

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

Widespread localisation of lncRNA to ribosomes: Distinguishing features and evidence for regulatory roles.

Juna Carlevaro-Fita^{1,2,3}, Anisa Rahim⁴, Roderic Guigo^{1,2,3}, Leah A. Vardy^{4,5,6}, Rory Johnson^{1,2,3,6}

1. Centre for Genomic Regulation (CRG), Dr. Aiguader 88, 08003 Barcelona, Spain.
2. Universitat Pompeu Fabra (UPF), Barcelona, Spain.
3. Institut Hospital del Mar d'Investigacions Mèdiques (IMIM), 08003 Barcelona, Spain.
4. A*STAR Institute of Medical Biology, 8A Biomedical Grove, Immunos, Singapore 138648.
5. School of Biological Sciences, Nanyang Technological University, 60 Nanyang Drive, Singapore.
6. Correspondence to RJ (rorj.johnson@crg.eu) or LV (leah.vardy@imb.a-star.edu.sg).

Keywords: Long noncoding RNA; lncRNA; ribosome; translation; cytoplasm; transposable element; antisense.

25 **Abstract**

26 The function of long noncoding RNAs (lncRNAs) depends on their location within the
27 cell. While most studies to date have concentrated on their nuclear roles in
28 transcriptional regulation, evidence is mounting that lncRNA also have cytoplasmic
29 roles. Here we comprehensively map the cytoplasmic and ribosomal lncRNA population
30 in a human cell. Three-quarters (74%) of lncRNAs are detected in the cytoplasm, the
31 majority of which (62%) preferentially cofractionate with polyribosomes. Ribosomal
32 lncRNA are highly expressed across tissues, under purifying evolutionary selection, and
33 have cytoplasmic-to-nuclear ratios comparable to mRNAs and consistent across cell
34 types. lncRNAs may be classified into three groups by their ribosomal interaction: non-
35 ribosomal cytoplasmic lncRNAs, and those associated with either heavy or light
36 polysomes. A number of mRNA-like features destin lncRNA for light polysomes,
37 including capping and 5'UTR length, but not cryptic open reading frames or
38 polyadenylation. Surprisingly, exonic retroviral sequences antagonise recruitment. In
39 contrast, it appears that lncRNAs are recruited to heavy polysomes through basepairing
40 to mRNAs. Finally, we show that the translation machinery actively degrades lncRNA.
41 We propose that light polysomal lncRNAs are translationally engaged, while heavy
42 polysomal lncRNAs are recruited indirectly. These findings point to extensive and
43 reciprocal regulatory interactions between lncRNA and the translation machinery.

44

45

46 **Introduction**

47 The past decade has witnessed the discovery of a tens of thousands of long non-
48 protein coding RNAs (lncRNAs) in our genome, with profound implications for our
49 understanding of molecular genetics, disease and evolution. Focus is now shifting to
50 understanding the function to these molecules. We reason that such function is likely to
51 be intimately linked to the location of lncRNA within the cell.

52 Following the first compelling discoveries of chromatin regulatory lncRNAs such
53 as XIST (1) and HOTAIR (2), a paradigm was established for lncRNAs as nuclear-
54 restricted, epigenetic regulatory molecules (3). However, it is not clear to what extent
55 this is true for the >10,000 lncRNAs that remain uncharacterised (4-7). Indeed growing
56 evidence points to lncRNA having diverse roles outside of the cell nucleus, including
57 regulation of microRNA activity (8), protein sequestration (9), and mRNA translation (10).

58 Somewhat paradoxically, cytoplasmic lncRNA has recently been reported to
59 interact with the ribosome. In footprinting experiments to map ribosome-bound
60 transcripts genome-wide, the Weissman group identified a considerable number of
61 lncRNAs directly engaged by the translation machinery (11), an observation
62 subsequently corroborated in an independent study (12). These transcripts do not
63 contain classical features of protein-coding sequence, and various analyses have
64 argued that these lncRNAs are not productively translated in most cases (13,14). It is
65 not yet clear whether ribosomal recruitment is a general property of all lncRNA in the
66 cell. If not, it is of interest to understand what features distinguish ribosomal lncRNAs.

67 The biological significance of ribosomal lncRNA remains unclear. Two principle
68 types of potential regulatory functions for ribosomal lncRNAs have been proposed:
69 either sequence-specific regulation of mRNA translation or general regulation of
70 ribosome function (15). lncRNA and mRNA arising from opposite genomic strands can
71 form stable RNA-RNA hybrids that are localised in ribosomes (10). Through such “cis-
72 antisense” interactions, lncRNA may specifically regulate stability and translation of their
73 mRNA partner (10), although this has not yet been demonstrated at a genome-wide
74 scale. The advent of ribosome footprinting technology has prompted the idea that
75 lncRNA may non-specifically regulate translation through direct binding by ribosomes
76 (15). Cryptic open reading frame (ORF) sequences within lncRNA may be recognised

77 by the ribosome, directly resulting in translational repression or else enabling
78 recruitment of regulatory proteins. Other more mundane scenarios are also possible:
79 ribosomes might be a default destination of all polyadenylated mRNA-like transcripts,
80 where they are recognised as non-coding and processed by one of various known
81 quality surveillance pathways.

82 In the present study we take these studies further by comprehensively mapping
83 the entire known cytoplasmic and ribosomal lncRNA population of a human cell line. We
84 show that the majority of cytoplasmic lncRNAs are robustly and verifiably associated
85 with ribosomes. We show evidence that lncRNAs can be divided into classes based on
86 subcellular location and distinguished by a variety of features. These classes likely
87 serve distinct regulatory roles in translation. Finally we show that the translation
88 machinery serves as the endpoint of the lncRNA life-cycle. We conclude that, rather
89 than being an exception, ribosomal recruitment is frequently the destination of
90 cytoplasmic lncRNAs.

91

92 **Results**

93

94 **Creating a high confidence lncRNA catalogue**

95 Our aim was to map the distribution of lncRNAs in the cytoplasm and on the
96 polysomes of human cells. A potential confounding factor in any analysis of ribosome-
97 bound RNAs is the possibility of misannotated protein-coding transcripts (16). These
98 represent a non-negligible fraction of lncRNA annotation, due to the technical
99 challenges of correctly identifying protein coding sequences with high sensitivity, as well
100 as biological factors: a number of annotated lncRNAs have subsequently been found to
101 encode peptides, including small “micropeptides”, which were overlooked by
102 conventional annotations (17,18).

103 We decided to implement the most stringent possible filtering to remove protein
104 coding transcripts from our analysis, even at the expense of omitting some genuine
105 non-coding transcripts. We first removed lncRNAs that could be unannotated
106 extensions of protein-coding genes or pseudogenes. Remaining genes were filtered
107 using a panel of methods for identifying protein coding sequence (Figure 1A and
108 Materials and Methods). Altogether 9057 lncRNA transcripts (61.9%), 6763 genes
109 (73.8%) were unanimously classified as non-coding - these we refer to as “filtered
110 lncRNAs” (Figure 1A). The remaining genes of uncertain protein coding status are
111 henceforth referred to as “potential protein coding RNAs” (4415 transcripts, 1878
112 genes). The complete sets of potential protein coding and filtered lncRNAs are available
113 in Supplementary Table S1.

114

115 **Mapping the cytoplasmic and ribosomal lncRNA population**

116 We sought to create a comprehensive map of cytoplasmic lncRNA localisation in
117 a human cell. We chose as a model the K562 human myelogenous leukaemia cell line.
118 Being an ENCODE Tier I cell, it has extensive transcriptomic, proteomic and
119 epigenomic data publically available (19). We subjected cytoplasmic cellular extracts to
120 polysome profiling, an ultracentrifugation method to identify ribosome-bound RNAs and
121 distinguish transcripts bound to single or multiple ribosomes (Figure 1B) (20). Extracts
122 were divided into three pools: “Heavy Polysomal”, corresponding to high molecular

123 weight complexes cofractioning with >6 ribosomes; “Light Polysomal”, cofractioning with
124 2-6 ribosomes; and low molecular weight complexes corresponding to non-translated,
125 cytoplasmic RNAs (Figure 1C). The latter contains free mRNAs found in the high peak
126 in fraction 1, the 40 and 60S ribosomal subunits (fractions 2 and 3) and mRNAs that are
127 bound by a single ribosome (fraction 4) - we define these as “Free Cytoplasmic”
128 throughout the paper. It is important to note that although this fraction includes some
129 RNAs bound by ribosomal subunits, or individual ribosomes, the majority of these are
130 not considered to be efficiently translated (20).

131 Custom microarrays probing the entire Gencode v7 long noncoding RNA
132 catalogue were used to analyse RNAs in the free cytoplasmic, light and heavy
133 polysome fractions, in addition to total input RNA (see Materials and Methods)(5).
134 Microarrays also contained probes targeting 2796 protein-coding genes. High positive
135 correlation was observed between microarray RNA concentration measurements and
136 RNA-sequencing of the same cells from ENCODE (Supplementary Figure S1)(19).
137 Correlation between microarray results and RNAseq measurements of cytoplasmic
138 RNA was higher than with either nuclear or whole-cell RNA from the same cells,
139 attesting to the purity of these cytoplasmic extracts (Supplementary Table S2). Using
140 stringent cutoffs we detected 10.6% of filtered lncRNA transcripts (962 transcripts,
141 representing 665 or 9.8% of genes) and 52.8% of mRNAs (1476) in K562 cytoplasm
142 (Figure 1D). An additional 292 transcripts (3.2%, representing 255 or 3.7% genes) were
143 detected only in the nucleus. Altogether, 1254 filtered lncRNA transcripts (13.9%,
144 representing 875 or 13.0% of genes) were detected.

145 We classified cytoplasmic lncRNAs according to their maximal ribosomal
146 association, resulting in 347 (37.6% of cytoplasmic lncRNA transcripts) Free
147 Cytoplasmic, 373 (40.4%) Light Polysomal, and 204 (22.1%) Heavy Polysomal
148 transcripts (Figure 1D). Altogether, 62.5% of lncRNA transcripts detected in the
149 cytoplasm have maximal detection in Light or Heavy Polysomal fractions. Two lines of
150 evidence support this classification approach. First, 75% (959/1287) of protein-coding
151 mRNAs are classified as Heavy Polysomal, consistent with their being actively
152 translated and in accordance with previous studies (Figure 1E)(20,21). Second, protein
153 abundance measurements show that Heavy Polysomal mRNAs are translated most

154 efficiently (Supplementary Table S3) (22). In contrast, potential protein-coding
155 transcripts had a similar global ribosome-association profile to filtered lncRNA,
156 suggesting that they are not translated efficiently and underlining the stringency of our
157 lncRNA filtering (Figure 1E). Ribosomal lncRNA are not apparently enriched for those
158 that produce small peptides (Supplementary Table S4).

159 Cytoplasmic and ribosomal localisation has previously been reported for a number of
160 lncRNA. To test the degree of agreement between these and our data, we examined
161 the 297 lncRNA transcripts (from 60 genes) from the lncRNA Database (23) that are
162 also present in the Gencode v7 annotation. SNHG5 (5) and Gas5 (9) were detected in
163 the cytoplasm and classified as Free Cytoplasmic transcripts, consistent with previous
164 reports. The snoRNA host Gas5 has previously been reported as associated with
165 ribosomes (24). Although we classified this gene as Free Cytoplasmic based on its
166 maximal detection, 11 out of 16 transcript isoforms of Gas5 were also clearly detected
167 in Light and Heavy polysomal fractions although with lower microarray probe intensities.
168 SNHG1 is another snoRNA host reported to be bound by ribosomes (25), which we
169 classify in the Light Polysomal fraction. For other known lncRNAs, we map their
170 subcellular location for the first time: GNAS-AS1 (Nespas) and MEM161B-AS1 are
171 specifically associated with the Light Polysomal fraction.

172

173 **Independent evidence for ribosomal interaction of lncRNA**

174 We next looked for additional evidence to support ribosomal interaction of
175 lncRNA. During ultracentrifugation, it is possible that lncRNAs associated with non-
176 ribosomal, high molecular weight complexes may co-sediment with polyribosomes and
177 thus represent false positives. To investigate this, we repeated polysome profiling on
178 cells treated with puromycin (puro), a drug that disrupts ribosomes, and profiled a set of
179 candidate transcripts by volume-normalised RT-PCR (Figure 2A, B). Bona fide
180 ribosome-bound transcripts are expected to relocalise to the free cytoplasmic fraction in
181 response to puromycin. Eleven out of 16 (69%) ribosomal lncRNAs were validated in
182 this way, similar to the 4/4 protein coding mRNAs tested. In contrast, 2/3 Free
183 Cytoplasmic lncRNAs we examined showed minimal response to puro treatment. Thus

184 in the majority of cases, cosedimentation reflects a physical interaction between lncRNA
185 and ribosomes.

186 We performed additional validation using fluorescence in situ hybridisation
187 (FISH) to visualise the localisation of lncRNA at subcellular resolution. We tested three
188 Light Polysomal lncRNAs (Figure 3). ENST0000504230 displays diffuse cytoplasmic
189 localisation and exclusion from nucleoli. In addition to cytoplasmic localisation, the
190 snoRNA precursor transcript ENST00000545440 (SNHG1) shows pronounced
191 concentrations around the periphery of the nucleus, likely to be endoplasmic reticulum,
192 and at three nuclear loci – possibly its site of transcription, given that the HeLa genome
193 is predominantly triploid (26). Finally, ENST00000545462 (previously described as
194 HEIH, a prognostic factor in hepatocellular carcinoma)(27), also has pronounced
195 staining in the nuclear periphery, as well as within the nucleolus. Thus, both PCR and
196 hybridisation methods support the interpretation from microarray data of ribosomal
197 recruitment of lncRNA.

