1  **LncATLAS database for subcellular localisation of long noncoding RNAs.**

2

3

4  David Mas-Ponte[1,2,6]

5  Joana Carlevaro-Fita[4,5,6]

6  Emilio Palumbo[1]

7  Toni Hermoso[1]

8  Roderic Guigo[1,2,3]

9  Rory Johnson[4,5]*

10

11

12  1. Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Dr.
13  Aiguader 88, Barcelona 08003, Spain
14  2. Universitat Pompeu Fabra (UPF), Barcelona, Spain.
15  3. Institut Hospital del Mar d'Investigacions Mèdiques (IMIM), Dr. Aiguader 88, 08003 Barcelona,
16  Catalonia, Spain.
17  4. Department of Clinical Research, University of Bern, 3008 Bern, Switzerland
18  5. Department of Medical Oncology, Inselspital, University Hospital and University of Bern, 3010
19  Bern, Switzerland
20  6. Equal contribution

21

22  **\* Correspondence to rory.johnson@dkf.unibe.ch**

23

24

25

26

27

28

29  **Keywords: long noncoding RNA; lncRNA; subcellular localisation; nucleus; cytoplasm;**
30  **chromatin.**

**Abstract**

**Background**

The subcellular localisation of long noncoding RNAs (lncRNAs) holds valuable clues to their molecular function. However, measuring localisation of newly-discovered lncRNAs involves time-consuming and costly experimental methods.

**Results**

We have created "LncATLAS", a comprehensive resource of lncRNA localisation in human cells based on RNA-sequencing datasets. Altogether, 6768 GENCODE-annotated lncRNAs are represented across various compartments of 15 cell lines. We introduce "Relative concentration index" (RCI) as a useful measure of localisation derived from ensemble RNAseq measurements. LncATLAS is accessible through an intuitive and informative webserver, from which lncRNAs of interest are accessed using identifiers or names. Localisation is presented across cell types and organelles, and may be compared to the distribution of all other genes. Publication-quality figures and raw data tables are automatically generated with each query, and the entire dataset is also available to download.

**Conclusions**

LncATLAS makes lncRNA subcellular localisation data available to the widest possible number of researchers. It is available at lncATLAS.crg.eu.

1  **<u>Introduction</u>**

2

3       The functions of long noncoding RNAs (lncRNAs) are intimately linked to location in the cell.
4  The first discovered lncRNAs tended to be located in the nucleus and chromatin, and epigenetically
5  regulate gene expression (Hutchinson et al. 2007; Mondal et al. 2010; Rinn et al. 2007; Tsai et al. 2010;
6  Whitehead, Pandey, and Kanduri 2009; Zhao et al. 2008). However, we now appreciate lncRNAs'
7  localisation, and molecular functions, to be highly diverse. There exists a substantial population of
8  lncRNAs in the cytoplasm (Carlevaro-Fita et al. 2016; van Heesch et al. 2014; Ulitsky and Bartel 2013),
9  with evidence for roles such as translation regulation (Schein et al. 2016; Yoon et al. 2012; Zucchelli et
10 al. 2016), miRNA decoys (Cesana et al. 2011), or protein trafficking (Aoki et al. 2010; Kino et al. 2010;
11 Willingham et al. 2005). Consequently, ascertaining nuclear-cytoplasmic localisation has become one
12 of the primary sources of evidence when investigating the molecular role of newly-discovered lncRNAs
13 (J. Chen et al. 2016; L.-L. Chen 2016; Hansji et al. 2016; Hutchinson et al. 2007; Ishizuka et al. 2014;
14 Ounzain et al. 2015).

15

16      The various methods to map RNA molecules in the cell operate with trade-offs in throughput,
17 convenience, and accuracy. Amongst the single-gene approaches, probably the most commonly used is
18 qRTPCR on RNA extracts of purified cellular compartments (Wang, Zhu, and Levy 2006). It yields
19 information on relative RNA concentrations between compartments, but not of absolute molecule
20 numbers per cell. Another method is fluorescence in situ hybridization (FISH), which can in principle
21 yield absolute counts of molecules at subcellular resolution (Dunagin et al. 2015; Raj et al. 2008).
22 However, FISH is time-consuming and low-throughput, and requires expensive reagents (Cabili et al.
23 2015). More recently, the ingenious in situ sequencing method, FISSEQ, has established high-
24 throughput subcellular RNA counting (Lee et al. 2015). But at present just one dataset is available, and
25 is restricted to several hundred highly-expressed lncRNAs (Lee et al. 2014).

26

27      The only method presently capable of whole-genome localisation mapping is subcellular RNA
28 sequencing (subcRNAseq). Here cells are fractionated, and extracted RNA sequenced (Djebali et al.
29 2012). SubcRNAseq yields high-throughput and quantitative data, although as with RTPCR approaches
30 mentioned above, the absolute counts of RNA molecules per cell are lost (Ulitsky and Bartel 2013).
31 Recently, large amounts of raw subcRNAseq data have become available, most notably from the
32 ENCODE consortium (Djebali et al. 2012; Dunham et al. 2012). These data remain under-utilized and
33 have not been made readily accessible.

34

35      In light of the growing use of RNA localisation to infer function of newly-discovered lncRNAs,
36 and the availability of large amounts of unprocessed subcRNAseq data, we have created a resource to
37 make lncRNA localisation data available to the broader scientific community. This resource,
38 "lncATLAS", enables non-expert users to rapidly access a rich variety of easily-interpreted data on their
39 lncRNA of interest.

