lncRNAs. Most validated lncRNAs were implicated in gene expression regulation, in pathways related to cellular pluripotency, differentiation and organ morphogenesis, suggesting that functional lncRNAs may be more abundant in embryonic development, rather than in adult organs.

**Results** Here, we perform a multi-dimensional comparative transcriptomics analysis, across five developmental time-points (two embryonic stages, newborn, adult and aged individuals), four organs (brain, kidney, liver and testes) and three species (mouse, rat and chicken). Overwhelmingly, lncRNAs are preferentially expressed in adult and aged testes, consistent with the presence of permissive transcription during spermatogenesis. LncRNAs are often differentially expressed among developmental stages and are less abundant in embryos and newborns compared to adult individuals, in agreement with a requirement for tighter expression control and less tolerance for noisy transcription early in development. However, lncRNAs expressed during embryonic development show increased levels of evolutionary conservation, both in terms of primary sequence and of expression patterns, and in particular at their promoter regions. We find that species-specific lncRNA transcription is frequent for enhancer-associated loci and occurs in parallel with expression pattern changes for neighboring protein-coding genes.

**Conclusions** We show that functionally constrained lncRNA loci are enriched in developing organ transcriptomes, and propose that many of these loci may function in an RNA-independent manner.

47 **Background**

48 Long non-coding RNAs (lncRNAs, loosely defined as transcripts that lack protein-coding potential, at

49 least 200 nucleotides long) are an excellent illustration of the ongoing conceptual tug-of-war between

50 biochemical activity and biological function (Graur et al. 2013; Doolittle 2018). The development of

51 sensitive transcriptome exploration techniques led to the identification of thousands of lncRNA loci in

52 vertebrates . While this ever wider class of transcripts includes well-studied lncRNAs with undisputed

53 biological roles, such as *Xist* (Brown et al. 1991) or *H19* (Brannan et al. 1990), experimental validations

54 are lacking for the great majority of lncRNAs and their functionality is controversial.

55

56 The first functional characterizations of individual lncRNAs forged the idea that these non-coding

57 transcripts are important contributors to gene expression regulatory networks. This has been

58 unequivocally proven for some lncRNAs, such as *Xist*, whose transcription and subsequent coating of

59 the X chromosome triggers a complex chain of molecular events leading to X inactivation in placental

60 mammals (Gendrel and Heard 2014). Additional proposed mechanisms for gene expression regulation

61 by lncRNAs included directing chromatin-modifying complexes at specific genomic locations, to control

62 gene expression in *trans* (Rinn et al. 2007); providing decoy targets for microRNAs (Cesana et al. 2011);

63 enhancing expression of neighboring genes through an RNA-dependent mechanism (Ørom et al. 2010).

64 These initial studies generally asserted that the biological function of lncRNA loci is directly carried out

65 by the transcribed RNA molecule. However, it rapidly became evident that in some cases the function

66 resides in the act of transcription at a given genomic location, rather than in the product of

67 transcription (Latos et al. 2012). In recent years, this view has gained ground, with several publications

68 showing that lncRNA transcripts are not required, and that instead biological functions are carried out

69 by other elements embedded in the lncRNA genomic loci (Bassett et al. 2014). For example, it was

70 recently shown that transcription of the *Linc-p21* gene, originally described as a *cis*-acting enhancer

71 lncRNA, is not needed to regulate neighboring gene expression (Groff et al. 2016). Genetic engineering

72 of multiple lncRNA loci in mouse likewise showed that lncRNA transcripts are dispensable, and that

73    gene expression regulation by lncRNA loci is instead achieved by the process of lncRNA transcription

74    and splicing, or by additional regulatory elements found in lncRNA promoters (Engreitz et al. 2016;

75    Anderson et al. 2016). Furthermore, some attempts to look for lncRNA function through genetic

76    engineering approaches showed that the tested lncRNA loci are altogether dispensable (Amândio et

77    al. 2016; Zakany et al. 2017; Goudarzi et al. 2019). These recent observations signal a paradigm shift in

78    lncRNA biology, as it is increasingly acknowledged that, even when phenotypic effects can be

79    unambiguously mapped to lncRNA loci, the underlying biological processes are not necessarily driven

80    by the lncRNA transcripts themselves.

81

82    Importantly, this new perspective on lncRNA biology had been predicted by evolutionary analyses,

83    which have long been used to evaluate the functionality of diverse genomic elements (Haerty and

84    Ponting 2014; Ulitsky 2016). Evolutionary studies of lncRNAs in vertebrates all agree that the extent of

85    selective constraint on lncRNA primary sequences is very low, though significantly above the genomic

86    background (Ponjavic et al. 2007; Kutter et al. 2012; Necsulea et al. 2014; Washietl et al. 2014; Hezroni

87    et al. 2015). These observations are compatible with the hypothesis that many of the lncRNAs detected

88    with sensitive transcriptomics techniques may be non-functional noise (Ponjavic et al. 2007), but may

89    also indicate that lncRNA functionality does not reside in the primary transcribed sequence. In

90    contrast, mammalian lncRNA promoters show higher levels of sequence conservation, similar to

91    protein-coding gene promoters, as expected if they carry out enhancer-like regulatory functions

92    independently of the transcribed RNA molecule. Moreover, it was previously reported that, in multi-

93    exonic lncRNAs, splicing signals are more conserved than the rest of the exonic sequence (Schüler et

94    al. 2014; Haerty and Ponting 2015), which is compatible with the recent finding that lncRNA splicing

95    can contribute to neighboring gene regulation (Engreitz et al. 2016). Thus, detailed evolutionary

96    analyses of lncRNA loci can bring important insights into their functionality, and can help to prioritize

97    candidates for experimental validation.

98

99    At present, comparative transcriptomics analyses in vertebrates agree that the extent of evolutionary

100   conservation of lncRNA sequences and expression patterns is very limited. However, these studies

101   were so far restricted to adult organ transcriptomes. In particular, it was shown that most known

102   vertebrate lncRNAs are active in adult testes and thus likely during spermatogenesis, a process

103   characterized by a permissive chromatin environment, which can promote non-functional

104   transcription (Soumillon et al. 2013). The resulting lncRNA datasets may thus be enriched in non-

105   functional transcripts. Additional lines of evidence suggest that the search for functional lncRNAs

106   should be extended beyond adult organ transcriptomes. For example, involvement in developmental

107   phenotypes was proposed for many experimentally-tested lncRNAs (Sauvageau et al. 2013; Ulitsky et

108   al. 2011; Grote et al. 2013), and an enrichment for developmental transcription factor binding was

109   reported for the promoters of highly conserved lncRNAs (Necsulea et al. 2014). These observations

110   motivated us to add a temporal dimension to comparative lncRNA transcriptomics studies. Therefore,

111   we characterize here the lncRNA transcriptomes of two model mammalian species (mouse and rat), in

112   four major organs (brain, kidney, liver and testes), across five developmental stages that cover the

113   entire lifespan of the individuals (including two embryonic stages, newborn, young adult and aged

114   individuals). To gain a deeper evolutionary perspective, we generate similar data for embryonic stages

115   of chicken somatic organs. We analyze the spatial and temporal expression patterns of protein-coding

116   and lncRNA genes, in conjunction with their evolutionary conservation. We find that, while lncRNAs

117   are overall poorly conserved among species in terms of primary sequence or expression patterns,

118   higher frequencies of evolutionarily constrained lncRNAs are observed in embryonic transcriptomes.

119   For many of these loci, biological function may be RNA-independent, as the highest levels of sequence

120   conservation are observed on promoter regions and on splice signals, rather than on lncRNA exonic

121   sequence. Our results are thus compatible with unconventional, RNA-independent functions for

122   lncRNA loci, in particular for those that are expressed during embryonic development.

5

123 **Results**

124 _Comparative transcriptomics across species, organs and developmental stages_

125 To study protein-coding and lncRNA expression patterns across both developmental and evolutionary

126 time, we generated RNA-seq data for mouse and rat, for four major organs (brain, kidney, liver and

127 testes) and five developmental time points, including two embryonic stages, newborn, young and aged

128 adult individuals (Figure 1A, Supplementary Table 1, Methods). The selected time points allow us to

129 obtain a broad view of major organ ontogenesis and to capture drastic physiological changes during

130 development (Theiler 1989). We chose to include in our study both young adult (8-10 weeks old) and

131 aged adult individuals (12 to 24 months old), to investigate transcriptomic changes that occur later in

132 life, thus completing our overview of the temporal patterns of gene expression variation. At the

133 earliest embryonic stage (day 13.5 post-conception for mouse, day 15 for rat), only three of the four

134 studied organs, with the exception of the testes, are well differentiated and large enough to be readily

135 dissected. Our experimental design for mouse and rat thus comprises 19 organ / developmental stage

136 combinations. Although most of our study relies on mouse-rat comparisons, to obtain a broader

137 evolutionary perspective we generated comparable RNA-seq data for the chicken, for the two earliest

138 developmental stages (Figure 1A, Supplementary Table 1). We obtained between 2 and 4 biological

139 replicates for each species/organ/developmental stage combination (Supplementary Table 1).

140 Additional RNA-seq samples from previous publications were included in the lncRNA annotation

141 pipeline, to increase detection sensitivity (Supplementary Table 2, Methods).

142

143 The organs and developmental stages included in our study differ greatly in terms of their cellular

144 composition diversity. To verify that our whole-organ RNA-seq data reflects cellular composition

145 heterogeneity, we assessed the expression patterns of cell population markers derived from single-

146 cell transcriptomics studies (Tabula Muris Consortium 2018; Green et al. 2018) in our samples (Figure

147 1B, Supplementary Table 3). This analysis confirms that our transcriptome collection reflects expected

148 developmental patterns. For example, mature oligodendrocyte cell markers are systematically highly

6

149   expressed in adult brain, while oligodendrocyte precursor markers are more highly expressed in the

150   earliest developmental stages (Figure 1B). Similarly, *Neurod6,* a gene involved in neuronal

151   differentiation (Kathleen Baxter et al. 2009), is preferentially expressed in embryonic and newborn

152   brain. Moreover, spermatogenesis-specific markers are enriched in adult but not in embryonic and

153   newborn testes, while markers for somatic cells (Leydig, Sertoli cells) are expressed earlier during

154   testes development (Figure 1B). Immune cell markers tend to be more broadly shared across organs

155   and developmental stages, but show strongest expression in the late embryo and newborn liver (Figure

156   1B), consistent with this organ's crucial role in establishing immunity (Nakagaki et al. 2018). In general,

157   adult organ transcriptomes contain higher numbers of expressed cell type-specific markers (Figure 1B).

158   However, as these genes were defined based on adult organ data, this observation may indicate that

159   cell sub-populations that are specific to embryonic organs are under-represented in this marker set,

160   rather than reflecting the true cellular diversity at different developmental stages.

161

162   We note that, in some cases, the cell type-specific markers predicted by single-cell transcriptomics

163   studies have seemingly unexpected expression patterns in our whole-organ RNA-seq collection. For

164   example, the expression of *Parvalbumin* (*Pvalb*), which was proposed as a marker for collecting duct

165   epithelial cells in the kidney (Tabula Muris Consortium 2018), is highest in the adult and aged brain,

166   for both mouse and rat (Figure 1B). Likewise, cellular retinoic acid binding protein 1 (*Crabp1*), which

167   was predominantly detected  in spermatogonia in a single-cell transcriptomics study of mouse testes

168   (Green et al. 2018), is preferentially expressed in mid-stage embryonic kidney in our samples (Figure

169   1B). These apparent discrepancies likely reflect the pleiotropic nature of genes, as well as the presence

170   of similar cell types across organs with distinct physiological functions (Arendt et al. 2016).

171

172   Finally, the genes proposed as markers for major cell types generally behave similarly in mouse and

173   rat, although some species-specific patterns can be observed, in particular for immune cell markers

174   (Figure 1B). Likewise, for those genes that had orthologues in the chicken, expression patterns are

7

175     generally similar among species, with higher between-species divergence for immunity-related genes

176     (Supplementary Figure 1). This observation confirms that the organs and developmental stages

177     selected for our integrative transcriptomics study are comparable across species.

178

179     Overall, these results indicate that our whole-organ transcriptomics collection provides a good

180     overview of the cell composition changes that occur during development, and enables meaningful

181     comparisons across species.

182     _Variations in transcriptome complexity among organs and developmental stages_

183     We next sought to assess transcriptome complexity in different organs across developmental stages,

184     for both protein-coding genes and lncRNAs. To predict lncRNAs, we used the RNA-seq data to

185     reconstruct gene models with StringTie (Pertea et al. 2015), building on existing genomic annotations

186     (Cunningham et al. 2019). We verified the protein-coding potential of newly annotated transcripts,

187     based on the codon substitution frequency score (Lin et al. 2007, 2011) and on sequence similarity

188     with known proteins, and we applied a stringent series of filters to reduce contaminations from un-

189     annotated protein-coding UTRs and other artefacts (Methods). We thus obtain a total of 18,858

190     candidate lncRNAs in the mouse, 20,159 in the rat and 5,496 in the chicken, including both newly-

191     annotated and previously known lncRNAs transcribed in our samples (Supplementary Dataset 1). We

192     note that many of these candidate lncRNAs are expressed at very low levels. When imposing a

193     minimum normalized expression level (transcript per million, or TPM) at least equal to 1, in at least

194     one sample, the numbers of candidate lncRNAs falls to 12,199, 15,319 and 2,892 in the mouse, rat and

195     chicken, respectively (Supplementary Datasets 2-3, Supplementary Table 4).

196

197     The differences in lncRNA content among species may reflect discrepancies in RNA-seq read coverage

198     and sample distribution, as well as genome sequence and annotation quality. To correct for the effect

199     of RNA-seq read coverage, we down-sampled the RNA-seq data to obtain the same number of uniquely

200     mapped reads for each organ/developmental stage combination within each species (Methods). After

8

201  this equalizing procedure, the number of detectable protein-coding genes (supported by at least 10

202  uniquely mapped reads) still shows broad variations among organs and developmental stages, with

203  the highest numbers of genes detected in the testes, for all time points (Figure 1C). Large numbers of

204  protein-coding genes (between 12,800 and 16,700) are detected in all samples. In contrast, for

205  lncRNAs, the pattern is much more striking: the young and aged adult testes express between 11,000

206  and 12,000 lncRNAs, in both mouse and rat, while in somatic organs and earlier developmental stages

207  we can detect only between 1,800 and 4,800 lncRNAs (Figure 1D). This observation is in agreement

208  with previous findings indicating that the particular chromatin environment of the adult testes, and in

209  particular of spermatogenesis-specific cell types, is extraordinarily permissive to transcription

210  (Soumillon et al. 2013). Interestingly, the numbers of protein-coding genes detectable in each organ

211  also varies among developmental stages. In young and aged adult individuals, the brain shows the

212  second-highest number of expressed protein-coding genes, after the testes, as previously observed

213  (Soumillon et al. 2013; Ramsköld et al. 2009). However, in embryonic and newborn samples, the kidney

214  expresses higher numbers of protein-coding genes than the brain (Figure 1C).

215  *Developmental expression patterns are well conserved among species for protein-coding genes*

216  Broad patterns of transcriptome evolution are already visible in our analyses of cell type specific

217  markers and of transcriptome complexity: individual gene expression profiles and numbers of

218  expressed genes are generally similar between mouse and rat, while more divergence is observed

219  between the two rodent species and the chicken (Figure 1B-D, Supplementary Figure 1).  To further

220  explore the evolution of developmental gene expression patterns, we performed a principal

221  component analysis (PCA) on normalized, log-transformed TPM values for 10,363 protein-coding

222  genes shared among the three species (Methods, Figure 2A). This analysis revealed that the main

223  source of gene expression variability among species, organs and developmental stages is the

224  distinction between adult and aged testes and the other samples, which are separated on the first PCA

225  axis (Figure 2A). In contrast, embryonic and newborn testes are grouped with kidney samples from

226  similar developmental stages, in agreement with the common developmental origin of the kidney and

9

227    the gonads (McMahon 2016). The first axis of the PCA, which explains 67% of the total expression

228    variance, also correlates with the developmental stage: samples derived from adult and aged

229    individuals have higher coordinates on this axis than embryonic and newborn samples, for mouse and

230    rat (Figure 2A). The second PCA axis (10% explained variability) mainly reflects the difference between

231    brain and the other organs (Figure 2A). While mouse and rat samples are generally undistinguishable,

232    the PCA confirms that there is considerably higher expression divergence between chicken and the

233    two rodent species (Figure 2A). However, differences among major organs are stronger than

234    differences among species, even at these broad evolutionary distances: brain samples all cluster

235    together, irrespective of the species of origin, and are clearly separated from kidney and liver samples

236    on the second PCA axis (Figure 2A). Interestingly, within the brain cluster, embryonic chicken samples

237    tend to be closer to adult and aged rodent brains than to embryonic or neonate samples (Figure 2A).

238

239    These broad patterns of gene expression variations among species, organs and developmental stages

240    are confirmed by a hierarchical clustering analysis based on Spearman's correlation coefficients

241    between pairs of samples (Figure 2B). The strongest clustering is observed for adult and aged testes

242    samples, followed by a robust grouping of brain samples, irrespective of the species (Figure 2B).

243

244    For the mouse and rat, we could delve deeper into the evolutionary conservation of gene expression

245    patterns, by asking whether variations among developmental stages are shared between species. We

246    used models from the DESeq2 (Love et al. 2014) package to detect differential gene expression among

247    developmental stages, independently for each species and organ (Supplementary Dataset 4,

248    Methods). As expected given the wide range of developmental stages that we sampled, the great

249    majority of protein-coding genes are significantly differentially expressed (FDR<0.01) among stages, in

250    each organ (Supplementary Dataset 4). We selected orthologous protein-coding genes that are

251    differentially expressed (DE) in both species, and used the K-means clustering algorithm to discover

252    broad patterns of variations among species and stages (Methods). In general, differentially expressed

253 genes show parallel patterns of variation among developmental stages in mouse and rat, for somatic

254 organs (Figure 2C, Supplementary Figure 2). Genes with shared patterns of variation among

255 developmental stages are enriched in organ-specific functional categories, such as nervous system

256 development and axon guidance for the first cluster of genes presented in Figure 2C, which have high

257 expression levels in the embryonic and newborn samples (Supplementary Dataset 4). While temporal

258 expression variations are generally conserved between species for brain, kidney and liver, almost 25%

259 of differentially expressed genes show different trends for mouse and rat in the testes (Supplementary

260 Figure 2). These sets of genes do not show any strong functional enrichment (Supplementary Dataset

261 4). This pattern confirms previous reports indicating that gene expression evolution is faster in the

262 adult testes (Brawand et al. 2011), and extends them by showing that patterns of variations among

263 developmental stages are often species-specific in the testes.

264 *Spatial and temporal expression pattern differences between protein-coding genes and lncRNAs*

265 We next compared spatial and temporal expression patterns between protein-coding genes and

266 lncRNAs. In agreement with previous findings -----, we show that lncRNAs are overwhelmingly

267 preferentially expressed in the testes (Figure 3A). Indeed, more than 68% of lncRNAs reach their

268 maximum expression level in this organ, compared to only approximately 32% of protein-coding genes,

269 for both mouse and rat (Figure 3A). Interestingly, more than 80% of lncRNAs are preferentially

270 expressed in young and aged adult samples, compared to only 62% of protein-coding genes (Figure

271 3B).