198

199 **Evidence for conserved function of ribosomal lncRNAs**

200 Purifying evolutionary selection represents powerful evidence of functionality. A
201 number of studies have shown that lncRNAs are under weak but non-neutral purifying
202 evolutionary selection (5,28,29). We sought to test if this holds true for cytoplasmic
203 lncRNAs, and in particular whether different classes of cytoplasmic lncRNA described
204 above might have experienced different strengths of selection. We extracted PhastCons
205 measures of exonic conservation and compared lncRNAs of distinct subcellular origins
206 (Figure 4). Ancestral repeats were treated as neutrally-evolving DNA for comparison. As
207 expected, protein coding exons have highly elevated conservation. Free Cytoplasmic,
208 Light Polysomal and nuclear lncRNAs exhibit similar rates of non-neutral evolution. In
209 addition, Heavy Polysomal lncRNAs contain a subset (~10%) of transcripts with
210 elevated conservation, second only to the potential protein coding transcripts, and
211 higher than other expressed lncRNAs ($P= 0.002$, $OR=2.40$ Fisher test, testing the top
212 10% of Heavy Polysomal lncRNA vs other lncRNAs pooled). Thus, cytoplasmic
213 lncRNAs experience purifying evolutionary selection consistent with conserved function.

214

215 **Ribosomal lncRNA are highly expressed and consistently localised across cell**
216 **types**

217 We next investigated the organismal and subcellular expression patterns of
218 lncRNA, in addition to their post-transcriptional processing. The steady state expression
219 levels of cytoplasmic lncRNA is similar across cytoplasmic classes in independent K562
220 whole cell RNAseq, similar to that of mRNAs and well above nuclear-specific lncRNA
221 (Figure 5A)(19). A similar trend is observed in human tissues: mean RPKM across
222 Human Body Map tissues for all three cytoplasmic classifications exceed nuclear RNA
223 ($P = 5e-5 / 2e-5 / 0.0009$ for Light Polysomal / Heavy Polysomal / Free Cytoplasmic vs
224 nuclear, Wilcox test) (Supplementary Figure S2).

225 lncRNAs have been reported to be more tissue specific than mRNAs (4,5).
226 Analysis of ubiquity, an inverse measure of tissue-specificity, of lncRNA in human
227 tissues was consistent with this (Figure 5B). Despite this similarity in expression
228 profiles, we find Heavy Polysomal lncRNAs to be significantly more ubiquitous in their
229 tissue expression profiles compared to other lncRNA classes ($P= 1.2e-4$, Fisher Test
230 Heavy Polysomal vs Nucleus), and essentially the same as mRNAs ($P=0.129$, Fisher
231 exact test).

232 Subcellular localisation of lncRNA reported by polysome profiling is consistent
233 with similar analysis using ENCODE RNAseq (19). Transcripts we report as ribosomal
234 or free cytoplasmic have significantly elevated cytoplasmic-nuclear ratios (Figure 5C)
235 ($P=0.027$ OR 2.4, $1.8e-07$ OR 4.4, 0.005 OR 2.6 for Heavy Polysomal, Light Polysomal
236 and Free Cytoplasmic vs Nuclear transcripts, Fisher exact test). Indeed, Light
237 Polysomal lncRNA have median cytoplasmic specificity that exceeds protein coding
238 mRNAs. Heavy Polysomal transcripts have a more nuclear distribution, suggesting that
239 while some transcripts are ribosomally bound, other copies are present in the nucleus.
240 We next asked whether the observed subcellular localisation of lncRNA in K562 is
241 conserved across other cell types (Figure 5D). Similar analysis on RNAseq from other
242 cell types showed Light Polysomal and Free Cytoplasmic transcripts tend to have high
243 cytoplasmic-nuclear distributions, often exceeding that of mRNAs, while Heavy
244 Polysomal has a more mixed distribution that nevertheless differs from nuclear-specific
245 transcripts. Protein-binding profiles of lncRNA yields a consistent picture, with lncRNA

246 tending to interact with proteins that localise to the same cellular compartment
247 (Supplementary Figure S3). In summary, lncRNA subcellular localisation is consistent
248 across cell types.

249

250 **mRNA-like 5' regions distinguish ribosomally-bound lncRNAs**

251 We next wished to identify factors that control the recruitment of lncRNA to
252 ribosomes. The most obvious candidate feature is the ORF, especially given that
253 lncRNAs contain abundant small ORF sequences that may be recognised by
254 ribosomes. In protein, ORF length influences the number of ribosomes that can
255 simultaneously bind, and hence the ribosomal fraction (compare mean sense ORF
256 length for heavy and light polysome mRNA in Supplementary Figure S4)(12). However
257 for lncRNA we could no evidence that ORFs determine ribosomal recruitment: neither
258 their total ORF coverage, nor their number of ORFs, nor the length of their longest ORF
259 is different from random sequence or correlates with ribosomal recruitment
260 (Supplementary Figure S5). Nor apparently does gross gene structure or GC content,
261 both clearly distinct between lncRNA and mRNA, appear to influence ribosomal
262 recruitment (Supplementary Figure S6, S7).

263 We hypothesised that factors known to influence mRNA recognition by
264 ribosomes may also apply to lncRNA. For mRNAs, a number of factors control the
265 scanning and engagement by ribosomes, including 3' polyadenylation, RNA structures
266 within the 5' UTR and 7-methylguanylate capping (30). To investigate whether
267 polyadenylation influences ribosomal recruitment, we estimated the efficiency of
268 polyadenylation of cytoplasmic and nuclear lncRNAs using ENCODE RNAseq on
269 polyA+ and polyA- nuclear RNA. Although mRNA are more polyadenylated than
270 lncRNA, we found no difference in polyadenylation efficiency between ribosomal and
271 non-ribosomal lncRNAs (Supplementary Figure S8). We recently showed that splicing
272 efficiency of lncRNAs is lower than mRNAs (31), but it does not distinguish ribosomal
273 lncRNAs from other types (Supplementary Figure S9).

274 We next looked at the role of the 5' end in ribosomal recruitment. Although
275 lncRNAs do not have identifiable ORFs and hence 5' UTRs, nevertheless they do
276 contain abundant short "pseudo-ORFs": random occurrences of in-frame start and stop

277 codons. We defined the “pseudo-5’UTR” to be the region upstream of the first AUG
278 trinucleotide of the lncRNA sequence. Although secondary structures in the 5’UTR have
279 been shown to strongly influence translation of mRNAs (32), there is no overall
280 difference in structural propensity between ribosomal and other lncRNAs
281 (Supplementary Figure S10). However, the length of pseudo-5’UTRs does distinguish
282 ribosomal from non-ribosomal lncRNA. Similar to protein-coding transcripts, Light
283 Polysomal lncRNA have significantly longer 5’UTR regions than expected by chance
284 (here estimated from the transcript’s reverse complement) (Figure 6A), while this effect
285 is essentially absent for other cytoplasmic lncRNAs. Thus long 5’UTR-like regions would
286 appear to contribute positively to ribosomal recognition of lncRNA.

287 Recognition of the 5’ methyl-guanosine cap is required for mRNA scanning by
288 the 40S ribosomal subunit. Using CAGE (cap analysis of gene expression) data (19),
289 we examined the relationship between the ribosomal recruitment of lncRNA and
290 capping using logistic regression. As shown in Figure 6B, there is a strong positive
291 relationship between capping and recruitment to the Light Polysomal Fraction. In
292 contrast, this relationship is negative for Free Cytoplasmic and Heavy Polysomal
293 recruitment. This data suggests that capping of lncRNA is a driver of ribosomal
294 recruitment, at least to the light polysomal fraction.

295

296 **Endogenous retroviral fragments are negatively correlated with ribosomal** 297 **recruitment**

298 There is growing evidence that transposable elements (TEs) contribute functional
299 sequence to lncRNA (33,34). Taking all TE classes together, we observed an excess of
300 TE-derived sequence within Free Cytoplasmic lncRNAs ($P=4e-14$, compared to
301 remaining detected filtered lncRNAs, Wilcoxon test) (Figure 7A). Potentially protein
302 coding transcripts are significantly depleted for TEs ($P=2e-16$, compared to all detected,
303 filtered lncRNAs, Wilcoxon test). Given that protein coding transcripts are strongly
304 depleted for TE insertions (35), this latter observation supports the idea that a subset of
305 potential protein coding transcripts do indeed encode functional protein.

306 We were curious whether there exist TEs whose presence correlates with the
307 subcellular localisation of their host transcript (Figure 7B) (Materials and Methods).

308 Thus we systematically tested the relationship between subcellular localisation and TE
309 class. We observed a relationship between the presence of Alu and transcript
310 expression in K562: Alu are enriched amongst detected compared to undetected filtered
311 lncRNAs ($P=6e-7$, Hypergeometric test), as recently described for human tissues (34).
312 TcMar.Tigger, although rare, show evidence for preferential enrichment in
313 polyribosomal lncRNAs ($P=9e-4$, Hypergeometric test). However the most obvious case
314 is for the class of ERVL-MaLR, which are approximately two-fold enriched in free
315 cytoplasmic lncRNAs compared to other expressed lncRNAs (Figure 7B). Closer
316 inspection revealed that this effect is not due to a single repeat type, but rather to
317 around a dozen subclasses of MST, MLT and THE endogenous retroelements (Figure
318 7C). We found no significant difference in the length of ERVL-MaLR insertions between
319 lncRNA classes (Supplementary Figure S11). Rather it is the relative proportion of
320 transcripts carrying an insertion that differs between groups. A selection of ERVL-MaLR
321 containing lncRNAs are shown in Figure 7D.

322 Enrichment of ERVL-MaLR class elements in Free Cytoplasmic lncRNAs
323 appears to be independent of cell type: using ribosome footprinting data from HeLa(36)
324 we observe that ERVL-MaLR class TEs are specifically depleted from ribosome-bound
325 lncRNAs (Figure 7E). Together these data suggest that endogenous retrovirus
326 fragments may influence lncRNA trafficking in the cell.

327

328 **Evidence for cis-antisense lncRNA-mRNA pairing in ribosomes**

329 Several reports exist describing hybridisation of lncRNA to mRNA through
330 complementary sequences, resulting in trafficking of the former to ribosomes. Antisense
331 complementarity between lncRNA and mRNA could take one of two forms: more
332 conventionally, the two transcripts may originate from opposite strands of the same
333 genomic locus, thus sharing complementary sequence regions (here “exonic antisense”,
334 also referred to as “cis-antisense”) (Figure 8A) (10). More recently, it was shown that
335 lincRNA-P21 contains regions of complementarity to mRNAs, through which they
336 hybridise and consequently localise together in the ribosome (37). Importantly, the
337 genes for lincRNA-P21 and its targets are located in distinct genomic loci – these we
338 define here as “trans-antisense” pairs.

339 We investigated whether either type of antisense may contribute to the observed
340 recruitment of lncRNA to the ribosomes. We first hypothesised that exonic antisense
341 lncRNAs would be more frequently localised in heavy polysomes, due to hybridisation to
342 their corresponding (actively translated) mRNA. We classified all lncRNA by their
343 genomic organisation with respect to protein-coding genes (5): intergenic (not
344 overlapping), exonic antisense, intronic antisense, or intronic same sense. Consistent
345 with our hypothesis, lncRNAs identified in heavy polysomes are significantly enriched
346 for exonic antisense transcripts compared to those in other cellular compartments
347 ($P=4.2e-5$, Fisher exact test) (Figure 8B). This finding is consistent with lncRNA / mRNA
348 hybrids existing in human ribosomes. An example of such a cis-antisense pair is shown
349 in Figure 8C. If this is the case, we would expect mRNAs bound by antisense heavy
350 polysomal lncRNA to be more highly expressed than others. Examining RNAseq
351 expression data we find this to be the case: mRNAs antisense to heavy polysomal
352 lncRNA are significantly more highly expressed than mRNAs antisense to other lncRNA
353 classes ($P=7e-4$, Wilcoxon test)(Figure 8D). In the course of this analysis, we also made
354 the incidental observation that intronic same sense lncRNAs tend to be nuclear-specific
355 (Supplementary Figure S12).

356 Trans-antisense hybridisation is another potential means by which lncRNA could
357 interact with mRNA and be recruited to ribosomes. In this model, the transcripts share
358 homology on opposite strands, but are not transcribed from the same genomic locus
359 (Figure 8A). Using a BLAST approach, we compiled all sense-antisense homology
360 relationships between intergenic lncRNA and mRNA (Figure 8E). As a control, we
361 performed the same operation with size-matched, randomised genomic regions instead
362 of lncRNA. This analysis resulted in two observations: first, lncRNA as a whole are more
363 likely to have trans-antisense homology to mRNA compared to random genomic
364 sequence ($P=1e-14$, Fisher test, comparing all lncRNA to all shuffled); second, this
365 tendency was observed with statistical significance in ribosomal lncRNA ($P=0.002$,
366 Fisher test for heavy and light lncRNA combined) but not in free cytoplasmic and
367 nuclear lncRNA. These findings were consistent across a range of different BLAST
368 settings. This data point to possible trans-antisense lncRNA-mRNA hybridisation as a
369 general regulatory mechanism of ribosomal lncRNA.

