

1 **Cancer LncRNA Census 2 (CLC2): an enhanced resource reveals clinical features of**
2 **cancer lncRNAs**

3

4 Authors: Adrienne Vancura (1,2), Andrés Lanzós (1,2,3), Núria Bosch (1,2,3), Mònica Torres
5 (1,3), Alejandro Hionides Gutierrez (1), Simon Haefliger (1), Rory Johnson (1,3,4,5*)

6

7 Affiliations:

8 ¹Department of Medical Oncology, Inselspital, Bern University Hospital, University of Bern,
9 Bern, Switzerland.

10 ²Graduate School of Cellular and Biomedical Sciences, University of Bern, Bern, Switzerland.

11 ³Department for BioMedical Research, University of Bern, Bern, Switzerland.

12 ⁴School of Biology and Environmental Science, University College Dublin, Dublin, Ireland.

13 ⁵Conway Institute of Biomedical and Biomolecular Research, University College Dublin,
14 Dublin, Ireland.

15

16 * Correspondence: rory.johnson@dbmr.unibe.ch

17

18

19 **Keywords:** long noncoding RNA; long-noncoding RNA; lncRNA, cancer, curation, *DGCR5*,
20 *CARMEN*, *CARMN*, *LINC00570*

1 **Abstract**

2 Long noncoding RNAs play key roles in cancer and are at the vanguard of precision
3 therapeutic development. These efforts depend on large and high-confidence
4 collections of cancer lncRNAs. Here we present the Cancer LncRNA Census 2
5 (CLC2): at 492 cancer lncRNAs, it is 4-fold greater than its predecessor, without
6 compromising on strict criteria of confident functional / genetic roles and inclusion in
7 the GENCODE annotation scheme. This increase was enabled by leveraging high-
8 throughput transposon insertional mutagenesis (TIM) screening data, yielding 95
9 novel cancer lncRNAs. CLC2 makes a valuable addition to existing collections: it is
10 amongst the largest, holds the greatest number of unique genes, and carries functional
11 labels (oncogene / tumour suppressor). Analysis of this dataset reveals that cancer
12 lncRNAs are impacted by germline variants, somatic mutations, and changes in
13 expression consistent with inferred disease functions. Furthermore, we show how
14 clinical / genomic features can be used to vet prospective gene sets from high-
15 throughput sources. The combination of size and quality makes CLC2 a foundation for
16 precision medicine, demonstrating cancer lncRNAs' evolutionary and clinical
17 significance.

1 **Introduction**

2 Tumours arise and grow via genetic and non-genetic changes that give rise to
3 widespread alterations gene expression programmes (1–3). The numerous dysregulated
4 genes may encode classical protein-coding mRNAs or non-protein coding RNAs, but it is likely
5 that just a subset of these actually functionally contribute to pathogenic cellular hallmarks. The
6 identification of such functional cancer genes is critical for the development of targeted cancer
7 therapies, as well as emerging methods to identify additional cancer genes. For protein-coding
8 genes (pc-genes), datasets such as the Cancer Gene Census (CGC) collect and organise
9 comprehensive gene collections according to defined criteria, and has proven invaluable for
10 scientific research and drug discovery (4).

11 The past decade has witnessed the discovery of numerous non-protein-coding RNA
12 genes in mammalian cells (5, 6). The most numerous but poorly understood produce long
13 noncoding RNAs (lncRNAs), defined as transcripts >200 nt in length with no detectable
14 protein-coding potential (7). Although their molecular mechanisms are highly diverse, many
15 lncRNAs have been shown to interact with other RNA molecules, proteins and DNA by
16 structural and sequence-specific interactions (8, 9). Most lncRNAs are clade- and species-
17 specific, but a subset display deeper evolutionary conservation in their gene structure (10)
18 and a handful have been demonstrated to have functions that were conserved across millions
19 of years of evolution (10, 11). The numbers of known lncRNA genes in human have grown
20 rapidly, and present catalogues range from 18,000 to ~100,000 (12), however just a tiny
21 fraction have been functionally characterized (13–16). As lncRNAs likely represent a huge yet
22 poorly understood component of cellular networks, understanding the clinical and therapeutic
23 significance of these numerous novel genes is a key contemporary challenge.

24 lncRNAs have been implicated in molecular processes governing tumorigenesis (17).
25 lncRNAs may promote or oppose cancer hallmarks (18). This fact, coupled to the emergence
26 of potent *in vivo* inhibitors in the form of antisense oligonucleotides (ASOs) (19), has given
27 rise to serious interest in lncRNAs as drug targets in cancer by both academia and pharma
28 (17, 20–22).

1 Initially, cancer lncRNAs were discovered by classical functional genomics workflows
2 employing microarray or RNA-seq expression profiling (23, 24). More recently, CRISPR-based
3 functional screening (25) and bioinformatic predictions (26–28) have also emerged as
4 powerful tools for novel cancer gene discovery. To assess their accuracy, these approaches
5 require accurate benchmarks in the form of curated databases of known cancer lncRNAs.

6 Any discussion of lncRNAs and cancer requires careful terminology. Experimental
7 evidence suggest that for some lncRNAs, it is a DNA element within the gene, in addition to
8 or instead of the RNA transcript, which mediates downstream gene regulation (29–31). This
9 introduces the need for meticulous assessment of the basis of each lncRNA gene’s
10 functionality. Furthermore, it has been shown that lncRNAs can exert strong phenotypic
11 effects in one cell background, but none in another (32). In the context of tumours, this means
12 that amongst the large numbers of differentially expressed lncRNAs (24), just a fraction are
13 likely to functionally contribute to a relevant cellular phenotype or cancer hallmark (20, 33–
14 36). Such genes, termed here “functional cancer lncRNAs”, are the focus of this study.
15 Remaining changing genes are non-functional “bystanders”, which are largely irrelevant in
16 understanding or inhibiting the molecular processes causing cancer and highlight the
17 importance of not assessing functionality evidence simply by expressional changes.

18 There are a number of excellent databases of cancer-associated lncRNAs:
19 lncRNADisease (37), CRlncRNA (38), EVLncRNAs (39) and Lnc2Cancer (40). These
20 principally employ labour-intensive manual curation, and rely extensively on differential
21 expression to identify candidates. On the other hand, these databases have not yet taken
22 advantage of recent high-confidence sources of functional cancer lncRNAs, such as high-
23 throughput functional screens (25, 41). For these reasons, existing annotations likely contain
24 unknown numbers of bystander lncRNAs, while omitting large numbers of *bona fide* functional
25 cancer lncRNAs. Thus, studies requiring high-confidence gene sets, including benchmarking
26 or drug discovery, call for a database focussed exclusively on functional cancer lncRNAs.

27 Here we address this need through the creation of the Cancer lncRNA Census 2
28 (CLC2). It not only extends our previous CLC dataset by several fold (42), but more

- 1 importantly, CLC2 takes a major step forward methodologically, by implementing an
- 2 automated curation component that utilises functional evolutionary conservation for the first
- 3 time. Using this data, we present a comprehensive analysis of the genomic and clinical
- 4 features of cancer lncRNAs.

1 **Results**

2 **Integrative, semi-automated cataloguing of cancer lncRNAs**

3 We sought to develop an improved map of lncRNAs with functional roles in either
4 promoting or opposing cancer hallmarks or tumorigenesis. Such a map should prioritise
5 lncRNAs with genuine causative roles, and exclude false-positive “bystanders”: genes whose
6 expression changes but play no functional role.

7 We began with conventional manual curation of lncRNAs from the scientific literature,
8 covering the period from January 2017 (directly after the end of the first CLC (42)) to the end
9 of December 2018. We continued to use stringent criteria for defining cancer lncRNAs: genes
10 must be annotated in GENCODE (here version 28), and cancer function must be
11 demonstrated either by functional *in vitro* or *in vivo* experiments, or germline or somatic
12 mutational evidence (see Methods) (Figure 1A). Altogether we collected 253 novel lncRNAs
13 in this way, which added to the original CLC amounts to 375 lncRNAs, hereafter denoted as
14 “literature lncRNAs” (Figure 1A).

15 We recently showed that some literature-curated lncRNAs were also targeted by
16 previously-overlooked mutations in published transposon insertion mutagenesis (TIM)
17 screens (42). We hypothesised that this insight could be extended to identify novel functional
18 cancer lncRNAs. Thus we developed a pipeline to automatically identify human lncRNAs by
19 orthology to a collection of TIM hits in mouse (41). In this way 123 lncRNAs were detected, of
20 which 102 were not already in the literature set. These were added to the CLC2, henceforth
21 denoted as “mutagenesis lncRNAs” (Figure 1B). This analysis is discussed in more detail in
22 the next section.

23 Pooled functional screens based on CRISPR-Cas9 loss-of-function have recently
24 emerged as a powerful means of identifying function cancer lncRNAs (25). However there has
25 been relatively little validation of the hits from such screens, and it is possible that they contain
26 substantial false positives (43, 44). Amongst the few datasets presently available, the most
27 comprehensive comes from a CRISPR-inhibition (CRISPRi) screen of ~16,000 lncRNAs in
28 seven human cell lines, with proliferation as a readout (45). Of the 499 hits identified, 322 are

1 annotated by GENCODE and hence could potentially be included in CLC2. These are
2 moderately enriched for known cancer lncRNAs from the literature search (Figure 1C). That
3 study independently validated 21 GENCODE-annotated hits, of which four (19%) were already
4 mentioned in the literature, and two (10%) were detected by TIM above. Given the uncertainty
5 over the true-positive rates of unvalidated screen hits, we opted for a conservative approach
6 and included the remaining 15 novel and independently-validated lncRNAs from this study
7 (“CRISPRi lncRNAs”) (Figure 1C).

8 Altogether, the resulting CLC2 set comprises 492 unique lncRNA genes, representing
9 a 4.0-fold increase over its predecessor. The entire CLC2 dataset is available in
10 Supplementary Table 1 and 2. Importantly, the dataset is fully annotated with evidence
11 information, affording users complete control over the particular subsets of lncRNAs
12 (literature, mutagenesis, CRISPRi) that they wish to include in their analyses.

13

14 **Automated annotation of human cancer lncRNAs via functional conservation**

15 We recently showed that transposon insertional mutagenesis (TIM) screens identify
16 cancer lncRNAs in mouse (42, 46), and that some of these overlapped previously-known
17 human cancer lncRNAs (Figure 2A). TIM screens identify “common insertion sites” (CIS),
18 where multiple transposon insertions at a particular genomic location have given rise to a
19 tumour, thereby implicating the underlying gene as an oncogene or tumour suppressor.

20 We here extend this strategy to identify new functional cancer lncRNAs, by developing
21 a new pipeline called CLIO-TIM (cancer lncRNA identification by orthology to TIM). Briefly,
22 CLIO-TIM uses chain alignments (47) to map mouse CIS to orthologous regions of the human
23 genome, and then identifies the most likely gene target (see Methods) (Figure 2B) (SUPP FIG
24 1B). Available CIS maps are based on a variety of identification methods, resulting in CIS with
25 a range of sizes, from 1 bp upwards. We opted to remove our previously conservative size
26 criterion (CIS = 1 bp), to now consider elements of any size resulting in 26,345 CIS (compared
27 to 2,806 previously (42)) (SUPP FIG 1A). This yields a 3-fold increase in sensitivity for true-
28 positive CGC genes (72% compared to 26.4% previously (42)) (SUPP FIG 1D).