272

273 As noted previously, between 59 and 82% of protein-coding genes are significantly differentially

274 expressed (DE) among developmental stages, at a false discovery rate (FDR) below 1%, in each organ

275 and species (Figure 3C, Supplementary Dataset 4). The proportions of DE lncRNAs are much lower in

276 somatic organs, between 18 and 40%, but are similar in the testes, around 75% (Figure 3C). However,

277 we suspected that this could be due to the low expression levels of this class of genes, as total read

278 counts are known to affect the sensitivity of DE analyses (Anders and Huber 2010). Indeed, as

11

279     previously observed, lncRNAs are expressed at much lower levels and in fewer organ/developmental

280     stage combinations than protein-coding genes (Supplementary Figure 3A-C). To control for this effect,

281     we down-sampled the read counts observed for protein-coding genes, bringing them to the same

282     average counts as lncRNAs but preserving relative gene abundance (Methods). Strikingly, when

283     performing the DE analysis on this dataset, we observe higher proportions of DE loci for lncRNAs

284     compared to protein-coding genes (Figure 3C). Moreover, the amplitude of expression variation

285     among developmental stages are more important for lncRNAs than for protein-coding genes

286     (Supplementary Figure 3D). This is expected given the lower lncRNA expression levels, which preclude

287     detecting subtle expression shifts among time points.  Finally, we observe that the developmental

288     stage with maximum expression is generally different between protein-coding genes and lncRNAs,

289     even when considering genes that are significantly DE among stages. For all organs, DE lncRNAs tend

290     to show highest expression levels in the young and aged adults, while DE protein-coding genes are

291     more homogeneously distributed among developmental stages (Figure 3D, Supplementary Figure 3E).

292

293     Similar conclusions are reached when performing DE analyses between consecutive time points

294     (Supplementary Dataset 4). For both protein-coding genes and lncRNAs, the strongest expression

295     changes are observed between newborn and young adult individuals. Almost 10,000 lncRNAs are

296     significantly up-regulated between newborn and young adult testes, confirming the strong enrichment

297     for lncRNAs during spermatogenesis (Supplementary Dataset 4). We note that, as expected, the lowest

298     numbers of DE genes are observed at the transition between young and aged adult organs. At this

299     time-point, we observe more changes for the rat than for the mouse, potentially due to a higher

300     proportion of immune cell infiltrates in rat aged organs. Genes associated with antigen processing and

301     presentation tend to be expressed at higher levels in aged adults than in young adults, for mouse

302     kidney, rat brain and liver (Supplementary Dataset 4).

303    *Stronger selective constraint on lncRNAs expressed earlier in development*

304    We next analyzed the patterns of long-term evolutionary sequence conservation for lncRNAs, in

305    conjunction with their spatio-temporal expression pattern (Supplementary Table 5). We used the

306    PhastCons score (Siepel et al. 2005) across placental mammals (Casper et al. 2018), to assess the level

307    of sequence conservation for various aspects of mouse lncRNAs: exonic sequences, promoter regions

308    (defined as 1 kb regions upstream of the transcription start site, masking any exonic sequence within

309    this region), splice sites (first and last two bases of the introns, for multi-exonic loci). As approximately

310    20% of lncRNAs overlap with exonic regions from other genes on the opposite strand (Supplementary

311    Dataset 1), we masked exonic sequences from other genes before computing sequence conservation

312    scores. We analyzed sets of protein-coding genes and lncRNAs that are expressed above noise levels

313    (TPM>=1, averaged across all replicates) in each organ / developmental stage combination.

314

315    For exonic sequences and splice site regions, the extent of sequence conservation is much lower for

316    lncRNAs than for protein-coding genes, irrespective of the organ and developmental stage in which

317    they are expressed (Figure 4A, C). In contrast, promoter sequence conservation levels are more

318    comparable between protein-coding genes and lncRNAs (Figure 4B). For all examined regions and for

319    both categories of genes, the spatio-temporal expression pattern is well correlated with the level of

320    sequence conservation. Globally, sequence conservation is higher for genes that are expressed earlier

321    in development than for genes expressed later in development, and is significantly higher for somatic

322    organs than for adult and aged testes (Figure 4). Interestingly, for genes that are highly expressed in

323    mid-stage embryonic brain and kidney samples, the levels of promoter sequence conservation are

324    higher for lncRNAs than for protein-coding genes (Figure 4B). We also observed that lncRNAs that are

325    transcribed from bidirectional promoters tend to have higher sequence conservation levels than other

326    lncRNAs (Supplementary Figure 4).

327

328    Finally, we asked whether the highest level of evolutionary sequence conservation is seen at exons,

329    promoter or splice site regions, for each lncRNA locus taken individually. We show that this pattern

330    also depends on the organs and the developmental stages where the lncRNAs are expressed: for loci

331    detected in somatic organs and in the developing testes, there is significantly higher conservation for

332    the promoter and the splice sites than for exonic regions (Supplementary Figure 4). However, for

333    lncRNAs that are highly transcribed in the adult and aged testes (which constitutes the great majority

334    of genes), this pattern is absent (Supplementary Figure 4).

335    *Detection of homologous lncRNAs across species*

336    Having investigated the patterns of long-term sequence conservation of mouse lncRNAs, we next

337    sought to assess the conservation of lncRNA repertoires in mouse, rat and chicken. We detected

338    lncRNA separately in each species, using only RNA-seq data and existing genome annotations, as

339    previously suggested (Hezroni et al. 2015). We then searched for putative 1-to-1 orthologous lncRNAs

340    between species using pre-computed whole-genome alignments as a guide (Methods), to increase the

341    sensitivity of orthologous gene detection in the presence of rapid sequence evolution (Washietl et al.

342    2014). The orthologous lncRNA detection procedure involves several steps, including the identification

343    of putative homologous (projected) loci across species, filtering to remove large-scale structural

344    changes in the loci and intersection with predicted loci in the target species (Methods). As illustrated

345    in Figure 5, for comparisons between rodents the extent of sequence divergence is low enough that

346    more than 90% of 12,199 high-confidence lncRNA loci (expressed at TPM>=1 in at least one sample)

347    are successfully projected from mouse to rat (Figure 5A, Supplementary Dataset 5). However, only 53%

348    of projected loci have even weak levels of detectable transcription in the target species (at least 10

349    uniquely mapped reads). Only 27% of mouse lncRNA loci have predicted 1-to-1 orthologues in the rat,

350    and only 18% are orthologous to confirmed lncRNA loci in the rat (Figure 5A, Supplementary Dataset

351    5). The 1,081 mouse lncRNAs that have non-lncRNAs orthologues in the rat are generally matched with

352    loci discarded because of low read coverage, minimum exonic length or distance to protein-coding

353 genes (Supplementary Dataset 5). Cases of lncRNA-protein-coding orthologues are rare at this

354 evolutionary distance (Supplementary Dataset 5), and they may stem from gene classification errors.

355

356 At larger evolutionary distances, the rate of sequence evolution is the main factor hampering detection

357 of orthologous lncRNAs. Only 1,940 (16%) of mouse high-confidence lncRNAs (TPM>=1) could be

358 projected onto the chicken genome, and after subsequent filters we detect only 56 mouse – chicken

359 lncRNA orthologues (Figure 5A, Supplementary Dataset 5). We note that our lncRNA detection power

360 is likely weaker for the chicken than for the rodents because of organ and developmental stage

361 sampling, although we did strive to include RNA-seq data from adult organs in the lncRNA detection

362 process (Methods, Supplementary Table 2).

363

364 Conserved lncRNAs differ from non-conserved lncRNAs in terms of expression patterns. While only

365 subtle differences can be observed when comparing mouse-rat orthologous lncRNAs to the mouse-

366 specific lncRNA set, lncRNAs that are conserved across mouse, rat and chicken are dramatically

367 enriched in somatic organs and early developmental stages (Figure 5B,C, Supplementary Table 6).

368 Although their expression patterns have a strong species-specific component, shared patterns of organ

369 specificity can be detected (Supplementary Figure 5).

370 *Global patterns of lncRNA expression across species, organs and developmental stages*

371 We next assessed the global patterns of expression variation across species, organs and developmental

372 stages, for predicted mouse – rat lncRNA orthologues (Supplementary Dataset 6). As for protein-coding

373 genes, the main source of variability in a PCA performed on lncRNA expression levels is the difference

374 between adult and aged testes and the other samples (Figure 6A, Supplementary Figure 6). However,

375 for lncRNAs samples cluster according to the species of origin already on the second factorial axis (10%

376 explained variance), thus confirming that lncRNA expression patterns evolve rapidly. Overall,

377 differences between organs and developmental stages are less striking for lncRNAs, compared to the

378 variation stemming from the species factor (Figure 6A, Supplementary Figure 6). This pattern is also

15

379    visible on a hierarchical clustering analysis (performed on distances derived from Spearman's

380    correlation coefficient): in contrast with what is observed for protein-coding genes, for lncRNAs

381    samples generally cluster by species, with the exception of adult and aged testes which are robustly

382    grouped.

383

384    The higher rates of lncRNA expression evolution are also visible when analyzing within-species

385    variations, through comparisons across biological replicates (Figure 7A). We sought to measure the

386    global extent of gene expression conservation, by contrasting between-species and within-species

387    variations. Briefly, we constructed an expression conservation index by dividing the between-species

388    and the within-species Spearman's correlation coefficient, computed on all genes from a category, for

389    a given organ / developmental stage combination (Methods). The resulting expression conservation

390    values are very high for protein-coding genes, in particular for the brain and the mid-stage embryonic

391    kidney. However, there is significant less conservation between species for the adult and aged testes

392    (Figure 7B). For lncRNAs, expression conservation values are much lower than those observed for

393    protein-coding genes, with strikingly low values for adult and aged testes (Figure 7C).

394    *Evolutionary divergence of individual lncRNA expression profiles*

395    Having established that, globally, lncRNA expression patterns evolve very rapidly, we next sought to

396    assess expression divergence at the individual gene level. We first asked whether temporal patterns

397    of expression variations are conserved across rodent species. We selected lncRNAs that are

398    significantly differentially expressed (FDR<0.01) across developmental stages, in both mouse and rat

399    (Supplementary Dataset 4), and grouped them into clusters (Methods). We observed that in general,

400    lncRNAs show consistent patterns of variation among developmental stage in mouse and rat, with a

401    few exceptions in the kidney and liver (Supplementary Figure 7). Interestingly, lncRNAs that are DE in

402    the testes only rarely show divergent profiles between species, in contrast with what is observed for

403    protein-coding genes, where 25% of genes have different temporal patterns for mouse and rat

16

404    (Supplementary Figures 2,7). Overwhelmingly, lncRNAs are more highly expressed in adult and aged

405    testes than in developing testes, in both mouse and rat.

406

407    To further quantify lncRNA expression profile differences among species, we measured the amount of

408    expression divergence as the Euclidean distance between relative expression profiles (average TPM

409    values across biological replicates, normalized by dividing by the sum of all values for a gene, for each

410    species), for mouse and rat orthologues (Methods, Supplementary Dataset 7, Supplementary Table 7).

411    The resulting expression divergence values correlate negatively with the average expression level

412    (Figure 8A), as expected. While the raw expression divergence values are significantly higher for

413    lncRNAs than for protein-coding genes (Figure 8B), this is largely due to the low lncRNA expression

414    levels. Indeed, the effect disappears when analyzing the residual expression divergence after

415    regressing the mean expression level (Figure 8C). For lncRNAs, we also observe a weak negative

416    correlation between expression divergence and the extent of exonic sequence conservation (Figure

417    8D). We measured the relative contribution of each organ/developmental stage to the expression

418    divergence estimate (Figure 8E). For both protein-coding genes and lncRNAs, by far the highest

419    contributors are the young adult and aged testes samples, which are responsible for almost 30% of the

420    lncRNA expression divergence (Figure 8E). This is visible in the expression patterns of the 2 protein-

421    coding and lncRNA genes with the highest residual expression divergence: the lncRNA expression

422    divergence is mostly due to changes in adult testes, while more complex expression pattern changes

423    seem to have occurred for the protein-coding genes (Supplementary Figure 8). The most divergent

424    protein-coding genes are enriched in functions related to immunity (Supplementary Dataset 7).

425    *Candidate species-specific lncRNAs*

426    We next sought to investigate the most extreme cases of expression divergence: situations where

427    expression can be robustly detected in one species, but not in the other one, despite the presence of

428    perfect sequence alignment (Methods). We selected lncRNA loci that were supported by at least 100

429    uniquely mapped reads in one species, with no reads detected in the predicted homologous region in

430    the other species. With this convention, we obtain 1,041 candidate mouse-specific and 1,646

431    candidate rat-specific loci (Supplementary Dataset 8). These lists include striking examples, such as the

432    region downstream of the *Fzd4* protein-coding gene, which contains a mouse-specific and a rat-specific

433    lncRNA candidate, each perfectly aligned in the other species (Supplementary Figure 9A). We could

434    not identify any differential transcription factor binding or transposable element enrichment in the

435    promoters of these species-specific lncRNAs (data not shown). Interestingly however, they are

436    increasingly associated with predicted expression enhancers (Supplementary Figure 10). While the

437    evolutionary and mechanistic origin of these lncRNAs is still mysterious, we could confirm that their

438    presence is associated with increased expression divergence in the neighboring genes. To test this, we

439    selected species-specific and orthologous lncRNAs that are transcribed from bidirectional promoters

440    shared with protein-coding genes, and evaluated the expression divergence of their protein-coding

441    neighbors (Supplementary Figure 9B,C). Though the difference is subtle, genes that are close to

442    species-specific lncRNAs have significantly higher expression divergence than the ones that have

443    conserved lncRNA neighbors, even after correcting for expression levels (Wilcoxon test, p-value < 10-

444    3). It thus seems that expression changes that led to the species-specific lncRNA transcription extend

445    beyond the lncRNA locus and affect the neighboring genes, as previously proposed (Kutter et al. 2012).

446    **Discussion**

447    *Assessing lncRNA functionality: current challenges and insights from evolutionary approaches*

448    More than a decade after the publication of the first genome-wide lncRNA datasets (Guttman et al.

449    2009; Khalil et al. 2009), the debate regarding their functionality is still not settled. While experimental

450    assessments of lncRNA functions are rapidly accumulating, they are lagging behind the exponential

451    increase of RNA sequencing datasets, each one revealing thousands of previously unreported

452    noncoding transcripts (Pertea et al. 2018). There is thus a need to define biologically relevant criteria

453    to prioritize lncRNAs for experimental investigation. Furthermore*, in vivo* tests of lncRNA functions

454    need to be carefully designed to account for ubiquitous confounding factors, such as the presence of

18

455    overlapping regulatory elements at lncRNA loci (Bassett et al. 2014). Another challenge is the fact that

456    some lncRNA loci undoubtedly have "unconventional" biological functions, that require for example

457    the presence of a transcription and splicing at a given genomic location, independently of the lncRNA

458    molecule that is produced (Latos et al. 2012; Engreitz et al. 2016).

459

460    Evolutionary approaches can provide important tools to assess biological functionality (Haerty and

461    Ponting 2014), and they have been already successfully applied to lncRNAs. Although only a few large-

462    scale comparative transcriptomics studies are available so far for vertebrate lncRNAs (Kutter et al.

463    2012; Washietl et al. 2014; Hezroni et al. 2015; Necsulea et al. 2014), they all agree that lncRNAs evolve

464    rapidly in terms of primary sequence, exon-intron structure and expression patterns, indicating that

465    there is little selective constraint and thus little functionality for these loci. However, these studies

466    have all focused on lncRNAs detected in adult organs. We hypothesized that lncRNAs expressed during

467    embryogenesis are enriched in functional loci, as suggested by the increasing number of lncRNAs with

468    proposed roles in development (Rinn et al. 2007; Sauvageau et al. 2013; Grote et al. 2013; Grote and

469    Herrmann 2015). To test this hypothesis, we performed a multi-dimensional comparative

470    transcriptomics analysis, following lncRNA and protein-coding gene expression patterns across

471    species, organs and developmental stages.

472    *Spatio-temporal lncRNA expression patterns*

473    Our first major observation is that lncRNAs are overwhelmingly detected in the adult and aged testes,

474    in agreement with previous data (Soumillon et al. 2013). Their relative depletion in embryonic and

475    newborn testes reinforces the association between lncRNA production and spermatogenesis, in accord

476    with the hypothesis that the particular chromatin environment during spermatogenesis is a driver for

477    promiscuous, non-functional transcription (Kaessmann 2010; Soumillon et al. 2013). Interestingly, we

478    show that lncRNAs are significantly differentially expressed among developmental stages, at least as

479    frequently as protein-coding genes, after correcting for their lower expression levels. However, in

480    contrast with protein-coding genes, the majority of lncRNAs reach their highest expression levels in

481 adult rather than in developing organs. As requirements for tight gene expression control are

482 undoubtedly higher during embryonic development (Ben-Tabou de-Leon and Davidson 2007), an

483 explanation for the relative lncRNA depletion in embryonic and newborn transcriptomes is that

484 transcriptional noise is more efficiently blocked during the early stages of development. Differences in

485 cellular composition heterogeneity may also be part of the explanation. Expression analyses of cell-

486 type specific markers suggest that adult and aged organ transcriptomes may be a mix of more diverse

487 cell types, notably including substantial immune cell infiltrates. A higher cell type diversity may explain

488 the increased abundance of lncRNAs in adult and aged organs, especially given that lncRNAs are

489 thought to be cell-type specific (Liu et al. 2016).

490 *Functionally constrained lncRNAs are enriched in developmental transcriptomes*

491 We show that, for those lncRNAs that are expressed above noise levels (TPM>=1) in somatic organs

492 and in the earlier developmental stages, there is a higher proportion of functionally constrained loci

493 than in testes-expressed lncRNAs. Strikingly, we find that the level of long-term sequence conservation

494 for lncRNA promoter regions is higher than the one observed for protein-coding promoters, when we

495 analyze genes that are robustly expressed (TPM>=1) in embryonic brain and kidney. Moreover, for

496 lncRNAs that are expressed in somatic organs and in the developing testes, there is significantly more

497 evolutionary constraint on promoter and splice site sequences than on exonic regions, while these

498 patterns are not seen for the bulk of lncRNAs, expressed in adult and aged testes. Thus, we show that

499 lncRNAs that are expressed in somatic organs and in the developing testes differ from those expressed

500 in the adult testes not only in terms of overall levels of sequence conservation, but also with respect

501 to the regions of the lncRNA loci that are under selective constraint.  We validate previous reports of

502 increased constraint on splicing regulatory regions in mammalian lncRNAs (Schüler et al. 2014; Haerty

503 and Ponting 2015), and we show that this pattern is specifically seen in lncRNAs that are expressed in

504 somatic organs and in the developing testes. These results are also in agreement with a series of recent

505 findings, suggesting that at many lncRNA loci, biological function may reside in the presence of

506 additional non-coding regulatory elements at the lncRNA promoter rather than in the production of a

507    specific transcript (Engreitz et al. 2016; Groff et al. 2016). Thus, while there is evidence for increased

508    functionality for those lncRNA loci that are detected in developmental transcriptomes or in adult

509    somatic organs, our sequence conservation analyses suggest that their biological functions may be

510    carried out in an RNA-independent manner, as exonic sequences are under less constraint than

511    promoter or splice site regions.

512    *Evolutionary divergence of spatio-temporal expression profiles for lncRNAs*

513    We previously established that lncRNA expression patterns evolve rapidly across species in adult

514    organs. Here, we show that this rapid evolution of lncRNA expression is not restricted to adult and

515    aged individuals, but is also true for embryonic and newborn developmental stages. Expression

516    patterns comparisons across species, organs and developmental stages are dominated by differences

517    between species for lncRNAs, while similarities between organs and developmental stages are

518    predominant for protein-coding genes, even across distantly related species. We assessed the extent

519    of expression level conservation by contrasting between-species and within-species expression

520    variations and we showed that lncRNAs have significantly lower levels of conservation than protein-

521    coding genes, for all organs and developmental stages. However, lncRNA expression is significantly

522    more conserved in somatic organs and in early embryonic stages than in the adult testes. Interestingly,

523    when we evaluate expression divergence individually for each orthologous gene pair, and when we

524    correct for the lower lncRNA expression levels, we find that lncRNAs are comparable with protein-

525    coding genes, on average. Nevertheless, lncRNAs show a broader distribution of expression divergence

526    levels than protein-coding genes, and these patterns are mainly driven by species-specific expression

527    in the adult testes.