40

41

42

43

1  **Results**
2
3  **A database of lncRNA localisation based on human RNAseq data**
4          In light of growing interest in lncRNA and their functions, we decided to create a resource for
5  accessing and visualizing lncRNA localisation within human cells. We collected data from the largest
6  dataset of subcRNAseq, produced by the ENCODE consortium (Djebali et al. 2012; Dunham et al.
7  2012). Raw RNAseq data from a panel of human cell lines were used to quantify the reference
8  GENCODE gene annotation (Derrien et al. 2012; Harrow et al. 2012). RNAseq experiments were
9  obtained for a total of 15 cell lines comprising 48 individual experiments (see Supplementary Table
10 S1). These cells originate from a wide diversity of adult and embryological organ sites, and comprise
11 both transformed and normal cells (Figure 1A). For each cell, cytoplasmic and nuclear data are
12 available, and for the majority of these, whole-cell data were also obtained. In addition, from a single
13 cell line, K562, sub-nuclear and sub-cytoplasmic compartment data are also available (Figure 1B;
14 Supplementary Table S1). Hereafter we refer to these as "compartments". PolyA+ RNA samples were
15 available for whole cell, cytoplasm, nucleus and sub-cytoplasmic compartments, and total RNA for sub-
16 nuclear compartments.
17
18
19 **Defining localisation from RNAseq data**
20         Throughout the present study, for practical reasons, we adopt a relative scheme to define and
21 quantify RNA localisation: the "Relative Concentration Index" (RCI). RCI is defined as the log2-
22 transformed ratio of FPKM (fragments per kilobase per million mapped) in two samples, for example
23 the cytoplasm and nucleus (Figure 1C). A similar approach has been used previously (Derrien et al.
24 2012; Ulitsky and Bartel 2013). It is worth commenting on exactly how these values should be
25 interpreted: RCI is the ratio of a transcript's concentration, per unit mass of sampled RNA, between
26 two compartments. Sampled RNA populations may be PolyA+ RNA or total RNA, and we are careful
27 to only compute RCI values within the same population.
28         The mass of RNA per compartment per cell is not equal, and typically not quantified prior to
29 RNAseq (Djebali et al. 2012). Therefore, without knowing the total mass of PolyA+ RNA in the nucleus
30 and cytoplasm of a single cell, we cannot make statements about the relative *number* of RNA transcripts
31 in cellular compartments of a single cell (Cabili et al. 2015).
32         We here briefly digress to contrast this approach with another possible way to define relative
33 subcellular localisation. Perhaps more obvious is a "molecular" definition, in terms of numbers of
34 molecules of a given RNA transcript $X$ in the compartments of a single cell. For example, if one cell
35 has 10 and 5 molecules of $X$ in the cytoplasm and nucleus, respectively, then its cytoplasmic/nuclear
36 localisation would be defined as 10/5=2. We define this measure as "Relative Molecules Index" (RMI).
37 Such information is, in principle, directly accessible from fluorescence-based techniques, and has been
38 calculated previously (Cabili et al. 2015; Lee et al. 2014). As mentioned above, our ignorance of the
39 total PolyA+ RNA mass of the cell lines used here, precludes the calculation of RMI in this study.
40
41 **Computing localisation across genes and cell types**
42         RCI was calculated for various selected pairs of cellular compartments (Figure 1D). In the
43 majority of cases, we calculated the cytosoplasmic/nuclear RCI – "CN-RCI" (Supplementary Table S2)
44 (see Materials and Methods). This is a measure of the relative concentration of an RNA sequence in the
45 cytoplasm, compared to the nucleus, in log2 units. For one cell type, K562, total RNA data from sub-
46 nuclear and PolyA+ RNA from sub-cytoplasmic compartments were also calculated, by reference to
47 total RNA from the nucleus or PolyA+ RNA from the cytoplasm, as appropriate (Supplementary Table
48 S3) (see Materials and Methods). Altogether this yielded localisation estimates in 20 compartment / cell

1     combinations (Figure 1D).

2         Where available, replicate data were used to assess reliability of RCI measurements (see
3     Supplementary Table S1 for information about availability of replicates). Silent and unreliable genes
4     were excluded from further consideration (see Materials and Methods for details in the filtering steps).
5     Between 3114 (max) and 582 (min) lncRNAs' CN-RCI localisation could be estimated per cell, after
6     filtering (Figure 1D,2A). Note that H1.hESC has a greater number of detected genes (4923 genes),
7     because biological replicates were not available for cytoplasm nor nucleus for this cell line. A total of
8     24,538 genes (17,770 mRNAs and 6,768 lncRNAs) were quantified in at least one cell type. Of these,
9     31 lncRNAs were detected in all samples tested (Figure 2B and Supplementary Table S4). LncRNAs
10     display a highly cell-type specific detection pattern, in contrast to mRNAs, as observed previously
11     (Figure 2B) (Cabili et al. 2015; Derrien et al. 2012; Guttman and Rinn 2012).

12         RCI data are consistent with known cytoplasmic-nuclear localisation tendencies of lncRNAs
13     and mRNAs. Amongst the top 15 most cytoplasmic measurements, 14 represent mRNAs (the remainder
14     is an annotation of uncertain biotype and may be protein-coding) (Supplementary Table S5). In contrast,
15     12 of the 15 most nuclear RCI values represent lncRNAs (Supplementary Table S6). The nuclear-
16     enriched X-chromosome inactivating transcript *XIST* occupies the top four positions (Brown et al. 1992;
17     Clemson et al. 1996). Manual inspection of several well-known lncRNAs showed that localisation
18     reported here tended to be consistent with literature reports (see next Section).

19         Detected lncRNA genes reported by lncATLAS cover a substantial fraction of the entries from
20     manually-curated and widely-used databases, such as lncRNAdb and LncRNADisease, and new
21     localization specific databases such as RNALocate (Zhang et al. 2016). Note that these databases
22     contain a mixture of GENCODE annotated and non-GENCODE annotated entries. 74, 128 and 150
23     genes from lncRNAdb, lncRNADisease and RNALocate respectively are detected in lncATLAS. These
24     numbers represent a 39%, 48% and 47% of the total number of human lncRNAs and a 63%, 73% and
25     90% of the total number of human GENCODE annotated lncRNAs of each database, respectively
26     (Figure 2C). LncRNAs from these databases that are not displayed in lncATLAS are either not detected
27     in any of the cell lines considered, or else do not belong to the GENCODE annotation.

28

### LncATLAS webserver for exploring localisation data

30         The lncATLAS dataset was compiled into a relational database that is searchable through a
31     webserver at lncATLAS.crg.eu. LncRNAs of interest are accessed using official gene names or
32     GENCODE gene identifiers. A maximum of three genes may be investigated simultaneously. Several
33     well-known lncRNAs with known localisation are also available for reference. Once a gene or genes
34     has been selected, a series of data interpretations are presented, summarized below. As examples, results
35     for *MALAT1* (nuclear localized) (Hutchinson et al. 2007) and *DANCR* (cytoplasm localized) (Cabili et
36     al. 2015; van Heesch et al. 2014) lncRNAs data are shown in Figure 3 and 4.

37         The following sections summarise the data presented to a user for their gene of interest (GOI).

38

39     *1.*     *Inspect the cytoplasmic-nuclear localisation of your gene of interest (GOI)*
40     Data are only displayed for selected gene(s).

41

42         **Plot 1: Cytoplasmic/Nuclear Localisation: RCI and expression values (all cell types)**: In
43     this basic summary, the Cytoplasmic/Nuclear RCI is shown as a bar plot across all available cell types.
44     Bars are coloured to reflect the expression level of the gene, as inferred from nuclear RNAseq. The
45     individual FPKM values, upon which RCI values are based, are displayed. When a gene is expressed
46     only in one compartment, RCI cannot be computed; then, dashed bars with expression values are shown
47     instead (Figure 3A).

48

*2. Inspect the cytoplasmic-nuclear localisation of your GOI within the distribution of all genes.*

  The aim of the second section is to understand, in terms of localisation, how the genes of interest behave relative to all other genes. Three different plots show CN-RCI values distribution for all lncRNAs and mRNAs, within which the location of the GOI is indicated.