1 Based on this expanded dataset, CLIO-TIM identified 16,430 orthologous regions in
2 human (hCIS) (Figure 2B) (SUPP FIG 1A). Altogether, 123 lncRNAs and 9,295 pc-genes were
3 identified as potential cancer genes. An example is the human-mouse orthologous lncRNA
4 locus shown in Figure 2B, comprising *Gm36495* in mouse and *LINC00570* in human. A CIS
5 lies upstream of the mouse gene's TSS, mapping to the first intron of the human orthologue.
6 *LINC00570* is an alternative identifier for ncRNA-a5 *cis*-acting lncRNA identified by Orom et
7 al. (48), that has not previously been associated with cancer or cell growth.

8 We expected that hCIS regions are enriched in known cancer genes. Consistent with
9 this, the 698 pc-genes from the COSMIC Cancer Gene Census (CGC) (4) (red in SUPP FIG
10 1D) are 155-fold enriched with hCIS over intergenic regions (light grey). Turning to lncRNAs,
11 the 375 literature lncRNAs are 19.5-fold enriched, supporting their disease relevance (Figure
12 2C). Thus, CLIO-TIM predictions are enriched in genuine protein-coding and lncRNA
13 functional cancer genes. Supporting its accuracy, the overall numbers of genes implicated by
14 CLIO-TIM agree with independent analysis in the CCGD database (SUPP FIG 1C).

15 An additional 209 hCIS fall in intergenic regions that are neither part of pc-genes or
16 lncRNAs, leading us to ask whether some may affect lncRNAs that are not annotated by
17 GENCODE (Figure 2C). To test this, we utilised the large set of cancer-associated lncRNAs
18 from miTranscriptome (24). 186 hCIS intersect 2167 miTranscriptome genes, making these
19 potentially novel non-annotated transcripts involved in cancer. Nevertheless, simulations
20 indicated that this rate of overlap was no greater than expected by random chance (see
21 Methods), making it unlikely that substantial numbers of undiscovered cancer lncRNAs remain
22 to be discovered in intergenic regions, at least with the datasets used here (SUPP FIG 1E).

23 In addition to known cancer lncRNAs, CLIO-TIM identifies 102 lncRNAs not previously
24 linked to cancer (FIG 2C, dark grey) with a 3.8-fold enrichment of insertions over intergenic
25 genome. As will be shown below, these lncRNAs bear clinical and genomic features of
26 functional cancer genes, and hence we decided to include them in CLC2. It should be noted,
27 however, that these "mutagenesis" lncRNAs are labelled and hence may be removed by end
28 users, as desired.

1 To experimentally test the principal that human orthologues of mouse cancer genes
2 have a conserved function, we selected *LINC00570*, identified by CLIO-TIM but never
3 previously been linked to cancer or cell proliferation. We asked whether *LINC00570* promotes
4 cell growth in transformed cells. We used RNA-sequencing data to search for cell models
5 where *LINC00570* is expressed, and identified robust expression in cervical carcinoma HeLa
6 cells (SUPP FIG 2A) and to a lesser extent in HCT116 colon carcinoma cells (SUPP FIG 2A).
7 We designed three distinct antisense oligonucleotides (ASOs) targeting the *LINC00570* intron
8 2 and 3 and exon 3 of the short isoform (intronic targeting ASOs are known to have
9 degradation efficiency comparable to exonic ones (49, 50)). Transfection of these ASOs led
10 to strong and reproducible decreases in steady state RNA levels in HeLa cells (Figure 2C).
11 This resulted in significant decreases in cell proliferation rates (Figure 2D, SUPP FIG 2B). We
12 observed a similar effect through CRISPRi-mediated inhibition of gene transcription by two
13 independent guide RNAs in HeLa (Figure 2D), and with the same ASOs in HCT116 cells
14 (SUPP FIG 2C and D). Therefore, *LINC00570* predicted by CLIO-TIM pipeline promotes
15 growth of human cancer cells, and is likely to have a deeply evolutionarily-conserved
16 tumorigenic activity.

17

18 **Enhanced cancer lncRNA catalogue integrating manual annotation, CRISPR screens** 19 **and functional conservation**

20 We next tallied the distinct lncRNAs in CLC3 and compared them with existing cancer
21 lncRNA databases. Figure 3A shows a breakdown of the composition of CLC2 in terms of
22 source, gene function and evidence strength. Where possible, the genes are given a functional
23 annotation, oncogene (og) or tumour suppressor (ts), according to evidence for promoting or
24 opposing cancer hallmarks. Oncogenes (n=275) quite considerably outnumber tumour
25 suppressors (n=95), although it is not clear whether this reflects genuine biology or an
26 ascertainment bias relating to scientific interest or technical issues. Smaller sets of lncRNAs
27 are associated with both functions, or have no functional information (those from TIM screens
28 where the functions of hits are ambiguous).

1 In terms of the quality of evidence sources, CLC2 represents a considerable
2 improvement over the original CLC. The fraction of lncRNAs with high quality *in vivo* evidence
3 (defined as functional validation in mouse models or mutagenesis analysis) now represent
4 66% compared to 24% previously (Figure 3A, SUPP FIG 3B). In total, the updated CLC2
5 comprises 33 cancer types (vs 29) and more lncRNAs are reported for every cancer subtype
6 (SUPP FIG 3A).

7 We were curious how much novelty the CLC2 gene set brought to the known universe
8 of cancer lncRNAs, as estimated from respected and longstanding cancer lncRNA collections
9 (Figure 3B). Considering only GENCODE-annotated genes, CLC2 with 492 is second only to
10 Lnc2Cancer (n= 512) in terms of size (40). However, Lnc2Cancer uses looser inclusion
11 criteria, including lncRNAs that are differentially expressed in tumours without additional
12 functional evidence. The three remaining databases are smaller (<200 genes). Importantly,
13 CLC2 holds the greatest number of unique genes, i.e. those that are not found in other
14 databases (n=225). These contain 118 literature-annotated cases, and also 95 novel
15 mutagenesis lncRNAs. Just 40 lncRNAs are common to all five databases (37–40). In
16 summary, CLC2 achieves large size without compromising on confidence, while also including
17 numerous new cancer lncRNAs for the first time.

18

19 **Unique genomic properties of CLC2 lncRNAs**

20 Cancer genes, both protein-coding and not, display elevated characteristics of
21 essentiality and clinical importance compared to other genes (4, 18, 51, 52). In order to confirm
22 their quality as a resource, we next asked whether CLC2 lncRNAs, and the mutagenesis
23 subset, display features expected for cancer genes.

24 In the following analyses, we compared gene features of selected lncRNAs to all other
25 lncRNAs. Comparison of gene sets can often be confounded by covariates such as gene
26 length or gene expression, therefore where appropriate we used control gene sets that were
27 matched to CLC2 by expression (denoted “nonCLCmatched”) (SUPP FIG 4A) and reported
28 findings correcting for gene length (SUPP FIG 4B).

1 Evolutionary conservation and steady-state expression are widely-used proxies for
2 gene function (53–55). Using the LnCompare tool (56), we find that the promoters and exons
3 of CLC2 genes display elevated evolutionary conservation in mammalian and vertebrate
4 phylogeny (Figure 4A) and elevated expression in cancer cell lines (Figure 4B). Strikingly we
5 observe a similar effect when considering the mutagenesis lncRNAs alone: their promoters
6 are significantly more conserved than expected by chance, and their expression is an order
7 of magnitude higher than other lncRNAs (Figure 4C and D).

8 Further, we found that CLC2 lncRNAs are enriched in repetitive elements (SUPP FIG
9 5A) and are more likely to house a small RNA gene, possibly indicating that some act as
10 precursor transcripts (SUPP FIG 5B). CLC2 lncRNAs also have non-random distributions of
11 gene biotypes, being depleted for intergenic class and enriched in divergent orientation to
12 other genes (SUPP FIG 5C).

13 In summary, CLC2 lncRNAs are significantly more conserved and more expressed
14 than expected by chance, pointing to biological function. Mutagenesis lncRNAs discovered by
15 the CLIO-TIM also carry these features, supports their designation as functional cancer
16 lncRNAs.

17

18 **CLC2 lncRNAs display consistent tumour expression changes and prognostic** 19 **properties**

20 Although gene expression was not a criterion for inclusion, we would expect that CLC2
21 lncRNAs' expression will be altered in tumours. Furthermore, one might expect that the nature
22 of this alteration should vary with disease function: oncogenes overexpressed, and tumour
23 suppressors downregulated.

24 To test this, we analysed TCGA RNA-sequencing (RNA-seq) data from 686 individual
25 tumours with matched healthy tissue (total n=1,372 analysed samples) in 20 different cancer
26 types (SUPP FIG 6A and B), and classified every gene as either differentially expressed (in at
27 least one cancer subtype, with log₂ Fold Change >1 and FDR <0.05) or not. We found that
28 CLC2 lncRNAs are 3.4-fold more likely to be differentially expressed compared to expression-

1 matched lncRNAs (Figure 5A). LncRNAs from each individual evidence source (literature,
2 mutagenesis, CRISPRi) behaved similarly, again supporting their inclusion. Similar effects
3 were found for pc-genes (SUPP FIG 7A).

4 Next, we asked whether the direction of expression change corresponds to gene
5 function. Indeed, oncogenes are enriched for overexpressed genes, whereas tumour
6 suppressors are enriched for down-regulated genes, supporting the functional labelling
7 scheme (Figure 5B).

8 Cancer genes' expression is often prognostic for patient survival. By correlating
9 expression to patient survival, we found that the expression of 392 CLC2 lncRNAs correlated
10 to patient survival in at least one cancer type (SUPP FIG 7C). When analysing the most
11 significant correlation of each CLC2 lncRNA compared to expression-matched nonCLC
12 lncRNAs, we find a weak but significant enrichment (SUPP FIG 7C), suggesting that CLC2
13 lncRNAs can be prognostic for patient survival.

14 In summary, gene expression characteristics of CLC2 genes, and subsets from
15 different evidence sources, support their functional labels as oncogenes and tumour
16 suppressors and is more broadly consistent with their important roles in tumorigenesis.

17

18 **CLC2 lncRNAs are enriched with cancer genetic mutations**

19 Cancer genes are characterized by a range of germline and somatic mutations that
20 lead to gain- or loss-of-function. It follows that cancer lncRNAs should be enriched with
21 germline single nucleotide polymorphisms that have been linked to cancer predisposition (57).
22 We obtained 5,331 germline cancer-associated single nucleotide polymorphisms (SNPs) from
23 genome-wide association studies (GWAS) (58) and mapped them to lncRNA and pc-gene
24 exons, calculating a density score that normalises for exon length (SUPP FIG 4B). As
25 expected, exons of known cancer pc-genes are >2-fold enriched in cancer SNPs (SUPP FIG
26 7B). When performing the same analysis with CLC2 lncRNAs, one observes an even more
27 pronounced enrichment of 4.0-fold when comparing to expression-matched nonCLC lncRNAs
28 (Figure 5C). Once again, the lncRNAs from each evidence source individually show

1 enrichment for cancer SNPs >2-fold (Figure 5C). Three mutagenesis lncRNAs, namely
2 *miR143HG/CARMN*, *LINC00511* and *LINC01488*, carry an exonic cancer SNP (Figure 5D).