528

529    Finally, we analyzed extreme cases of expression divergence between species, namely situations

530    where transcription can be robustly detected in one species but not in the other, despite the presence

531    of good sequence conservation. We identify more than a thousand candidate species-specific lncRNAs,

532    in both mouse and rat. Interestingly, we observe that candidate mouse-specific lncRNAs are more

21

533    frequently transcribed from enhancers than lncRNAs conserved between mouse and rat. This

534    observation is consistent with previous reports that enhancers and enhancer-associated lncRNAs

535    evolve rapidly (Villar et al. 2015; Marques et al. 2013). The genetic basis of these extreme transcription

536    pattern changes is still not elucidated, and deserves further detailed investigations. Nevertheless, we

537    show that these lncRNA expression patterns do not occur in an isolated manner. When such species-

538    specific transcription was detected at protein-coding genes bidirectional promoters, the neighboring

539    protein-coding genes also showed increased expression divergence, compared to genes that are

540    transcribed from conserved lncRNA promoters. This observation is compatible with previous reports

541    that lncRNA turnover is associated with changes in neighboring gene expression levels (Kutter et al.

542    2012). While lncRNAs changes may be directly affecting gene expression, it is also possible that a

543    common mechanism affects both lncRNAs and protein-coding genes transcribed from bidirectional

544    promoters.

545    **Conclusions**

546    Our comparative transcriptomics approach confirms the established finding that lncRNAs repertoires,

547    sequences and expression patterns evolve rapidly across species, and shows that the accelerated rates

548    of lncRNA evolution are also seen in developmental transcriptomes. These observations are consistent

549    with the hypothesis that the majority of lncRNAs (or at least of those detected with sensitive

550    transcriptome sequencing approaches, in particular in the adult testes) may be non-functional.

551    However, we are able to modulate this conclusion, by showing that there are increased levels of

552    functional constraint on lncRNAs expressed during embryonic development, in particular in the

553    developing brain and kidney. These increased levels of constraint apply to all analyzed aspects of

554    lncRNAs, including sequence conservation for exons, promoter and splice sites, but also expression

555    pattern conservation. For many of these loci, biological function may be RNA-independent, as the

556    highest levels of selective constraint are observed on promoter regions and on splice signals, rather

557    than on lncRNA exonic sequences. Our results are thus compatible with unconventional, RNA-

558    independent functions for lncRNAs expressed during embryonic development.

22

559

**Methods**

*Biological sample collection*

We collected samples from three species (mouse C57BL/6J strain, rat Wistar strain and chicken White Leghorn strain), four organs (brain, kidney, liver and testes) and five developmental stages (including two embryonic stages, newborn, young and aged adult individuals). We sampled the following stages in the mouse: embryonic day post-conception (dpc) 13.5 (E13.5 dpc, hereafter mid-stage embryo); E17 to E17.5 dpc (late embryo); post-natal day 1 to 2 (newborn); young adult (8-10 weeks old); aged adult (24 months old). For the rat, we sampled the following stages: E15 dpc (mid-stage embryo); E18.5 to E19 dpc (late embryo); post-natal day 1 to 2 (newborn); young adult (8-10 weeks old); aged adult (24 months, with the exception of kidney samples and two of four liver samples, derived from 12 months old individuals). The embryonic and neonatal developmental stages were selected for maximum comparability based on Carnegie stage criteria (Theiler 1989). For chicken, we collected samples from Hamburger-Hamilton stages 31 and 36, hereafter termed mid-stage and late embryo. We selected these two stages for comparability with the two embryonic stages in mouse and rat (Hamburger and Hamilton 1951). In general, each sample corresponds to one individual, except for mouse and rat mid-stage embryonic kidney, for which tissue from several embryos was pooled prior to RNA extraction. For adult and aged organs, multiple tissue pieces from the same individual were pooled and homogenized prior to RNA extraction. For brain dissection, we sampled the cerebral cortex. For mouse and rat samples, with the exception of the mid-stage embryonic kidney, individuals were genotyped and males were selected for RNA extraction. Between two and four biological replicates were obtained for each species/organ/stage combination, amounting to 97 samples in total (Supplementary Table 1).

*RNA-seq library preparation and sequencing*

We performed RNA extractions using RNeasy Plus Mini kit from Qiagen. RNA quality was assessed using the Agilent 2100 Bioanalyzer. Sequencing libraries were produced using the Illumina TruSeq

24

584    stranded mRNA protocol with polyA selection, and sequenced as 101 base pairs (bp) single-end reads,

585    at the Genomics Platform of iGE3 and the University of Geneva (https://ige3.genomics.unige.ch/).

*Additional RNA-seq data*

587    To improve detection power for lowly expressed lncRNAs, we complemented our RNA-seq collection

588    with samples generated with the same technology for Brown Norway rat adult organs (Cortez et al.

589    2014). We added data generated by the Chickspress project (http://geneatlas.arl.arizona.edu/) for

590    adult chicken (red jungle fowl strain UCD001) organs, as well as for embryonic chicken (White Leghorn)

591    organs from two publications (Uebbing et al. 2015; Ayers et al. 2013). Almost all samples were strand-

592    specific, except the chicken adult organs and early embryonic testes. As the data were not perfectly

593    comparable with our own in terms of library preparation and animal strains, the additional rat and

594    chicken samples were only used to increase lncRNA detection sensitivity.

*RNA-seq data processing*

596    We used HISAT2  (Kim et al. 2015) release 2.0.5 to align the RNA-seq data on reference genomes. The

597    genome sequences (assembly versions mm10/GRCm38, rn6/Rnor_6.0 and galGal5/Gallus_gallus-5.0)

598    were downloaded from the Ensembl database (Cunningham et al. 2019). Genome indexes were built

599    using only genome sequence information. To improve detection sensitivity, at the alignment step we

600    provided known splice junction coordinates extracted from Ensembl. We set the maximum intron

601    length for splice junction detection at 1 million base pairs (Mb). To verify the strandedness of the RNA-

602    seq data, we analyzed spliced reads that spanned introns with canonical (GT-AG or GC-AG) splice sites

603    and compared the strand inferred based on the splice site with the one assigned based on the library

604    preparation protocol (Supplementary Table 1). Finally, to estimate the mappability of each genomic

605    region, we generated error-free artificial RNA-seq reads (single-end, 101 bp long, with 5 bp distance

606    between consecutive read starts) from the genome sequence and realigned them to the genome with

607    the same HISAT2 parameters. Regions for which the corresponding reads could be aligned

608    unambiguously were considered "mappable"; the remaining regions were said to be "unmappable".

25

609    *Transcript assembly and filtering*

610    We assembled transcripts for each sample using StringTie (Pertea et al. 2015), release 1.3.5, based on

611    read alignments obtained with HISAT2. We provided genome annotations from Ensembl release 94 as

612    a guide for transcript assembly. We filtered Ensembl annotations to remove transcripts that spanned

613    a genomic length above 2.5 Mb. For protein-coding genes, we kept only protein-coding transcripts,

614    discarding isoforms annotated as "retained_intron", "processed_transcript" etc. We set the minimum

615    exonic length at 150 bp, the minimum anchor length for splice junctions at 8bp and the minimum

616    isoform fraction at 0.05. We compared the resulting assembled transcripts with Ensembl annotations

617    and we discarded read-through transcripts, defined as overlapping with multiple multi-exonic

618    Ensembl-annotated genes. For strand-specific samples, we discarded transcripts for which the ratio of

619    sense to antisense unique read coverage was below 0.01. We discarded multi-exonic transcripts that

620    were not supported by splice junctions with correctly assigned strands. The filtered transcripts

621    obtained for each sample were assembled into a single dataset *per* species using the merge option in

622    StringTie. For increased sensitivity, we removed the minimum FPKM and TPM thresholds for transcript

623    inclusion. We constructed a combined annotation dataset, starting with Ensembl annotations, to

624    which we added newly-assembled transcripts that had no exonic overlap with Ensembl genes. We also

625    included newly-annotated isoforms for known genes if they had exonic overlap with exactly one

626    Ensembl gene, thus discarding potential read-through transcripts or gene fusions.

627    *Protein-coding potential of assembled transcripts*

628    To determine whether the newly assembled transcripts were protein-coding or non-coding, we mainly

629    relied on the codon substitution frequency (CSF) score (Lin et al. 2007). As in a previous publication

630    (Necsulea et al. 2014) we scanned whole genome alignments and computed CSF scores in 75 bp sliding

631    windows moving with a 3 bp step. We used pre-computed alignments downloaded from the UCSC

632    Genome Browser (Casper et al. 2018), including the alignment between the mouse genome and 59

633    other vertebrates (for mouse classification), between the human genome and 99 other vertebrates

634    (for rat and chicken classification) and between the rat genome and 19 other vertebrates (for rat

635    classification). For each window, we computed the score in each of the 6 possible reading frames and

636    extracted the maximum score for each strand. We considered that transcripts are protein-coding if

637    they overlapped with positive CSF scores on at least 150 bp. As positive CSF scores may also appear on

638    the antisense strand of protein-coding regions due to the partial strand-symmetry of the genetic code,

639    in this analysis we considered only exonic regions that did not overlap with other genes. In addition,

640    we searched for sequence similarity between assembled transcripts and known protein sequences

641    from the SwissProt 2017_04 (The UniProt Consortium 2017) and Pfam 31.0 (El-Gebali et al. 2019)

642    databases. We kept only SwissProt entries with confidence scores 1, 2 or 3 and we used the Pfam-A

643    curated section of Pfam. We searched for sequence similarity using the blastx utility in the BLAST+

644    2.8.1 package (Camacho et al. 2009; Altschul et al. 1990), keeping hits with maximum e-value 1e-3 and

645    minimum protein sequence identity 40%, on repeat-masked cDNA sequences. We considered that

646    transcripts were protein-coding if they overlapped with blastx hits over at least 150 bp. Genes were

647    said to be protein-coding if at least one of their isoforms was classified as protein-coding, based on

648    either the CSF score or on sequence similarity with known proteins.

649    *Long non-coding RNA selection*

650    To construct a reliable lncRNA dataset, we selected newly-annotated genes classified as non-coding

651    based on both the CSF score and on sequence similarity with known proteins and protein domains, as

652    well as Ensembl-annotated genes with non-coding biotypes ("lincRNA", "processed_transcript",

653    "antisense", "TEC", "macro_lncRNA", "bidirectional_promoter_lncRNA", "sense_intronic"). For newly

654    detected genes, we applied several additional filters: we required a minimum exonic length

655    (corresponding to the union of all annotated isoforms) of at least 200 bp for multi-exonic loci and of

656    at least 500 bp for mono-exonic loci; we eliminated genes that overlapped for more than 5% of their

657    exonic length with unmappable regions; we kept only loci that were classified as intergenic and at least

658    5 kb away from Ensembl-annotated protein-coding genes on the same strand; for multi-exonic loci, we

659    required that all splice junctions be supported by reads with correct strand assignment (cf. above). For

27

660    both *de novo* and Ensembl annotations, we removed transcribed loci that overlapped on at least 50%

661    of their length with retrotransposed gene copies, annotated by the UCSC Genome Browser and from

662    a previous publication (Carelli et al. 2016); we discarded loci that overlapped with UCSC-annotated

663    tRNA genes and with RNA-type elements from RepeatMasker (Smit et al. 2003) on at least 25% of their

664    length. We kept loci supported by at least 10 uniquely mapped RNA-seq reads and for which a ratio of

665    sense to antisense transcription of at least 1% was observed in at least one sample.

666    *Gene expression estimation*

667    We computed the number of uniquely mapping reads unambiguously attributed to each gene using

668    the Rsubread package in R (Liao et al. 2019), discarding reads that overlapped with multiple genes. We

669    also estimated read counts and TPM (transcript *per* million) values *per* gene using Kallisto (Bray et al.

670    2016). To approach absolute expression levels estimates, for better comparisons across samples, we

671    further normalized TPM values using a scaling approach (Brawand et al. 2011). Briefly, we ranked the

672    genes in each sample according to their TPM values, we computed the variance of the ranks across all

673    samples for each gene, and we identified the 100 least-varying genes, found within the inter-quartile

674    range (25%-75%) in terms of average expression levels across samples. We derived normalization

675    coefficients for each sample such that the median of the 100 least-varying genes be identical across

676    samples. We then used these coefficients to normalize TPM values for each sample. We excluded

677    mitochondrial genes from expression estimations and analyses, as these genes are highly expressed

678    and can be variable across samples.

679    *Differential expression analyses*

680    We used the DESeq2 (Love et al. 2014)(Smedley et al. 2009)(74)(75) package release 1.22.2 in R release

681    3.5.0 (R Core Team 2018) to test for differential expression across developmental stages, separately

682    for each organ and species. We analyzed both protein-coding genes and lncRNAs, selected according

683    to the criteria described above. We first performed a global differential expression analysis, using the

684    likelihood ratio test to contrast a model including an effect of the developmental stage against the null

28

685    hypothesis of homogeneous expression level across all developmental stages. This analysis was

686    performed on all annotated protein-coding and lncRNA genes for each species, as well as on 1-to-1

687    orthologous genes for mouse and rat. In addition, we down-sampled the numbers of reads assigned

688    to protein-coding genes to obtain identical average numbers of reads for protein-coding genes and

689    lncRNAs. We also contrasted consecutive developmental stages, for each species and organ. For each

690    test, we also computed the expression fold change based on average TPM values for each

691    developmental stage/organ combination.

692    *Expression specificity index*

693    We used the previously proposed tissue specificity index (Liao et al. 2006) to measure gene expression

694    specificity across organs and developmental stages, provided by the formula: tau = sum $(1 - r_i)/(n-1)$,

695    where $r_i$ represents the ratio between the expression level in sample i and the maximum expression

696    level across samples, and n represents the total number of samples. We computed this index on

697    normalized TPM values, averaged across all replicates for a given species / organ / developmental stage

698    combination (Supplementary Dataset 3).

699    *Homologous lncRNA family prediction*

700    We used existing whole-genome alignments as a guide to predict homologous lncRNAs across species,

701    as previously proposed (Washietl et al. 2014). We first constructed for each gene the union of its exon

702    coordinates across all isoforms, hereafter termed "exon blocks". We projected exon block coordinates

703    between pairs of species using the liftOver utility and whole-genome alignments generated with blastz

704    (http://www.bx.psu.edu/miller_lab/), available through the UCSC Genome Browser (Casper et al.

705    2018). To increase detection sensitivity, for the initial liftOver projection we required only that 10% of

706    the reference bases remap on the target genome. Projections were then filtered, retaining only cases

707    where the size ratio between the projected and the reference region was between 0.33 and 3 for

708    mouse and rat (0.2 and 5 for comparisons involving chicken). To exclude recent lineage-specific

709    duplications, regions with ambiguous or split liftOver projections were discarded. For genes where

710    multiple exon blocks could be projected across species, we defined the consensus chromosome and

711    strand in the target genome and discarded projected exon blocks that did not match this consensus.

712    We then evaluated the order of the projected exon blocks on the target genes, to identify potential

713    internal rearrangements. If internal rearrangements were due to the position of a single projected

714    exon block, the conflicting exon block was discarded; otherwise, the entire projected gene was

715    eliminated. As the projected reference gene coordinates could overlap with multiple genes in the

716    target genome, we constructed gene clusters based on the overlap between projected exon block

717    coordinates and target annotations, using a single-link clustering approach. We then realigned entire

718    genomic loci for each pair of reference-target genes found within a cluster, using lastz

719    (http://www.bx.psu.edu/miller_lab/) and the threaded blockset aligner (Blanchette et al. 2004). Using

720    this alignment, we computed the percentage of exonic sequences aligned without gaps and the

721    percentage of identical exonic sequence, for each pair of reference-target genes. We then extracted

722    the best hit in the target genome for each gene in the reference genome based on the percentage of

723    identical exonic sequence, requiring that the ratio between the maximum percent identity and the

724    percent identity of the second-best hit be above 1.1. Reciprocal best hits were considered to be 1-to-

725    1 orthologous loci between pairs of species. For analyses across all three species, we constructed

726    clusters of reciprocal best hits from pairwise species comparisons, using a single-link clustering

727    approach. Resulting clusters with more than 1 representative *per* species were discarded. To examine

728    the validity of our procedure for homologous gene family prediction, we compared the resulting gene

729    families with predictions from the Ensembl Compara pipeline (Herrero et al. 2016), extracted from

730    Ensembl release 94, for protein-coding genes.

731    *Sequence evolution*

732    We evaluated long-term evolutionary sequence conservation based on PhastCons (Siepel et al. 2005)

733    scores, computed for the mouse genome using either a placental mammal or a vertebrate multiple

734    species alignment available from the UCSC Genome Browser (Casper et al. 2018). We computed

735    average PhastCons scores on exonic sequences (excluding exonic regions overlapping with other

30

736     genes), promoter regions (defined as 1 kb immediately upstream of the transcription start site) and

737     splice sites (defined as the first two and last two bases of each intron). For genes with multiple

738     promoters, we computed the average score across all promoters.

739     *Gene expression evolution*

740     We computed global and *per*-gene expression level conservation between mouse and rat, for 1-to-1

741     orthologous genes. We first measured gene expression conservation for protein-coding genes and

742     lncRNAs as a class. For each organ/developmental stage, we computed the expression level correlation

743     between mouse and rat average TPM levels, across all orthologous pairs. We also computed the

744     correlation between individuals within the same species; for organ/stages with more than two

745     biological replicates we computed the average correlation coefficient across all possible pairs of

746     individuals. We then evaluated the global extent of gene expression conservation through the ratio of

747     the between-species correlation coefficient to the average within-species correlation coefficient.

748     Spearman's rank correlation coefficients were used in all cases. We obtained 95% confidence intervals

749     for expression conservation measures through a bootstrap procedure, resampling 100 times the same

750     number of genes with replacement. In addition to this global measure of expression conservation, we

751     estimated the extent of between-species expression divergence *per* gene by computing Euclidean

752     distances between relative expression profiles for each species. The relative expression profiles were

753     derived from TPM values *per* organ/developmental stage, averaged across biological replicates,

754     divided by the sum of all average TPM values.

755     *Statistical analyses and graphical representations*

756     All statistical analyses and graphical representations were done with R (R Core Team 2018), version

757     3.5.0. We performed principal component analyses using the ade4 library (Dray and Dufour 2007) and

758     hierarchical clustering of gene expression matrices using the hclust function in the stats package in R,

759     on pairwise Euclidean distances. For all analyses involving multiple statistical tests, false discovery

760     rates were computed with the Benjamini-Hochberg procedure (Benjamini and Hochberg 1995). 95%

761 confidence intervals for median values of distributions were computed with the following formula:

762 median +/- 1.57 x IQR/sqrt(N), where IQR is the inter-quartile range, sqrt denotes the square root and

763 N the number of points.

**Availability of data and materials**

765 The raw and processed RNA-seq data were submitted to the NCBI Gene Expression Omnibus (GEO),

766 under accession number GSE108348. Additional processed files and all scripts used to analyze the

767 data are available at the address: ftp://pbil.univ-lyon1.fr/pub/datasets/Darbellay_LncEvoDevo.

768

**Author contributions**

770 FD performed organ dissections, RNA extractions, quality control, prepared samples for sequencing

771 and contributed to study design and manuscript preparation. AN designed the study, performed

772 computational analyses and wrote the manuscript. All authors read and approved the final manuscript.

785

786 **Figure legends**

787 **Figure 1. Transcriptome complexity across organs and developmental stages.**

788 **A.** Experimental design. The developmental stages selected for mouse, rat and chicken are marked on

789 a horizontal axis. Organs sampled for each species and developmental stage are shown below.

790 Abbreviations: br, brain; kd, kidney; lv, liver; ts, testes.

791 **B.** Expression of cell type-specific markers derived from single-cell experiments (full list provided in

792 Supplementary Table 3), in our mouse and rat RNA-seq samples. The heatmap represents centered

793 and scaled log2-transformed TPM levels (z-score). Developmental stages are indicated by numeric

794 labels, 1 to 5. Average levels across biological replicates are shown. Species are color-coded, shown

795 below the heatmap.

796 **C.** Number of protein-coding genes supported by at least 10 uniquely mapped reads in each sample,

797 after read resampling to homogenize coverage (Methods).

798 **D.** Number of lncRNAs supported by at least 10 uniquely mapped reads in each sample, after read

799 resampling to homogenize coverage.