370

371 **Degradation of lncRNA by the ribosome**

372 We were next curious whether recruitment to ribosomes had any effect on
373 lncRNA. It was proposed by Chew et al(38) that lncRNA at ribosomes are subject to
374 nonsense mediated decay (NMD), and indeed one report does exist of ribosome-
375 dependent degradation of a snoRNA host gene (35). We tested whether blocking
376 translation has any outcome on the stability of ribosomal lncRNA identified here. Using
377 the same candidate genes tested previously, we tested whether interfering with
378 ribosomal function through drug-induced stalling (cyclohexamide) influenced lncRNA
379 stability (Figure 9). In a number of cases we observed elongation-dependent
380 degradation of lncRNA, often decreasing lncRNA amounts by several fold over 6 hours.
381 This effect was highly heterogeneous, with other transcripts unaffected by
382 cyclohexamide treatment. Thus, ribosomal recruitment leads to degradation of many
383 lncRNAs.

384

385 **Discussion**

386 In order to gain clues as to lncRNA function at a global level, we have
387 comprehensively mapped the ribosomal and cytoplasmic lncRNA populations of a
388 human cell. The very substantial populations of lncRNA we discover in these fractions is
389 at odds with existing paradigms of lncRNA function as principally nuclear molecules. We
390 must now consider the possibility that lncRNA play more diverse roles outside the
391 nucleus, including translational control, cellular metabolism or signal transduction.

392 One key challenge in the study of cytoplasmic lncRNAs is to rule out the
393 possibility that they encode a cryptic, unannotated protein product. This question has
394 been discussed in excellent reviews elsewhere (15), and has not yet been satisfactorily
395 resolved. Indeed, it is likely that an extensive “grey zone” of transcripts with weak
396 protein coding potential exists (and indeed may form the substrate for novel protein
397 evolution (39)). It is also plausible that some or many transcripts do exist that function
398 both as protein-coding and noncoding transcripts, although apart from the archetypal
399 SRA1 (40), few concrete examples have so far been presented (41,42). In this study we
400 took great pains to filter any transcripts with even minimal probability of encoding
401 protein, in the process collecting many weakly coding (“potential protein coding”)
402 transcripts that may be of rich scientific interest in future. We describe a set of 1867
403 annotated lncRNAs that have varying degrees of evidence for encoding protein. These
404 transcripts have intriguing characteristics intermediate between coding and non-coding
405 RNA: they are under higher evolutionary selection than lncRNA, are depleted for
406 transposable elements, and tend to be cytoplasmically enriched – similar to mRNAs. In
407 contrast, they have ribosomal association profiles and expression ubiquity similar to
408 lncRNA. Finally, their gene structures and expression levels are intermediate between
409 coding and noncoding sequences. It will be fascinating in future to study whether these
410 transcripts represent an intermediate timepoint in the evolution of either new proteins
411 from non-coding sequences, or the evolution of non-coding RNAs from formerly coding
412 transcripts.

413 We cannot rule out the possibility that some of our filtered lncRNA produce a
414 peptide product. Recent studies have revealed the potential for large volumes of
415 unrecognised protein coding capacity in mammalian genomes, either as small peptides

416 (43) or non-canonical ORF translation (44,45). However, even if these lncRNAs give
417 rise to peptides, it does not necessarily follow that all are functional: occasional
418 nonsense translation of ribosomally-localised lncRNA may occur with no functional
419 consequences - “translational noise”. Future intensive mass spectrometry studies of
420 short peptides, such as that carried out by Slavoff et al (18), will hopefully allow us to
421 further improve lncRNA annotations.

422 We present several lines of evidence that ribosomal lncRNAs are a large
423 functional gene class that genuinely interacts with the translation machinery: (1)
424 ribosomal lncRNAs are puromycin sensitive; (2) fluorescence in situ hybridisation
425 indicates their cytoplasmic localisation; (3) they have elevated cytoplasmic-nuclear
426 ratios by independent ENCODE RNAseq data across diverse cell types; (4) their
427 sequence is under similar or even elevated evolutionary selection compared to nuclear
428 lncRNAs. Thus these transcripts are functional and appear to have a regulated and
429 consistent subcellular localisation.

430 Polysome profiling appears to distinguish lncRNAs with distinct properties. We
431 have attempted to rather crudely classify transcripts according to their fraction of
432 maximum detection, but most transcripts are detected at varying concentrations in all
433 fractions. Nevertheless, through this classification we have managed to discover
434 features that distinguish lncRNAs and laid a foundation for predicting lncRNA localisation
435 *de novo*. Similarly, a recent study discovered an RNA motif that predicts and appears to
436 confer cytoplasmic localisation (46). We find that lncRNAs localised in the Light
437 Polysomal fraction tend to have mRNA-like 5' features, more specifically a non-
438 randomly long pseudo-5' UTR length and the presence of a cap structure. This is
439 consistent with the importance of 5' recognition in the initiation of translation. Other
440 mRNA-like features such as polyadenylation, GC content or open reading frames do not
441 appear to affect ribosomal interaction at a global level. In contrast, repetitive sequence
442 features, and particularly human endogenous retrovirus fragments, are negatively
443 associated with ribosomal recruitment. This is perhaps to be expected, given that
444 mRNAs tend to be depleted of such repeats, in contrast to lncRNAs (35). The
445 mechanism by which hERV prevent lncRNAs from ribosomal recruitment remains to be
446 ascertained, although we proposed recently that such fragments may interact with

447 protein complexes that could antagonise ribosomal binding (33). In summary, these
448 findings represent a starting point for discovering features that distinguish lncRNA
449 classes and may eventually lead to useful models for predicting such classes.

450 Light Polysomal and Heavy Polysomal lncRNA appear to represent functionally
451 distinct classes of lncRNA. It is generally considered that heavy polysomes tend to be
452 actively translating, while light polysomes represent more weakly translated
453 messengers, and this is supported by our mRNA data. We interpret the lncRNAs in the
454 Light Polysomal fraction to be engaged by two or a few ribosomes that are not in the
455 process of translating an mRNA. The above features of 5' processing distinguish Light
456 Polysomal transcripts clearly from Free Cytoplasmic transcripts, but not Heavy
457 Polysomal. The latter tend to be more nuclear and more strongly evolutionarily
458 conserved. Some clues to the origin of these differences may be gleaned from the
459 observation that cis-antisense transcripts are enriched in the heavy fraction. Cis-
460 antisense transcripts have been studied for a number of years, and cases have been
461 described where the antisense lncRNA hybridises with its sense mRNA and
462 accompanies it to the ribosome (10). Thus we might posit that lncRNA in heavy
463 polysomes are involved in active translational processes and include transcripts that
464 exist as hybrids with their sense mRNA partner. Such recognition is sequence specific,
465 and we may guess that this localisation occurs indirectly: the lncRNA is recruited
466 through its binding to a translated mRNA, and not directly engaged by ribosomes. In
467 contrast, given their mRNA-like 5' features, we propose that Light Polysomal lncRNA
468 include cases that are directly engaged by ribosomes, resulting in non-specific
469 translational repression and/or lncRNA degradation(15). This model is outlined in Figure
470 10.

471 Although it is tempting to propose that ribosomal lncRNA regulate protein
472 translation, we must also seriously consider an alternative possibility: that the ribosome
473 represents the default endpoint of the lncRNA lifecycle, and it is rather the non-
474 ribosomal cytoplasmic transcripts that are exceptional. Indeed, it is perhaps not
475 surprising that these mRNA-like transcripts - capped, polyadenylated and 100-10,000 nt
476 long – should be recognised by the cell and trafficked accordingly. We here show
477 evidence that, at least for a subset of transcripts, the result of ribosomal recruitment is

478 degradation. That is, the translation machinery is also responsible for lncRNA
479 clearance, and that the regulatory relationship between lncRNA and the translational
480 machinery is reciprocal.
481

482 **Materials and Methods**

483

484 **Polysome fractionation**

485 For polysome fractionations, 20 million K562 cells were incubated with 100 ug/mL of
486 cycloheximide (Sigma, Cat C4859) for 10 min. Cell pellets were resuspended in 200ul
487 RSB buffer (20 mM Tris-HCl, pH 7.4, 20 mM NaCl, 30 mM MgCl₂, 200ug/mL
488 cycloheximide, 0.2mg/mL heparin (Sigma, Cat No. H4787), 1000 unit/mL RNasin), then
489 lysed with an equal volume of Lysis Buffer (1X RSB, 1 % Triton X-100, 2% Tween-20,
490 200ug/ul heparin) with or without 1% Na deoxycholate. Following incubation on ice for
491 10 min, extracts were centrifuged at 13,000 x g for 3 min to remove the nuclei.
492 Supernatants were further centrifuged at 13,000 x g for 8 min at 4°C. Equal OD units
493 were loaded onto 10% to 50% linear sucrose gradients (prepared in 10 mM Tris-HCl pH
494 7.4, 75 mM KCl and 1.5mM MgCl₂), and centrifuged at 36,000 rpm for 90 min at 8° C in
495 a SW41 rotor (Beckman Coulter). Twelve fractions were collected from the top of the
496 gradient using a piston gradient fractionator (BioComp Instruments). A UV-M II monitor
497 (BIORAD) was used to measure the absorbance at 254 nm. 110ul of 10% SDS and 12
498 uL of proteinase K (10 mg/mL Invitrogen) was added to each 1ml fraction and incubated
499 for 30 min at 42°C. Fractions 1-5, 6-8 and 9-11 were pooled corresponding to groups
500 Free Cytoplasmic (Free / Monosomal), LP (Light Polysome) and HP (Heavy Polysome),
501 respectively. For puromycin-treated samples, cells were incubated in 100ug/ml
502 puromycin for 15 minutes prior to processing and puromycin was used in place of
503 cyclohexamide in all the buffers.

504 Unfractionated cytoplasmic RNA and pooled polysomal RNAs were purified using
505 phenol chloroform isoamyl extraction followed by LiCl precipitation to remove the
506 heparin. The integrity of the samples was monitored by Bioanalyzer. For qRT-PCR
507 analysis equal volumes of RNA were used to synthesise cDNA using the Superscript III
508 Reverse Transcriptase (Invitrogen) according to manufacturer's instructions. Two
509 bacterial spike-in RNAs, Dap and Thr were added before RNA purification to equal
510 volumes of each polysomal RNA pool. Gene specific primers were used with SYBR
511 Green for qRT-PCR on an ABI PRISM 7900 Sequence Detection Systems. Candidate
512 CT values were normalized to the spike in controls Dap and Thr that were present at

513 equal concentrations per pool. Relative RNA levels are presented as a percentage of
514 the RNA present in each pool with 100% RNA calculated as the sum of the FM, LP and
515 HP pools.

516

517 **Microarray Design**

518 This study was carried out using Agilent custom gene expression microarrays, in the
519 8x60k format with 60mer probes. Probes were designed using eArray software with
520 standard settings: Base composition methodology / 60mer / 4 probes per target / sense
521 probes / best probe methodology / 3' bias. Probes were designed for 14700 transcripts
522 from the entire Gencode v7 lncRNA catalogue, in addition to 26 known lncRNAs from
523 www.lncrnadb.org (23) and 90 randomly-selected protein-coding housekeeping genes.
524 The array was then filled with probes targeting 2796 randomly-selected protein-coding
525 gene probes. Microarray design details are available from the Gencode website
526 (http://www.gencodegenes.org/lncrna_microarray.html).

527

528 **Microarray Hybridization and Probe Quantification**

529 100 ng of total RNA was labeled using Low Input Quick Amp Labeling kit (Agilent 5190-
530 2305) following manufacturer instructions. mRNA was reverse transcribed in the
531 presence of T7-oligo-dT primer to produce cDNA. cDNA was then in vitro transcribed
532 with T7 RNA polymerase in the presence of Cy3-CTP to produce labeled cRNA. The
533 labeled cRNA was hybridized to the Agilent SurePrint G3 gene expression 8x60K
534 microarray according to the manufacturer's protocol. The arrays were washed, and
535 scanned on an Agilent G2565CA microarray scanner at 100% PMT and 3um resolution.
536 Intensity data was extracted using the Feature Extraction software (Agilent).

537 Raw data was taken from the Feature Extraction output files and was corrected
538 for background noise using the normexp method(47). To assure comparability across
539 samples we used quantile normalization(48). All statistical analyses were performed
540 with the Bioconductor project (<http://www.bioconductor.org/>) in the R statistical
541 environment (<http://cran.r-project.org/>) (49).

542

543 **Preparation of filtered lncRNA gene catalogues**

544 We first filtered the former set to remove any transcripts that potentially result from
545 misannotated extensions or isoforms of protein-coding genes or pseudogenes. Any
546 gene was discarded that has at least one transcript fulfilling one of the following
547 conditions: overlapping on the same strand a Gencode v18 annotated pseudogene,
548 overlapping on the same strand an exon of a protein-coding mRNA, or lying within 5 kb
549 and on the same strand as an Gencode v18 protein-coding transcript or pseudogene
550 (1169 transcripts, 521 genes). This resulted in a dataset of 13,472 lncRNA transcripts
551 (8641 genes). Next, genes having at least one transcript predicted as protein coding by
552 at least one method, were classified as “potential protein coding RNAs” (4415
553 transcripts, 1878 genes), while the remainder were classified as “filtered lncRNAs”. The
554 four filtering methods used were: 1) PhyloCSF, a comparative genomics method based
555 on phylogenetic conservation across species (50). The analysis was performed using
556 29 mammalian nucleotide sequence alignments and assessing the three sense frames.
557 The alignment of each transcript was extracted from stitch gene blocks given a set of
558 exons from Galaxy(51). Transcripts with score >95 were classified as potential protein
559 coding, following the work of Sun et al (52). 2) Coding Potential Assessment Tool
560 (CPAT)(53), using the score threshold of 0.364 described by the authors. 3) Coding
561 Potential Calculator (CPC), a support vector machine-based classifier based on six
562 biological sequence features, using a cutoff of 1 (54). 4) Peptides: We used
563 experimental mass spectrometry tag mappings from Pinstripe to identify any transcripts
564 producing peptides (55). Any transcript having an exonic, same strand tag mapping
565 were designated as “potential protein coding”. Collectively, sequence filters reduced the
566 pool of analyzed transcripts to 9057 transcripts (6763 genes). The full table of
567 classification data for all Gencode v7 lncRNA is available in Supplementary Table S1.