3 Cancer genes are also frequently the subject of large-scale somatic mutations, or copy
4 number variants (CNVs). Using a collection of CNV data from LncVar (59), we calculated the
5 gene-span length-normalized coverage of lncRNAs by CNVs. CLC2 lncRNAs are enriched for
6 CNVs compared to all lncRNAs (Figure 5E).

7 All information of the lncRNAs in the CLC2 with the corresponding cancer function,
8 evidence level, analysis method and cancer types can be found in the Supplementary Table
9 1. The Supplementary Table 2 can be used to filter lncRNAs based on their reported cancer
10 associated functionalities.

11 In summary, CLC2 lncRNAs and their subsets display germline and somatic mutational
12 patterns consistent with known oncogenes and tumour suppressors.

1 **Discussion**

2 We have presented the Cancer LncRNA Census 2, an expanded collection of lncRNAs
3 with functional roles in cancer. CLC2 is distinguished from other resources by several key
4 features. All its constituent lncRNAs have strong evidence for functional cancer roles (and not
5 merely differential expression), providing for lowest possible false positive rates. All CLC2
6 lncRNAs are included in the gold-standard GENCODE annotation, permitting smooth
7 interoperability with almost all public genomics projects and resources (12). The majority of
8 CLC2 entries are accompanied by functional labels (oncogene / tumour suppressor), enabling
9 one to link function to other observable features. Finally, we utilise transposon insertional
10 mutagenesis (TIM) datasets for the first time to discover 102 “mutagenesis” lncRNAs, of which
11 95 are completely novel. In spite of strict inclusion criteria, CLC2 is amongst the largest
12 available cancer lncRNA collections. Most striking, is that it contains the greatest number of
13 “unique” lncRNAs, not found in other resources. Overall, CLC2 makes a valuable addition to
14 the present landscape of cancer lncRNA resources.

15 A key novelty of CLC2 is its use of automated gene curation based on functional
16 evolutionary conservation, as inferred from TIM. This responds to the challenge from the rapid
17 growth of scientific literature, which makes manual curation increasingly impractical. Other
18 high-throughput / automated methods like CRISPR pooled screening, text mining and
19 machine learning will also be important, although it will be necessary to vet the quality of such
20 predictions prior to inclusion. Here we showed one way approach for this, by assessing the
21 TIM gene set across a range of genomic and clinical features. The fact that the “mutagenesis”
22 lncRNA set display rates of (i) nucleotide conservation, (ii) expression, (iii) tumour differential
23 expression, (iv) germline cancer polymorphisms and (v) tumour mutations similar to that of
24 gold-standard literature curated lncRNAs, coupled to thorough experimental validation of one
25 novel prediction (*LINC00570*), is powerful support for TIM and functional evolutionary
26 conservation as means for new cancer lncRNA discovery.

27 It might be argued that hits from TIM sites could be false positives that act via DNA
28 elements (for example, enhancers) that, by coincidence, overlap a non-functional lncRNA.

1 While certainly likely to occur in some cases, it would nevertheless appear unlikely to explain
2 the majority, in light of the features listed above, plus the observation that TIM sites are highly
3 enriched in independently-validated literature-curated lncRNAs (which act via RNA) including
4 *NEAT1*, *LINC-PINT* and *PVT1* (42). In spite of this, we recognise that some colleagues may
5 ascribe lower confidence to novel “mutagenesis” lncRNAs in CLC2. For this reason, the CLC2
6 data table is organised to facilitate filtering by source, enabling users to extract only the 375
7 literature-supported cases, or indeed any other subset based on source, evidence or function
8 as desired.

9 Apart from its usefulness as a resource, this study has enabled some important
10 conceptual insights. Firstly, we have replicated our previous finding that cancer lncRNAs are
11 distinguished by signatures of functionality, as inferred from evolutionary nucleotide
12 conservation and expression. These features were originally linked to protein-coding cancer
13 genes (51, 52), but are also utilised as markers for lncRNA functionality (42, 60). Moreover,
14 we extended this approach to clinical features, by showing that curated cancer lncRNAs are
15 dramatically more likely to be differentially expressed in tumours, suffer copy number
16 alteration, or carry a germline predisposition SNP. In the latter case, this rate even exceeds
17 cancer driver protein-coding genes. We also could demonstrate that changes in gene
18 expression in tumours are linked to function: oncogenes tend to be overexpressed, while
19 tumour-suppressors tend to be repressed. Finally, the demonstration that cancer lncRNAs can
20 be predicted on the basis of orthology to a TIM hit in mouse, lends powerful support to the
21 notion that there is widespread functional evolutionary conservation of lncRNAs in networks
22 related to cell growth and transformation.

23 *LINC00570* is a new functional cancer lncRNA predicted by CLIO-TIM. The gene was
24 previously discovered by Orom and colleagues, as a *cis*-activating enhancer-like RNA named
25 *ncRNA-a5* (48). That and a subsequent study showed that perturbation by siRNA transfection
26 affects the expression of the nearby pc-gene *ROCK2* in HeLa. However, these studies did not
27 investigate the effect on cell proliferation. We here show by means of two independent

1 perturbations, that *LINC00570* promotes proliferation of HeLa and HCT116 cells. These
2 findings make *LINC00570* a potential therapeutic target for follow up.

3 Intriguingly, amongst the novel mutagenesis lncRNAs identified by CLIO-TIM are
4 genes previously linked to other diseases. *miR143HG/CARMEN1* (*CARMN*) was shown to
5 regulate cardiac specification and differentiation in mouse and human hearts (61). In addition
6 to being a TIM target, *CARMEN1* also contains a germline cancer SNP correlating to the risk
7 of developing lung cancer (62), adding further weight to the notion that it also plays a role in
8 oncogenesis. Similarly, *DGCR5*, is located in the DiGeorge critical locus and has been linked
9 to neurodevelopment and neurodegeneration (63), and was recently implicated as a tumour
10 suppressor in prostate cancer (64). These results raise the possibility that developmental
11 lncRNAs can also play roles in cancer.

12 In summary, CLC2 establishes a new benchmark for cancer lncRNA resources. We
13 hope this dataset will enable a wide range of studies, from bioinformatic identification of new
14 disease genes, to developing a new generation of cancer therapeutics with anti-lncRNA ASOs
15 (65).

1 **Material and Methods**

2 **Gene curation**

3 If not stated otherwise, GENCODE v28 gene IDs (gencode.v28.annotation.gtf) were used.

4 **Literature search.** PubMed was searched for publications linking lncRNA and cancer using
5 keywords: long noncoding RNA cancer, lncRNA cancer. Additional inclusion criteria consisted
6 of GENCODE annotation, reported cancer subtype and cancer functionality
7 (oncogene/tumour suppressor). The manual curation and assigning evidence levels to each
8 lncRNA was performed exactly as previously (42) and included reports until December 2018.

9 **CLIO-TIM.** From the CCGD website (<http://ccgd-starrlab.oit.umn.edu/about.php>, May 2018
10 (41)) a table with all CIS elements was downloaded. These mouse genomic regions
11 (mm10) were converted to homologous regions in the human genome assembly hg38 using
12 the LiftOver tool (<https://genome.ucsc.edu/cgi-bin/hgLiftOver>). Settings: original Genome was
13 Mouse GRCm38/mm10 to New Genome Human GRCh38/hg38, minMatch was 0.1 and
14 minBlocks 0.1. For insertion sites intersecting several lncRNA genes, all the genes were
15 reported. IntersectBed from bedtools was used to align human insertion sites to GENCODE
16 IDs by intersecting at least 1nt and assigned to protein-coding or lncRNA gene families.
17 Insertion sites aligning to protein-coding and lncRNA genes were always assigned to protein-
18 coding genes. If insertion sites overlap multiple ENSGs, all genes are reported. Insertion sites
19 not aligning to protein-coding or lncRNAs genes were added to the intergenic region.

20 CCGD human Entrez gene results were converted to GENCODE IDs using the “Entrez gene
21 ids” Metadata file from <https://www.genecodegenes.org/human/> to compare CLIO-TIM results
22 with CCGD results for each gene set.

23 **MiTranscriptome data for evaluating intergenic insertion sites.** The cancer associated
24 MiTranscriptome IDs (24) previously used in Bergada et al. (66) were intersected with
25 intergenic insertion sites using IntersectBed. With ShuffleBed the intergenic insertions were
26 randomly shuffled 1000x and assigned to MiTranscriptome IDs.

27 **CRISPRi.** We used the Supplementary Table 1 from the 2017 Liu et al. paper (45) to extract
28 ENST IDs and gene names which are then converted to GENCODE IDs to match each guide

1 (LH identifier in the screen). From Supplementary Table S4 from the 2017 Liu et al. paper
2 (Liu_et_al_aah7111-TableS4) (45) we extracted genes with “hit” (validated as a hit in the
3 screen), “LH” (unique identifiers correlating to a gene in the screen) and “lncRNA” (referring
4 to a lncRNA gene and to exclude lncRNA hits close to a protein-coding gene (“Neighbor hit”))
5 resulting in 499 hits. Of these, 322 hits contain a GENCODE IDs and were used for enrichment
6 analysis, tested by one-sided Fisher’s test.

7 We included n=21 CRISPRi genes to the CLC2 from the Supplementary Figure 8A from the
8 2017 Liu et al. paper (45) , the tested cancer cell line and the effect of the CRISPRi on the
9 growth phenotype (either promoting (tumor suppressor) or inhibiting (oncogene)) of each
10 lncRNA was reported.

11 **Cancer gene sets.** For downstream analysis protein-coding (pc) genes (GENCODE IDs) are
12 grouped in cancer-associated pc-genes (CGC genes) and non cancer-associated pc-genes
13 (nonCGC n=19,174). The TSV file containing the CGC data was downloaded from
14 <https://cancer.sanger.ac.uk/census> with 700 ENSGs with 698 ENSG IDs detected in
15 GENCODE v28 of which 696 are unique (CGC n=696). The same is done for lncRNAs, into
16 CLC2 (n=492) and nonCLC genes (n= 15,314).

17 **Matched expression analysis.** Based on an in house script used for Survival analysis
18 (section below), TCGA survival expression data for each GENCODE ID is reported and the
19 average FPKM across all tumor samples is calculated. The count distribution of nonCGC and
20 nonCLC gene expression to CGC and CLC2 expression, respectively, is matched using the
21 matchDistribution.pl script (<https://github.com/julienlag/matchDistribution>).

22 **Cancer lncRNA databases.** The tested databases were first filtered for lncRNAs in the
23 GENCODE v28 long noncoding annotation (n=15,767).

24 **Lnc2cancer** GENCODE IDs from datatable ([http://www.bio-](http://www.bio-bigdata.com/lnc2cancer/download.html)
25 [bigdata.com/lnc2cancer/download.html](http://www.bio-bigdata.com/lnc2cancer/download.html)) were evaluated (n=512) (40).

26 **CRlncRNA** gene names from (<http://crlnc.xtbq.ac.cn/download/>) were converted to
27 GENCODE IDs (n=146) (38).

1 **EVLncRNAs** gene names (<http://biophy.dzu.edu.cn/EVLncRNAs/>) were converted to
2 GENCODE IDs (n=187) (39).

3 **lncRNADisease** gene names from
4 (<http://www.rnanut.net/lncrnadisease/index.php/home/info/download>) and only cancer-
5 associated transcripts (carcinoma, lymphoma, cancer, leukemia, tumor, glioma, sarcoma,
6 blastoma, astrocytoma, melanoma, meningioma) were extracted. Names were converted to
7 GENCODE IDs (n=137) (37).