  **Plot 2: Cytoplasmic/Nuclear Localisation: RCI distribution (all cell types):** To put RCI values in context, their percentile rank within the distribution of all lncRNAs is indicated (ranks relative to lowest value). Data are shown for all cell types (Figure 3B).

  **Plot 3: Cytoplasmic/Nuclear Localisation: RCI distribution (individual cell type):** The same data are shown as for Plot 2, but in the form of a density plot. The User must here specify a single cell type. When genes are not classed as "Detected", RCI cannot be computed and no data are shown (Figure 3C).

  **Plot 4: Cytoplasmic/Nuclear Localisation: Comparison with expression (individual cell type):** As for Plot 3, gene values are shown in the context of all other genes in the same cell, but here also indicating whole-cell expression values. As before, the data are shown for a single cell type chosen by the User, and plots are only generated for cells where RCI values are "Detected" (Figure 3D).

*3. Inspect the localisation of your GOI at sub-compartment level.*

  The final section gives information about enrichment in the cytoplasmic and nuclear subcompartments of K562 cells. As in the previous section, RCI values for the genes of interest are indicated in the context of full lncRNA and mRNA distributions.

  **Plot 5: Sub-cytoplasmic and Sub-nuclear Localisation in K562 Cells:** Here data are shown for sub-nuclear and sub-cytoplasmic compartments K562 cell line. As in Plots 2 and 3, distributions across all detected genes are shown (Figure 4).

  In the examples shown, the differences in localisation of *MALAT1* and *DANCR* are clear. Their cytoplasmic-nuclear localisations are highly divergent (Figure 3 and 4), and broadly consistent across all the cell lines observed. The difference in localisation is observed even in cells where their overall expression level is similar (eg HUVEC, Figure 3D).

  All figures may be downloaded as publication-quality files in pdf format. Similarly, in the *Get Raw Data* tab, the underlying RCI and raw expression values for selected genes may be accessed as a batch query. Furthermore, the entire set of data tables for lncATLAS may be downloaded in the same tab using the *Download All raw data* button.

1 **Discussion**

2        Subcellular localisation provides important clues to the molecular function of novel lncRNAs.

3 LncATLAS is designed to make such data available to the largest number of researchers. To our

4 knowledge, only one other database of lncRNA localisation exists: RNALocate (Zhang et al. 2016).

5 RNALocate contains manually-curated localisation classifications across multiple species. Despite

6 focussing on a single species (human), due to the limited availability of subcellular RNAseq data,

7 LncATLAS has two key advantages: it is quantitative, and it is based on standard GENCODE

8 annotations, the *de facto* official annotation for both protein-coding and lncRNA genes (Derrien et al.

9 2012). These features boost the usefulness of lncATLAS data for other research groups and ensure its

10 integration with diverse other genomics datasets. Future subcellular RNAseq data from other cell types,

11 or other species, will be integrated as they become available.

12

1 **Materials and Methods**
2
3 **Data source**
4 Cytoplasmic and nuclear PolyA+ RNAseq data from 15 different cell lines were obtained from
5 ENCODE (Djebali et al. 2012). (ENCODE RNAseq data in BAM format were obtained from ENCODE
6 Data Coordination Centre (DCC) in September 2016 -
7 **https://www.encodeproject.org/matrix/?type=Experiment).** For most cell lines, whole-cell data
8 were also obtained (exceptions being HT1080, NCI.H460, SK.MEL.5 and SK.N.DZ). A full list of
9 processed RNAseq libraries are available in Supplementary Table S1.
10
11
12 **Data processing**
13 Data were mapped to human genome assembly GRCh38 using STAR software (Dobin et al.
14 2013) and quantified with RSEM (Li and Dewey 2011) for all GENCODE v24 gene quantification,
15 within the GRAPE analysis pipeline (Harrow et al. 2012; Knowles et al. 2013). Data consisted of two
16 independent biological replicates per cell line and fraction (exceptions being H1.hESC cytoplasm and
17 nucleus and NCI.H460 cytoplasm for which only one replicate was available) (see Supplementary Table
18 S1 for a full list of source datasets). For sub-cytoplasmic RCI, instead of using poly+ cytoplasmic
19 sample coming from Gingeras lab (used for CNRCI), we used the corresponding sample from the lab
20 where the sub-cytoplasmic fractionation was done (Lécuyer lab). This is not considered as an additional
21 biological replicate.
22 Throughout, RNAseq data are processed at the level of genes, rather than transcripts. From the
23 whole GENCODE v24 annotation, genes contained in the "Long non-coding RNA gene annotation"
24 define our lncRNA set of genes. Protein coding gene set is defined by GENCODE biotype
25 "protein_coding" (see Supplementary Table S7).
26 In order to remove genes with high variability between replicates, genes with >2-fold difference
27 between replicates are labelled "unreliable" and excluded from further analysis. This cutoff was not
28 possible for the samples, mentioned above, for which replicate experiments were not available.
29
30 **Cytoplasmic-nuclear relative concentration index (CN-RCI)**
31 For cytoplasmic/nuclear localisation, PolyA+ RNA data were used. At this stage, all genes are
32 defined in one of four categories in each cell line: (1) Detected: genes with non-zero and "reliable"
33 values in both cellular compartments; (2) Compartment-specific: genes considered "reliable" in both
34 compartments but expressed >0 FPKM only in one; (3) Silent: both compartments have FPKM=0; (4)
35 Unreliable: genes that are unreliable in at least one compartment.
36 Genes classed as "Detected" were retained, and their localisation was computed as the C/N
37 Relative Concentration Index (CN-RCI) thus:
38

39 $$RCI^{pA+}_{C/N} = log2(\frac{Cytosoplasmic\ Expression\ (FPKM)}{Nuclear\ Expression\ (FPKM)})$$

40
41 Only detected genes are shown in plots of lncATLAS, with the exception of compartment-
42 specific genes in Plot 1. For this group of genes, CN-RCI value is not available in the plot and the bar
43 only indicates the tendency of the gene towards nucleus or cytoplasm. Colour and FPKM are shown
44 normally, to indicate the level of expression.
45
46 **Sub-cytoplasmic and sub-nuclear fractions for K562 (subN RCI, subC RCI)**
47 RNAseq data from subnuclear fractions (chromatin, nucleolus and nucleoplasm) and from sub-

cytoplasmic fractions (membrane and insoluble fraction) were retrieved and processed as above. These data are only available for K562 cells, and in the case of subnuclear samples correspond to total RNA (not PolyA+ selected RNA). Data were processed and RCI calculated as above, with the only differences being: (1) the RCI was calculated with reference to the nuclear fraction for sub-nuclear compartments, and cytoplasmic fraction for sub-cytoplasmic compartments; (2) for subnuclear compartments total RNA samples were used, instead of PolyA+.