568

569 **Microarray probe filtering**

570 lncRNA transcripts were considered to be present in a sample when at least three out
571 of four microarray probes were reliable and not absent. The expression intensity value
572 for “present” transcripts was computed as the median of its present probes. Protein
573 coding genes were considered “present” if at least one probe was reliable and not
574 absent, and the intensity value was that of one of the present probes, chosen randomly.

575 Variances in probe intensity values within probesets were significantly different when
576 comparing all probesets from present transcripts in a sample (Levene's test). To avoid
577 non representative intensity values, 5% of transcripts (for each sample) with highest
578 probeset variance were removed from our dataset. Applying these filters we define as
579 cytoplasmically detected 962 filtered lncRNAs (665 genes), 906 potential protein-coding
580 transcripts (382 genes) and 1476 protein-coding genes that are detected in K562
581 cytoplasm.

582

583 **RNAseq correlation analysis**

584 ENCODE RNA-sequencing quantifications (Gencode v10 annotation) from cytoplasmic
585 fraction of K562 cells was used to check correlation with microarray data. Correlation
586 was calculated only with transcripts present in both ENCODE data (considered as
587 present when RPKM bio-replicates mean different to 0 and IDR < 0.1) and microarray
588 data.

589

590 **Classification of array transcripts**

591 From the polysome profiling analysis, detected lncRNAs and mRNAs were classified
592 according to the microarray sample (condition) where they displayed the highest
593 transcript-level signal. Thus, present transcripts were classified into Heavy Polysomal
594 (Heavy P.), Light Polysomal (Light P.) and Free Cytoplasmic transcripts (Free C.)
595 transcripts. The remaining protein coding genes, which were not present in any
596 microarray condition were considered not present. Remaining filtered lncRNA
597 transcripts were subsequently checked in ENCODE K562 nucleus RNAseq. Those
598 detected (defined as RPKM bio-replicates mean > 0 and IDR < 0.1) were classified as
599 nuclear specific transcripts. Remaining transcripts, which are not present in cytoplasm
600 nor in the nucleus are considered not present (NotPre).

601

602 **Cytoplasmic-nuclear localisation using RNAseq data**

603 Cytoplasmic and nuclear RNAseq data from six different cell lines (K562, HeLa, NHEK,
604 HepG2, GM12878, HUVEC) were obtained from ENCODE (19). For each cell line we
605 calculated cytoplasmic-nuclear RPKM ratios for transcripts detected in both that cell line

606 and K562. RPKM was calculated as the mean of two available technical replicates, and
607 only transcripts with mean > 0 RPKMs and IDR < 0.1 were considered present. We
608 calculated log₂ ratios of cytoplasmic expression versus nuclear expression (RPKM
609 units) for those transcripts present in both nucleus and cytoplasm.

610

611 **Tissue Expression Analysis**

612 We extracted tissue expression values for 16 human tissues from Human Body Map
613 (HBM) RNAseq data, downloaded from ArrayExpress under accession number E-
614 MTAB-513. These data were used to quantify Gencode v10 transcripts using the
615 GRAPE pipeline(56). Transcripts were defined as ubiquitous if they had >0.1 RPKM
616 expression in all 16 tissues.

617

618 **Transposable element analysis**

619 The 2013 version of RepeatMasker human genomic repetitive element annotations
620 were downloaded from UCSC Genome Browser, and was converted to BED format.
621 Using the tool IntersectBED, we calculated (1) the number of instances of intersection,
622 and (2) the number of nucleotides of overlap, between each lncRNA transcript and each
623 transposable element. This analysis was carried out for both transposable element
624 types, and transposable element classes.

625

626 **ORF analysis**

627 We mapped all possible canonical open reading frames (ORFs) in each of six frames in
628 lncRNA and protein coding transcripts from Gencode. If more than one start codon is in
629 frame with a stop codon, only the start codon for the longest ORF is considered.

630

631 **CAGE analysis of lncRNA capping**

632 5' cap analysis was performed on cap analysis gene expression (CAGE) tags from
633 ENCODE (19) for K562 cytoplasmic poly+ RNA and mapped these tags to the
634 microarray region comprising between 100nt before and after transcription start sites of
635 lncRNA. In order to assess the relationship between cytoplasmic class and capping, we
636 compared CAGE tag presence to fractional occupancy in each class. The latter was

637 calculated by subtracting input cytoplasmic log2 microarray expression intensity values
638 from each of the three polysome profiling fractions intensity values (Free C., Light P. or
639 Heavy P.). We divided transcripts into (log2) fraction occupancy bins from -2 to 2 at 0.5
640 bins. Transcripts with values outside this range were pooled into the last corresponding
641 bin. Logistic regression was performed to assess the relationship between CAGE tag
642 presence and occupancy.

643

644 **BLAST analysis of trans-homology between lncRNA and mRNA**

645 Gencode v7 transcript-level FASTA files of mRNA (Gene type “protein coding”) and
646 lncRNA were downloaded from Gencode. Two control sets analogous to lncRNA were
647 also collected and processed in exactly the same way: first, Bedtools “shuffle” tool was
648 used to extract random regions identical in size to the lncRNAs. Second, the introns of
649 each lncRNA were concatenated, then a fragment of sequence identical in size to the
650 mature lncRNA was extracted at a random location within this sequence. All sequences
651 were repeat-masked using RepeatMasker with “sensitive” and “human” settings. A
652 BLAST library was created using default settings with the mRNA sequences. lncRNA
653 and control sequences were BLASTed against this library with maximum expectation
654 value of 20.

655

656 **RNA Stability Assay**

657 K562 cells were incubated with or without cyclohexamide (100ug/ml) for three hours
658 followed by treatment with actinomycin D (5ug/ml) for 6 hours. RNA samples were taken
659 at 0 and 6 hours following actinomycin D treatment to assess the stability of the RNA in
660 the absence of transcription. RNA was purified using Trizol and Qiagen RNeasy
661 columns. 1ug of RNA was used to make cDNA using RevertAid H Minus reverse
662 transcriptase. Luminaris Color HiGreen High ROX qPCR master mix was used with
663 gene specific primers for qRT-PCR on an ABI PRISM 7900 Sequence Detection
664 Systems. Expression levels were normalised to the housekeeping gene GAPDH by the
665 delta-delta Ct method.

666

667 **Acknowledgements**

668

669 We thank members of the Guigo lab and the CRG Bioinformatics and Genomics
670 Programme for many ideas and discussions, particularly Marta Melé, Joao Curado, and
671 Marc Friedlaender, in addition to Fatima Gebauer (CRG, Gene Regulation Stem Cells
672 and Cancer Programme). We would particularly like to thank Ferran Reverter for
673 invaluable statistical advice. Thomas Derrien (University of Rennes) and Giovanni
674 Bussotti (EBI) helped with evolutionary conservation analysis. We thank the CRG
675 Genomics Core Facility, particularly Anna Ferrer, Maria Aguilar, Sarah Bonnin and
676 Manuela Hummel, for array hybridisation and analysis. The following CRG colleagues
677 generously helped with FISH experiments: Francois Le Dily, Maria Sanz, Carme Arnan,
678 Maria Teresa Zomeño, as well as Timmo Zimmermann and Raquel García from CRG
679 Microscopy Facility.

680 **References**

681

- 682 1. Brown, C.J., Ballabio, A., Rupert, J.L., Lafreniere, R.G., Grompe, M., Tonlorenzi, R. and
683 Willard, H.F. (1991) A gene from the region of the human X inactivation centre is
684 expressed exclusively from the inactive X chromosome. *Nature*, **349**, 38-44.
- 685 2. Rinn, J.L., Kertesz, M., Wang, J.K., Squazzo, S.L., Xu, X., Bruggmann, S.A., Goodnough,
686 L.H., Helms, J.A., Farnham, P.J., Segal, E. *et al.* (2007) Functional demarcation of active
687 and silent chromatin domains in human HOX loci by noncoding RNAs. *Cell*, **129**, 1311-
688 1323.
- 689 3. Khalil, A.M., Guttman, M., Huarte, M., Garber, M., Raj, A., Rivea Morales, D., Thomas,
690 K., Presser, A., Bernstein, B.E., van Oudenaarden, A. *et al.* (2009) Many human large
691 intergenic noncoding RNAs associate with chromatin-modifying complexes and affect
692 gene expression. *Proceedings of the National Academy of Sciences of the United States*
693 *of America*, **106**, 11667-11672.
- 694 4. Cabili, M.N., Trapnell, C., Goff, L., Koziol, M., Tazon-Vega, B., Regev, A. and Rinn, J.L.
695 (2011) Integrative annotation of human large intergenic noncoding RNAs reveals global
696 properties and specific subclasses. *Genes & development*, **25**, 1915-1927.
- 697 5. Derrien, T., Johnson, R., Bussotti, G., Tanzer, A., Djebali, S., Tilgner, H., Guernec, G.,
698 Martin, D., Merkel, A., Knowles, D.G. *et al.* (2012) The GENCODE v7 catalog of human
699 long noncoding RNAs: analysis of their gene structure, evolution, and expression.
700 *Genome research*, **22**, 1775-1789.
- 701 6. Hangauer, M.J., Vaughn, I.W. and McManus, M.T. (2013) Pervasive Transcription of the
702 Human Genome Produces Thousands of Previously Unidentified Long Intergenic
703 Noncoding RNAs. *PLoS genetics*, **9**, e1003569.
- 704 7. Managadze, D., Lobkovsky, A.E., Wolf, Y.I., Shabalina, S.A., Rogozin, I.B. and Koonin,
705 E.V. (2013) The vast, conserved mammalian lincRNome. *PLoS computational biology*,
706 **9**, e1002917.
- 707 8. Cesana, M., Cacchiarelli, D., Legnini, I., Santini, T., Sthandier, O., Chinappi, M.,
708 Tramontano, A. and Bozzoni, I. (2011) A long noncoding RNA controls muscle
709 differentiation by functioning as a competing endogenous RNA. *Cell*, **147**, 358-369.
- 710 9. Kino, T., Hurt, D.E., Ichijo, T., Nader, N. and Chrousos, G.P. (2010) Noncoding RNA
711 gas5 is a growth arrest- and starvation-associated repressor of the glucocorticoid
712 receptor. *Science signaling*, **3**, ra8.
- 713 10. Carrieri, C., Cimatti, L., Biagioli, M., Beugnet, A., Zucchelli, S., Fedele, S., Pesce, E.,
714 Ferrer, I., Collavin, L., Santoro, C. *et al.* (2012) Long non-coding antisense RNA controls
715 Uchl1 translation through an embedded SINEB2 repeat. *Nature*, **491**, 454-457.
- 716 11. Ingolia, N.T., Lareau, L.F. and Weissman, J.S. (2011) Ribosome profiling of mouse
717 embryonic stem cells reveals the complexity and dynamics of mammalian proteomes.
718 *Cell*, **147**, 789-802.
- 719 12. van Heesch, S., van Iterson, M., Jacobi, J., Boymans, S., Essers, P.B., de Bruijn, E.,
720 Hao, W., Macinnes, A.W., Cuppen, E. and Simonis, M. (2014) Extensive localization of
721 long noncoding RNAs to the cytosol and mono- and polyribosomal complexes. *Genome*
722 *biology*, **15**, R6.
- 723 13. Banfai, B., Jia, H., Khatun, J., Wood, E., Risk, B., Gundling, W.E., Jr., Kundaje, A.,
724 Gunawardena, H.P., Yu, Y., Xie, L. *et al.* (2012) Long noncoding RNAs are rarely
725 translated in two human cell lines. *Genome research*, **22**, 1646-1657.
- 726 14. Guttman, M., Russell, P., Ingolia, N.T., Weissman, J.S. and Lander, E.S. (2013)
727 Ribosome profiling provides evidence that large noncoding RNAs do not encode
728 proteins. *Cell*, **154**, 240-251.