800

801 **Figure 2. Protein-coding gene expression is conserved across organs and developmental stages.**

802 **A.** First factorial map of a principal component analysis, performed on log2-transformed TPM values,

803 for 10,363 protein-coding genes with orthologues in mouse, rat and chicken. Colors represent different

804 organs and  developmental stages, point shapes represent different species.

805 **B.** Hierarchical clustering, performed on a distance matrix derived from Spearman correlations

806 between pairs of samples, for 10,363 protein-coding genes with orthologues in mouse, rat and chicken.

807 Organ and developmental stages are color-coded, shown below the heatmap. Species of origin is color-

808 coded, shown on the right. Sample clustering is shown on the left.

809 **C.** Expression profiles of protein-coding genes that are significantly differentially expressed (FDR<0.01)

810 among developmental stages, for both mouse and rat, in the brain. TPM values were averaged across

811 replicates and normalized by dividing by the maximum value, for each species. The resulting relative

33

812     expression profiles were combined across species and clustered with the K-means algorithm. The

813     average profiles of the genes belonging to each cluster are shown. Gray lines represent profiles of

814     individual genes from a cluster.

815

816     **Figure 3. Different expression patterns for protein-coding genes and lncRNAs.**

817     **A.** Distribution of the organ in which maximum expression is observed, for protein-coding genes (pc)

818     and lncRNAs (lnc), for mouse, rat and chicken. Organs are color-coded, shown above the plot.

819     **B.** Distribution of the developmental stage in which maximum expression is observed, for protein-

820     coding genes and lncRNAs, for mouse, rat and chicken. Developmental stages are color-coded, shown

821     above the plot.

822     **C.** Percentage of protein-coding and lncRNA genes that are significantly (FDR<0.01) DE among

823     developmental stages, with respect to the total number of genes tested for each organ. Left panel:

824     differential expression analysis performed with all RNA-seq reads. Right panel: differential expression

825     analysis performed after down-sampling read counts for protein-coding genes, to match those of

826     lncRNAs (Methods).

827     **D.** Distribution of the developmental stage in which maximum expression is observed, for protein-

828     coding genes and lncRNAs that are significantly DE (FDR<0.01) in each organ, for the mouse. The

829     percentages are computed with respect to the total number of DE genes in each organ and each gene

830     class.

831

832     **Figure 4. Increased levels of long-term sequence conservation for lncRNAs expressed early in**

833     **development.**

834     **A.** Distribution of the PhastCons sequence conservation score for protein-coding and lncRNAs exonic

835     regions, for subsets of genes expressed above noise levels (TPM>=1) in each organ and developmental

836     stage. We used precomputed PhastCons score for placental mammals, downloaded from the UCSC

837     Genome Browser. Exonic regions that overlap with exons from other genes were masked. Dots

34

838    represent median values, vertical bars represent 95% confidence intervals. Numbers of analyzed genes

839    are provided in Supplementary Table 4.

840    **B.** Same as A, for promoter regions (1kb upstream of transcription start sites). Exonic sequences were

841    masked before assessing conservation.

842    **C.** Same as B, for splice sites (first and last two bases of each intron).

843

844    **Figure 5. Orthologous lncRNA families for mouse, rat and chicken.**

845    **A.** Number of mouse protein-coding genes and lncRNAs in different classes of evolutionary

846    conservation. From left to right: all loci (with TPM>=1 in at least one mouse sample), loci with

847    conserved sequence in the rat, loci for which transcription could be detected (at least 10 unique reads)

848    in predicted orthologous locus in the rat, loci with predicted 1-to-1 orthologues, loci for which the

849    predicted orthologue belonged to the same class (protein-coding or lncRNA) in the rat, loci with

850    conserved sequence in the chicken, loci for which transcription could be detected (at least 10 unique

851    reads) in predicted orthologous locus in the chicken, loci with predicted 1-to-1 orthologues, loci for

852    which the predicted orthologue belonged to the same class (protein-coding or lncRNA) in the chicken.

853    We analyze 17,868 protein-coding genes and 12,199 candidate lncRNAs with an expression level (TPM)

854    >=1 in at least one mouse sample.

855    **B.** Distribution of the organ in which maximum expression is observed, for mouse protein-coding and

856    lncRNA genes that have no orthologues in the rat or chicken, for genes with orthologues in the rat and

857    for genes with orthologues in chicken.

858    **C.** Same as B, for the distribution of the developmental stage in which maximum expression is

859    observed.

860

861    **Figure 6. Global comparison of lncRNA expression patterns across species.**

862    **A.** First factorial map of a principal component analysis, performed on log2-transformed TPM values,

863    for 2,754 orthologous mouse and rat lncRNAs expressed above noise levels (TPM >= 1) in at least one

35

864    mouse or rat sample. Colors represent different organs and  developmental stages, point types

865    represent species.

866    **B.** Hierarchical clustering, performed on a distance matrix derived from Spearman correlations

867    between pairs of samples, for 2,754 orthologous mouse and rat lncRNAs. Organ and developmental

868    stages are shown below the heatmap. Species of origin is shown on the right. Sample clustering is

869    shown on the left.

870

871    **Figure 7. Global estimates of expression conservation across organs and developmental stages.**

872    **A.** Example of between-species and within-species variation of expression levels, for protein-coding

873    genes (left) and lncRNAs (right), for orthologous genes between mouse and rat, for the mid-stage

874    embryonic brain. Spearman's correlation coefficients (rho) are shown above each plot. We show a

875    smoothed color density representation of the scatterplots, obtained through a (2D) kernel density

876    estimate (smoothScatter function in R).

877    **B.** Expression conservation index, defined as the ratio of the between-species and the within-species

878    expression level correlation coefficients, for protein-coding genes, for each organ and developmental

879    stage. The vertical segments represent minimum and maximum values obtained from 100 bootstrap

880    replicates. We analyzed 14,919 pairs of orthologous genes between mouse and rat, with TPM >= 1 in

881    at least one sample.

882    **C.** Same as B, for lncRNAs. We analyzed 2,754 orthologous mouse and rat lncRNAs with TPM >= 1 in at

883    least one mouse or rat sample.

884

885    **Figure 8. Per-gene estimates of expression pattern divergence between species.**

886    **A.** Relationship between the per-gene expression divergence measure (Euclidean distance of relative

887    expression profiles among organs/stages, between mouse and rat), and the average expression values

888    (log2-transformed TPM) across all mouse and rat samples. We show a smoothed color density

889     representation of the scatterplots, obtained through a (2D) kernel density estimate (smoothScatter

890     function in R). Red line: linear regression.

891     **B.** Distribution of the expression divergence value for all protein-coding and lncRNA genes with

892     predicted 1-to-1 orthologues in mouse and rat.

893     **C.** Distribution of the residual expression divergence values, after regressing the average expression

894     level, for protein-coding genes and lncRNAs.

895     **D.** Relationship between expression divergence and exonic sequence conservation (% exonic sequence

896     aligned without gaps between mouse and rat), for protein-coding genes and lncRNAs.

897     **E.** Average contribution of each organ/developmental stage combination to expression divergence, for

898     protein-coding genes and lncRNAs.

899

900     **Supplementary Figure 1. Expression patterns of cell-type specific markers in mouse, rat and chicken**

901     **samples.**

902     **A.** Expression of cell type-specific markers derived from single-cell experiments (full list provided in

903     Supplementary Table 3), in our mouse, rat and chicken RNA-seq samples. The heatmap represents

904     centered and scaled log2-transformed TPM levels (z-score). Developmental stages are indicated by

905     numeric labels, 1 to 5. Average levels across biological replicates are shown. We show only organs and

906     developmental stages that were sampled in all three species, for genes with 1-to-1 orthologues.

907

908     **Supplementary Figure 2. Conservation of developmental expression profiles between mouse and**

909     **rat, for protein-coding genes.**

910     **A.** Expression profiles of orthologous protein-coding genes that are significantly differentially

911     expressed (FDR<0.01) among developmental stages, for both mouse and rat, in the kidney. TPM values

912     were averaged across replicates and normalized by dividing by the maximum, for each species. The

913     resulting relative expression profiles were combined across species and clustered with the K-means

37

914    algorithm. The average profiles of the genes belonging to each cluster are shown. Gray lines represent

915    profiles of individual genes from a cluster. Numbers of genes in each cluster are shown in the plot.

916    **B.** Same as A, for the liver.

917    **C.** Same as A, for the testes. For this organ, we searched for only 4 clusters with the K-means algorithm.

918

919    **Supplementary Figure 3. Protein-coding genes and lncRNA expression patterns.**

920    **A.** Distribution of the maximum expression level (log2-transformed TPM values), for protein-coding

921    genes (red) and lncRNAs (blue), for mouse, rat and chicken. We show only genes that are expressed

922    above noise levels (TPM >= 1) in at least one sample.

923    **B.** Distribution of the expression specificity index (Methods) for protein-coding genes and lncRNAs, in

924    the mouse. Genes were divided into 5 expression bins, based on their maximum expression level across

925    samples.

926    **C.** Same as B, for the rat.

927    **D.** Distribution of the ratio between the minimum and the maximum TPM value across developmental

928    stages, for genes that are significantly differentially expressed among stages for each organ and

929    species. Lower values indicate stronger expression changes among developmental stages.

930    **E.** Distribution of the developmental stage in which maximum expression is observed, for protein-

931    coding genes and lncRNAs that are significantly DE (FDR<0.01) in each organ, for the rat. The

932    percentages are computed with respect to the total number of DE genes in each organ and each gene

933    class.

934

935    **Supplementary Figure 4. Estimates of long-term sequence conservation scores for different regions**

936    **of lncRNAs loci.**

937    **A.** Distribution of the difference between the exonic PhastCons score and the promoter PhastCons

938    score, for mouse lncRNAs that are expressed above noise levels (TPM>=1) in each organ and

939    developmental stage. Precomputed PhastCons score for placental mammals were provided by the

940    UCSC Genome Browser. Exonic regions that overlap with other genes were masked. Dots represent

941    median values, vertical bars represent 95% confidence intervals. Numbers of analyzed genes are

942    provided in Supplementary Table 4.

943    **B.** Same as A, for the difference between exonic and splice site PhastCons score.

944    **C.** Distribution of the promoter sequence conservation score, for all lncRNAs, for lncRNAs that have

945    bidirectional promoters and for lncRNAs that overlap with Encode-annotated enhancers, in the mouse.

946

947    **Supplementary Figure 5. Expression patterns of 30 lncRNAs conserved in mouse, rat and chicken.**

948    **A.** Heatmap of the centered and scaled expression values (log2-transformed TPM), for 30 lncRNAs that

949    are shared across mouse, rat and chicken. For comparability with chicken, we show only mid-stage and

950    late embryo samples for mouse and rat, for somatic organs. The annotation source is shown on the

951    right: gray rectangles indicate a newly-annotated gene. Organs and developmental stages are depicted

952    by color rectangles below the heatmap. The list of genes used for this analysis is provided in

953    Supplementary Table 6.

954

955    **Supplementary Figure 6. Main sources of expression pattern variability for protein-coding genes and**

956    **lncRNAs.**

957    **A.** Coordinates on the first five axes of the principal component analysis, for mouse and rat orthologous

958    protein-coding genes. Points represent individual samples. Organs are color-coded and developmental

959    stages are distinguished by point types. Mouse (m, filled dots) and rat (r, unfilled dots) samples are

960    shown on separate x-axis positions.

961    **B.** Same as A, for lncRNAs.

962

963    **Supplementary Figure 7. Conservation of developmental expression profiles between mouse and**

964    **rat, for lncRNAs.**

965   **A.** Expression patterns of orthologous lncRNAs that are significantly differentially expressed

966   (FDR<0.01) among developmental stages, for both mouse and rat, in the brain. TPM values were

967   averaged across replicates and normalized by dividing by the maximum value for each species. The

968   resulting relative expression profiles were combined across species and clustered with the K-means

969   algorithm. The average profiles of the genes belonging to each cluster are shown. Gray lines represent

970   profiles of individual genes from a cluster.

971   **B.** Same as A, for the kidney.

972   **C.** Same as A, for the liver.

973   **D.** Same as A, for the testes. For this organ, we searched for only 4 clusters with the K-means algorithm.

974

975   **Supplementary Figure 8. Examples of genes with high expression pattern divergence between**

976   **mouse and rat.**

977   **A.** Examples of average expression profile in mouse and rat, for the top 2 most-divergent protein-

978   coding and lncRNA genes.

979

980   **Supplementary Figure 9. Candidate species-specific lncRNAs.**

981   **A.** Genomic localization and RNA-seq read coverage of a candidate mouse-specific lncRNA, situated

982   downstream of the *Fzd4* gene. RNA-seq data is shown for young and aged adult kidney.

983   **B.** Distribution of the raw expression divergence for protein-coding genes that are transcribed from

984   the same bidirectional promoters as lncRNAs with 1-to-1 orthologues in mouse and rat (black), or as

985   candidate species-specific lncRNAs (red).

986   **C.** Same as A, after correcting for the average expression level of the protein-coding genes.

987

988   **Supplementary Figure 10. Genomic and expression characteristics of candidate species-specific**

989   **lncRNAs.**

990     **A.** Percentage of mouse lncRNAs for which the predicted transcription start site is found within 1kb

991     of an Encode-annotated enhancer. LncRNAs are divided into loci with predicted 1-to-1 orthologues in

992     the rat (1-1, dark blue) and mouse-specific lncRNAs (light blue). LncRNAs are further separated into

993     newly-annotated (new) or previously known (Ensembl).

994     **B.** Same as A, for the percentage of multi-exonic loci, for mouse and rat.

995     **C.** Same as A, for the percentage of loci that have predicted bidirectional promoters, for mouse and

996     rat.

997     **D.** Distribution of the organ in which maximum expression is observed, for mouse and rat lncRNAs.

998     LncRNAs are divided into loci with predicted 1-to-1 orthologues (1-1) and species-specific lncRNAs (sp).

999     **E.** Same as D, for the distribution of the developmental stage in which maximum expression is

1000    observed.

1001

1002    **Supplementary Table List.**

1003    **Supplementary Table 1.** List of RNA-seq samples generated specifically for this project, and used for

1004    all downstream expression analyses.

1005    **Supplementary Table 2.** List of additional, previously published RNA-seq samples, included in the

1006    lncRNA detection pipeline.

1007    **Supplementary Table 3.** Cell-type markers for the four organs analyzed here, derived from single-cell

1008    transcriptomics analyses.

1009    **Supplementary Table 4.** Numbers of protein-coding genes and lncRNAs that have an average TPM

1010    expression level of at least 1 in each organ / developmental stage combination, for each species.

1011    **Supplementary Table 5.** Sequence conservation scores (average PhastCons scores), for exons, introns,

1012    promoters and splice sites, for mouse protein-coding genes and lncRNAs.

1013    **Supplementary Table 6.** List of 30 lncRNAs that are predicted to be 1-to-1 orthologues in mouse, rat

1014    and chicken.

1015    **Supplementary Table 7.** Expression pattern and sequence conservation scores for protein-coding

1016    genes and lncRNAs, for mouse and rat 1-to-1 orthologues.

1017

1018    **Supplementary Dataset List.**

1019    **Supplementary Dataset 1.** Complete gene annotations for mouse, rat and chicken.

1020    **Supplementary Dataset 2.** Gene expression levels (raw and normalized TPM values, unique read

1021    counts).

1022    **Supplementary Dataset 3.** Expression patterns (average across replicates, samples with maximum

1023    expression) and expression specificity indexes.

1024    **Supplementary Dataset 4.** Results of the differential expression analyses across all developmental

1025    stages, or between consecutive developmental stages, for each organ and each species.

1026    **Supplementary Dataset 5.** Predicted orthologous gene families and sequence conservation statistics.

1027    **Supplementary Dataset 6.** Raw and normalized expression values (TPM) for orthologous protein-

1028    coding and lncRNA families.

1029    **Supplementary Dataset 7.** Expression pattern divergence for mouse and rat orthologous genes.

1030    **Supplementary Dataset 8.** Lists of candidate species-specific lncRNAs.

1031

1032

**References**

1034 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol*
1035       **215**: 403–410.

1036 Amândio AR, Necsulea A, Joye E, Mascrez B, Duboule D. 2016. Hotair is dispensable for mouse
1037       development. *PLoS Genet* **12**: e1006232.

1038 Anders S, Huber W. 2010. Differential expression analysis for sequence count data. *Genome Biol* **11**:
1039       R106.

1040 Anderson KM, Anderson DM, McAnally JR, Shelton JM, Bassel-Duby R, Olson EN. 2016. Transcription
1041       of the non-coding RNA upperhand controls Hand2 expression and heart development. *Nature*
1042       **539**: 433–436.

1043 Arendt D, Musser JM, Baker CVH, Bergman A, Cepko C, Erwin DH, Pavlicev M, Schlosser G, Widder S,
1044       Laubichler MD, et al. 2016. The origin and evolution of cell types. *Nat Rev Genet* **17**: 744–757.

1045 Ayers KL, Davidson NM, Demiyah D, Roeszler KN, Grützner F, Sinclair AH, Oshlack A, Smith CA. 2013.
1046       RNA sequencing reveals sexually dimorphic gene expression before gonadal differentiation in
1047       chicken and allows comprehensive annotation of the W-chromosome. *Genome Biol* **14**.

1048 Bassett AR, Akhtar A, Barlow DP, Bird AP, Brockdorff N, Duboule D, Ephrussi A, Ferguson-Smith AC,
1049       Gingeras TR, Haerty W, et al. 2014. Considerations when investigating lncRNA function in vivo.
1050       *eLife* **3**: e03058.

1051 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful approach
1052       to multiple testing. *J Roy Stat Soc B* **57**: 289+300.

1053 Ben-Tabou de-Leon S, Davidson EH. 2007. Gene regulation: gene control network in development.
1054       *Annu Rev Biophys Biomol Struct* **36**: 191.

1055 Blanchette M, Kent WJ, Riemer C, Elnitski L, Smit AFA, Roskin KM, Baertsch R, Rosenbloom K, Clawson
1056       H, Green ED, et al. 2004. Aligning multiple genomic sequences with the threaded blockset
1057       aligner. *Genome Res* **14**: 708–715.

1058 Brannan CI, Dees EC, Ingram RS, Tilghman SM. 1990. The product of the H19 gene may function as an
1059       RNA. *Mol Cell Biol* **10**: 28–36.

1060 Brawand D, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, Weier M, Liechti A, Aximu-Petri A,
1061       Kircher M, et al. 2011. The evolution of gene expression levels in mammalian organs. *Nature*
1062       **478**: 343–348.

1063 Bray NL, Pimentel H, Melsted P, Pachter L. 2016. Near-optimal probabilistic RNA-seq quantification.
1064       *Nat Biotechnol* **34**: 525–527.

1065 Brown CJ, Ballabio A, Rupert JL, Lafreniere RG, Grompe M, Tonlorenzi R, Willard HF. 1991. A gene from
1066       the region of the human X inactivation centre is expressed exclusively from the inactive X
1067       chromosome. *Nature* **349**: 38–44.

1068 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+:
1069       architecture and applications. *BMC Bioinformatics* **10**: 421.

1070  Carelli FN, Hayakawa T, Go Y, Imai H, Warnefors M, Kaessmann H. 2016. The life history of retrocopies
1071         illuminates the evolution of new mammalian genes. *Genome Res* **26**: 301–314.

1072  Casper J, Zweig AS, Villarreal C, Tyner C, Speir ML, Rosenbloom KR, Raney BJ, Lee CM, Lee BT, Karolchik
1073         D, et al. 2018. The UCSC Genome Browser database: 2018 update. *Nucleic Acids Res* **46**: D762–
1074         D769.

1075  Cesana M, Cacchiarelli D, Legnini I, Santini T, Sthandier O, Chinappi M, Tramontano A, Bozzoni I. 2011.
1076         A long noncoding RNA controls muscle differentiation by functioning as a competing
1077         endogenous RNA. *Cell* **147**: 358–369.