**Database design**

LncATLAS is a relational database implemented in MySQL (http://www.mysql.com) and designed through the official MySQL WorkBench tool for Linux. The Entity-Relationship (ER) diagram extracted summarizes its structure (Supplementary Figure S1). The tables are hierarchically organized from general information of the samples to the expression value per gene that is stored in the expression table.

**Web-tool Implementation**

LncATLAS is constructed using the Shiny R package (version 0.13.2) (http://www.rstudio.com/shiny/). The database is connected to the application itself via the R package RMySQL (version 0.10.9). Other packages used in the implementation are the *ggplot2* package (version 2.1.0) and *dplyr* (version 0.5.0) used to build custom plots and manipulate the data.

## Acknowledgements

## References

Aoki, Kazuma et al. 2010. "A Thymus-Specific Noncoding RNA, Thy-ncR1, Is a Cytoplasmic Riboregulator of MFAP4 mRNA in Immature T-Cell Lines." *BMC Molecular Biology* 11(1):99.

Brown, C. J. et al. 1992. "The Human XIST Gene: Analysis of a 17 Kb Inactive X-Specific RNA That Contains Conserved Repeats and Is Highly Localized within the Nucleus." *Cell* 71(3):527–42.

Cabili, Moran N. et al. 2015. "Localization and Abundance Analysis of Human lncRNAs at Single-Cell and Single-Molecule Resolution." *Genome Biology* 16(1):20.

Carlevaro-Fita, Joana, Anisa Rahim, Roderic Guigó, Leah A. Vardy, and Rory Johnson. 2016. "Cytoplasmic Long Noncoding RNAs Are Frequently Bound to and Degraded at Ribosomes in Human Cells." *RNA (New York, N.Y.)* 22(6):867–82.

Cesana, Marcella et al. 2011. "A Long Noncoding RNA Controls Muscle Differentiation by Functioning as a Competing Endogenous RNA." *Cell* 147(2):358–69.

Chen, G. et al. 2013. "LncRNADisease: A Database for Long-Non-Coding RNA-Associated Diseases." *Nucleic Acids Research* 41(D1):D983–86.

Chen, J. et al. 2016. "The Role and Possible Mechanism of lncRNA U90926 in Modulating 3T3-L1 Preadipocyte Differentiation." *International Journal of Obesity*.

Chen, Ling-Ling. 2016. "Linking Long Noncoding RNA Localization and Function." *Trends in Biochemical Sciences*.

Clemson, C. M., J. A. McNeil, H. F. Willard, and J. B. Lawrence. 1996. "XIST RNA Paints the Inactive X Chromosome at Interphase: Evidence for a Novel RNA Involved in Nuclear/chromosome Structure." *The Journal of Cell Biology* 132(3):259–75.

Derrien, Thomas et al. 2012. "The GENCODE v7 Catalog of Human Long Noncoding RNAs: Analysis of Their Gene Structure, Evolution, and Expression." *Genome Research* 22(9):1775–89.

Djebali, Sarah et al. 2012. "Landscape of Transcription in Human Cells." *Nature* 489(7414):101–8.

Dobin, A. et al. 2013. "STAR: Ultrafast Universal RNA-Seq Aligner." *Bioinformatics* 29(1):15–21.

Dunagin, Margaret, Moran N. Cabili, John Rinn, and Arjun Raj. 2015. "Visualization of lncRNA by Single-Molecule Fluorescence In Situ Hybridization." Pp. 3–19 in.

Dunham, Ian et al. 2012. "An Integrated Encyclopedia of DNA Elements in the Human Genome." *Nature* 489(7414):57–74.

Guttman, Mitchell and John L. Rinn. 2012. "Modular Regulatory Principles of Large Non-Coding RNAs." *Nature* 482(7385):339–46.

Hansji, Herah et al. 2016. "ZFAS1: A Long Noncoding RNA Associated with Ribosomes in Breast Cancer Cells." *Biology Direct* 11(1):62.

Harrow, J. et al. 2012. "GENCODE: The Reference Human Genome Annotation for The ENCODE Project." *Genome Research* 22(9):1760–74.

van Heesch, Sebastiaan et al. 2014. "Extensive Localization of Long Noncoding RNAs to the Cytosol and Mono- and Polyribosomal Complexes." *Genome Biology* 15(1):R6.

Hutchinson, John N. et al. 2007. "A Screen for Nuclear Transcripts Identifies Two Linked Noncoding RNAs Associated with SC35 Splicing Domains." *BMC Genomics* 8(1):39.

Ishizuka, Akira, Yuko Hasegawa, Kentaro Ishida, Kaori Yanaka, and Shinichi Nakagawa. 2014. "Formation of Nuclear Bodies by the lncRNA Gomafu-Associating Proteins Celf3 and SF1." *Genes to Cells : Devoted to Molecular & Cellular Mechanisms* 19(9):704–21.

Kino, T., D. E. Hurt, T. Ichijo, N. Nader, and G. P. Chrousos. 2010. "Noncoding RNA Gas5 Is a Growth Arrest- and Starvation-Associated Repressor of the Glucocorticoid Receptor." *Science Signaling* 3(107):ra8-ra8.

Knowles, David G., Maik Röder, Angelika Merkel, and Roderic Guigó. 2013. "Grape RNA-Seq Analysis Pipeline Environment." *Bioinformatics (Oxford, England)* 29(5):614–21.

Lee, J. H. et al. 2014. "Highly Multiplexed Subcellular RNA Sequencing in Situ." *Science* 343(6177):1360–63.

Lee, Je Hyuk et al. 2015. "Fluorescent in Situ Sequencing (FISSEQ) of RNA for Gene Expression Profiling in Intact Cells and Tissues." *Nature Protocols* 10(3):442–58.

Li, Bo and Colin N. Dewey. 2011. "RSEM: Accurate Transcript Quantification from RNA-Seq Data

with or without a Reference Genome." *BMC Bioinformatics* 12(1):323.

Mondal, T., M. Rasmussen, G. K. Pandey, A. Isaksson, and C. Kanduri. 2010. "Characterization of the RNA Content of Chromatin." *Genome Research* 20(7):899–907.

Ounzain, Samir et al. 2015. "CARMEN, a Human Super Enhancer-Associated Long Noncoding RNA Controlling Cardiac Specification, Differentiation and Homeostasis." *Journal of Molecular and Cellular Cardiology* 89(Pt A):98–112.

Quek, Xiu Cheng et al. 2015. "lncRNAdb v2.0: Expanding the Reference Database for Functional Long Noncoding RNAs." *Nucleic Acids Research* 43(Database issue):D168-73.