- 729 15. Ulitsky, I. and Bartel, D.P. (2013) lincRNAs: genomics, evolution, and mechanisms. *Cell*,
730 **154**, 26-46.
- 731 16. Dinger, M.E., Pang, K.C., Mercer, T.R. and Mattick, J.S. (2008) Differentiating protein-
732 coding and noncoding RNA: challenges and ambiguities. *PLoS computational biology*, **4**,
733 e1000176.
- 734 17. Kondo, T., Plaza, S., Zanet, J., Benrabah, E., Valenti, P., Hashimoto, Y., Kobayashi, S.,
735 Payre, F. and Kageyama, Y. (2010) Small peptides switch the transcriptional activity of
736 Shavenbaby during *Drosophila* embryogenesis. *Science*, **329**, 336-339.
- 737 18. Slavoff, S.A., Mitchell, A.J., Schwaid, A.G., Cabili, M.N., Ma, J., Levin, J.Z., Karger, A.D.,
738 Budnik, B.A., Rinn, J.L. and Saghatelian, A. (2013) Peptidomic discovery of short open
739 reading frame-encoded peptides in human cells. *Nature chemical biology*, **9**, 59-64.
- 740 19. Djebali, S., Davis, C.A., Merkel, A., Dobin, A., Lassmann, T., Mortazavi, A., Tanzer, A.,
741 Lagarde, J., Lin, W., Schlesinger, F. *et al.* (2012) Landscape of transcription in human
742 cells. *Nature*, **489**, 101-108.
- 743 20. Zhang, D., Zhao, T., Ang, H.S., Chong, P., Saiki, R., Igarashi, K., Yang, H. and Vardy,
744 L.A. (2012) AMD1 is essential for ESC self-renewal and is translationally down-regulated
745 on differentiation to neural precursor cells. *Genes & development*, **26**, 461-473.
- 746 21. Beilharz, T.H. and Preiss, T. (2004) Translational profiling: the genome-wide measure of
747 the nascent proteome. *Briefings in functional genomics & proteomics*, **3**, 103-111.
- 748 22. Khatun, J., Yu, Y., Wrobel, J.A., Risk, B.A., Gunawardena, H.P., Secret, A., Spitzer,
749 W.J., Xie, L., Wang, L., Chen, X. *et al.* (2013) Whole human genome proteogenomic
750 mapping for ENCODE cell line data: identifying protein-coding regions. *BMC genomics*,
751 **14**, 141.
- 752 23. Amaral, P.P., Clark, M.B., Gascoigne, D.K., Dinger, M.E. and Mattick, J.S. (2011)
753 lncRNADB: a reference database for long noncoding RNAs. *Nucleic acids research*, **39**,
754 D146-151.
- 755 24. Smith, C.M. and Steitz, J.A. (1998) Classification of gas5 as a multi-small-nucleolar-RNA
756 (snoRNA) host gene and a member of the 5'-terminal oligopyrimidine gene family
757 reveals common features of snoRNA host genes. *Molecular and cellular biology*, **18**,
758 6897-6909.
- 759 25. Pelczar, P. and Filipowicz, W. (1998) The host gene for intronic U17 small nucleolar
760 RNAs in mammals has no protein-coding potential and is a member of the 5'-terminal
761 oligopyrimidine gene family. *Molecular and cellular biology*, **18**, 4509-4518.
- 762 26. Adey, A., Burton, J.N., Kitzman, J.O., Hiatt, J.B., Lewis, A.P., Martin, B.K., Qiu, R., Lee,
763 C. and Shendure, J. (2013) The haplotype-resolved genome and epigenome of the
764 aneuploid HeLa cancer cell line. *Nature*, **500**, 207-211.
- 765 27. Yang, F., Zhang, L., Huo, X.S., Yuan, J.H., Xu, D., Yuan, S.X., Zhu, N., Zhou, W.P.,
766 Yang, G.S., Wang, Y.Z. *et al.* (2011) Long noncoding RNA high expression in
767 hepatocellular carcinoma facilitates tumor growth through enhancer of zeste homolog 2
768 in humans. *Hepatology*, **54**, 1679-1689.
- 769 28. Guttman, M., Amit, I., Garber, M., French, C., Lin, M.F., Feldser, D., Huarte, M., Zuk, O.,
770 Carey, B.W., Cassady, J.P. *et al.* (2009) Chromatin signature reveals over a thousand
771 highly conserved large non-coding RNAs in mammals. *Nature*, **458**, 223-227.
- 772 29. Ponjavic, J., Ponting, C.P. and Lunter, G. (2007) Functionality or transcriptional noise?
773 Evidence for selection within long noncoding RNAs. *Genome research*, **17**, 556-565.
- 774 30. Jackson, R.J., Hellen, C.U. and Pestova, T.V. (2010) The mechanism of eukaryotic
775 translation initiation and principles of its regulation. *Nature reviews. Molecular cell*
776 *biology*, **11**, 113-127.
- 777 31. Tilgner, H., Knowles, D.G., Johnson, R., Davis, C.A., Chakraborty, S., Djebali, S.,
778 Curado, J., Snyder, M., Gingeras, T.R. and Guigo, R. (2012) Deep sequencing of

- 779 subcellular RNA fractions shows splicing to be predominantly co-transcriptional in the
780 human genome but inefficient for lncRNAs. *Genome research*, **22**, 1616-1625.
- 781 32. Kudla, G., Murray, A.W., Tollervey, D. and Plotkin, J.B. (2009) Coding-sequence
782 determinants of gene expression in *Escherichia coli*. *Science*, **324**, 255-258.
- 783 33. Johnson, R. and Guigo, R. (2014) The RIDL hypothesis: transposable elements as
784 functional domains of long noncoding RNAs. *RNA*, **20**, 959-976.
- 785 34. Kelley, D. and Rinn, J. (2012) Transposable elements reveal a stem cell-specific class of
786 long noncoding RNAs. *Genome biology*, **13**, R107.
- 787 35. Kapusta, A., Kronenberg, Z., Lynch, V.J., Zhuo, X., Ramsay, L., Bourque, G., Yandell,
788 M. and Feschotte, C. (2013) Transposable elements are major contributors to the origin,
789 diversification, and regulation of vertebrate long noncoding RNAs. *PLoS genetics*, **9**,
790 e1003470.
- 791 36. Ingolia, N.T., Brar, G.A., Rouskin, S., McGeachy, A.M. and Weissman, J.S. (2012) The
792 ribosome profiling strategy for monitoring translation in vivo by deep sequencing of
793 ribosome-protected mRNA fragments. *Nature protocols*, **7**, 1534-1550.
- 794 37. Yoon, J.H., Abdelmohsen, K., Srikantan, S., Yang, X., Martindale, J.L., De, S., Huarte,
795 M., Zhan, M., Becker, K.G. and Gorospe, M. (2012) lincRNA-p21 suppresses target
796 mRNA translation. *Molecular cell*, **47**, 648-655.
- 797 38. Chew, G.L., Pauli, A., Rinn, J.L., Regev, A., Schier, A.F. and Valen, E. (2013) Ribosome
798 profiling reveals resemblance between long non-coding RNAs and 5' leaders of coding
799 RNAs. *Development*, **140**, 2828-2834.
- 800 39. Knowles, D.G. and McLysaght, A. (2009) Recent de novo origin of human protein-coding
801 genes. *Genome research*, **19**, 1752-1759.
- 802 40. Lanz, R.B., McKenna, N.J., Onate, S.A., Albrecht, U., Wong, J., Tsai, S.Y., Tsai, M.J.
803 and O'Malley, B.W. (1999) A steroid receptor coactivator, SRA, functions as an RNA and
804 is present in an SRC-1 complex. *Cell*, **97**, 17-27.
- 805 41. Ulveling, D., Francastel, C. and Hube, F. (2011) Identification of potentially new
806 bifunctional RNA based on genome-wide data-mining of alternative splicing events.
807 *Biochimie*, **93**, 2024-2027.
- 808 42. Marques, A.C., Tan, J., Lee, S., Kong, L., Heger, A. and Ponting, C.P. (2012) Evidence
809 for conserved post-transcriptional roles of unitary pseudogenes and for frequent
810 bifunctionality of mRNAs. *Genome biology*, **13**, R102.
- 811 43. Ma, J., Ward, C.C., Jungreis, I., Slavoff, S.A., Schwaid, A.G., Neveu, J., Budnik, B.A.,
812 Kellis, M. and Saghatelian, A. (2014) Discovery of human sORF-encoded polypeptides
813 (SEPs) in cell lines and tissue. *Journal of proteome research*, **13**, 1757-1765.
- 814 44. Zu, T., Gibbens, B., Doty, N.S., Gomes-Pereira, M., Huguet, A., Stone, M.D., Margolis,
815 J., Peterson, M., Markowski, T.W., Ingram, M.A. *et al.* (2011) Non-ATG-initiated
816 translation directed by microsatellite expansions. *Proceedings of the National Academy
817 of Sciences of the United States of America*, **108**, 260-265.
- 818 45. Vanderperre, B., Lucier, J.F., Bissonnette, C., Motard, J., Tremblay, G., Vanderperre, S.,
819 Wisztorski, M., Salzet, M., Boisvert, F.M. and Roucou, X. (2013) Direct detection of
820 alternative open reading frames translation products in human significantly expands the
821 proteome. *PloS one*, **8**, e70698.
- 822 46. Zhang, B., Gunawardane, L., Niazi, F., Jahanbani, F., Chen, X. and Valadkhan, S.
823 (2014) A novel RNA motif mediates the strict nuclear localization of a long noncoding
824 RNA. *Molecular and cellular biology*, **34**, 2318-2329.
- 825 47. Ritchie, M.E., Silver, J., Oshlack, A., Holmes, M., Diyagama, D., Holloway, A. and
826 Smyth, G.K. (2007) A comparison of background correction methods for two-colour
827 microarrays. *Bioinformatics*, **23**, 2700-2707.
- 828 48. Bolstad, B. (2001) Probe Level Quantile Normalization of High Density Oligonucleotide
829 Array Data.

- 830 49. Gentleman, R.C., Carey, V.J., Bates, D.M., Bolstad, B., Dettling, M., Dudoit, S., Ellis, B.,
831 Gautier, L., Ge, Y., Gentry, J. *et al.* (2004) Bioconductor: open software development for
832 computational biology and bioinformatics. *Genome biology*, **5**, R80.
- 833 50. Lin, M.F., Jungreis, I. and Kellis, M. (2011) PhyloCSF: a comparative genomics method
834 to distinguish protein coding and non-coding regions. *Bioinformatics*, **27**, i275-282.
- 835 51. Goecks, J., Nekrutenko, A. and Taylor, J. (2010) Galaxy: a comprehensive approach for
836 supporting accessible, reproducible, and transparent computational research in the life
837 sciences. *Genome biology*, **11**, R86.
- 838 52. Sun, K., Chen, X., Jiang, P., Song, X., Wang, H. and Sun, H. (2013) iSeeRNA:
839 identification of long intergenic non-coding RNA transcripts from transcriptome
840 sequencing data. *BMC genomics*, **14 Suppl 2**, S7.
- 841 53. Wang, L., Park, H.J., Dasari, S., Wang, S., Kocher, J.P. and Li, W. (2013) CPAT:
842 Coding-Potential Assessment Tool using an alignment-free logistic regression model.
843 *Nucleic acids research*, **41**, e74.
- 844 54. Kong, L., Zhang, Y., Ye, Z.Q., Liu, X.Q., Zhao, S.Q., Wei, L. and Gao, G. (2007) CPC:
845 assess the protein-coding potential of transcripts using sequence features and support
846 vector machine. *Nucleic acids research*, **35**, W345-349.
- 847 55. Gascoigne, D.K., Cheetham, S.W., Cattenoz, P.B., Clark, M.B., Amaral, P.P., Taft, R.J.,
848 Wilhelm, D., Dinger, M.E. and Mattick, J.S. (2012) Pinstripe: a suite of programs for
849 integrating transcriptomic and proteomic datasets identifies novel proteins and improves
850 differentiation of protein-coding and non-coding genes. *Bioinformatics*, **28**, 3042-3050.
- 851 56. Knowles, D.G., Roder, M., Merkel, A. and Guigo, R. (2013) Grape RNA-Seq analysis
852 pipeline environment. *Bioinformatics*, **29**, 614-621.
- 853 57. Li, J.H., Liu, S., Zhou, H., Qu, L.H. and Yang, J.H. (2014) starBase v2.0: decoding
854 miRNA-ceRNA, miRNA-ncRNA and protein-RNA interaction networks from large-scale
855 CLIP-Seq data. *Nucleic acids research*, **42**, D92-97.
- 856 58. Lorenz, R., Bernhart, S.H., Honer Zu Siederdisen, C., Tafer, H., Flamm, C., Stadler,
857 P.F. and Hofacker, I.L. (2011) ViennaRNA Package 2.0. *Algorithms for molecular*
858 *biology : AMB*, **6**, 26.

861

862 **Figure Legends**

863

864 **Figure1: Discovery and classification of ribosome-associated lncRNAs.**

865 (A) Numbers of Gencode v7 lncRNA genes filtered by protein-coding prediction
866 methods. Genes (and all their constituent transcripts) having at least one transcript
867 identified as protein-coding by at least one method were designated “Potential Protein
868 Coding”. Remaining genes with no evidence for protein coding potential were defined as
869 “Filtered lncRNAs”. (B) Outline of the subcellular mapping of K562 lncRNA by
870 polysome profiling and microarray hybridisation. (C) Definition of the pooled fractions
871 from sucrose ultracentrifugation used in this study. (D) Summary of the numbers of
872 genes and transcripts classified in subcellular fractions. (E) Heatmaps show the relative
873 microarray intensity measured for each RNA sample. The colour scale runs from blue
874 (low detection) through white to red (high detection). “Protein coding” refers to the 2796
875 probes for protein-coding mRNAs included on the microarray, “Known lncRNAs” are
876 those filtered transcripts that also belong to the lncRNAdb database(23).

877

878 **Figure 2: Validation of selected ribosome-associated lncRNA candidates.**

879 (A) We individually validated nine predicted ribosome-associated lncRNAs in
880 independent ribosome profile experiments. In each case, two replicate experiments
881 each were carried out with control K562 (red) and cells treated with puromycin (blue),
882 for three distinct RNA fractions: (from left to right) free cytoplasmic, light polysomal,
883 heavy polysomal. RNA levels are normalized to absolute levels of an RNA spiked into
884 equal volumes of RNA sample. The first four panels represent protein coding mRNAs.
885 Transcript IDs and classifications are shown above each panel. (B) Genomic map of
886 ENST00000423918, a ribosome-associated transcript validated in this way.

887

888 **Figure 3: Fluorescence in situ hybridisation of ribosomal lncRNA in HeLa.**

889 Panel: DAPI staining of DNA; Middle: FISH probe; Right: merged. The actively
890 translated housekeeping mRNA GAPDH was tested as a positive control for
891 cytoplasmic localisation.