8

9 **Features of CLC2 genes**

10 **Genomic classification.** The genomic classification was performed as previously (42) using
11 an in house script ([https://github.com/gold-](https://github.com/gold-lab/shared_scripts/tree/master/lncRNA.annotator)
12 [lab/shared_scripts/tree/master/lncRNA.annotator](https://github.com/gold-lab/shared_scripts/tree/master/lncRNA.annotator)).

13 **Small RNA analysis.** For this analysis “snoRNA”, “snRNA”, “miRNA” and “miscRNA”
14 coordinates were extracted from GENCODE v28 annotation file and intersected with the
15 genomic region of the genes (intronic and exonic regions).

16 **Repeat elements.** In total 452 CLC2 lncRNAs compared to 1693 expression-matched
17 nonCLC lncRNAs using the LnCompare Categorical analysis
18 (<http://www.rnanut.net/lnccompare/>) (56).

19 **Feature analysis.** In total 452 CLC2 lncRNAs and 120 mutagenesis lncRNAs are compared
20 to the GENCODE v24 reference using LnCompare (<http://www.rnanut.net/lnccompare/>) (56).

21

22 **Cancer characteristic analysis**

23 **Differential gene expression analysis (DEA).** was performed using TCGA data and
24 TCGAbiolinks. Analysis was performed as reported in manual for matching tumor and normal
25 tissue samples using the HTseq analysis pipeline as described previously.

26 (<https://www.bioconductor.org/packages/devel/bioc/vignettes/TCGAbiolinks/inst/doc/analysis>
27 [.html](https://www.bioconductor.org/packages/devel/bioc/vignettes/TCGAbiolinks/inst/doc/analysis.html)) (67). For this analysis only matched samples were used and the TCGA data was
28 presorted for tumor tissue samples (TP with 01 in sample name) and solid tissue normal (NT

1 with 11 in sample name). Settings used for DEA analysis: $fdr.cut = 0.05$, $logFC.cut = 1$ for
2 DGE output between matched TP and NT samples for 20 cancer types. CLC2 cancer types
3 had to be converted to TCGA cancer types (Supp Fig 6A) Cancer types and number of
4 samples used in the analysis can be found in Supp Fig 6B. DEA enrichment analysis tested
5 with one-sided Fisher's test. For each CLC2 gene reported as true oncogene (n=275) or tumor
6 suppressor (n=95), hence where no double function is reported (n=22), the positive and
7 negative fold change (FC) values were counted and compared to expression-matched lncRNA
8 genes found in the DEA.

9

10 **Survival analysis.** An inhouse script for extracting TCGA survival data was used to generate
11 p values correlating to survival for each gene. Expression and clinical data from 33 cohorts
12 from TCGA with the "TCGAbiolinks" R package
13 (<https://bioconductor.org/packages/release/bioc/html/TCGAbiolinks.html>) were downloaded
14 (67). P value and Hazard ratio were calculated with the Cox proportional hazards regression
15 model from "Survival" R package ([https://cran.r-](https://cran.r-project.org/web/packages/survival/survival.pdf)
16 [project.org/web/packages/survival/survival.pdf](https://cran.r-project.org/web/packages/survival/survival.pdf)). All scripts were adapted from here
17 (<https://www.biostars.org/p/153013/>) and are available upon request. For downstream
18 analysis, only groups with at least 20 patient samples in high or low expression group were
19 used. The plot comprises only the most significant cancer survival p value per gene and was
20 assessed by the Komnogorow-Smirnow-Test (ks-test).

21

22 **Cancer-associated SNP analysis.** SNP data linked to tumor/cancer/tumour were extracted
23 from the GWAS page (<https://www.ebi.ac.uk/gwas/docs/file-downloads>) (n=5,331) and
24 intersected with the whole exon body of the genes. SNPs were intersected to the transcript
25 bed file and plotted per nt in each subset (SNP/nt y axis) and tested using one-sided Fisher's
26 test.

27 **CNV analysis.** Human CNV in lncRNAs downloaded from
28 <http://bioinfo.ibp.ac.cn/LncVar/download.php> (59). NONCODE IDs were converted to

1 GENCODE IDs using [NONCODEv5_hg38.lncAndGene.bed.gz](https://www.encodeproject.org/track/GENCODEv5_hg38.lncAndGene.bed.gz). CLC2 and nonCLC ENSGs
2 were matched to NONHSAT IDs with a significant pvalue (0.05, n=733) in the LncVAR table
3 and tested using one-sided Fisher's test.

4 **Code availability.** Custom code are available from the corresponding author upon request.

5

6 **In vitro validation**

7 **Cell culture.** HeLa and HCT116 were cultured on Dulbecco's Modified Eagles Medium
8 (DMEM) (Sigma-Aldrich, D5671) supplemented with 10% Fetal Bovine Serum (FBS)
9 (ThermoFisher Scientific, 10500064), 1% L-Glutamine (ThermoFisher Scientific, 25030024),
10 1% Penicillin-Streptomycin (ThermoFisher Scientific, 15140122). Cells were grown at 37°C
11 and 5% CO₂ and passaged every two days at 1:5 dilution.

12 **Generation of Cas9 stable cell lines.** HeLa cells were infected with lentivirus carrying the
13 Cas9-BFP (blue fluorescent protein) vector (Addgene 52962). HCT116 were transfected with
14 the same vector using Lipofectamine 2000 (ThermoFisher Scientific, 11668019). Both cell
15 types were selected with blasticidin (4ug/ml) for at least five days and selected for BFP-
16 positive cells twice by fluorescence activated cell sorting.

17 **CRISPR inhibition sgRNA pair design and cloning.** sgRNA pairs targeting *LINC00570*
18 were designed using GPP sgRNA designer (<https://portals.broadinstitute.org/gpp/>). The
19 sgRNA pairs were manually selected from the output list and cloned into the pGECKO
20 backbone (CRISPRi.1: 5' *GTTACTTCCAACGTACCATG* 3', CRISPRi.2: 5'
21 *CCTGTACCCCATGGTACGT* 3') (Addgene 78534; (68))

22 **Antisense LNA GapmeR design.** Antisense LNA GapmeR Control (5'
23 *AACACGTCTATACGC* 3') and three Antisense LNA GapmeR Standard targeting *LINC00570*
24 (LNA1: 5' *GGAAATTGCTCTGATG* 3', LNA2: 5' *GATTGGCATTGGGATA* 3', LNA3: 5'
25 *GAAGTGGCCTGAGAAA* 3') were designed and purchased at Qiagen.

26 **RT-qPCR.** For each time point total RNA was extracted and reverse transcribed (Promega).
27 Transcript levels of *LINC00570* (FP: 5' *TAGGAGTGCTGGAGACTGAG* 3', RP: 5'
28 *GTCGCCATCTTGGTTGTCTG* 3') and housekeeping gene *HPRT1* (FP: 5'

1 ATGACCAGTCAACAGGGGACAT 3', RP: 5' CAACACTTCGTGGGGTCCTTTTCA 3') were
2 measured using GoTaq qPCR Master Mix (Promega, A6001) on a TaqMan Viia 7 Real-Time
3 PCR System. Data were normalized using the $\Delta\Delta C_t$ method (69)).

4 **Gene-specific RT-PCR and cDNA amplification.** From the extracted total RNA, we
5 performed a gene specific reverse transcription using the reverse primers for *LINC00570* and
6 *HPRT1* to enrich for their cDNA. Presence or absence of transcript was detected by a regular
7 PCR using GoTaq® G2 DNA Polymerase (Promega, M7841) from 100ng cDNA and
8 visualized on an agarose gel.

9 **Viability assay.** HeLa and HCT116 cells were transfected with Antisense LNA GapmeRs at
10 a concentration of 50nM using Lipofectamine 2000 (Thermofisher) according to
11 manufacturer's protocol. One day after, transfected cells were plated in a white, flat 96-well
12 plate (3000 cells/well). Viability was measured in technical replicates using CellTiter-Glo 2D
13 Kit (Promega) following manufacturer's recommendations at 0, 24, 48, 72 hours after seeding.
14 Luminescence was detected with Tecan Reader Infinite 200. Statistical significance calculated
15 by t-test.

16 For CRISPR inhibition experiments, HeLa-Cas9 and HCT116-Cas9 cells were transfected
17 with control sgRNA plasmid and two *LINC00570* targeting plasmids. Cells were selected with
18 puromycin (2ug/ml) for 48h. Viability assay was performed as previously described.

19

20

1 **Figure Legends**

2 **Figure 1: Functional cancer lncRNAs from three sources are integrated in the CLC2.**

3 **A)** Literature curation with four criteria are used to define “literature lncRNAs”. **B)** Transposon
4 insertion mutagenesis screens identify “mutagenesis lncRNAs”. **C)** Validated hits from
5 CRISPRi proliferation screens are denoted “CRISPRi lncRNAs”. Statistical significance
6 calculated by one-sided Fisher’s test.

7

8 **Figure 2: The CLIO-TIM pipeline identifies human cancer lncRNAs via functional**
9 **evolutionary conservation.**

10 **A)** Overview of transposon insertional mutagenesis (TIM) method for identifying functional
11 cancer genes. Engineered transposons carry bidirectional cassettes capable of either blocking
12 or upregulating gene transcription, depending on orientation. Transposons are introduced into
13 a population of cells, where they integrate at random genomic sites. The cells are injected into
14 a mouse. In some cells, transposons will land in and perturb expression of a cancer gene
15 (either tumour suppressor or oncogene), giving rise to a tumour. DNA of tumour cells is
16 sequenced to identify the exact location of the transposon insertion. Clusters of such insertions
17 are termed Common Insertion Sites (CIS). **B)** (Left) Schematic of the CLIO-TIM pipeline used
18 here to identify human cancer genes using mouse CIS. (Right) An example of a CLIO-TIM
19 predicted cancer lncRNA. **C)** The density of hCIS sites, normalised by gene length, in indicated
20 classes of lncRNAs. Statistical significance calculated by one-sided Fisher’s test. **D)** Upper
21 panels: Expression of *LINC00570* RNA in response to inhibition by CRISPRi (left) or ASOs
22 (right). Lower panels: Measured populations of the same cells over time. Statistical
23 significance calculated by Student’s *t*-test.

24

25 **Figure 3: An overview of the CLC2 database and comparison with other lncRNA**
26 **databases.**

27 **A)** The CLC2 database broken down by source, function and evidence type. **B)** Comparison
28 of CLC2 to other leading cancer lncRNA databases.

1 **Figure 4: Features of functionality in CLC2 and mutagenesis lncRNAs.**

2 In each panel, two gene sets are compared: the test set (either all CLC2 genes, or
3 mutagenesis genes alone), and the set of all other lncRNAs (GENCODE v24). Y-axis: Log₂
4 fold difference between the means of gene sets. X-axis: false-discovery rate adjusted
5 significance, calculated by Wilcoxon test. **A)** Evolutionary conservation for all CLC2,
6 calculated by PhastCons. **B)** Expression of all CLC2 in cell lines. **C)** Evolutionary conservation
7 for mutagenesis lncRNAs, calculated by PhastCons. **D)** Expression of mutagenesis lncRNAs
8 in cell lines. For (A) and (C), “Promoter mean” and “Exon mean” indicate mean PhastCons
9 scores (7-vertebrate alignment) for those features, while “Exon-coverage” indicates percent
10 coverage by PhastCons elements. Promoters are defined as a window of 200 nt centered on
11 the transcription start site.