Raj, Arjun, Patrick van den Bogaard, Scott A. Rifkin, Alexander van Oudenaarden, and Sanjay Tyagi. 2008. "Imaging Individual mRNA Molecules Using Multiple Singly Labeled Probes." *Nature Methods* 5(10):877–79.

Rinn, John L. et al. 2007. "Functional Demarcation of Active and Silent Chromatin Domains in Human HOX Loci by Noncoding RNAs." *Cell* 129(7):1311–23.

Schein, Aleks et al. 2016. "Identification of Antisense Long Noncoding RNAs That Function as SINEUPs in Human Cells." *Scientific Reports* 6:33605.

Tsai, M. C. et al. 2010. "Long Noncoding RNA as Modular Scaffold of Histone Modification Complexes." *Science* 329(5992):689–93.

Ulitsky, Igor and David P. Bartel. 2013. "lincRNAs: Genomics, Evolution, and Mechanisms." *Cell* 154(1):26–46.

Wang, Yaming, Wei Zhu, and David E. Levy. 2006. "Nuclear and Cytoplasmic mRNA Quantification by SYBR Green Based Real-Time RT-PCR." *Methods* 39(4):356–62.

Whitehead, Joanne, Gaurav Kumar Pandey, and Chandrasekhar Kanduri. 2009. "Regulation of the Mammalian Epigenome by Long Noncoding RNAs." *Biochimica et Biophysica Acta (BBA) - General Subjects* 1790(9):936–47.

Willingham, A. T. et al. 2005. "A Strategy for Probing the Function of Noncoding RNAs Finds a Repressor of NFAT." *Science (New York, N.Y.)* 309(5740):1570–73.

Yoon, Je-Hyun et al. 2012. "LincRNA-p21 Suppresses Target mRNA Translation." *Molecular Cell* 47(4):648–55.

Zhang, Ting et al. 2016. "RNALocate: A Resource for RNA Subcellular Localizations." *Nucleic Acids Research* 35(D1):D810–14.

Zhao, J., B. K. Sun, J. A. Erwin, J. J. Song, and J. T. Lee. 2008. "Polycomb Proteins Targeted by a Short Repeat RNA to the Mouse X Chromosome." *Science* 322(5902):750–56.

Zucchelli, Silvia, Laura Patrucco, Francesca Persichetti, Stefano Gustincich, and Diego Cotella. 2016. "Engineering Translation in Mammalian Cell Factories to Increase Protein Yield: The Unexpected Use of Long Non-Coding SINEUP RNAs." *Computational and Structural Biotechnology Journal* 14:404–10.
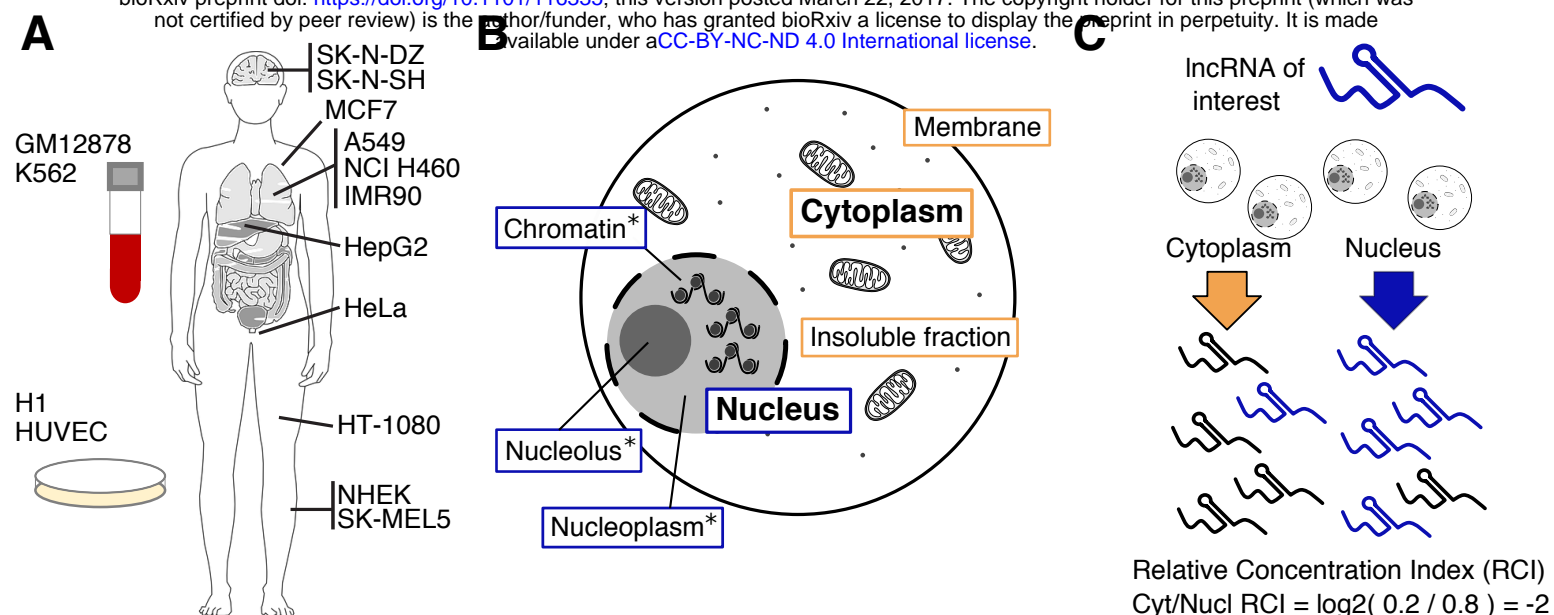
1    **Figure Legends**

2

3    **Figure 1.** Overview of lncATLAS data. A) Cell lines available in lncATLAS, indicating their

4    approximate origin. B) Cellular compartments available. * Compartments with only total RNA samples

5    available. C) The relative concentration index (RCI), in this case calculated for the cytoplasm and

6    nucleus (CN-RCI). The RCI can be thought of as the log-ratio, between two compartments, of the

7    concentration of a given RNA molecule per unit mass of RNA. D) Overview of the cell lines and cellular

8    compartments available for lncRNA RCI calculations. "Number of genes" indicates the number of

9    lncRNA for which the RCI could be calculated in the corresponding cell line (see Materials and

10   Methods for details). Sub-C RCI and Sub-N RCI correspond to sub-cytoplasmic RCI and sub-nuclear

11   RCI respectively.

12

13   **Figure 2.** Summary of Cytoplasmic/Nuclear detected genes by lncATLAS.  A) Number

14   of genes analysed in lncATLAS. (i) Detected genes: present a reliable expression value in both

15   compartments, cytoplasm and nucleus. Other categories comprise (ii) Compartment specific genes:

16   present a reliable non-zero value in one compartment and zero in the other, (iii) unreliable genes: genes

17   that did not pass the 2fold cutoff (see Materials and Methods, data processing) and (iv) silent genes:

18   not expressed in any compartment (see Materials and Methods, CN-RCI). *No biological

19   replicates were available for cytoplasm nor nucleus. **No biological replicates were available

20   for cytoplasm. B) Histogram showing how many genes are detected in a determined number

21   of cell lines. LncRNA genes in blue and mRNAs in red. C) Coverage by lncATLAS of widely-

22   used, manually-curated lncRNA databases, lncRNAdb (Quek et al. 2015) and lncRNAdisease (Chen et

23   al. 2013).  The new RNAlocate database (Zhang et al. 2016) is also shown for comparison. Note that

24   these databases hold a mixture of GENCODE annotated and non-GENCODE annotated lncRNAs. The

25   barplot displays the total number of human lncRNA genes in each database (whole bar). Bars are

26   colored to represent, from the total number of human lncRNAs in a database, the number of genes that:

27   (green) are displayed in lncATLAs (the percentage numbers indicate the proportion that this fraction

28   represents), (shaded red) are part of GENCODE annotation but are not detected in any of the 15 cell

29   lines, and therefore are not displayed in lncATLAs, (red) are not present in GENCODE annotation and

30   could not be considered in our database.

31

32   **Figure 3.** Subcellular localisation plots displayed by lncATLAS. MALAT1 and DANCR genes

33   are selected as examples of nuclear and cytoplasmic lncRNAs, respectively. A) Bars representing CN-

34   RCI values for the selected genes across all cell lines. Expression values (FPKMs) for the genes of

35   interest are shown for both compartments (cytoplasm on top of the bar, nucleus on the bottom). Bars

36   are colored by their absolute nuclear expression. B) Boxplot showing CN-RCI values distribution of all

37   lncRNAs (blue) and mRNAs (orange) for each cell line ("n" indicates total number of genes, "m"

38   median of CN-RCI values for lncRNAs and mRNAs separately). LncRNAs of interest are located in

39   the distribution and a percentage indicates their percentile rank within the distribution of all lncRNAs

40   (ranks relative to lowest value). C) Same than in the previous plot but in this case distribution is shown

41   as a density plot and only for a particular cell type, HUVEC. Again, genes of interest are located in the

42   distribution and their percentile rank (relative to lowest value) and RCI are indicated. D) Contour plot

43   showing lncRNA and mRNA populations as a function of CN-RCI values and whole cell expression

44   (log10(FPKMs)). LncRNAs of interest are specifically displayed together with their whole cell
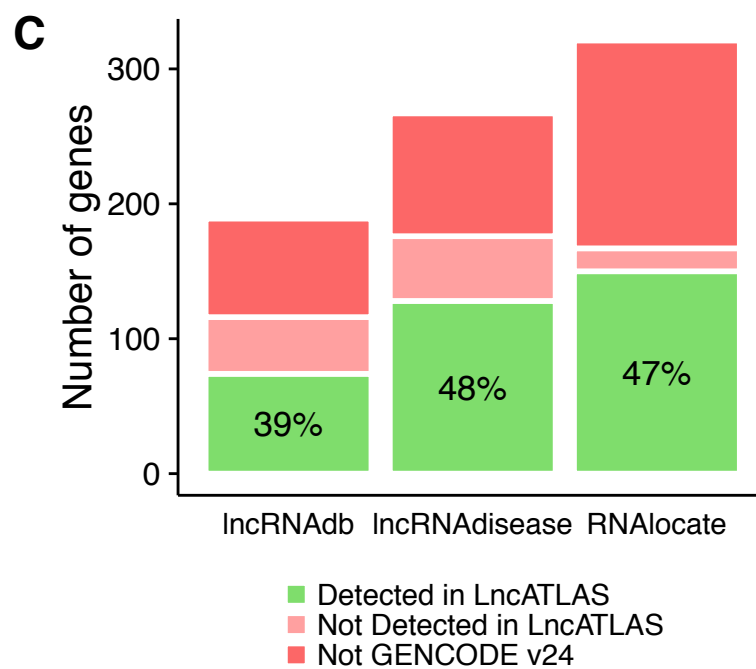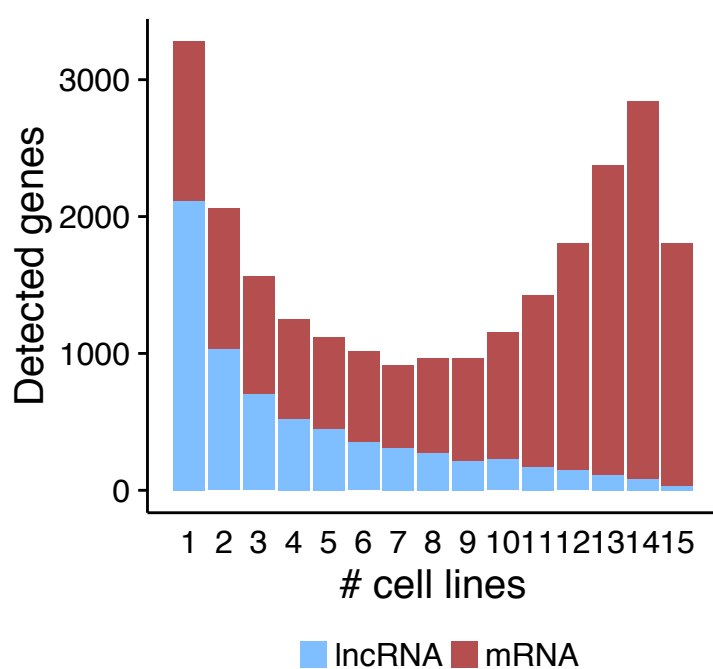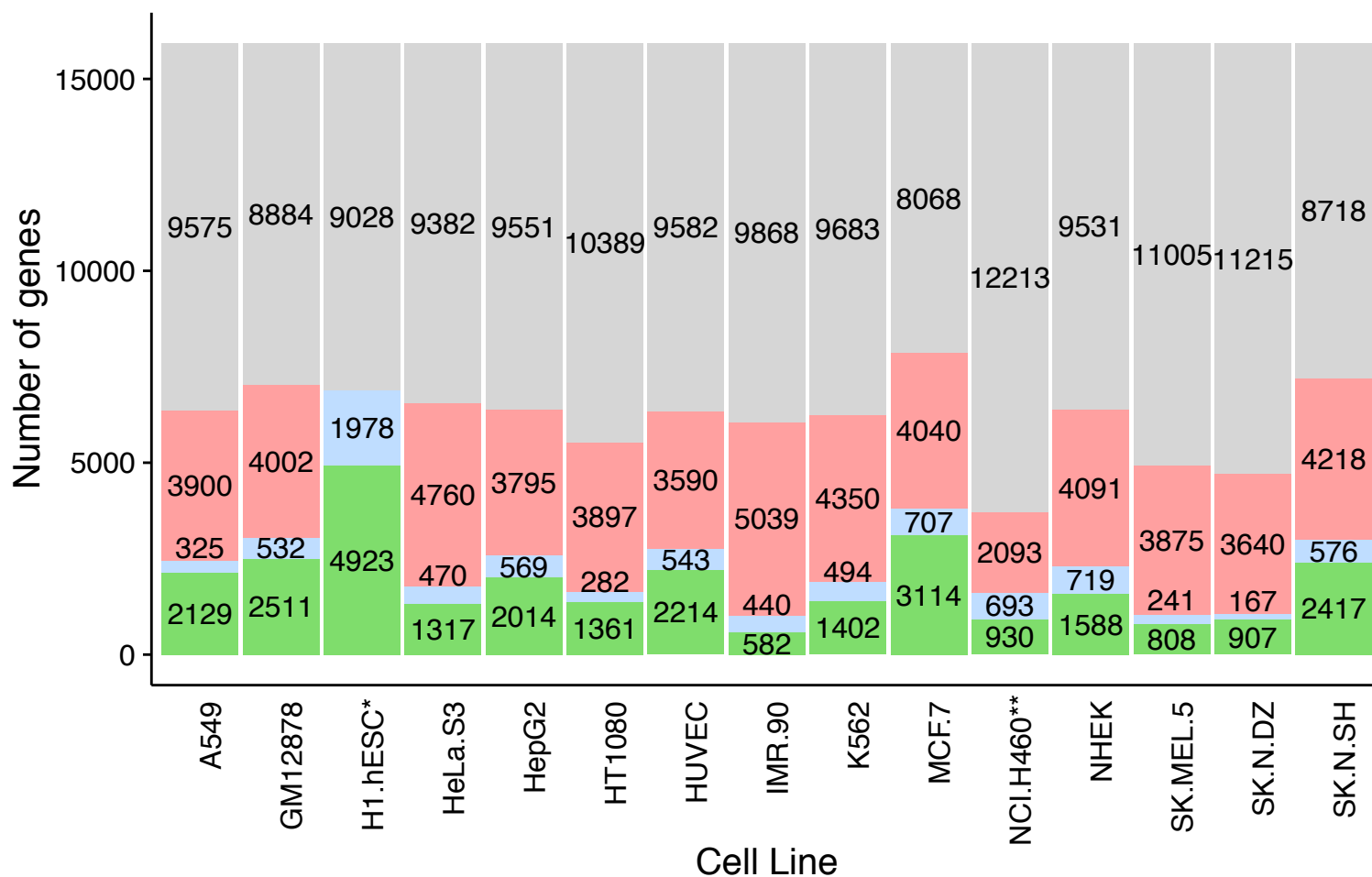
45   expression and RCI.