892

893 **Figure 4: Ribosomal and cytoplasmic lncRNA are under purifying selection.**

894 Cumulative distribution of the mean PhastCons nucleotide-level conservation for the
895 exons of the indicated transcript classes. PhastCons scores for ancestral repeats
896 regions are also included to represent neutral evolutionary rates.

897

898 **Figure 5: Intra- and Sub-cellular expression of ribosomal lncRNAs.**

899 (A) Expression in K562 whole cell by RNAseq. (B) Percent of transcripts having
900 ubiquitous expression across human tissues defined by Human Body Map RNAseq. (B)
901 Percent of ubiquitously expressed transcripts in each class. (C) Log2
902 cytoplasmic/nuclear RPKM ratios calculated from ENCODE RNAseq for indicated RNAs
903 in K562 (whole cell, polyA+). For potential protein coding transcripts and mRNAs, data
904 is only shown for detected transcripts. Numbers indicate median value. (D) Subcellular
905 localisation of lncRNA amongst different cell lines. Colours reflect median cytoplasmic /
906 nuclear RPKM values.

907

908 **Figure 6: Ribosomal lncRNAs have mRNA-like 5' ends.**

909 (A) The pseudo 5' UTR was defined to be the distance from the start to the first AUG
910 trinucleotide (top row). As a control, we calculated the same measure on the antisense
911 strand (bottom row). Shown is the distribution of these lengths for each set of transcripts
912 - protein coding mRNA (left), followed by cytoplasmic lncRNA classes. The red line
913 indicates the mean value. P-values are for comparison of sense and antisense
914 distributions using the Kolmogorov-Smirnov test. (B) Capping efficiency positively
915 correlates with Light Polysome localisation of lncRNA. We defined every transcript to be
916 capped if it has a K562 cytoplasmic polyA+ CAGE tag within 100bp upstream or
917 downstream of its transcription start site. lncRNA were binned according to their
918 relative enrichment in each of the three cytoplasmic fractions (x axis). In each bin, the
919 percent of capped transcripts is shown in the y axis. Logistic regression was used to
920 assess the relationship between these variables.

921

922 **Figure 7: Transposable element composition of lncRNAs.** (A) The fraction of each
923 transcript covered by annotated repeat sequence from RepeatMasker. (B) The heatmap

924 shows the normalised frequency of insertion for RepeatMasker-defined classes, ie the
925 number of insertions per class divided by the length of each transcript. (C) As in (B) but
926 showing data only for MLT-type repeats. (D) The repeat composition of a selection of
927 Free Cytoplasmic, MLT-containing lncRNAs. The direction of the arrows indicates the
928 annotated strand of the repeat with respect to the lncRNA. The colours represent the
929 repeat class. (E) As in (B), except showing data for HeLa derived from ribosome
930 footprinting experiments(36). For practical reasons, the lncRNAs are divided into three
931 classes, see Materials and Methods for more details.

932

933 **Figure 8: Cis- and Trans-antisense lncRNA-mRNA pairs and ribosomal**
934 **recruitment.**

935 (A) Cartoon illustrating the definition of cis- and trans-antisense lncRNA-mRNA pairs.
936 Red boxes indicate regions of opposite-strand homology. (B) The percent of each
937 subcellular lncRNA class defined as exonic-antisense (cis-antisense) to a protein coding
938 gene. (C) Example of a cis-antisense lncRNA-mRNA pair. ENST00000529247 (forward
939 strand) is a heavy polysomal lncRNA transcribed antisense to the EEF1D gene (reverse
940 strand), encoding a subunit of the translation elongation factor 1 complex. (D) Whole
941 cell K562 polyA+ steady state levels of mRNAs that are antisense to lncRNA in the
942 indicated subcellular classes. (E) The percent of lncRNA (blue bars) or size-matched
943 random genomic fragments, having an antisense trans homology match to at least one
944 mRNA.

945

946 **Figure 9: Changes in lncRNA stability in response to ribosome stalling.**

947 Bars represent mean detection in cells treated with cyclohexamide (CHX) for 6 hours
948 and control cells (0h). Experiments were performed with three biological replicates. Bars
949 show mean and standard deviation. Statistical significance was calculated by one-sided
950 t-test (* P<0.05, ** P<0.01, *** P<0.001)

951

952 **Figure 10: A model of lncRNA targeting within the cytoplasm.**

953

954 **Supplementary Data Files**

955

956 **Table S1: Gencode v7 lncRNA classification.** Rows represent lncRNA transcript from
957 Gencode v7.

958 Columns:

959 TransID: ENST ID for transcripts.

960 GeneID: ENSG ID for the corresponding gene.

961 Chr: Chromosome.

962 Trans_Start: Transcript start position.

963 Trans_End: Transcript end position.

964 Strand

965 Cellular_Localization: Classification of the transcripts into 5 different categories: 1:
966 Present in cytoplasm (from polysome profiling experiment, K562 cell line); 2: Present in
967 nucleus (from ENCODE nucleus RNAseq data, K562 cell line), but not in cytoplasm
968 (polysome profiling experiment); 3: Not present in cytoplasm (polysome profiling
969 experiment) nor in nucleus (ENCODE nucleus RNAseq data, K562 cell line); 4:
970 Transcripts classified as potential protein coding transcripts; 5: Discarded transcripts.

971 See Materials and Methods for details.

972 FreeC_intensity: Log2 intensity value for Free C. condition (NA if not present in this
973 condition).

974 LightP_intensity: Log2 intensity value for light P. condition (NA if not present in this
975 condition).

976 HeavyP_intensity: Log2 intensity value for Heavy P. condition (NA if not present in this
977 condition).

978 WholeC_intensity: Log2 intensity value for whole cytoplasmic fraction (NA if not present
979 in this condition).

980 Ribosomal_Classification: Classification for transcripts present in the cytoplasm: 1: Free
981 C.; 2: Light P.; 3: Heavy P.

982 CPAT: CPAT score.

983 PhyloCSF: PhyloCSF score.

984 CPC: CPC score.

985 MS: information about presence of peptides in Mass Spectrometry analysis for this
986 transcripts: 0: No peptide associated; 1: Peptide associated.

987

988 **Table S2: Correlation of gene expression quantification between microarray K562**
989 **cytoplasmic measurements, and ENCODE RNAseq data from K562 cellular**
990 **fractions (19).**

991

992 **Table S3: Heavy Polysome mRNAs are most actively translated.** Shown are the
993 numbers of ENCODE K562 mass spectrometry tags originating from ribosome-profiled
994 mRNAs.

995

996 **Table S4: Small peptides originating from lncRNA.** Shown are the numbers of
997 known small peptides discovered by mass spectrometry that map to Gencode v7
998 lncRNA (43).

999

1000 **Figure S1: Comparison of lncRNA microarray and RNAseq quantifications.**
1001 **Steady state values for K562 cytoplasmic RNA were analysed.** RNAseq data was
1002 obtained from ENCODE. Only transcripts detected in both experiments are shown.

1003

1004 **Figure S2: Mean tissue expression across 16 human tissues from Human Body**
1005 **Map.** In the cases of potential protein coding and protein coding transcripts, data is only
1006 shown for those transcripts detected in K562.

1007

1008 **Figure S3: Protein interactions related to subcellular compartmentalisation.**
1009 Heatmap depicts lncRNA genes that interact with the indicated proteins, as defined by
1010 the Starbase database (57). Interactions of “Low stringency” were used in all cases. The
1011 colour scale indicates the percent difference of the actual to expected number of
1012 overlaps. The rows show lncRNA gene sets assigned to the four subcellular lncRNA
1013 categories, and the columns represent various proteins for which CLIPseq binding sites
1014 were analysed. Arrows indicate the reported subcellular localisation of the protein,
1015 identified by manual curation of the literature. We found a number of cases where the

1016 localisation of lncRNA corresponded with the known distribution of the protein to which
1017 they are bound: Nuclear-associated lncRNAs showed elevated binding to nuclear-
1018 localised proteins, including hnRNPC ($P=0.007$, Fisher's exact test), U2AF65 ($P=0.002$),
1019 and eIFAIII ($P=0.0002$). In contrast, lncRNAs bound by the cytoplasmic acting IGFBP1
1020 were significantly enriched in the Light Polysomal and free cytoplasmic fractions
1021 ($P=0.007$). Light Polysomal lncRNAs are enriched for binding by the TAF15 protein
1022 ($P=0.033$). A general depletion of Heavy Polysomal lncRNA was observed in the protein
1023 binding data.

1024

1025 **Figure S4: Association of ORF length with polysome density for protein coding**
1026 **transcripts.** Shown are histograms for % coverage of transcripts by their longest ORF.
1027 Top row: sense strand ORFs; Bottom row: antisense strand ORFs (control). Data are
1028 shown for mRNAs included in microarray design and classified by ribosomal occupancy.
1029 P values compare sense / antisense distributions in each case. Red line indicates mean
1030 ORF coverage percentage. Note the difference in mean value between Heavy and Light
1031 Polysomal means.

1032

1033 **Figure S5: Association of ORF length with polysome density for lncRNA**
1034 **transcripts.** Shown are histograms for % coverage of transcripts by their longest ORF.
1035 Top row: sense strand ORFs; Bottom row: antisense strand ORFs (control). Data are
1036 shown for all mRNAs together, for comparison. P values compare sense / antisense
1037 distributions in each case. Red line indicates mean ORF coverage percentage. Note the
1038 lack of difference in mean value between Heavy and Light Polysomal means.

1039

1040 **Figure S6: Gene structure characteristics of lncRNAs.** (A) Exon length distributions.
1041 (B) Intron length distributions. (C) Exon number per transcript. (D) Mature (processed)
1042 transcript length.

1043

1044 **Figure S7: GC content of coding and noncoding transcripts.**

1045 **Figure S8: Polyadenylation of lncRNAs.** Using ENCODE data (19), we calculated the
1046 ratio of RPKM for PolyA+ / PolyA- K562 cytoplasmic RNAseq. RPKM values were
1047 averaged across the two available technical replicates, and only transcripts with non-
1048 zero mean values in both RNA samples retained. No statistically significant differences
1049 were found between Free Cytoplasmic lncRNAs and either of the ribosomal groups
1050 using either the Kolmogorov-Smirnov or Wilcoxon tests.

1051 **Figure S9: Splicing efficiency of lncRNAs.** Using ENCODE data, we calculated
1052 RPKM values separately for the exons and introns of all lncRNAs. Shown are the log10
1053 ratios of exon/intron values for all sets of transcripts. No statistically significant
1054 differences were found between Free Cytoplasmic lncRNAs and either of the ribosomal
1055 groups using either the Kolmogorov-Smirnov or Wilcoxon tests.

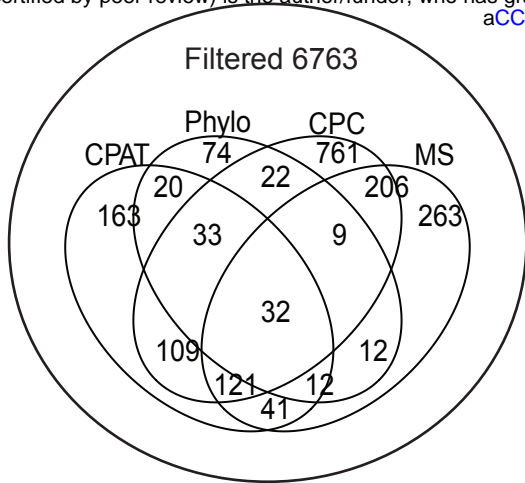
1056 **Figure S10: Comparison of 5' RNA folding energy.** Using the Vienna RNAfold
1057 programme (58) with default settings, we estimated the free energy of folding of the first
1058 50nt of lncRNA and mRNA. While mRNA have more stable folding on average than
1059 expressed lncRNA ($P=2.2e-16$, Wilcoxon test), we could find no difference between
1060 either Heavy Polysomal ($P=0.8$) or Light Polysomal ($P=0.7$) and Free Cytoplasmic
1061 lncRNAs.

1062 **Figure S11: ERVL-MaLR insertion length distributions.**

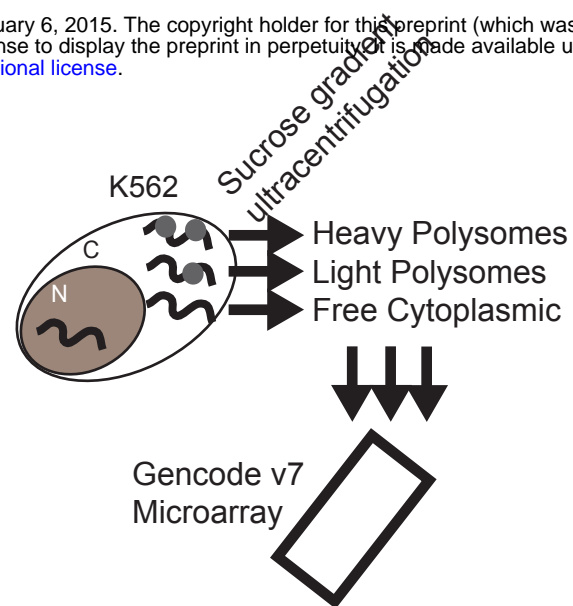
1063 **Figure S12: The association between sense-intronic lncRNAs and nuclear**
1064 **localisation.** Shown is the percent of transcripts in each subcellular category that are
1065 located within the intron of a same-strand protein coding gene.