1078  Cortez D, Marin R, Toledo-Flores D, Froidevaux L, Liechti A, Waters PD, Grützner F, Kaessmann H. 2014.
1079         Origins and functional evolution of Y chromosomes across mammals. *Nature* **508**: 488–93.

1080  Cunningham F, Achuthan P, Akanni W, Allen J, Amode MR, Armean IM, Bennett R, Bhai J, Billis K, Boddu
1081         S, et al. 2019. Ensembl 2019. *Nucleic Acids Res* **47**: D745–D751.

1082  Doolittle WF. 2018. We simply cannot go on being so vague about "function." *Genome Biol* **19**: 223.

1083  Dray S, Dufour AB. 2007. The ade4 package: implementing the duality diagram for ecologists. *J Stat*
1084         *Softw* **22**: 1–20.

1085  El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, Qureshi M, Richardson LJ, Salazar GA,
1086         Smart A, et al. 2019. The Pfam protein families database in 2019. *Nucleic Acids Res* **47**: D427–
1087         D432.

1088  Engreitz JM, Haines JE, Perez EM, Munson G, Chen J, Kane M, McDonel PE, Guttman M, Lander ES.
1089         2016. Local regulation of gene expression by lncRNA promoters, transcription and splicing.
1090         *Nature* **539**: 452–455.

1091  Gendrel A-V, Heard E. 2014. Noncoding RNAs and epigenetic mechanisms during X-chromosome
1092         inactivation. *Annu Rev Cell Dev Biol* **30**: 561–580.

1093  Goudarzi M, Berg K, Pieper LM, Schier AF. 2019. Individual long non-coding RNAs have no overt
1094         functions in zebrafish embryogenesis, viability and fertility. *eLife* **8**.

1095  Graur D, Zheng Y, Price N, Azevedo RBR, Zufall RA, Elhaik E. 2013. On the immortality of television sets:
1096         "function" in the human genome according to the evolution-free gospel of ENCODE. *Genome*
1097         *Biol Evol* **5**: 578–590.

1098  Green CD, Ma Q, Manske GL, Shami AN, Zheng X, Marini S, Moritz L, Sultan C, Gurczynski SJ, Moore BB,
1099         et al. 2018. A comprehensive roadmap of murine spermatogenesis defined by single-cell RNA-
1100         Seq. *Dev Cell* **46**: 651-667.e10.

1101  Groff AF, Sanchez-Gomez DB, Soruco MML, Gerhardinger C, Barutcu AR, Li E, Elcavage L, Plana O,
1102         Sanchez LV, Lee JC, et al. 2016. In vivo characterization of Linc-p21 reveals functional cis-
1103         regulatory DNA elements. *Cell Rep* **16**: 2178–2186.

1104  Grote P, Herrmann BG. 2015. Long noncoding RNAs in organogenesis: making the difference. *Trends*
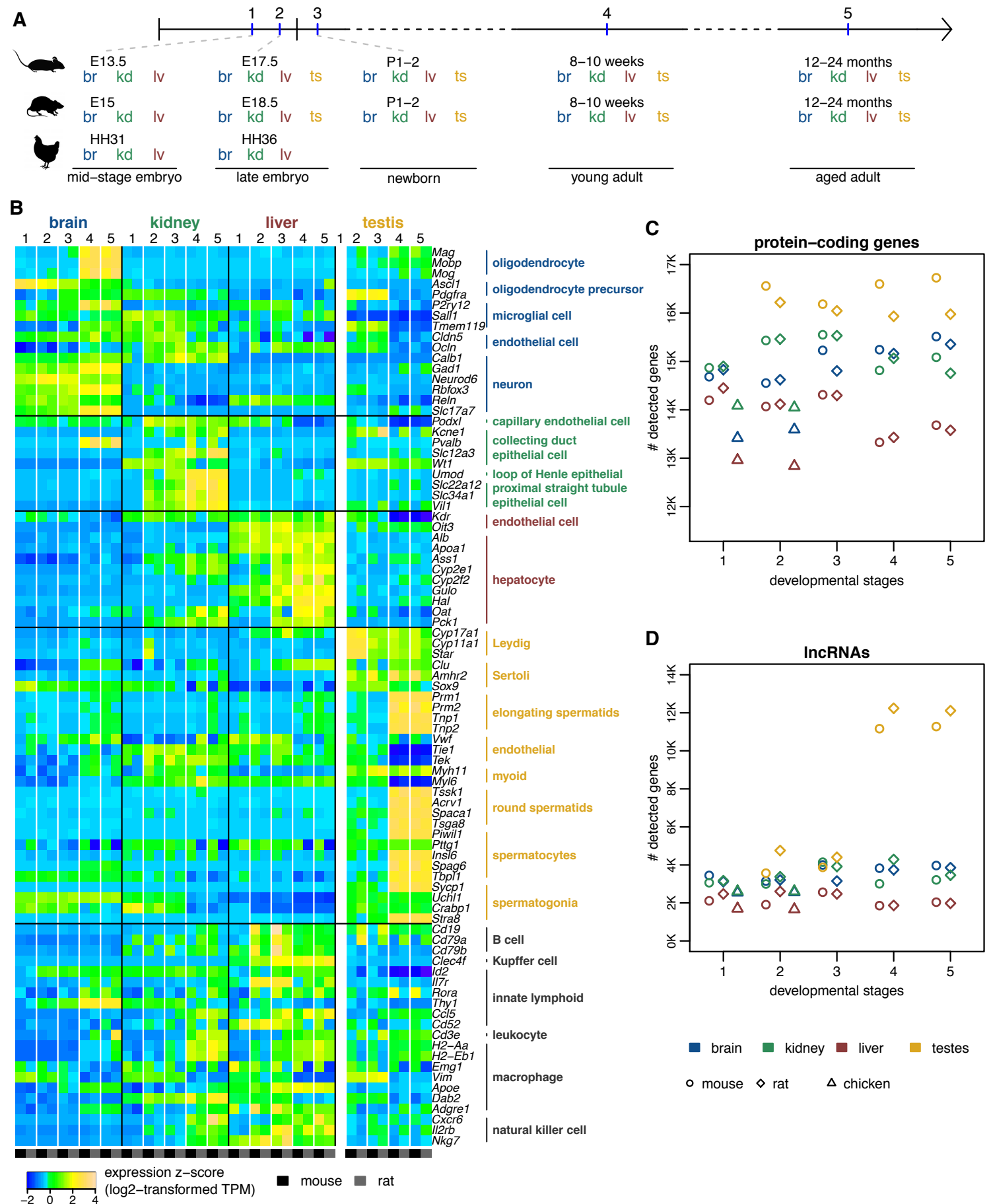1105         *Genet TIG* **31**: 329–335.

44

1106    Grote P, Wittler L, Hendrix D, Koch F, Währisch S, Beisaw A, Macura K, Bläss G, Kellis M, Werber M, et
1107         al. 2013. The tissue-specific lncRNA Fendrr is an essential regulator of heart and body wall
1108         development in the mouse. *Dev Cell* **24**: 206–214.

1109    Guttman M, Amit I, Garber M, French C, Lin MF, Feldser D, Huarte M, Zuk O, Carey BW, Cassady JP, et
1110         al. 2009. Chromatin signature reveals over a thousand highly conserved large non-coding RNAs
1111         in mammals. *Nature* **458**: 223–227.

1112    Haerty W, Ponting CP. 2014. No gene in the genome makes sense except in the light of evolution. *Annu*
1113         *Rev Genomics Hum Genet* **15**: 71–92.

1114    Haerty W, Ponting CP. 2015. Unexpected selection to retain high GC content and splicing enhancers
1115         within exons of multiexonic lncRNA loci. *RNA N Y N* **21**: 333–346.

1116    Hamburger V, Hamilton HL. 1951. A series of normal stages in the development of the chick embryo. *J*
1117         *Morphol* **88**: 49–92.

1118    Herrero J, Muffato M, Beal K, Fitzgerald S, Gordon L, Pignatelli M, Vilella AJ, Searle SMJ, Amode R, Brent
1119         S, et al. 2016. Ensembl comparative genomics resources. *Database J Biol Databases Curation*
1120         **2016**.

1121    Hezroni H, Koppstein D, Schwartz MG, Avrutin A, Bartel DP, Ulitsky I. 2015. Principles of long noncoding
1122         RNA evolution derived from direct comparison of transcriptomes in 17 species. *Cell Rep* **11**:
1123         1110–1122.

1124    Iyer MK, Niknafs YS, Malik R, Singhal U, Sahu A, Hosono Y, Barrette TR, Prensner JR, Evans JR, Zhao S,
1125         et al. 2015. The landscape of long noncoding RNAs in the human transcriptome. *Nat Genet* **47**:
1126         199–208.

1127    Kaessmann H. 2010. Origins, evolution, and phenotypic impact of new genes. *Genome Res* **20**: 1313–
1128         1326.

1129    Kathleen Baxter K, Uittenbogaard M, Yoon J, Chiaramello A. 2009. The neurogenic basic helix-loop-
1130         helix transcription factor NeuroD6 concomitantly increases mitochondrial mass and regulates
1131         cytoskeletal organization in the early stages of neuronal differentiation. *ASN Neuro* **1**.

1132    Khalil AM, Guttman M, Huarte M, Garber M, Raj A, Rivea Morales D, Thomas K, Presser A, Bernstein
1133         BE, van Oudenaarden A, et al. 2009. Many human large intergenic noncoding RNAs associate
1134         with chromatin-modifying complexes and affect gene expression. *Proc Natl Acad Sci U S A* **106**:
1135         11667–11672.

1136    Kim D, Langmead B, Salzberg SL. 2015. HISAT: a fast spliced aligner with low memory requirements.
1137         *Nat Methods* **12**: 357–60.

1138    Kutter C, Watt S, Stefflova K, Wilson MD, Goncalves A, Ponting CP, Odom DT, Marques AC. 2012. Rapid
1139         turnover of long noncoding RNAs and the evolution of gene expression. *PLoS Genet* **8**:
1140         e1002841.

1141    Latos PA, Pauler FM, Koerner MV, Şenergin HB, Hudson QJ, Stocsits RR, Allhoff W, Stricker SH, Klement
1142         RM, Warczok KE, et al. 2012. Airn transcriptional overlap, but not its lncRNA products, induces
1143         imprinted Igf2r silencing. *Science* **338**: 1469–1472.

1144  Liao B-Y, Scott NM, Zhang J. 2006. Impacts of gene essentiality, expression pattern, and gene
1145      compactness on the evolutionary rate of mammalian proteins. *Mol Biol Evol* **23**: 2072–2080.

1146  Liao Y, Smyth GK, Shi W. 2019. The R package Rsubread is easier, faster, cheaper and better for
1147      alignment and quantification of RNA sequencing reads. *Nucleic Acids Res*.

1148  Lin MF, Carlson JW, Crosby MA, Matthews BB, Yu C, Park S, Wan KH, Schroeder AJ, Gramates LS, St
1149      Pierre SE, et al. 2007. Revisiting the protein-coding gene catalog of Drosophila melanogaster
1150      using 12 fly genomes. *Genome Res* **17**: 1823–1836.

1151  Lin MF, Jungreis I, Kellis M. 2011. PhyloCSF: a comparative genomics method to distinguish protein
1152      coding and non-coding regions. *Bioinforma Oxf Engl* **27**: i275-282.

1153  Liu SJ, Nowakowski TJ, Pollen AA, Lui JH, Horlbeck MA, Attenello FJ, He D, Weissman JS, Kriegstein AR,
1154      Diaz AA, et al. 2016. Single-cell analysis of long non-coding RNAs in the developing human
1155      neocortex. *Genome Biol* **17**: 67.

1156  Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for RNA-seq
1157      data with DESeq2. *Genome Biol* **15**: 550.

1158  Marques AC, Hughes J, Graham B, Kowalczyk MS, Higgs DR, Ponting CP. 2013. Chromatin signatures at
1159      transcriptional start sites separate two equally populated yet distinct classes of intergenic long
1160      noncoding RNAs. *Genome Biol* **14**: R131.

1161  McMahon AP. 2016. Development of the Mammalian Kidney. *Curr Top Dev Biol* **117**: 31–64.

1162  Nakagaki BN, Mafra K, de Carvalho É, Lopes ME, Carvalho-Gontijo R, de Castro-Oliveira HM, Campolina-
1163      Silva GH, de Miranda CDM, Antunes MM, Silva ACC, et al. 2018. Immune and metabolic shifts
1164      during neonatal development reprogram liver identity and function. *J Hepatol* **69**: 1294–1307.

1165  Necsulea A, Soumillon M, Warnefors M, Liechti A, Daish T, Grutzner F, Kaessmann H. 2014. The
1166      evolution of lncRNA repertoires and expression patterns in tetrapods. *Nature* **505**: 635–640.

1167  Ørom UA, Derrien T, Beringer M, Gumireddy K, Gardini A, Bussotti G, Lai F, Zytnicki M, Notredame C,
1168      Huang Q, et al. 2010. Long noncoding RNAs with enhancer-like function in human cells. *Cell*
1169      **143**: 46–58.

1170  Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL. 2015. StringTie enables
1171      improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol* **33**: 290–5.

1172  Pertea M, Shumate A, Pertea G, Varabyou A, Breitwieser FP, Chang Y-C, Madugundu AK, Pandey A,
1173      Salzberg SL. 2018. CHESS: a new human gene catalog curated from thousands of large-scale
1174      RNA sequencing experiments reveals extensive transcriptional noise. *Genome Biol* **19**: 208.

1175  Ponjavic J, Ponting CP, Lunter G. 2007. Functionality or transcriptional noise? Evidence for selection
1176      within long noncoding RNAs. *Genome Res* **17**: 556–65.

1177  R Core Team. 2018. *R: A Language and Environment for Statistical Computing*. https://www.R-
1178      project.org/.

1179  Ramsköld D, Wang ET, Burge CB, Sandberg R. 2009. An abundance of ubiquitously expressed genes
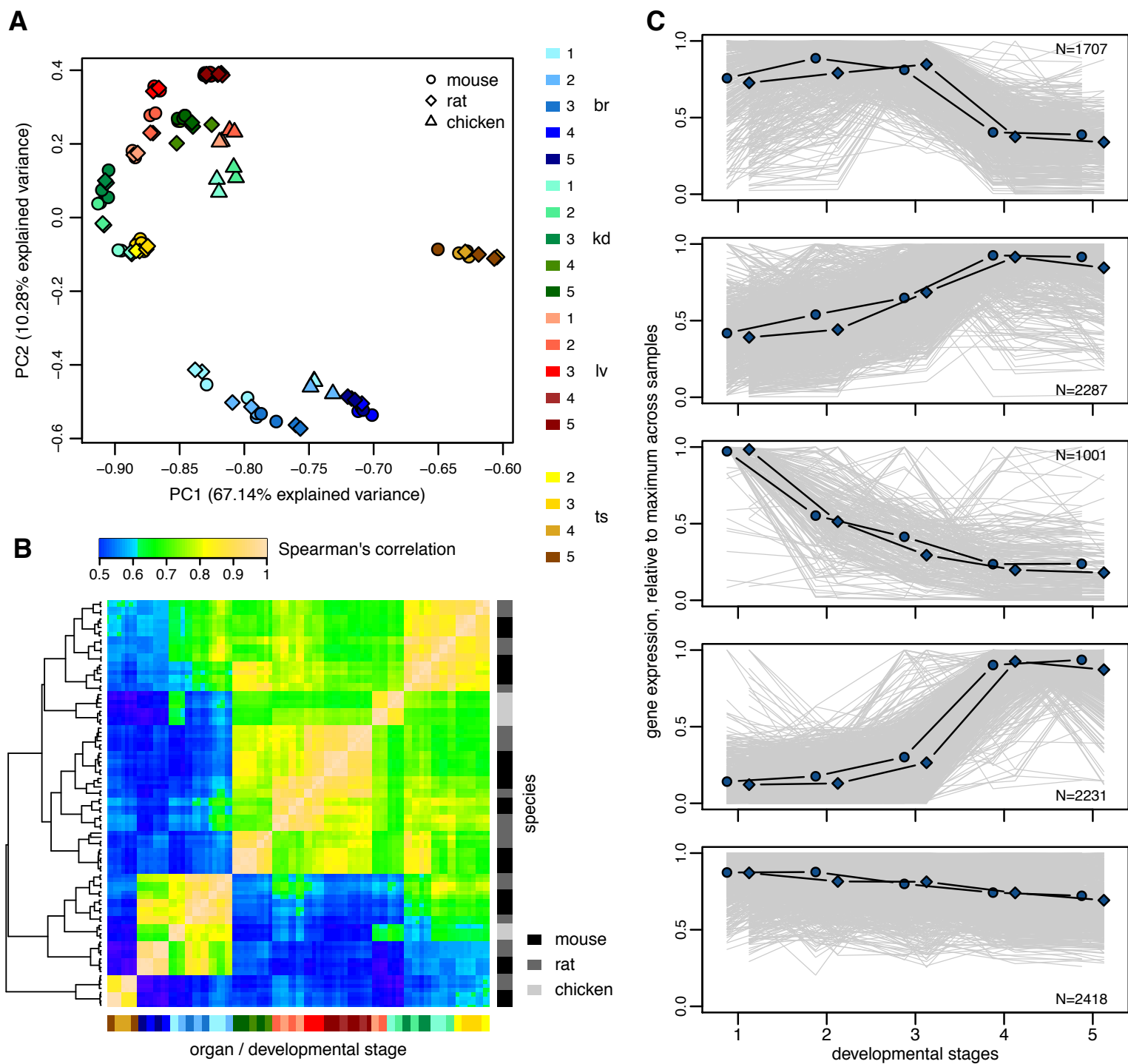1180      revealed by tissue transcriptome sequence data. *PLoS Comput Biol* **5**: e1000598.

1181    Rinn JL, Kertesz M, Wang JK, Squazzo SL, Xu X, Brugmann SA, Goodnough LH, Helms JA, Farnham PJ,
1182        Segal E, et al. 2007. Functional demarcation of active and silent chromatin domains in human
1183        HOX loci by noncoding RNAs. *Cell* **129**: 1311–1323.

1184    Sauvageau M, Goff LA, Lodato S, Bonev B, Groff AF, Gerhardinger C, Sanchez-Gomez DB, Hacisuleyman
1185        E, Li E, Spence M, et al. 2013. Multiple knockout mouse models reveal lincRNAs are required
1186        for life and brain development. *eLife* **2**: e01749.

1187    Schüler A, Ghanbarian AT, Hurst LD. 2014. Purifying selection on splice-related motifs, not expression
1188        level nor RNA folding, explains nearly all constraint on human lincRNAs. *Mol Biol Evol* **31**: 3164–
1189        3183.

1190    Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J, Hillier LW,
1191        Richards S, et al. 2005. Evolutionarily conserved elements in vertebrate, insect, worm, and
1192        yeast genomes. *Genome Res* **15**: 1034–50.

1193    Smedley D, Haider S, Ballester B, Holland R, London D, Thorisson G, Kasprzyk A. 2009. BioMart--
1194        biological queries made easy. *BMC Genomics* **10**: 22.

1195    Smit AF., Hubley R, Green P. 2003. *RepeatMasker Open-4.0.* http://www.repeatmasker.org.

1196    Soumillon M, Necsulea A, Weier M, Brawand D, Zhang X, Gu H, Barthès P, Kokkinaki M, Nef S, Gnirke
1197        A, et al. 2013. Cellular source and mechanisms of high transcriptome complexity in the
1198        mammalian testis. *Cell Rep* **3**: 2179–2190.

1199    Tabula Muris Consortium. 2018. Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris.
1200        *Nature* **562**: 367–372.

1201    Theiler K. 1989. *The house mouse: atlas of embryonic development*. Springer-Verlag, Berlin Heidelberg
1202        https://www.springer.com/la/book/9783642884207.

1203    The UniProt Consortium. 2017. UniProt: the universal protein knowledgebase. *Nucleic Acids Res* **45**:
1204        D158–D169.