46

47   **Figure 4.** Sub-compartment data as displayed in LncATLAS. For K562 cells subnuclear and

1  subcytoplasmic fractions were available and RCI for all lncRNAs and mRNAs was computed (see
2  Materials and Methods). Boxplot shows the distribution of these values for each subnuclear and
3  subcytoplasmic compartment ("n" indicates total number of genes in a distribution). Percentile rank
4  (relative to lowest value) of each gene of interest is displayed to contextualize the relative enrichment
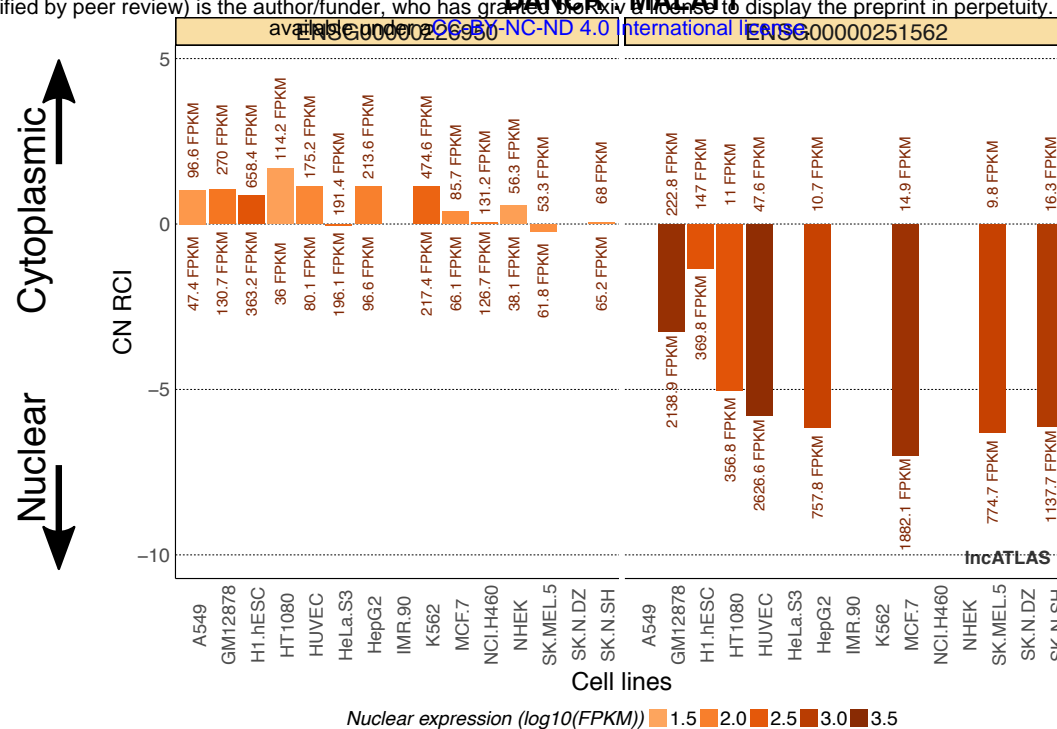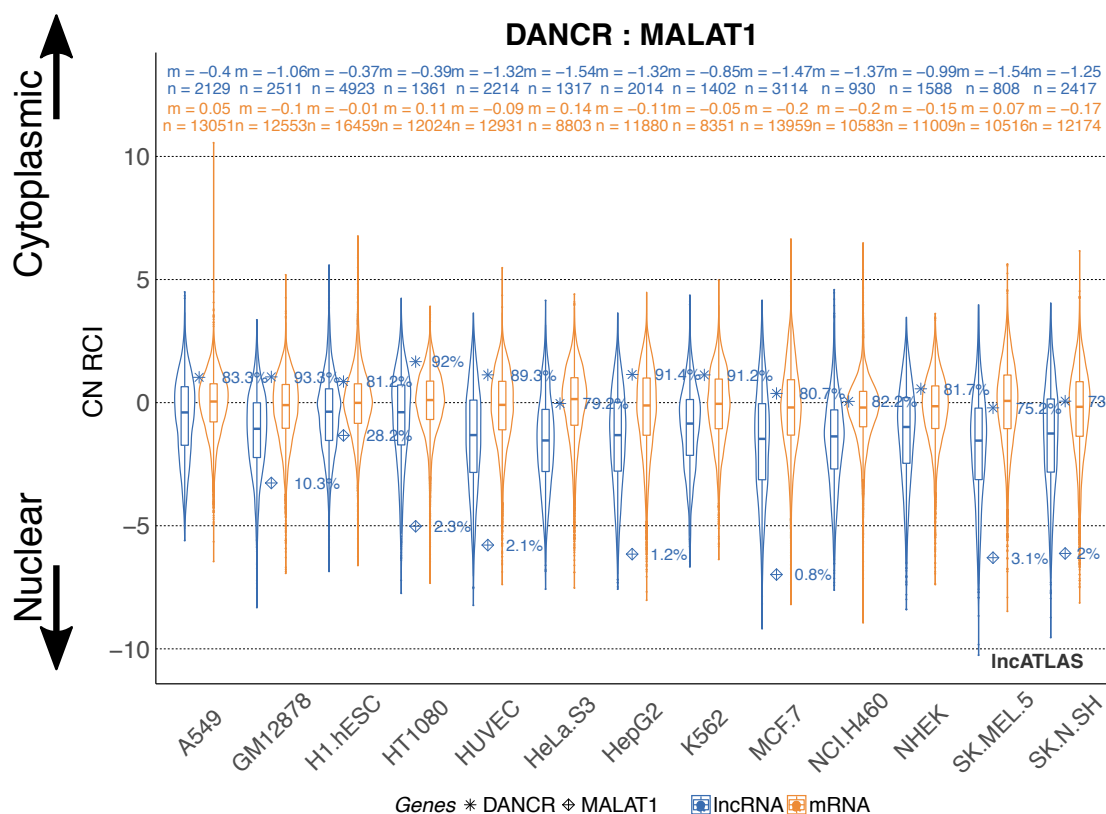5  of these genes in a subcompartment compared to the rest of lncRNAs and mRNAs.
6
7