12

13 **Figure 5: Clinical features of CLC2 lncRNAs.**

14 **A)** The percent of indicated genes that are significantly differentially expressed in at least one
15 tumour type from the TCGA. Statistical significance calculated by one-sided Fisher’s test. **B)**
16 Here, only differentially expressed genes from (A) are considered. lncRNAs with both tumour
17 suppressor and oncogene labels are excluded. Remaining lncRNAs are divided by those that
18 are up- or down-regulated (positive or negative fold change). Statistical significance calculated
19 by one-sided Fisher’s test. **C)** The density of germline cancer-associated SNPs is displayed.
20 Only SNPs falling in gene exons are counted, and are normalised to the total length of those
21 exons. Statistical significance calculated by one-sided Fisher’s test. **D)** Examples of
22 mutagenesis lncRNAs with an exonic cancer SNP. **E)** Length-normalised overlap rate of copy
23 number variants (CNVs) in lncRNA gene span. Statistical significance calculated by one-sided
24 Fisher’s test.

25

26

27

28

1 **Supplementary Data**

2 **Supplementary Table 1:** Excel table with all CLC2 and concertype and evidence level

3 **Supplementary Table 2:** Excel table with all CLC2 ENSG with cancer functionality

4

5 **Supplementary Figures**

6 **SUPP. Figure 1:** Insertion analysis

7 **A)** All insertion sizes after Liftover compared to GENCODE v28 gene length. Number of input
8 CIS elements in mouse compared to hCIS elements in human after Liftover. **B)** Assign genes
9 to GENCODE v28 genes and gene families. C) CCGD reported genes (dark) and CLIO-TIM
10 reported genes (fade) for each gene class. **D)** Genes with insertion categorized in gene types.
11 Statistical significance calculated by one-sided Fisher's test. **E)** Assign intergenic regions to
12 MiTranscriptome IDs and compare to shuffled hCIS overlaid with MiTranscriptome IDs.
13 Example of one insertion site with MiTranscriptome ID.

14

15 **SUPP. Figure 2:** *LINC00570* insertion candidate characteristics.

16 **A)** ENCODE expression data of *LINC00570* in HeLa (blue) and HCT116 cells (black). **B)** Cell
17 proliferation of HeLa cells treated with ASO negative control and ASO 3 at Day 1, 2 and 3. **C)**
18 Proliferation assay for HCT116 cells with ASO control and the three ASO targeting
19 *LINC00570*. Statistical significance calculated by Student's *t*-test. **D)** RT-PCR of *LINC00570*
20 and *HPRT1* from HCT116 cells to check for expression.

21

22 **SUPP. Figure 3:** comparison of first CLC and CLC2.

23 **A)** CLC2 genes are detected in 33 cancer types and compared to 29 in the first CLC. CLC
24 reported n=122 literature lncRNAs whereas CLC2 comprises 492 genes from 3 different
25 analysis. **B)** Comparing evidence levels of genes from the initial CLC with the CLC2.

26

27

1 **SUPP. Figure 4:** CLC2 expression and gene length bias.

2 **A)** CLC2 genes (turquoise) are higher expressed than nonCLC genes (grey), same for CGC
3 genes (red) compared to nonCGC genes (orange). Expression-matched CLC2 (blue) and
4 CGC (yellow) were generated and match the expression of the CLC2 and CGC, respectively.

5 **B)** CLC2 genes (turquoise) with increased exon and whole gene body length when compared
6 to expression-matched (blue) and all other lncRNAs (grey).

7

8 **SUPP. Figure 5:** CLC2 gene characteristics.

9 **A)** CLC2 genes are enriched for $\frac{2}{3}$ of the analyzed repeat element families when compared
10 to expression-matched nonCLC genes. Statistical significance calculated by hypergeometric
11 test (highly significant ****= <0.0001). **B)** CGC and CLC2 genes are enriched for small RNAs
12 compared to expression-matched nonCGC and nonCLC, respectively. In the bar graph we
13 report the fraction of genes of each dataset with (dark color) or without (light color) small RNA
14 encoded in the genomic region. Statistical significance calculated by one-sided Fisher's test.

15 **C)** Genomic classification of CLC2, expression-matched nonCLC and nonCLC genes.
16 Statistical significance calculated by two-sided Fisher's test (*= <0.05).

17

18 **SUPP. Figure 6:** TCGA cancer types for differential expression analysis.

19 **A)** CLC2 cancer types corresponding to TCGA cancer types. **B)** Samples for each TCGA
20 cancer type analyzed for differential expression analysis.

21

22 **SUPP. Figure 7:** Cancer characteristics for all analyzed gene types.

23 **A)** Differentially expressed genes enriched in cancer-associated gene families (CGC and
24 CLC2). Statistical significance calculated by one-sided Fisher's test. **B)** exonic cancer SNPs
25 enriched in cancer-associated gene families (CGC and CLC2). Statistical significance
26 calculated by one-sided Fisher's test. **C)** survival analysis comparing most significant p-value
27 for each lncRNA in the CLC2 compared to expression-matched lncRNAs. Statistical
28 significance calculated by ks-test.

1 **Acknowledgements**

2 We gratefully acknowledge administrative support from Ana Radovanovic and Silvia
3 Roesselet (DBMR, University of Bern). We also acknowledge Joana Carlevaro-Fita (EPFL,
4 Lausanne) and Judith Bergada (University of Zurich) for the helpful advice and discussions
5 as well as Roberta Esposito, Panagiotis Chouvardas, Hugo Guillen Ramirez and the other
6 members of the Laboratory for Genomics of LncRNA and Disease for their valuable input.

7

8 **Author contribution**

9 RJ conceived the project. RJ, AV, AH performed manual annotation of CLC2. AV performed
10 the feature analysis, evolutionary analysis, mutation analysis, differential expression, GWAS
11 SNP, CNV analysis and data integration. AL performed survival analysis. NB performed the
12 ASO and CRISPRi KD experiments. AV, NB and MT performed the qPCR experiments. RJ,
13 AV, AL, NB, MT and SH drafted the manuscript and prepared the figures and supplementary
14 material. All authors read and approved the final draft.

15

16 **Conflict of interest**

17 The authors declare that they have no competing interests.

18

19 **The Paper Explained**

20 Problem: Cancer is one of the leading causes of death worldwide. The development of
21 effective therapies depends on creating collections of known cancer genes. These can
22 comprise not only conventional protein coding genes, but also more recently discovered genes
23 like long noncoding RNAs (lncRNAs). LncRNAs are considered highly promising therapeutic
24 targets, however the relatively poor state of knowledge, and the lack of high quality cancer
25 lncRNA collections, represents a significant hurdle to developing lncRNA therapies.

26 Results: To address the need for collections of cancer lncRNAs, we have developed the
27 Cancer lncRNA Census 2 (CLC2). CLC2 consists of 492 cancer lncRNAs functionally
28 validated in 33 cancer subtypes. CLC2 is the first catalogue to incorporate automatic screen

1 data from mice, and is shown to be superior to existing collections across several criteria. We
2 show that CLC2 lncRNAs enriched for cancer associated mutations and tend to be
3 differentially expressed in tumours.

4 Impact: CLC2 is a critical resource for future development of cancer therapies targeting
5 lncRNAs. Analysis of these genes has provided new insights into their biological and clinical
6 properties.

7

8

9 **Ethics approval and consent to participate**

10 Not applicable.

11

12 **Consent for publication**

13 Not applicable.

14

15 **Availability of data and materials**

16 Information on CIS elements for mouse and human lncRNAs reported in this publication are
17 available in the Supplementary Table 1 and the code is available from the corresponding
18 author on request.

19

20

21 **Funding**

22 This work was funded by the Swiss National Science Foundation through the National Center
23 of Competence in Research (NCCR) "RNA & Disease", by the Medical Faculty of the
24 University and University Hospital of Bern, by the Helmut Horten Stiftung and Krebsliga
25 Schweiz (4534-08-2018).

1 References

- 2 1. Hanahan,D. and Weinberg,R.A. (2011) Hallmarks of Cancer: The Next Generation. *Cell*,
- 3 **144**, 646–674.
- 4 2. Yates,L.R. and Campbell,P.J. (2012) Evolution of the cancer genome. *Nat. Rev. Genet.*,
- 5 10.1038/nrg3317.
- 6 3. Calabrese,C., Davidson,N.R., Demircioğlu,D., Fonseca,N.A., He,Y., Kahles,A.,
- 7 Lehmann,K. Van, Liu,F., Shiraishi,Y., Soulette,C.M., *et al.* (2020) Genomic basis for
- 8 RNA alterations in cancer. *Nature*, 10.1038/s41586-020-1970-0.
- 9 4. Sondka,Z., Bamford,S., Cole,C.G., Ward,S.A., Dunham,I. and Forbes,S.A. (2018) The
- 10 COSMIC Cancer Gene Census: describing genetic dysfunction across all human
- 11 cancers. *Nat. Rev. Cancer*, 10.1038/s41568-018-0060-1.
- 12 5. Guttman,M., Amit,I., Garber,M., French,C., Lin,M.F., Feldser,D., Huarte,M., Zuk,O.,
- 13 Carey,B.W., Cassady,J.P., *et al.* (2009) Chromatin signature reveals over a thousand
- 14 highly conserved large non-coding RNAs in mammals. *Nature*, 10.1038/nature07672.
- 15 6. Uszczynska-Ratajczak,B., Lagarde,J., Frankish,A., Guigó,R. and Johnson,R. (2018)
- 16 Towards a complete map of the human long non-coding RNA transcriptome. *Nat. Rev.*
- 17 *Genet.*, 10.1038/s41576-018-0017-y.
- 18 7. Derrien,T., Johnson,R., Bussotti,G., Tanzer,A., Djebali,S., Tilgner,H., Guernec,G.,
- 19 Martin,D., Merkel,A., Knowles,D.G., *et al.* (2012) The GENCODE v7 catalog of human
- 20 long noncoding RNAs: Analysis of their gene structure, evolution, and expression.
- 21 *Genome Res.*, 10.1101/gr.132159.111.
- 22 8. Guttman,M. and Rinn,J.L. (2012) Modular regulatory principles of large non-coding RNAs.
- 23 *Nature*, 10.1038/nature10887.
- 24 9. Johnson,R. and Guigó,R. (2014) The RIDL hypothesis: Transposable elements as
- 25 functional domains of long noncoding RNAs. *RNA*, 10.1261/rna.044560.114.
- 26 10. Ulitsky,I., Shkumatava,A., Jan,C.H., Sive,H. and Bartel,D.P. (2011) Conserved function
- 27 of lincRNAs in vertebrate embryonic development despite rapid sequence evolution.
- 28 *Cell*, 10.1016/j.cell.2011.11.055.