1066

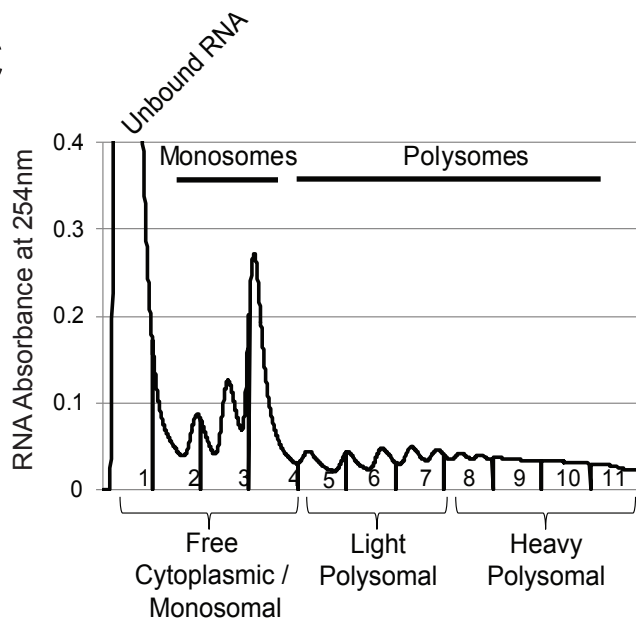
A



B



C



D

	Transcripts	% cytoplasmic transcripts	Genes
Heavy Polysome	204	22.1	177
Light Polysome	373	40.4	248
Free Cytoplasmic	347	37.5	248
Nucleus	292		255
Not Present	7754		6033
Potential Coding	4415		1878
Discarded	1169		521

E

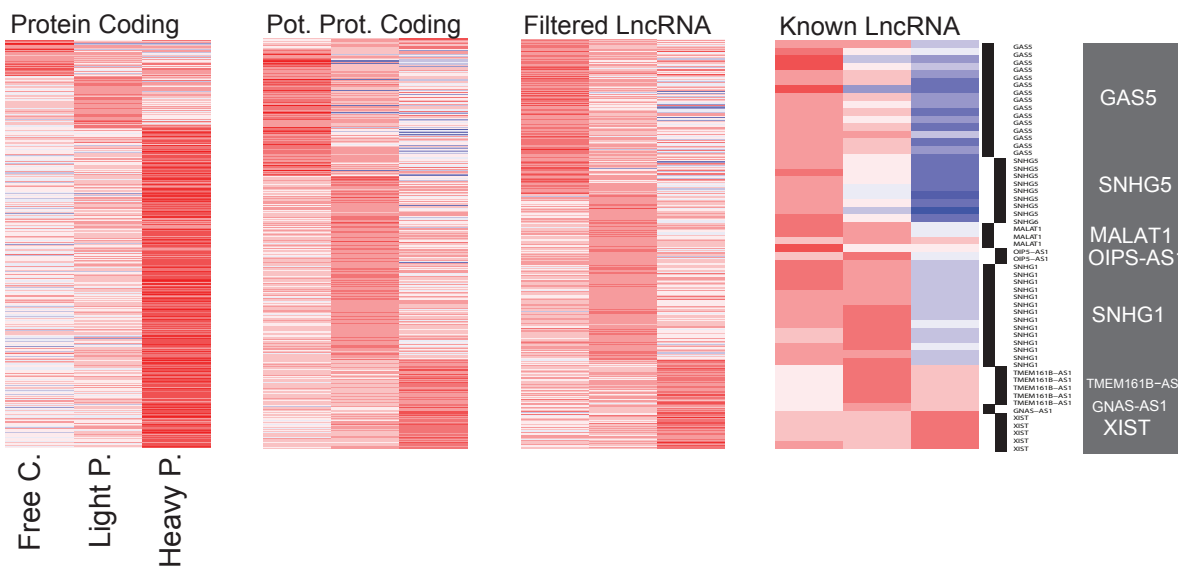
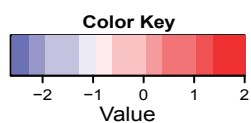


Figure 3

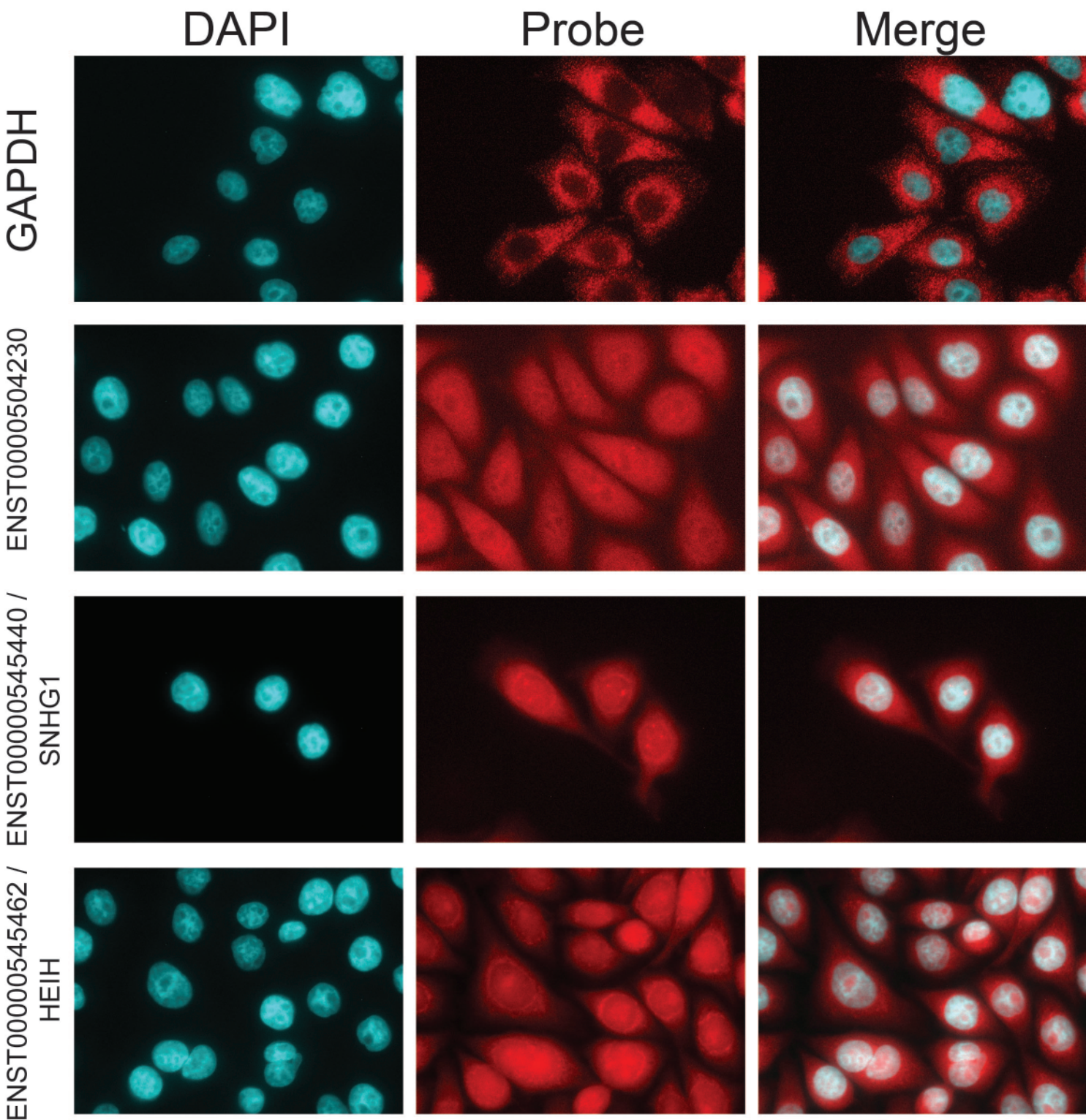


Figure 5

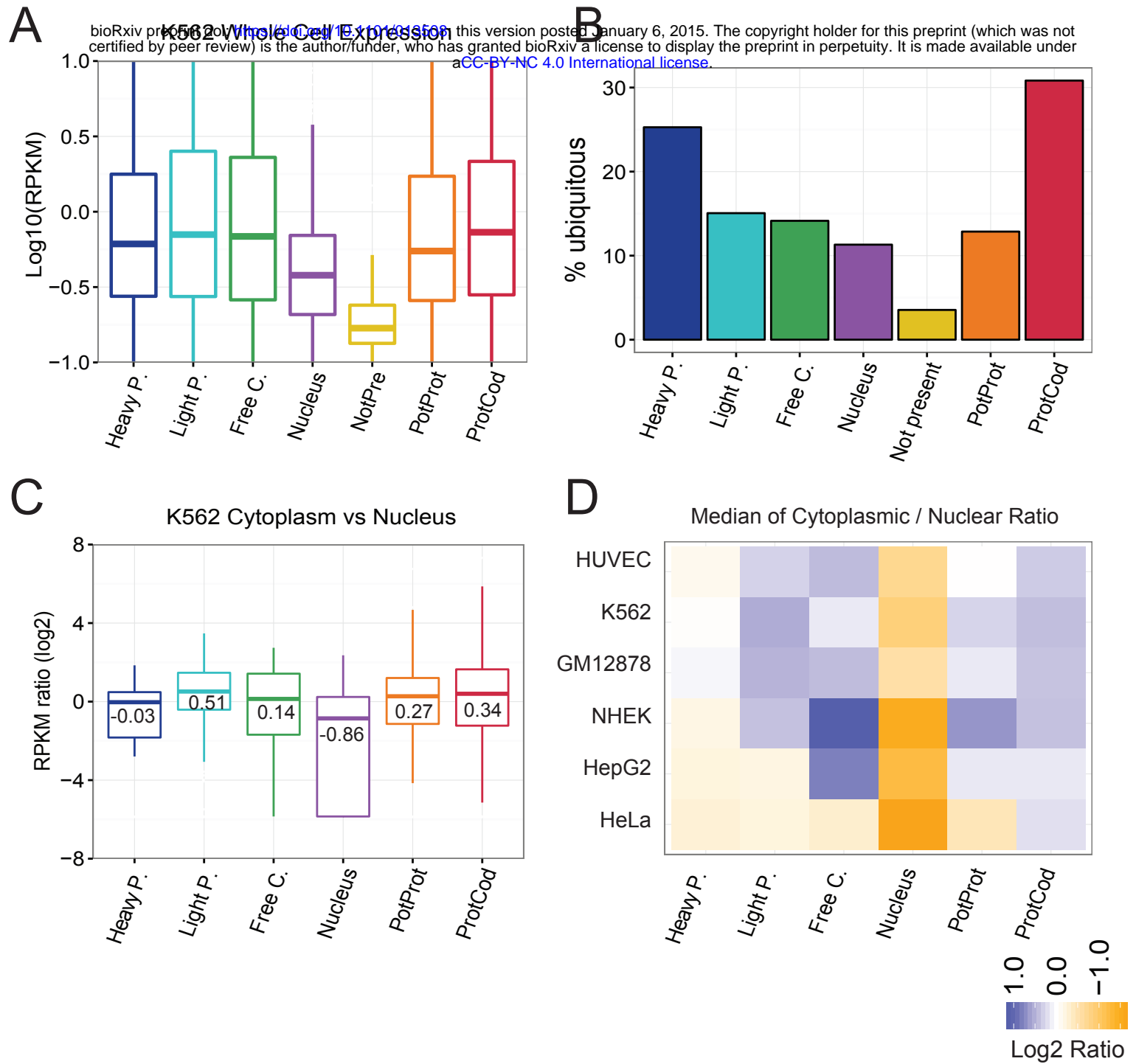
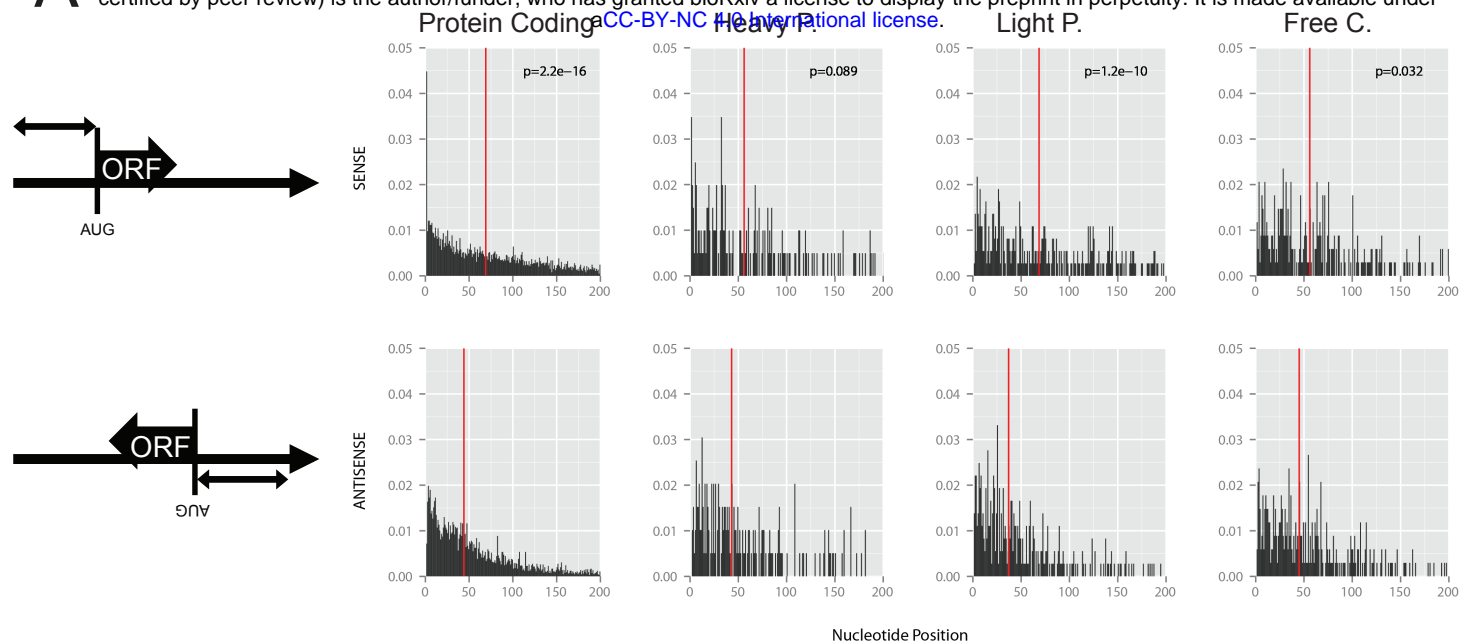
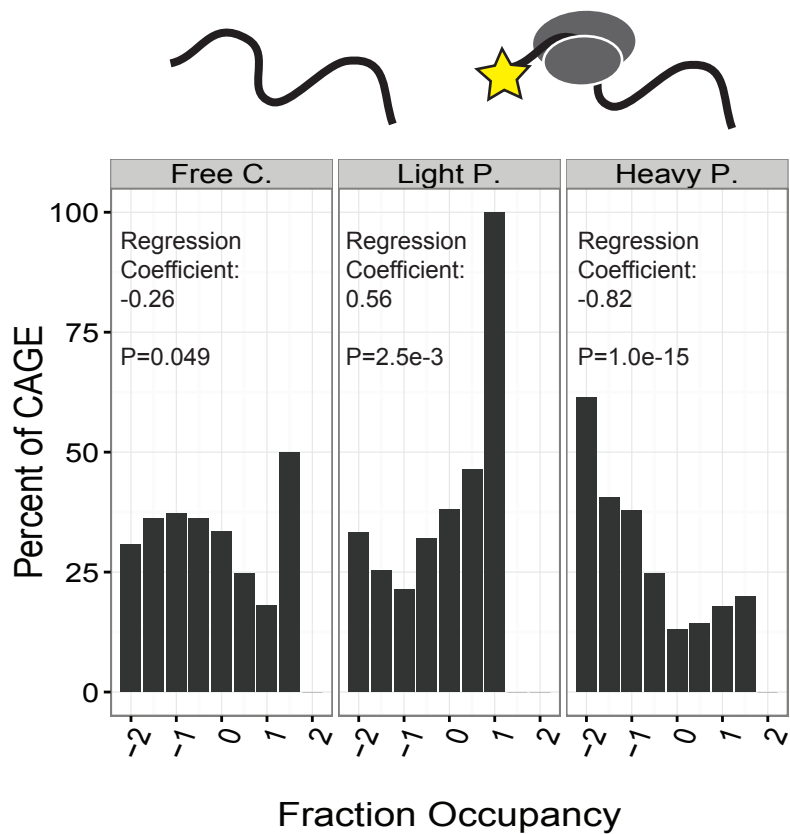