1205    Uebbing S, Konzer A, Xu L, Backström N, Brunström B, Bergquist J, Ellegren H. 2015. Quantitative mass
1206        spectrometry reveals partial translational regulation for dosage compensation in chicken. *Mol
1207        Biol Evol* **32**: 2716–2725.

1208    Ulitsky I. 2016. Evolution to the rescue: using comparative genomics to understand long non-coding
1209        RNAs. *Nat Rev Genet* **17**: 601–614.

1210    Ulitsky I, Shkumatava A, Jan CH, Sive H, Bartel DP. 2011. Conserved function of lincRNAs in vertebrate
1211        embryonic development despite rapid sequence evolution. *Cell* **147**: 1537–1550.

1212    Villar D, Berthelot C, Aldridge S, Rayner TF, Lukk M, Pignatelli M, Park TJ, Deaville R, Erichsen JT,
1213        Jasinska AJ, et al. 2015. Enhancer evolution across 20 mammalian species. *Cell* **160**: 554–566.

1214    Washietl S, Kellis M, Garber M. 2014. Evolutionary dynamics and tissue specificity of human long
1215        noncoding RNAs in six mammals. *Genome Res* **24**: 616–28.

1216    Zakany J, Darbellay F, Mascrez B, Necsulea A, Duboule D. 2017. Control of growth and gut maturation
1217        by HoxD genes and the associated lncRNA Haglr. *Proc Natl Acad Sci U S A* **114**: E9290–E9299.
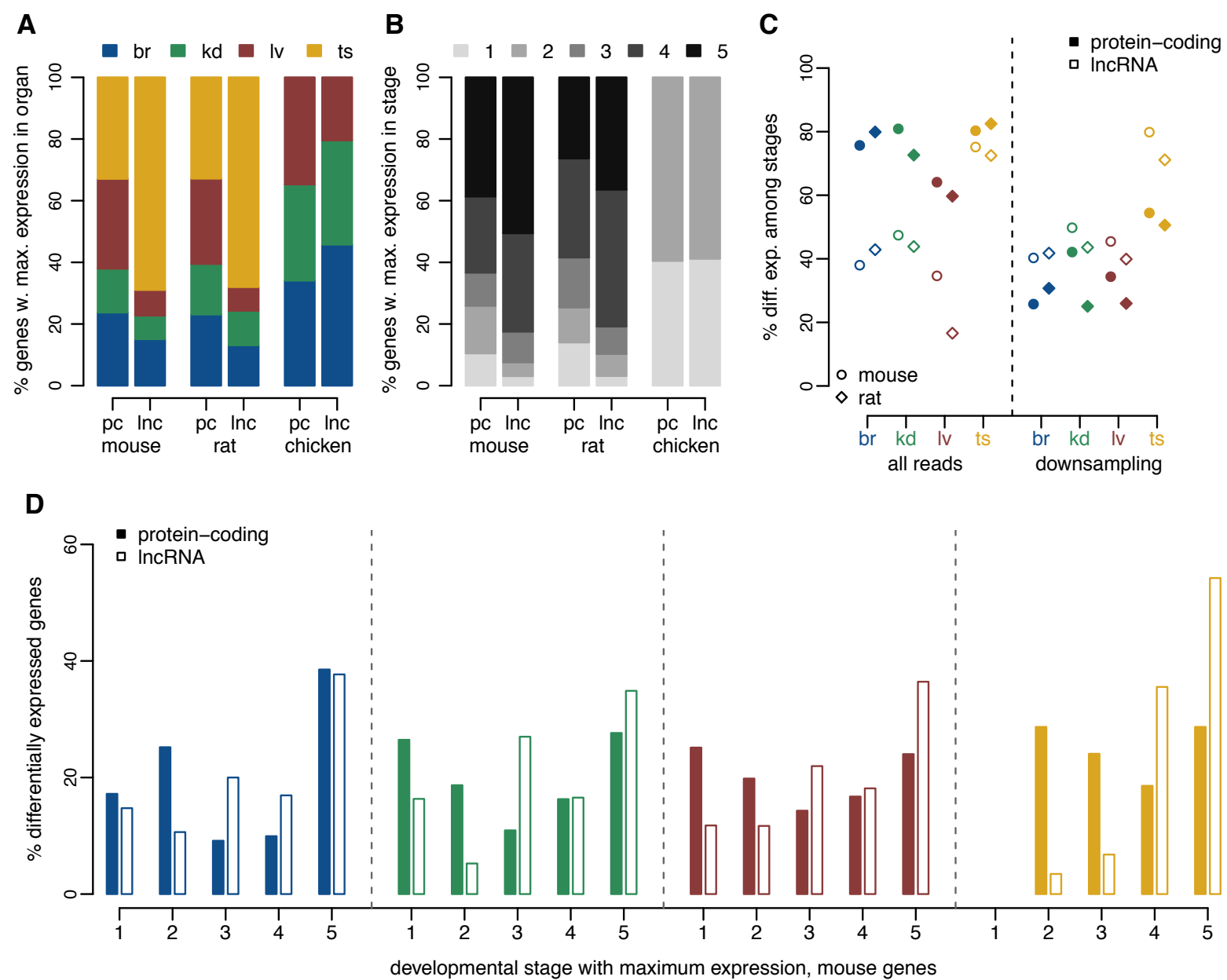
1218