- 1 11. Marín-Béjar,O., Mas,A.M., González,J., Martínez,D., Athie,A., Morales,X., Galduroz,M.,
2 Raimondi,I., Grossi,E., Guo,S., *et al.* (2017) The human lncRNA LINC-PINT inhibits
3 tumor cell invasion through a highly conserved sequence element. *Genome Biol.*,
4 10.1186/s13059-017-1331-y.
- 5 12. Frankish,A., Diekhans,M., Ferreira,A.M., Johnson,R., Jungreis,I., Loveland,J.,
6 Mudge,J.M., Sisu,C., Wright,J., Armstrong,J., *et al.* (2019) GENCODE reference
7 annotation for the human and mouse genomes. *Nucleic Acids Res.*,
8 10.1093/nar/gky955.
- 9 13. Kopp,F. and Mendell,J.T. (2018) Functional Classification and Experimental Dissection
10 of Long Noncoding RNAs. *Cell*, 10.1016/j.cell.2018.01.011.
- 11 14. Ulitsky,I. and Bartel,D.P. (2013) XlincRNAs: Genomics, evolution, and mechanisms.
12 *Cell*, 10.1016/j.cell.2013.06.020.
- 13 15. Ma,L., Cao,J., Liu,L., Du,Q., Li,Z., Zou,D., Bajic,V.B. and Zhang,Z. (2019) Lncbook: A
14 curated knowledgebase of human long non-coding rnas. *Nucleic Acids Res.*,
15 10.1093/nar/gky960.
- 16 16. Quek,X.C., Thomson,D.W., Maag,J.L.V., Bartonicek,N., Signal,B., Clark,M.B.,
17 Gloss,B.S. and Dinger,M.E. (2015) lncRNADB v2.0: Expanding the reference database
18 for functional long noncoding RNAs. *Nucleic Acids Res.*, 10.1093/nar/gku988.
- 19 17. Slack,F.J. and Chinnaiyan,A.M. (2019) The Role of Non-coding RNAs in Oncology. *Cell*,
20 10.1016/j.cell.2019.10.017.
- 21 18. Du,Z., Fei,T., Verhaak,R.G.W., Su,Z., Zhang,Y., Brown,M., Chen,Y. and Liu,X.S. (2013)
22 Integrative genomic analyses reveal clinically relevant long noncoding RNAs in human
23 cancer. *Nat. Struct. Mol. Biol.*, 10.1038/nsmb.2591.
- 24 19. Dias,N. and Stein,C.A. (2002) Antisense oligonucleotides: Basic concepts and
25 mechanisms. *Mol. Cancer Ther.*
- 26 20. Gutschner,T., Hämmerle,M., Eißmann,M., Hsu,J., Kim,Y., Hung,G., Revenko,A.,
27 Arun,G., Stentrup,M., Groß,M., *et al.* (2013) The noncoding RNA MALAT1 is a critical
28 regulator of the metastasis phenotype of lung cancer cells. *Cancer Res.*, 10.1158/0008-

- 1 5472.CAN-12-2850.
- 2 21. Wahlestedt,C. (2013) Targeting long non-coding RNA to therapeutically upregulate gene
3 expression. *Nat. Rev. Drug Discov.*, 10.1038/nrd4018.
- 4 22. Kaczmarek,J.C., Kowalski,P.S. and Anderson,D.G. (2017) Advances in the delivery of
5 RNA therapeutics: From concept to clinical reality. *Genome Med.*, 10.1186/s13073-017-
6 0450-0.
- 7 23. Huarte,M., Guttman,M., Feldser,D., Garber,M., Koziol,M.J., Kenzelmann-Broz,D.,
8 Khalil,A.M., Zuk,O., Amit,I., Rabani,M., *et al.* (2010) A Large Intergenic Noncoding RNA
9 Induced by p53 Mediates Global Gene Repression in the p53 Response. *Cell*, **142**,
10 409–419.
- 11 24. Iyer,M.K., Niknafs,Y.S., Malik,R., Singhal,U., Sahu,A., Hosono,Y., Barrette,T.R.,
12 Prensner,J.R., Evans,J.R., Zhao,S., *et al.* (2015) The landscape of long noncoding
13 RNAs in the human transcriptome. *Nat. Genet.*, 10.1038/ng.3192.
- 14 25. Esposito,R., Bosch,N., Lanzós,A., Polidori,T., Pulido-Quetglas,C. and Johnson,R. (2019)
15 Hacking the Cancer Genome: Profiling Therapeutically Actionable Long Non-coding
16 RNAs Using CRISPR-Cas9 Screening. *Cancer Cell*, 10.1016/j.ccell.2019.01.019.
- 17 26. Lanzós,A., Carlevaro-Fita,J., Mularoni,L., Reverter,F., Palumbo,E., Guigó,R. and
18 Johnson,R. (2017) Discovery of Cancer Driver Long Noncoding RNAs across 1112
19 Tumour Genomes: New Candidates and Distinguishing Features. *Sci. Rep.*,
20 10.1038/srep41544.
- 21 27. Mularoni,L., Sabarinathan,R., Deu-Pons,J., Gonzalez-Perez,A. and López-Bigas,N.
22 (2016) OncodriveFML: A general framework to identify coding and non-coding regions
23 with cancer driver mutations. *Genome Biol.*, 10.1186/s13059-016-0994-0.
- 24 28. Rheinbay,E., Nielsen,M.M., Abascal,F., Wala,J.A., Shapira,O., Tiao,G., Hornshøj,H.,
25 Hess,J.M., Juul,R.I., Lin,Z., *et al.* (2020) Analyses of non-coding somatic drivers in
26 2,658 cancer whole genomes. *Nature*, 10.1038/s41586-020-1965-x.
- 27 29. Engreitz,J.M., Haines,J.E., Perez,E.M., Munson,G., Chen,J., Kane,M., McDonel,P.E.,
28 Guttman,M. and Lander,E.S. (2016) Local regulation of gene expression by lncRNA

- 1 promoters, transcription and splicing. *Nature*, 10.1038/nature20149.
- 2 30. Yin,Y., Yan,P., Lu,J., Song,G., Zhu,Y., Li,Z., Zhao,Y., Shen,B., Huang,X., Zhu,H., *et al.*
3 (2015) Opposing roles for the lncRNA haunt and its genomic locus in regulating HOXA
4 gene activation during embryonic stem cell differentiation. *Cell Stem Cell*,
5 10.1016/j.stem.2015.03.007.
- 6 31. Groff,A.F., Sanchez-Gomez,D.B., Soruco,M.M.L., Gerhardinger,C., Barutcu,A.R., Li,E.,
7 Elcavage,L., Plana,O., Sanchez,L. V., Lee,J.C., *et al.* (2016) In Vivo Characterization of
8 Linc-p21 Reveals Functional cis-Regulatory DNA Elements. *Cell Rep.*,
9 10.1016/j.celrep.2016.07.050.
- 10 32. John Liu,S., Malatesta,M., Lien,B. V., Saha,P., Thombare,S.S., Hong,S.J., Pedraza,L.,
11 Koontz,M., Seo,K., Horlbeck,M.A., *et al.* (2020) CRISPRi-based radiation modifier
12 screen identifies long non-coding RNA therapeutic targets in glioma. *Genome Biol.*,
13 10.1186/s13059-020-01995-4.
- 14 33. Hosono,Y., Niknafs,Y.S., Prensner,J.R., Iyer,M.K., Dhanasekaran,S.M., Mehra,R.,
15 Pitchiaya,S., Tien,J., Escara-Wilke,J., Poliakov,A., *et al.* (2017) Oncogenic Role of
16 THOR, a Conserved Cancer/Testis Long Non-coding RNA. *Cell*,
17 10.1016/j.cell.2017.11.040.
- 18 34. Lee,S., Kopp,F., Chang,T.C., Sataluri,A., Chen,B., Sivakumar,S., Yu,H., Xie,Y. and
19 Mendell,J.T. (2016) Noncoding RNA NORAD Regulates Genomic Stability by
20 Sequestering PUMILIO Proteins. *Cell*, 10.1016/j.cell.2015.12.017.
- 21 35. Leucci,E., Vendramin,R., Spinazzi,M., Laurette,P., Fiers,M., Wouters,J., Radaelli,E.,
22 Eyckerman,S., Leonelli,C., Vanderheyden,K., *et al.* (2016) Melanoma addiction to the
23 long non-coding RNA SAMMSON. *Nature*, 10.1038/nature17161.
- 24 36. Munschauer,M., Nguyen,C.T., Sirokman,K., Hartigan,C.R., Hogstrom,L., Engreitz,J.M.,
25 Ulirsch,J.C., Fulco,C.P., Subramanian,V., Chen,J., *et al.* (2018) The NORAD lncRNA
26 assembles a topoisomerase complex critical for genome stability. *Nature*,
27 10.1038/s41586-018-0453-z.
- 28 37. Bao,Z., Yang,Z., Huang,Z., Zhou,Y., Cui,Q. and Dong,D. (2019) LncRNADisease 2.0: An

- 1 updated database of long non-coding RNA-associated diseases. *Nucleic Acids Res.*,
2 10.1093/nar/gky905.
- 3 38. Wang,J., Zhang,X., Chen,W., Li,J. and Liu,C. (2018) CRlncRNA: A manually curated
4 database of cancer-related long non-coding RNAs with experimental proof of functions
5 on clinicopathological and molecular features. *BMC Med. Genomics*, 10.1186/s12920-
6 018-0430-2.
- 7 39. Zhou,B., Zhao,H., Yu,J., Guo,C., Dou,X., Song,F., Hu,G., Cao,Z., Qu,Y., Yang,Y., *et al.*
8 (2018) EVLncRNAs: A manually curated database for long non-coding RNAs validated
9 by low-throughput experiments. *Nucleic Acids Res.*, 10.1093/nar/gkx677.
- 10 40. Gao,Y., Wang,P., Wang,Y., Ma,X., Zhi,H., Zhou,D., Li,X., Fang,Y., Shen,W., Xu,Y., *et al.*
11 (2019) Lnc2Cancer v2.0: Updated database of experimentally supported long non-
12 coding RNAs in human cancers. *Nucleic Acids Res.*, 10.1093/nar/gky1096.
- 13 41. Abbott,K.L., Nyre,E.T., Abrahante,J., Ho,Y.Y., Vogel,R.I. and Starr,T.K. (2015) The
14 candidate cancer gene database: A database of cancer driver genes from forward
15 genetic screens in mice. *Nucleic Acids Res.*, 10.1093/nar/gku770.
- 16 42. Carlevaro-Fita,J., Lanzós,A., Feuerbach,L., Hong,C., Mas-Ponte,D., Pedersen,J.S.,
17 Abascal,F., Amin,S.B., Bader,G.D., Barenboim,J., *et al.* (2020) Cancer LncRNA Census
18 reveals evidence for deep functional conservation of long noncoding RNAs in
19 tumorigenesis. *Commun. Biol.*, 10.1038/s42003-019-0741-7.
- 20 43. Bergadà-Pijuan,J., Pulido-Quetglas,C., Vancura,A. and Johnson,R. (2020) CASPR, an
21 analysis pipeline for single and paired guide RNA CRISPR screens, reveals optimal
22 target selection for long non-coding RNAs. *Bioinformatics*,
23 10.1093/bioinformatics/btz811.
- 24 44. Goyal,A., Myacheva,K., Groß,M., Klingenberg,M., Duran Arqué,B. and Diederichs,S.
25 (2017) Challenges of CRISPR/Cas9 applications for long non-coding RNA genes.
26 *Nucleic Acids Res.*, 10.1093/nar/gkw883.
- 27 45. Liu,S.J., Horlbeck,M.A., Cho,S.W., Birk,H.S., Malatesta,M., He,D., Attenello,F.J.,
28 Villalta,J.E., Cho,M.Y., Chen,Y., *et al.* (2017) CRISPRi-based genome-scale