Figure 6

A bioRxiv preprint doi: <https://doi.org/10.1101/013508>; this version posted January 6, 2015. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.

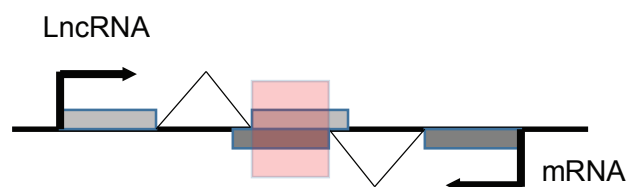


B

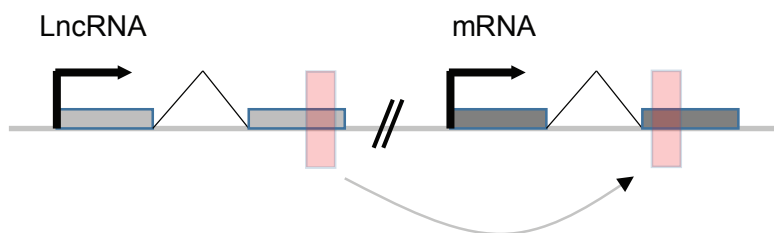


A

Exonic ("cis") antisense

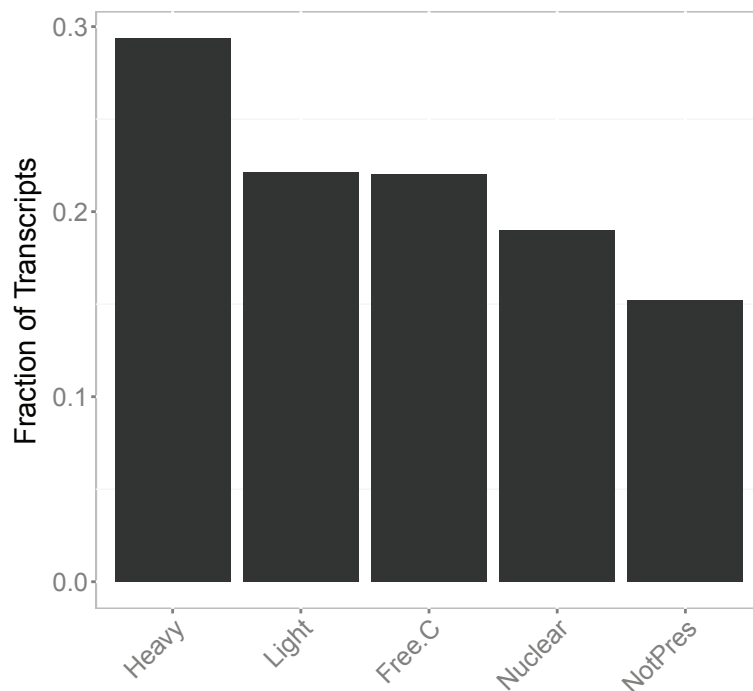


Trans antisense

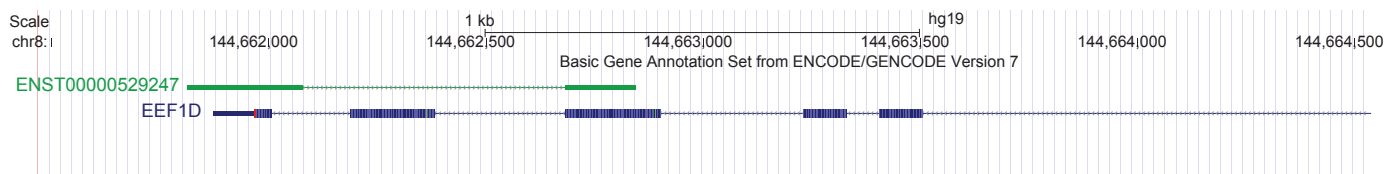


B

Exonic Antisense

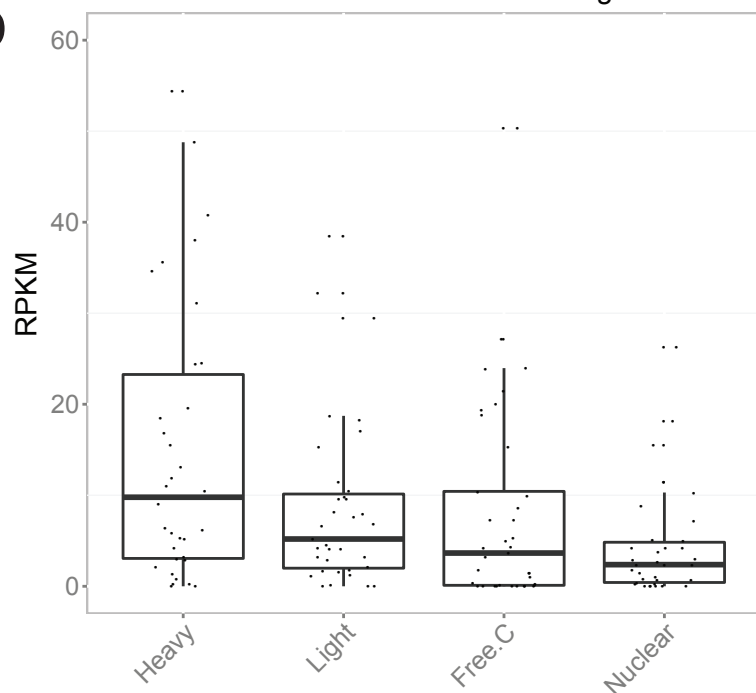


C



D

Exonic Antisense Protein Coding Genes



E

Trans Antisense Homology

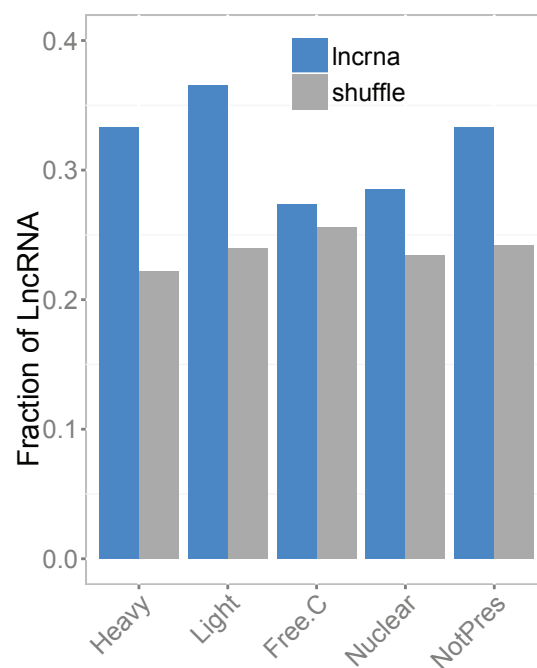


Figure 9

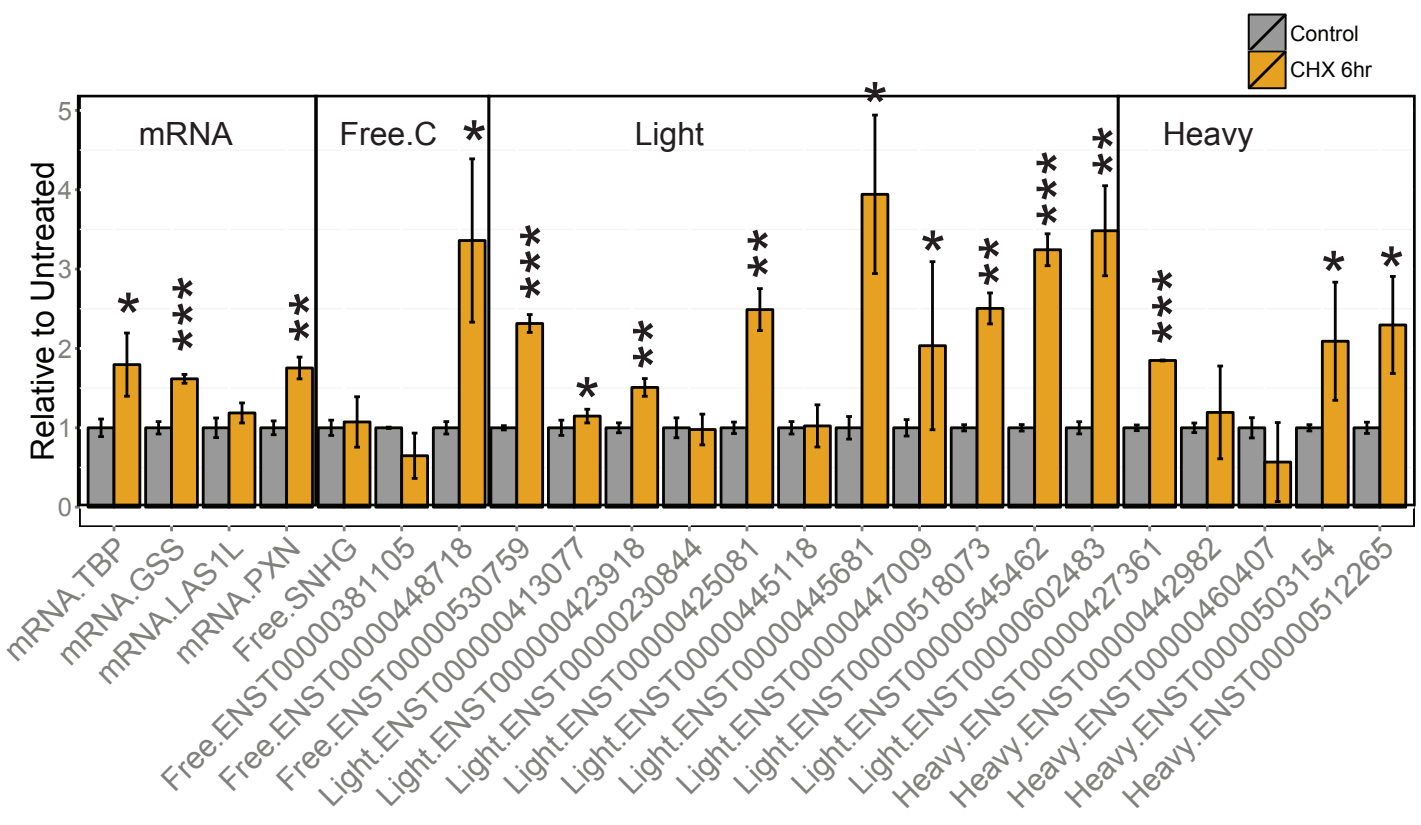


Figure 10