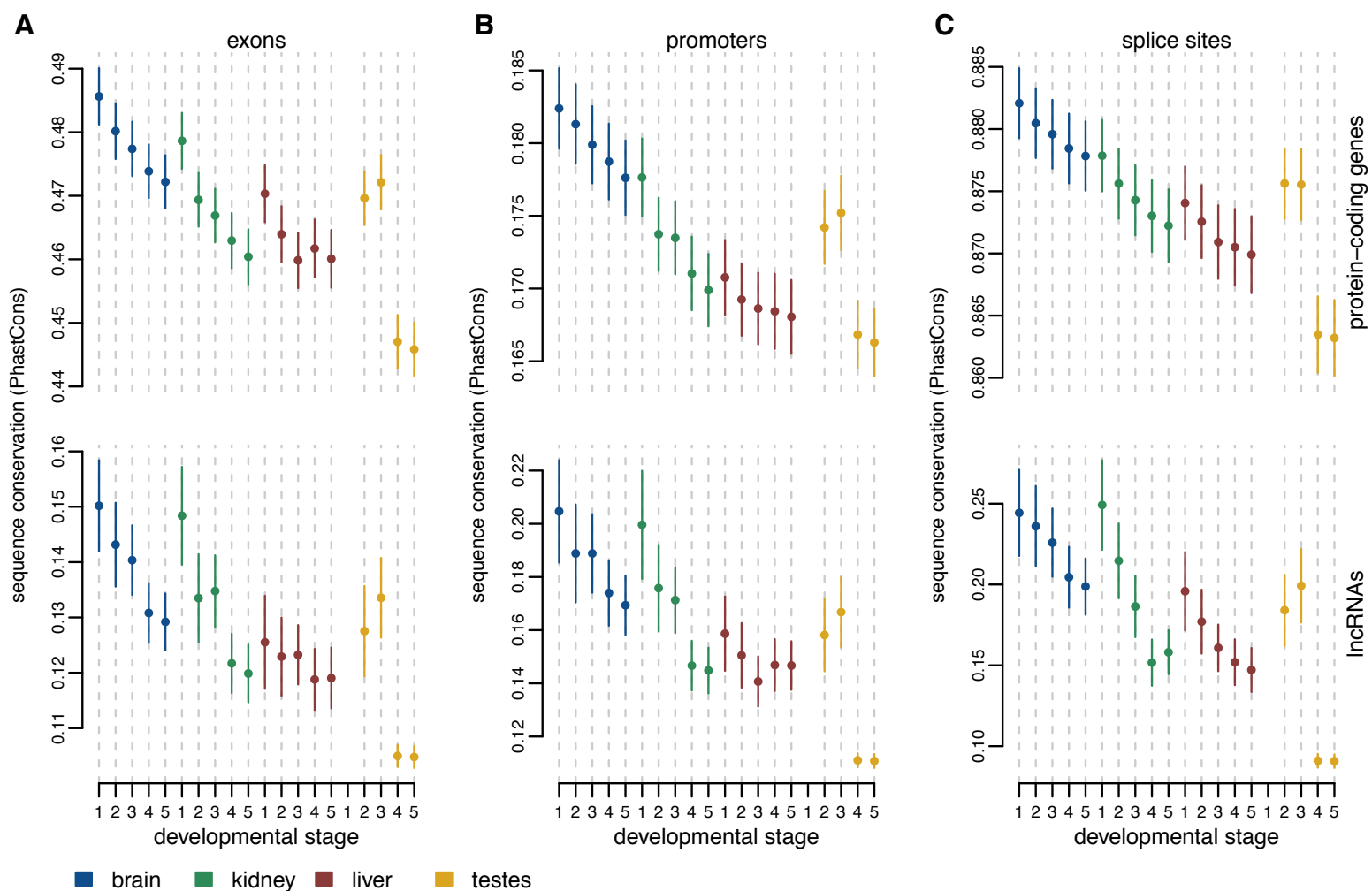Darbellay and Necsulea, Figure 1

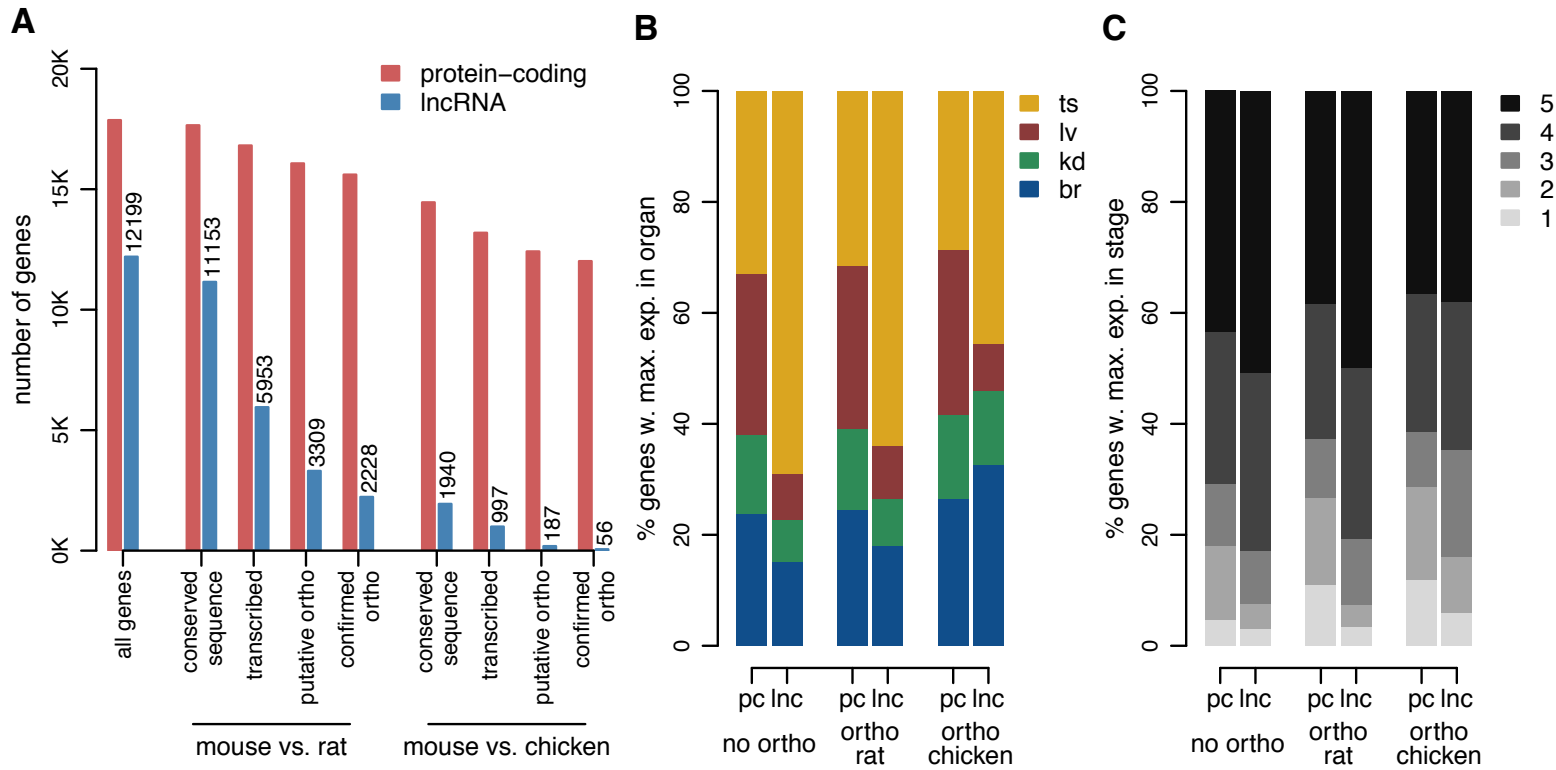Darbellay and Necsulea, Figure 2
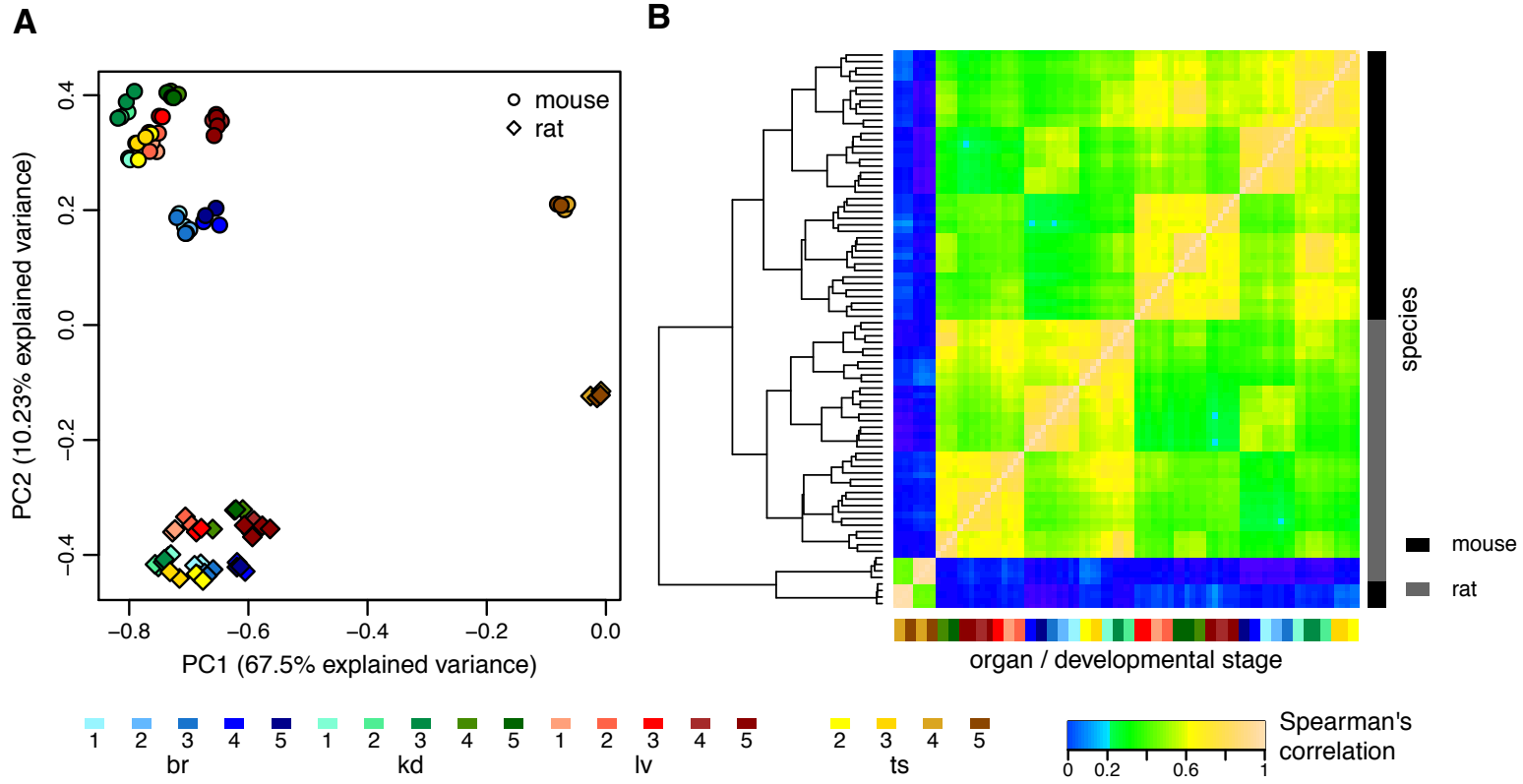
Darbellay and Necsulea, Figure 3
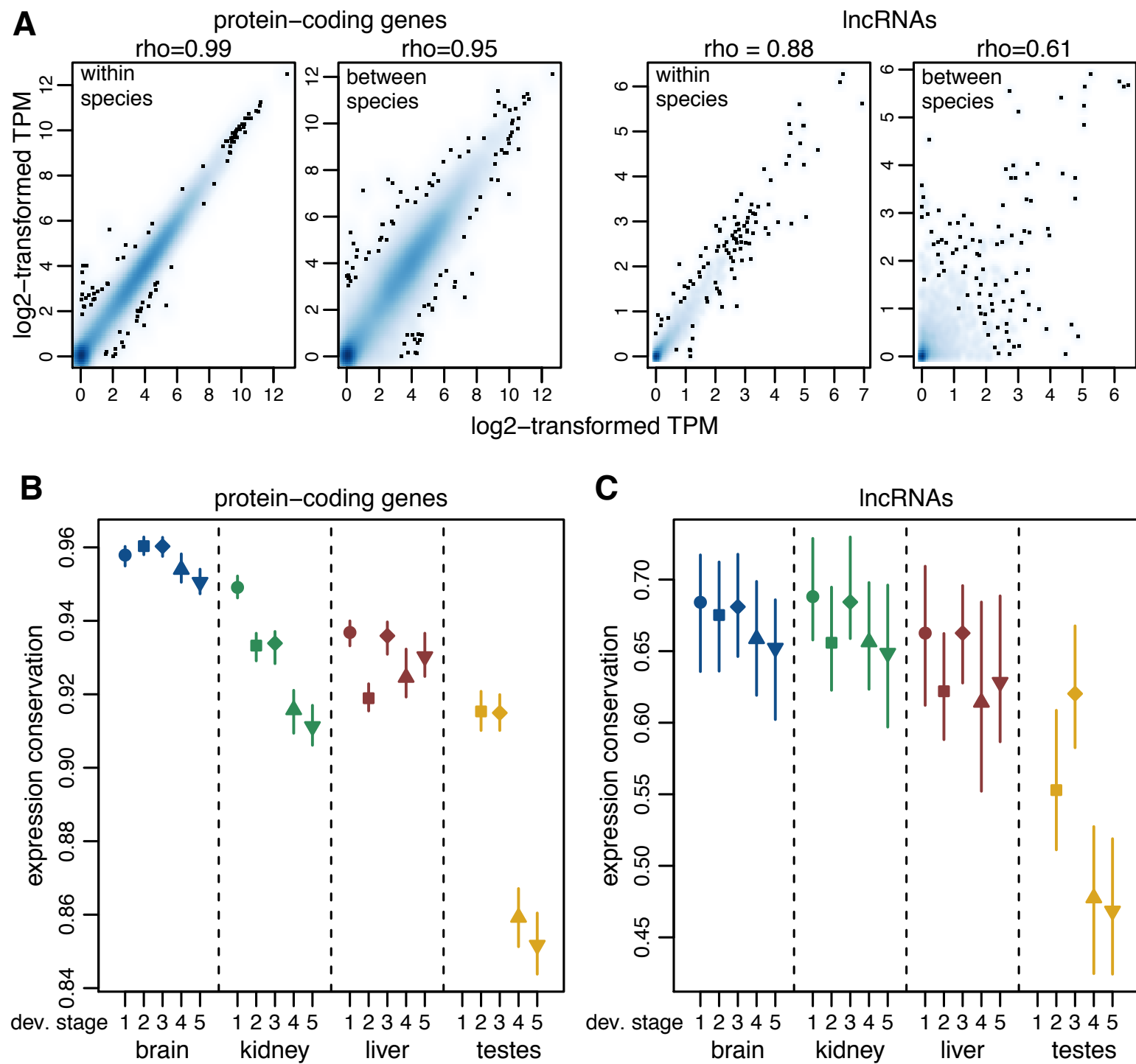
Darbellay and Necsulea, Figure 4
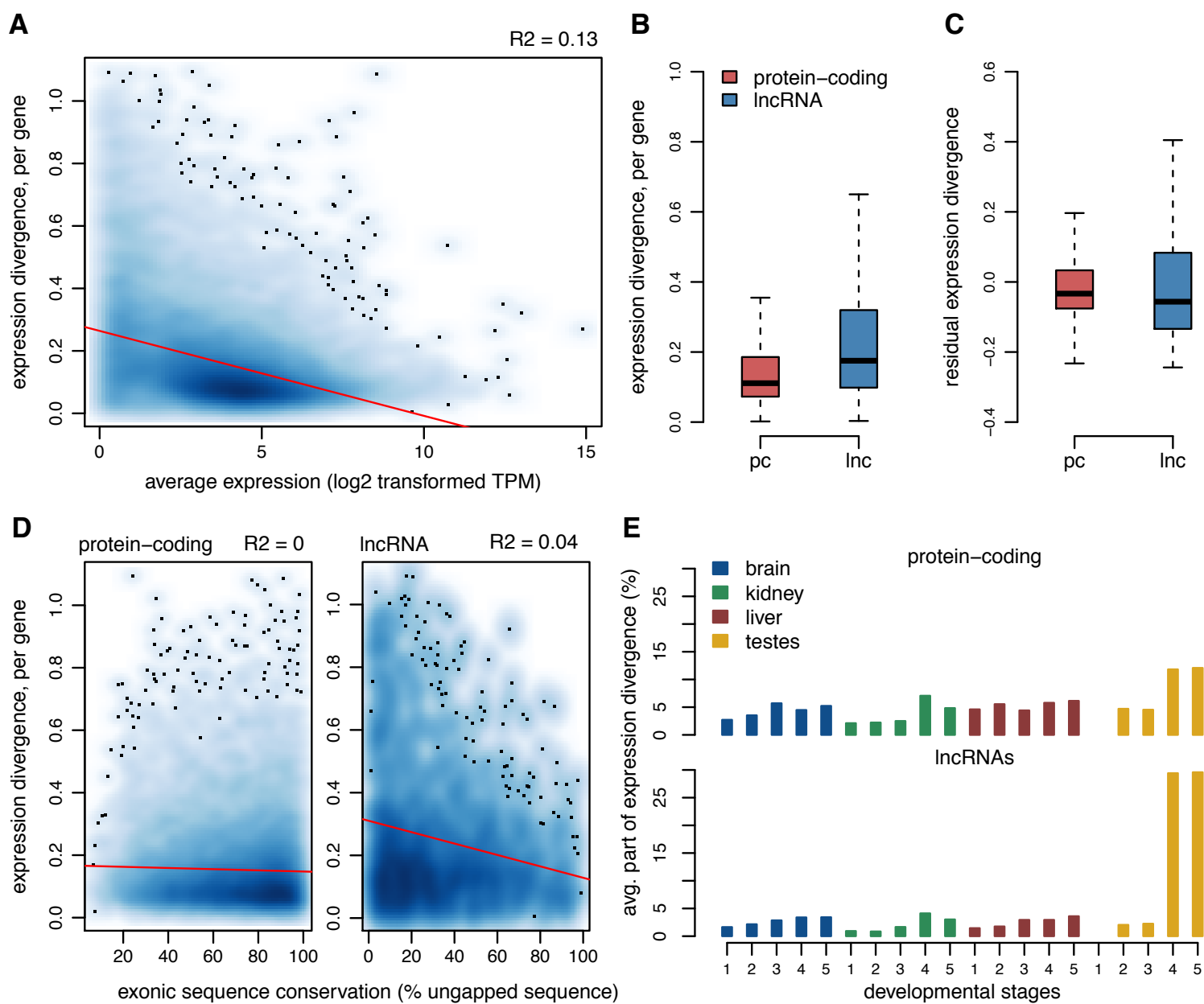
Darbellay and Necsulea, Figure 5
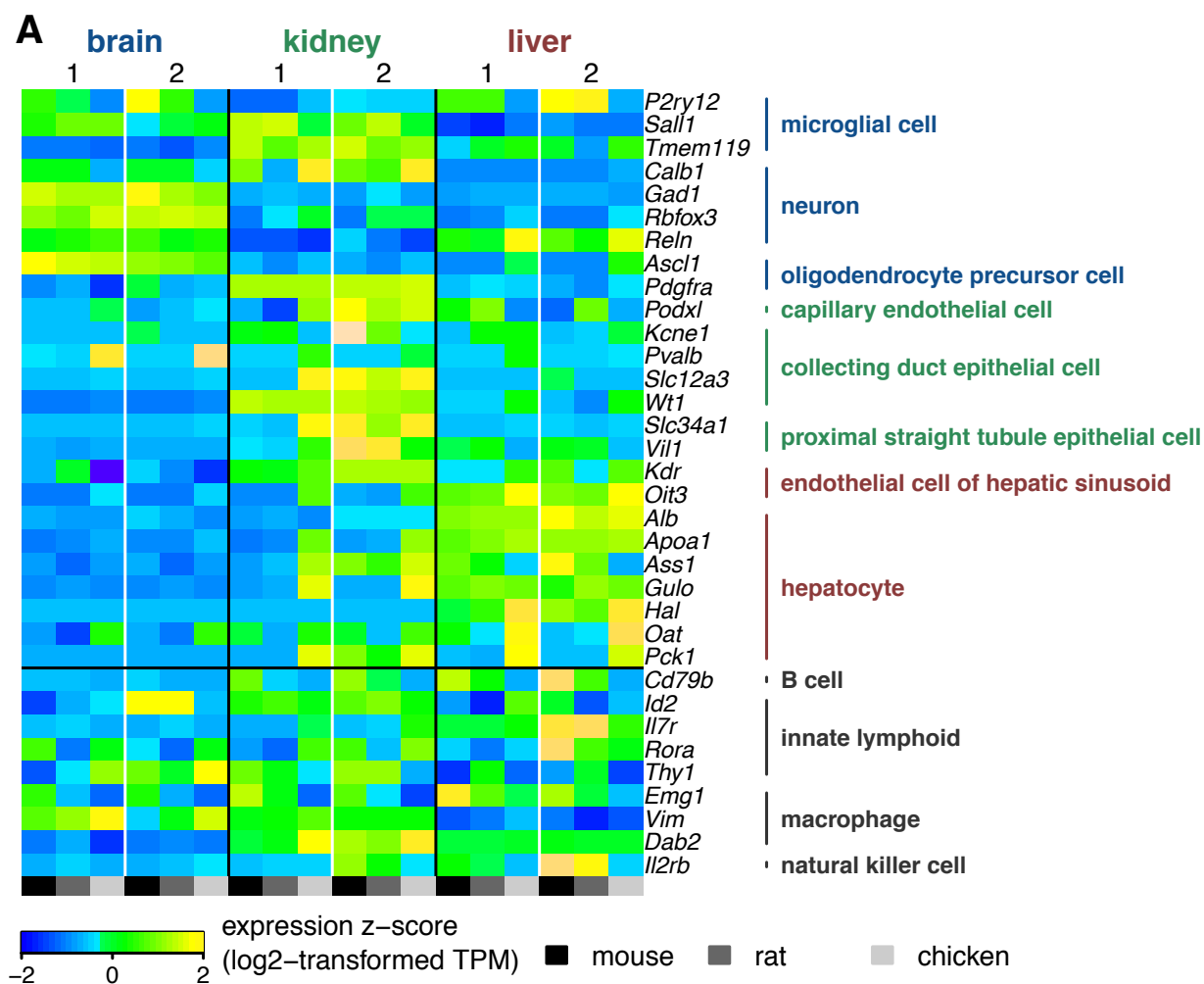
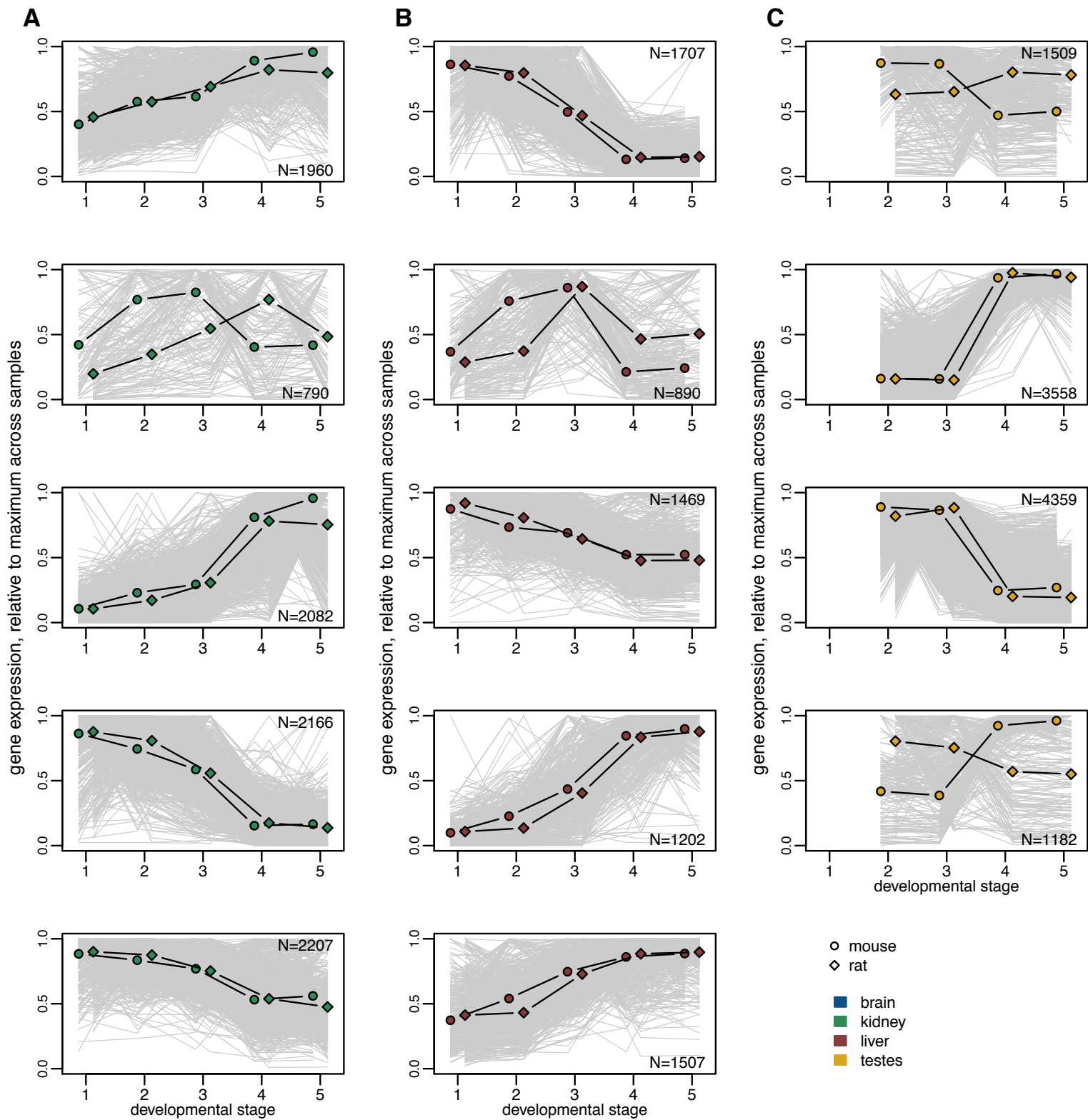Darbellay and Necsulea, Figure 6

Darbellay and Necsulea, Figure 7

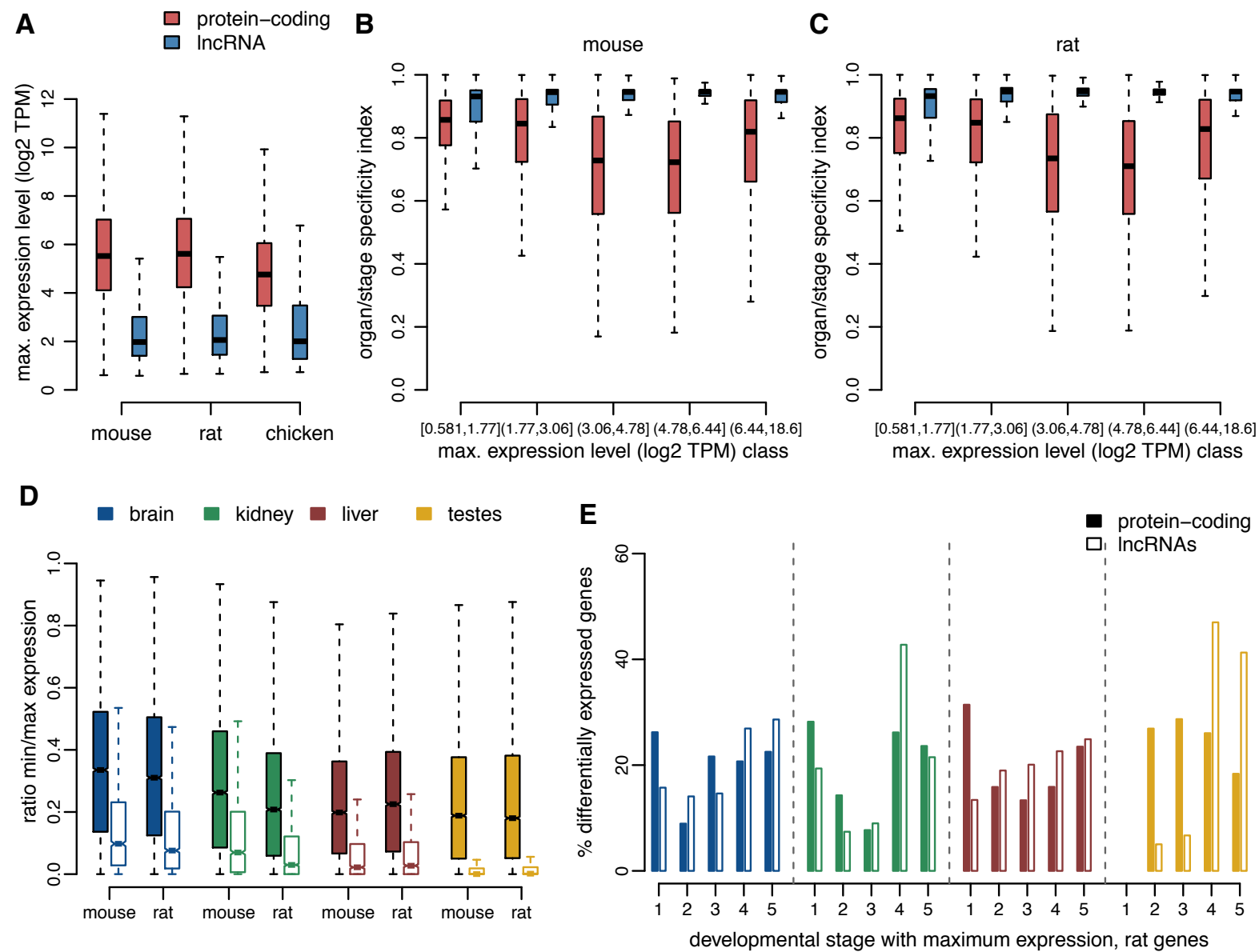Darbellay and Necsulea, Figure 8
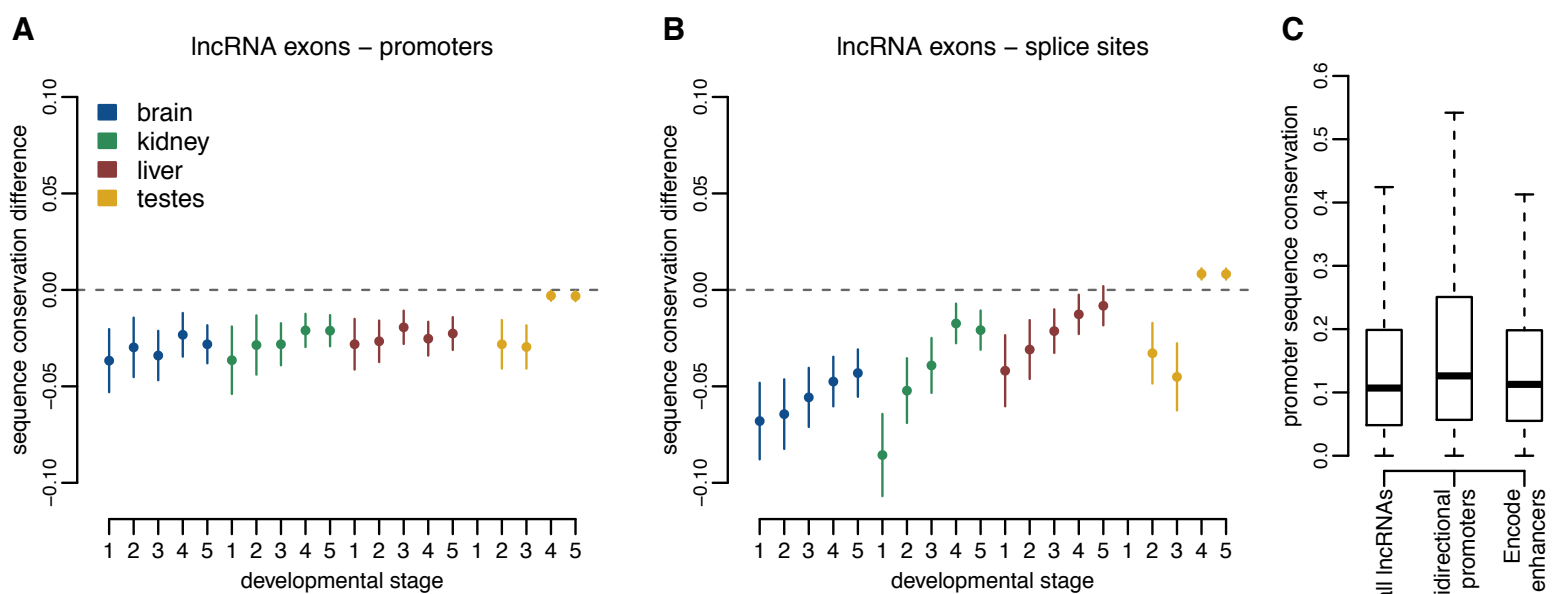
# Darbellay and Necsulea, Supplementary Figure 1