- 1 identification of functional long noncoding RNA loci in human cells. *Science* (80-.),
2 10.1126/science.aah71111.
- 3 46. Copeland,N.G. and Jenkins,N.A. (2010) Harnessing transposons for cancer gene
4 discovery. *Nat. Rev. Cancer*, 10.1038/nrc2916.
- 5 47. Siepel,A., Bejerano,G., Pedersen,J.S., Hinrichs,A.S., Hou,M., Rosenbloom,K.,
6 Clawson,H., Spieth,J., Hillier,L.D.W., Richards,S., *et al.* (2005) Evolutionarily conserved
7 elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.*,
8 10.1101/gr.3715005.
- 9 48. Ørom,U.A., Derrien,T., Beringer,M., Gumireddy,K., Gardini,A., Bussotti,G., Lai,F.,
10 Zytnicki,M., Notredame,C., Huang,Q., *et al.* (2010) Long noncoding RNAs with
11 enhancer-like function in human cells. *Cell*, 10.1016/j.cell.2010.09.001.
- 12 49. Liang,X.H., Sun,H., Nichols,J.G. and Crooke,S.T. (2017) RNase H1-Dependent
13 Antisense Oligonucleotides Are Robustly Active in Directing RNA Cleavage in Both the
14 Cytoplasm and the Nucleus. *Mol. Ther.*, 10.1016/j.ymthe.2017.06.002.
- 15 50. Kamola,P.J., Kitson,J.D.A., Turner,G., Maratou,K., Eriksson,S., Panjwani,A.,
16 Warnock,L.C., Douillard Guilloux,G.A., Moores,K., Koppe,E.L., *et al.* (2015) In silico and
17 in vitro evaluation of exonic and intronic off-target effects form a critical element of
18 therapeutic ASO gapmer optimization. *Nucleic Acids Res.*, 10.1093/nar/gkv857.
- 19 51. Furney,S.J., Madden,S.F., Kisiel,T.A., Higgins,D.G. and Lopez-Bigas,N. (2008) Distinct
20 patterns in the regulation and evolution of human cancer genes. *In Silico Biol.*
- 21 52. Furney,S.J., Higgins,D.G., Ouzounis,C.A. and López-Bigas,N. (2006) Structural and
22 functional properties of genes involved in human cancer. *BMC Genomics*,
23 10.1186/1471-2164-7-3.
- 24 53. Lee,D., Redfern,O. and Orengo,C. (2007) Predicting protein function from sequence and
25 structure. *Nat. Rev. Mol. Cell Biol.*, 10.1038/nrm2281.
- 26 54. Wilson,C.A., Kreychman,J. and Gerstein,M. (2000) Assessing annotation transfer for
27 genomics: Quantifying the relations between protein sequence, structure and function
28 through traditional and probabilistic scores. *J. Mol. Biol.*, 10.1006/jmbi.2000.3550.

- 1 55. Pennacchio,L.A., Ahituv,N., Moses,A.M., Prabhakar,S., Nobrega,M.A., Shoukry,M.,
2 Minovitsky,S., Dubchak,I., Holt,A., Lewis,K.D., *et al.* (2006) In vivo enhancer analysis of
3 human conserved non-coding sequences. *Nature*, 10.1038/nature05295.
- 4 56. Carlevaro-Fita,J., Liu,L., Zhou,Y., Zhang,S., Chouvardas,P., Johnson,R. and Li,J. (2019)
5 LnCompare: gene set feature analysis for human long non-coding RNAs. *Nucleic Acids*
6 *Res.*, 10.1093/nar/gkz410.
- 7 57. Deng,N., Zhou,H., Fan,H. and Yuan,Y. (2017) Single nucleotide polymorphisms and
8 cancer susceptibility. *Oncotarget*, 10.18632/oncotarget.22372.
- 9 58. Buniello,A., Macarthur,J.A.L., Cerezo,M., Harris,L.W., Hayhurst,J., Malangone,C.,
10 McMahon,A., Morales,J., Mountjoy,E., Sollis,E., *et al.* (2019) The NHGRI-EBI GWAS
11 Catalog of published genome-wide association studies, targeted arrays and summary
12 statistics 2019. *Nucleic Acids Res.*, 10.1093/nar/gky1120.
- 13 59. Chen,X., Hao,Y., Cui,Y., Fan,Z., He,S., Luo,J. and Chen,R. (2017) LncVar: A database
14 of genetic variation associated with long non-coding genes. *Bioinformatics*,
15 10.1093/bioinformatics/btw581.
- 16 60. Perry,R.B.T. and Ulitsky,I. (2016) The functions of long noncoding RNAs in development
17 and stem cells. *Dev.*, 10.1242/dev.140962.
- 18 61. Ounzain,S., Micheletti,R., Arnan,C., Plaisance,I., Cecchi,D., Schroen,B., Reverter,F.,
19 Alexanian,M., Gonzales,C., Ng,S.Y., *et al.* (2015) CARMEN, a human super enhancer-
20 associated long noncoding RNA controlling cardiac specification, differentiation and
21 homeostasis. *J. Mol. Cell. Cardiol.*, 10.1016/j.yjmcc.2015.09.016.
- 22 62. Park,S.L., Carmella,S.G., Chen,M., Patel,Y., Stram,D.O., Haiman,C.A., Le Marchand,L.
23 and Hecht,S.S. (2015) Mercapturic acids derived from the toxicants acrolein and
24 crotonaldehyde in the urine of cigarette smokers from five ethnic groups with differing
25 risks for lung cancer. *PLoS One*, 10.1371/journal.pone.0124841.
- 26 63. Johnson,R., Teh,C.H.L., Jia,H., Vanisri,R.R., Pandey,T., Lu,Z.H., Buckley,N.J.,
27 Stanton,L.W. and Lipovich,L. (2009) Regulation of neural macroRNAs by the
28 transcriptional repressor REST. *RNA*, 10.1261/rna.1127009.

- 1 64. Li,B., Guo,Z., Liang,Q., Zhou,H., Luo,Y., He,S. and Lin,Z. (2019) LncRNA DGCR5 Up-
2 regulates TGF- β 1, increases cancer cell stemness and predicts survival of prostate
3 cancer patients. *Cancer Manag. Res.*, 10.2147/CMAR.S231112.
- 4 65. Amodio,N., Stamato,M.A., Juli,G., Morelli,E., Fulciniti,M., Manzoni,M., Taiana,E.,
5 Agnelli,L., Cantafio,M.E.G., Romeo,E., *et al.* (2018) Drugging the lncRNA MALAT1 via
6 LNA gapmer ASO inhibits gene expression of proteasome subunits and triggers anti-
7 multiple myeloma activity. *Leukemia*, 10.1038/s41375-018-0067-3.
- 8 66. Bergadà-Pijuan,J., Pulido-Quetglas,C., Vancura,A. and Johnson,R. (2019) CASPR, an
9 analysis pipeline for single and paired guide RNA CRISPR screens, reveals optimal
10 target selection for long noncoding RNAs. *Bioinformatics*,
11 10.1093/bioinformatics/btz811.
- 12 67. Colaprico,A., Silva,T.C., Olsen,C., Garofano,L., Cava,C., Garolini,D., Sabedot,T.S.,
13 Malta,T.M., Pagnotta,S.M., Castiglioni,I., *et al.* (2016) TCGAbiolinks: An R/Bioconductor
14 package for integrative analysis of TCGA data. *Nucleic Acids Res.*,
15 10.1093/nar/gkv1507.
- 16 68. Aparicio-Prat,E., Arnan,C., Sala,I., Bosch,N., Guigó,R. and Johnson,R. (2015) DECKO:
17 Single-oligo, dual-CRISPR deletion of genomic elements including long non-coding
18 RNAs. *BMC Genomics*, 10.1186/s12864-015-2086-z.
- 19 69. Schmittgen,T.D. and Livak,K.J. (2008) Analyzing real-time PCR data by the comparative
20 CT method. *Nat. Protoc.*, 10.1038/nprot.2008.73.