**A**

GM12878
K562

SK-N-DZ
SK-N-SH
MCF7
A549
NCI H460
IMR90
HepG2
HeLa

H1
HUVEC

HT-1080

NHEK
SK-MEL5

**B**

Membrane

**Cytoplasm**

Chromatin*

Insoluble fraction

**Nucleus**

Nucleolus*

Nucleoplasm*

**C**

lncRNA of
interest

Cytoplasm        Nucleus

Relative Concentration Index (RCI)
Cyt/Nucl RCI = log2( 0.2 / 0.8 ) = -2

**D**

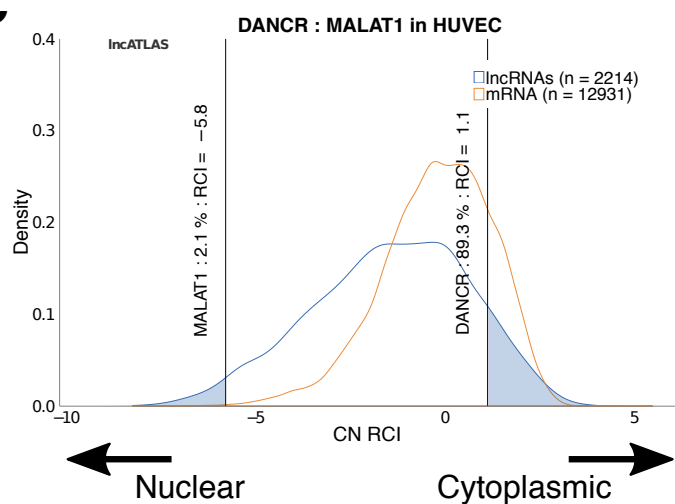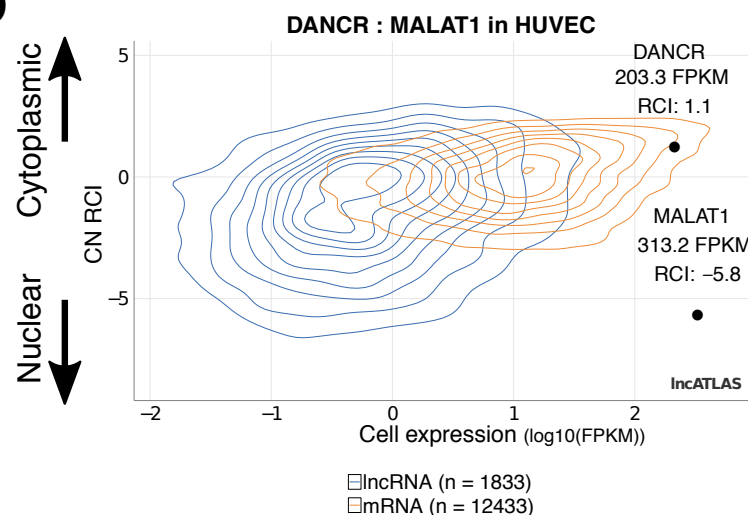| RCI | Cellular compartments | | Cell line | LncRNA genes |
|---|---|---|---|---|
| Cytoplasmic / Nuclear RCI (CN-RCI) | Cytoplasm / Nucleus | | K562 | 1402 |
| | | | A549 | 2129 |
| | | | GM12878 | 2511 |
| | | | H1.hESC | 4923 |
| | | | HeLa.S3 | 1317 |
| | | | HepG2 | 2014 |
| | | | HT1080 | 1361 |
| | | | HUVEC | 2214 |
| | | | IMR.90 | 582 |
| | | | MCF.7 | 3114 |
| | | | NCI.H460 | 930 |
| | | | NHEK | 1588 |
| | | | SK.MEL.5 | 808 |
| | | | SK.N.DZ | 907 |
| | | | SK.N.SH | 2417 |
| Sub - Compartments RCI | SubN- RCI | Chromatin / Nucleus | K562 | 1511 |
| | | Nucleolus / Nucleus | K562 | 2282 |
| | | Nucleoplasm / Nucleus | K562 | 2306 |
| | Sub C- RCI | Cell Membrane / Cytoplasm | K562 | 2324 |
| | | Insoluble Fraction / Cytoplasm | K562 | 2126 |

**A**



**B**



**C**

**A** DANCR : MALAT1

ENSG00000226950 — ENSG00000251562

**B** DANCR : MALAT1

**C** DANCR : MALAT1 in HUVEC

**D** DANCR : MALAT1 in HUVEC

DANCR : MALAT1