Darbellay and Necsulea, Supplementary Figure 2
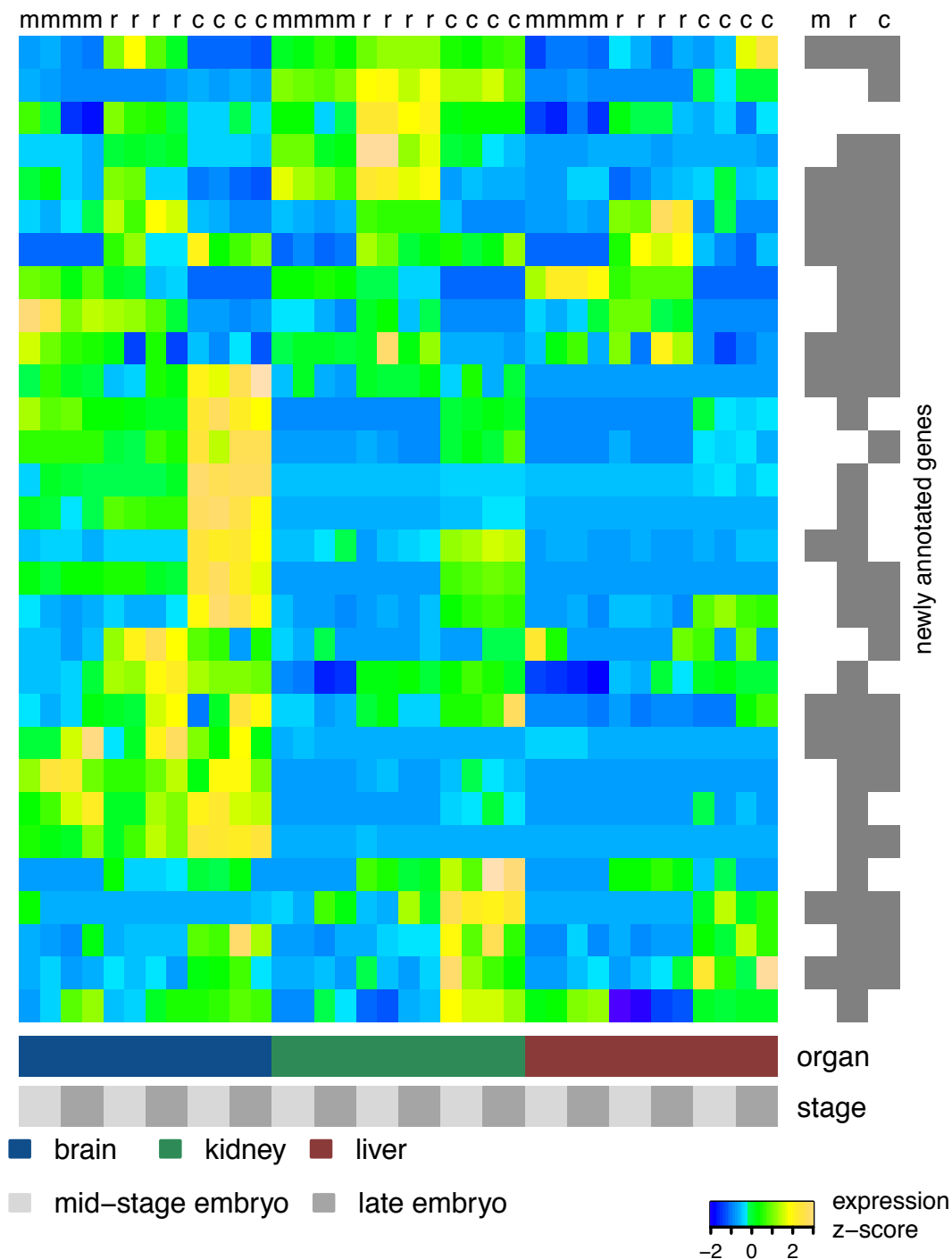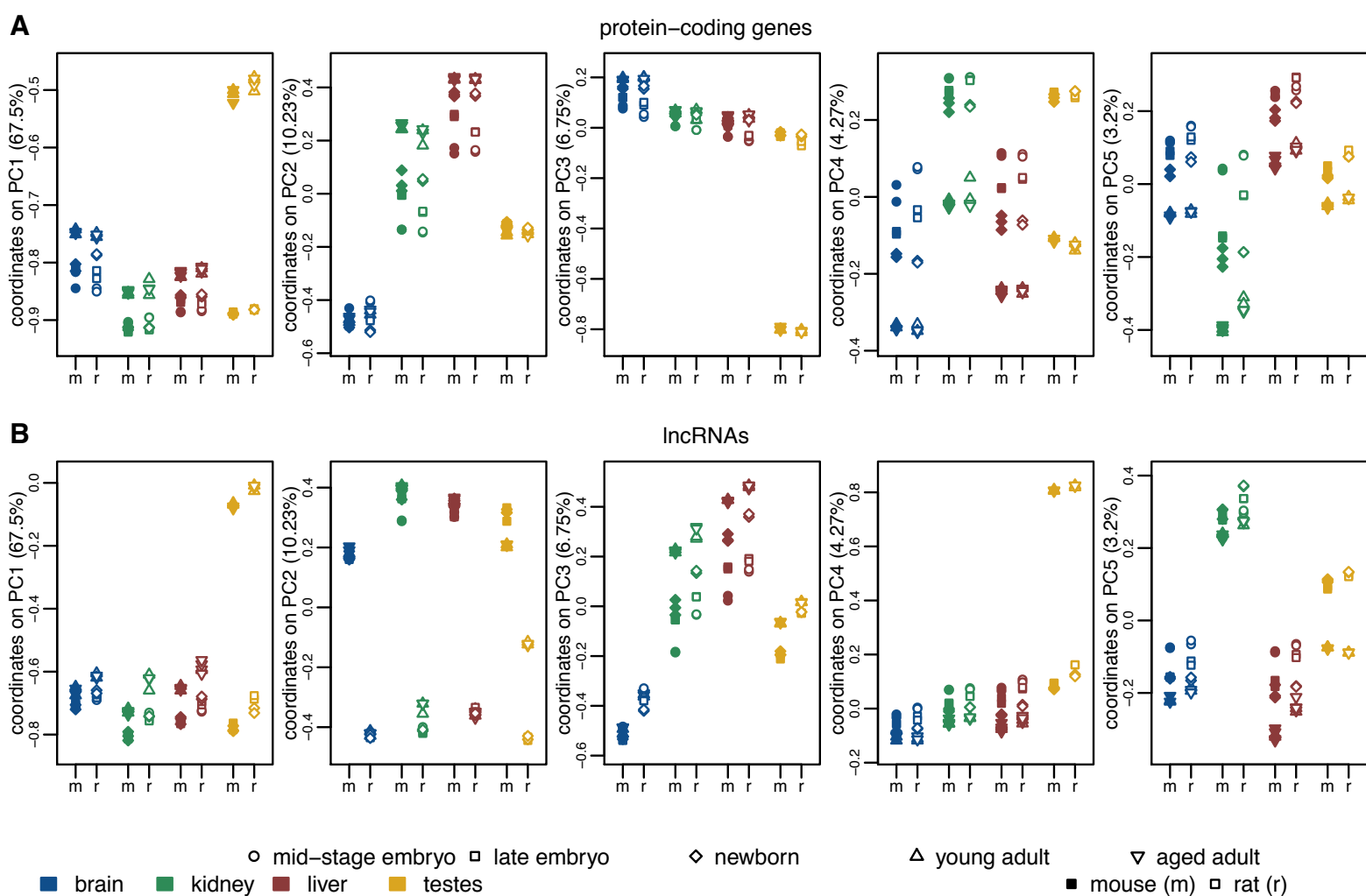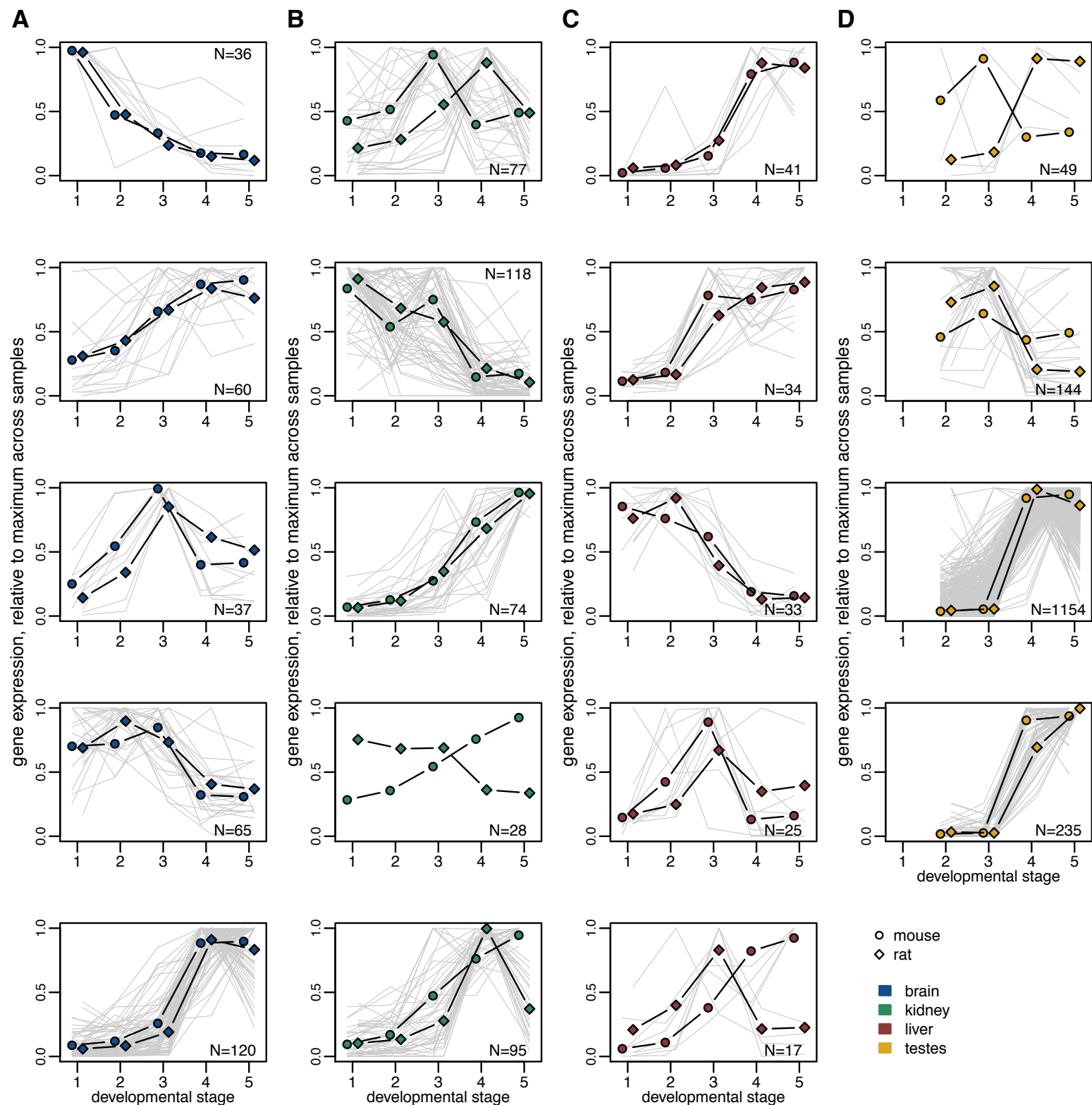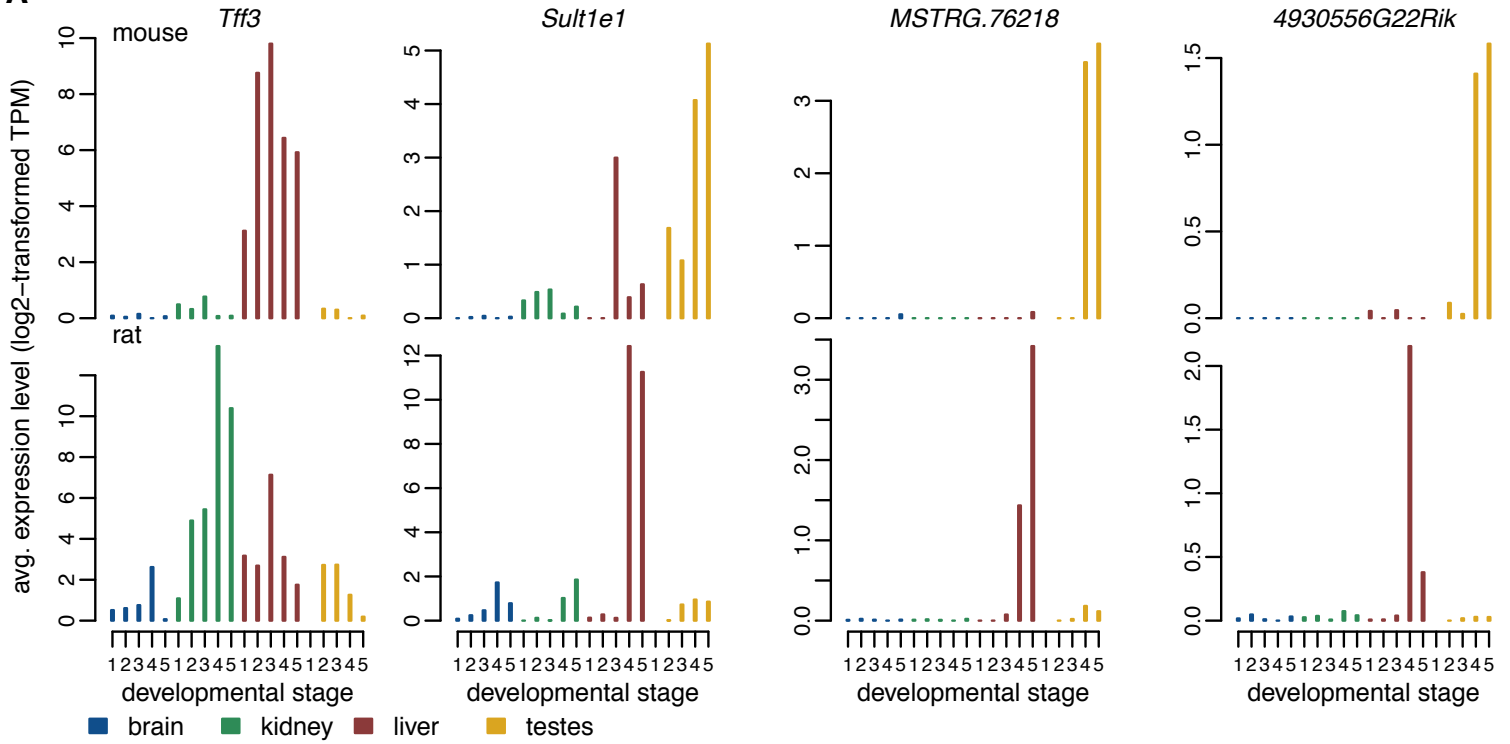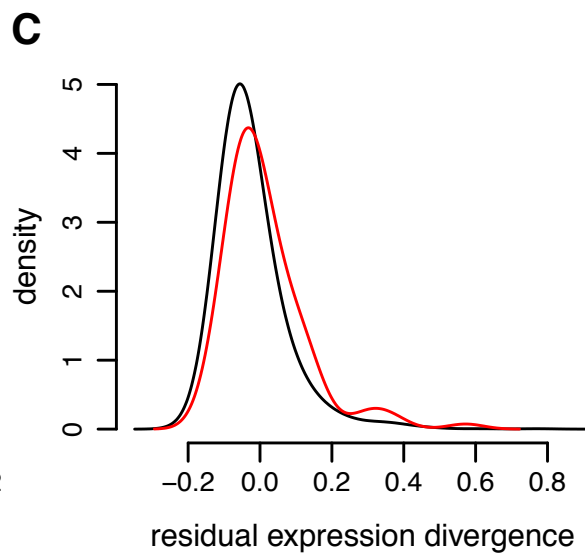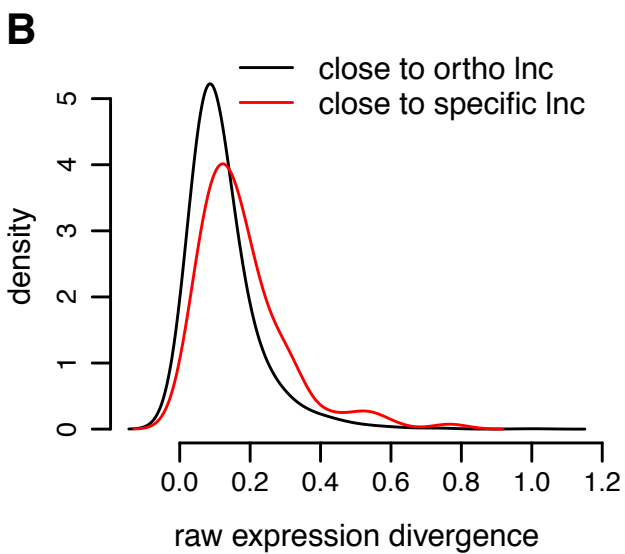
Darbellay and Necsulea, Supplementary Figure 3

Darbellay and Necsulea, Supplementary Figure 4

Darbellay and Necsulea, Supplementary Figure 5

**A**

Darbellay and Necsulea, Supplementary Figure 6
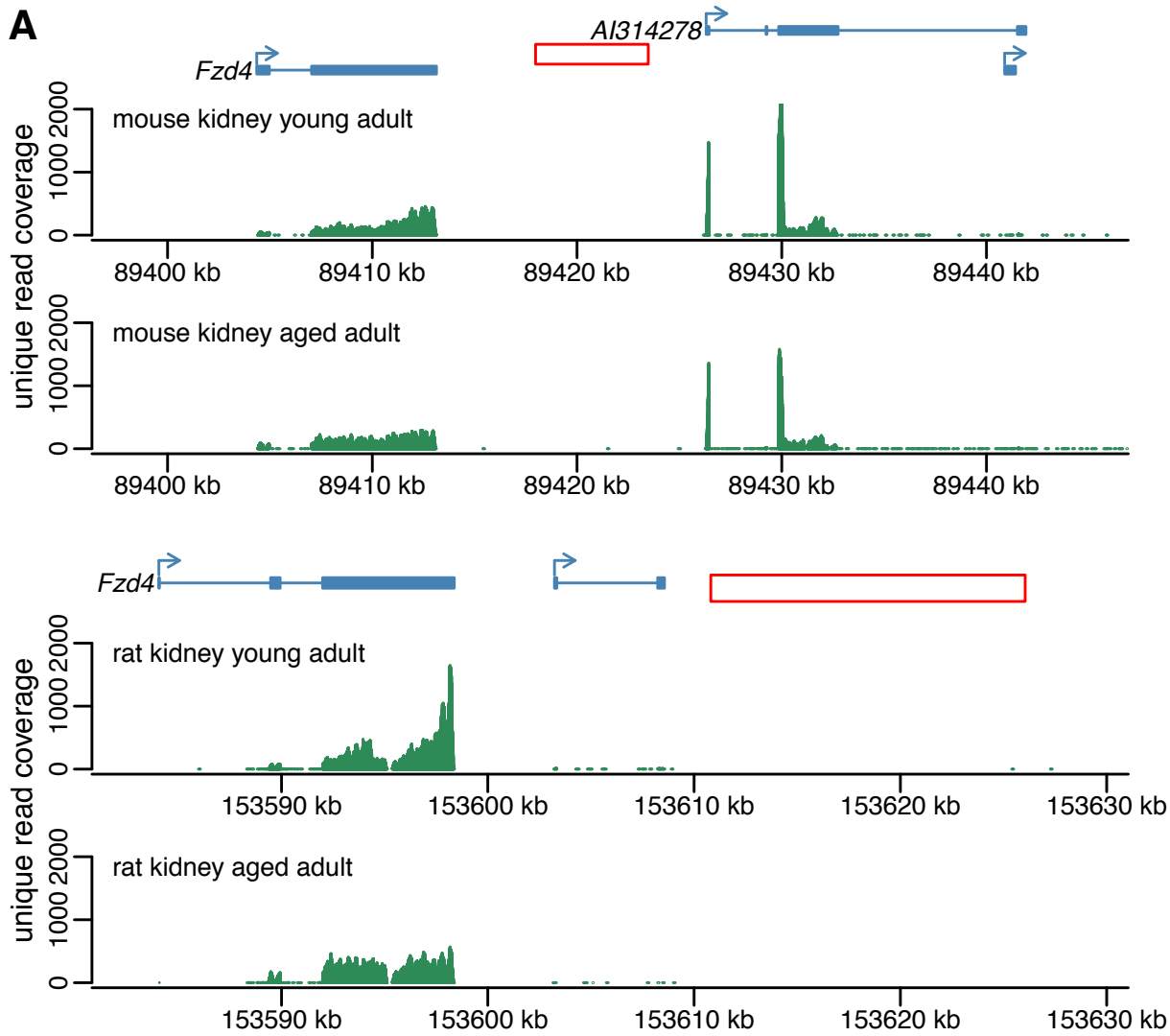
Darbellay and Necsulea, Supplementary Figure 7

Darbellay and Necsulea, Supplementary Figure 8

**A**

Darbellay and Necsulea, Supplementary Figure 9

Darbellay and Necsulea, Supplementary Figure 10