21

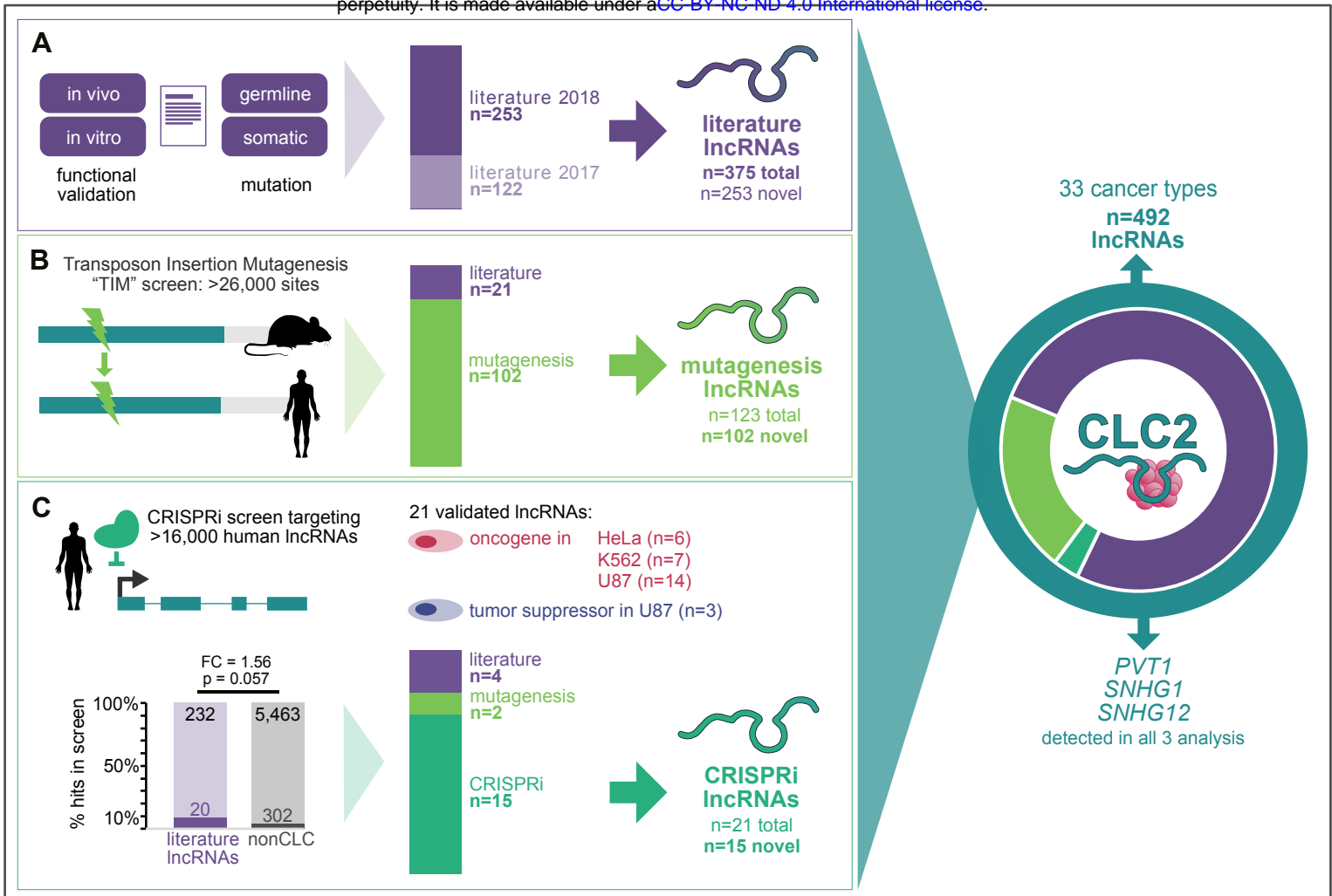


Figure 1

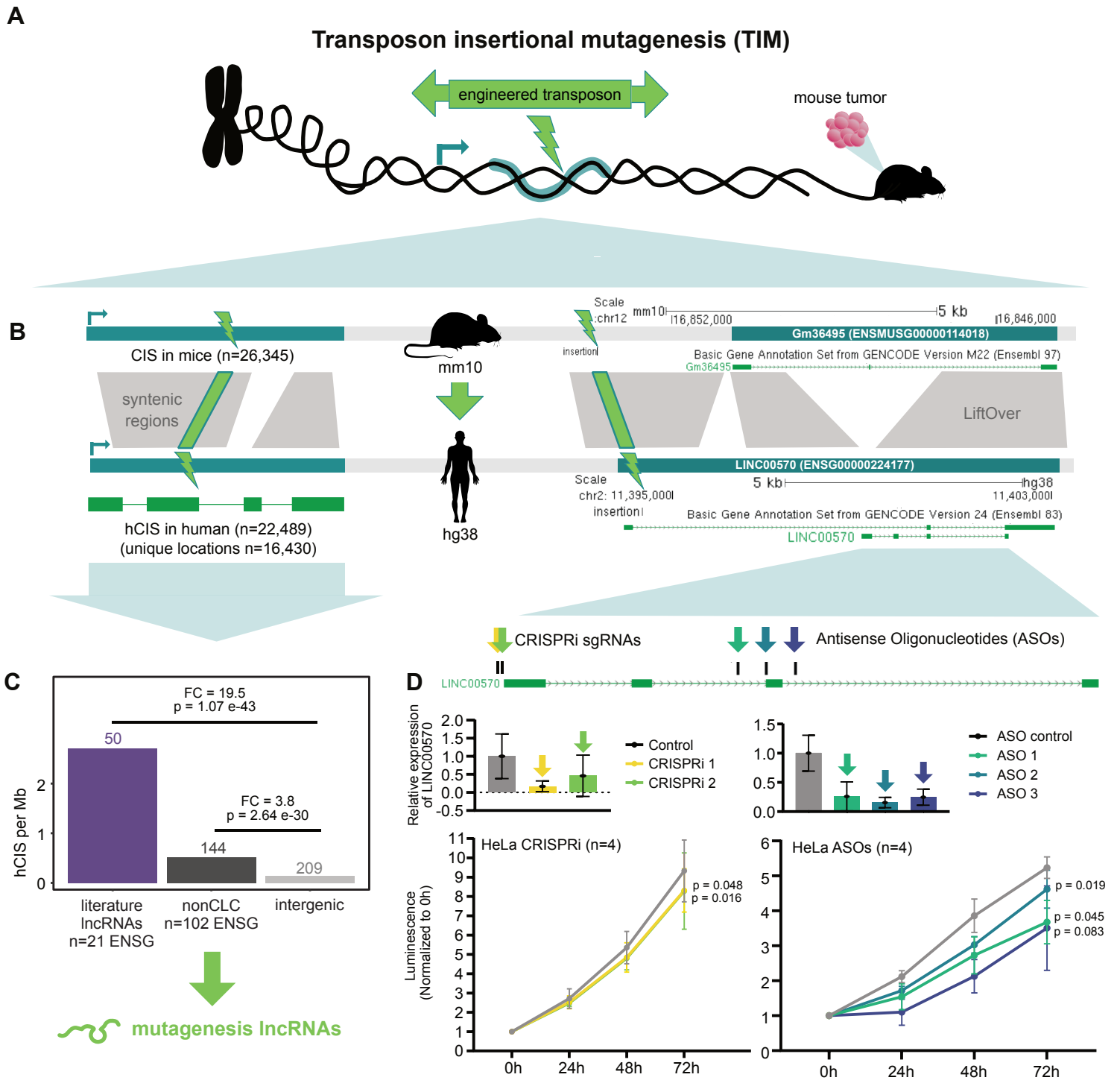


Figure 2

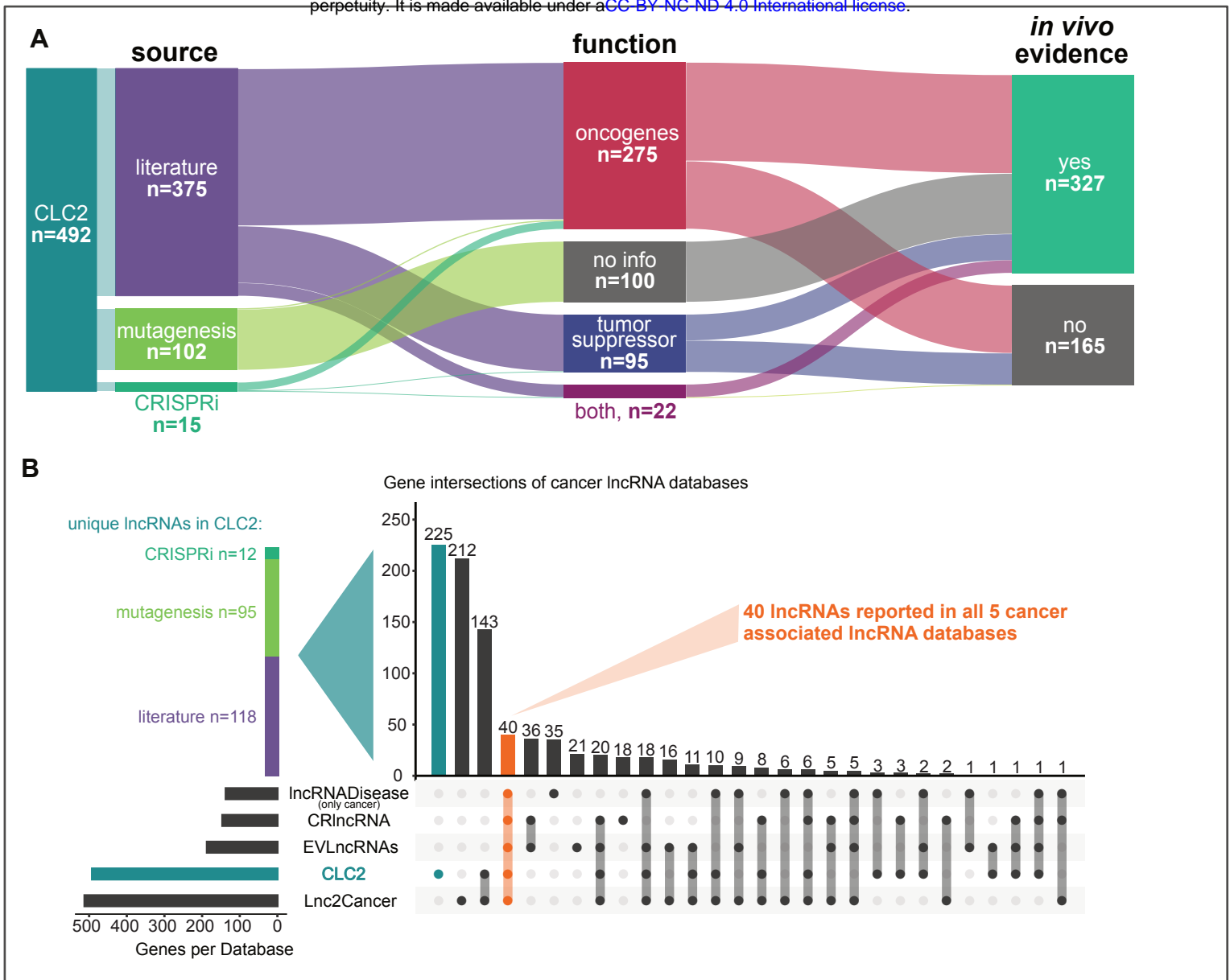


Figure 3

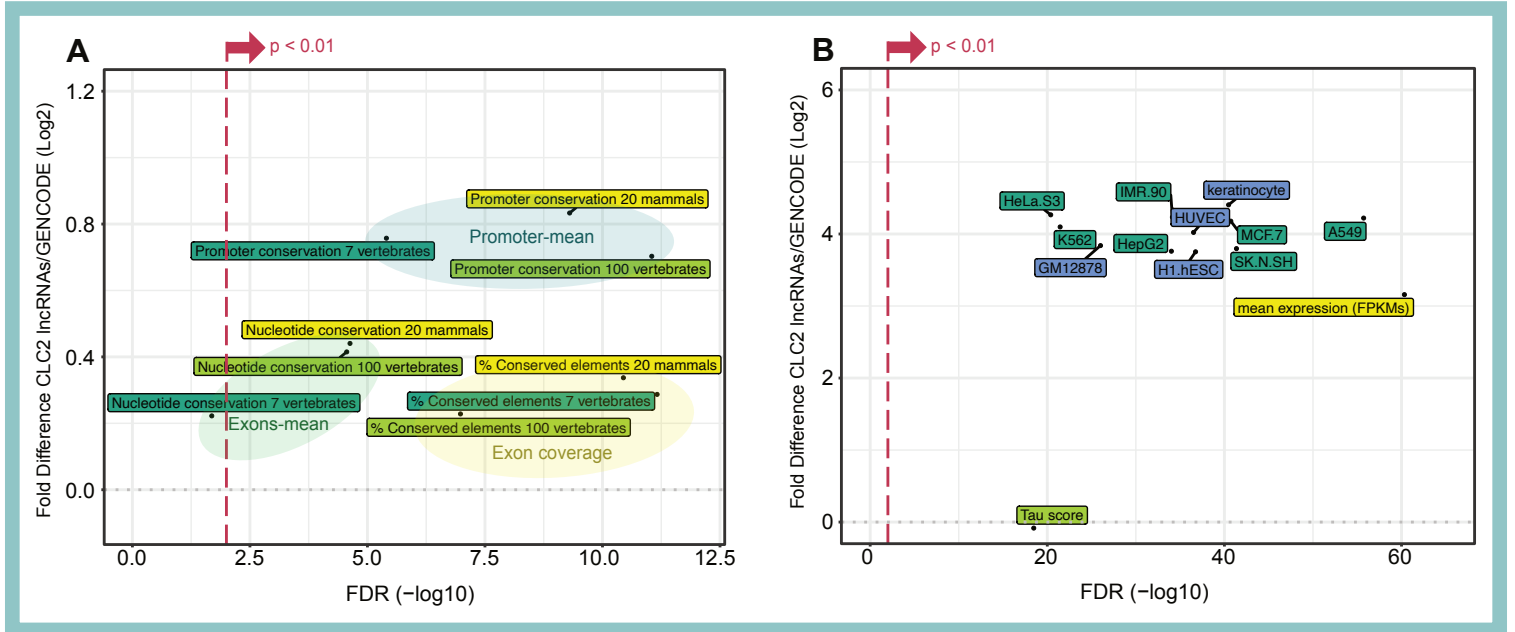
Conservation

- conservation mammals (n=20)
- conservation vertebrates (n=100)
- conservation vertebrates (n=7)

Expression

- expression in human tissue (HBM samples)
- tissue specificity in human tissue (HBM samples)
- whole cell expression (cancer cell line)
- whole cell expression (non cancer cell line)

CLC2 lncRNAs n=492 compared to GENCODE lncRNAs (n= 15,827)



mutagenesis lncRNAs n=123 compared to GENCODE lncRNAs (n= 15,827)

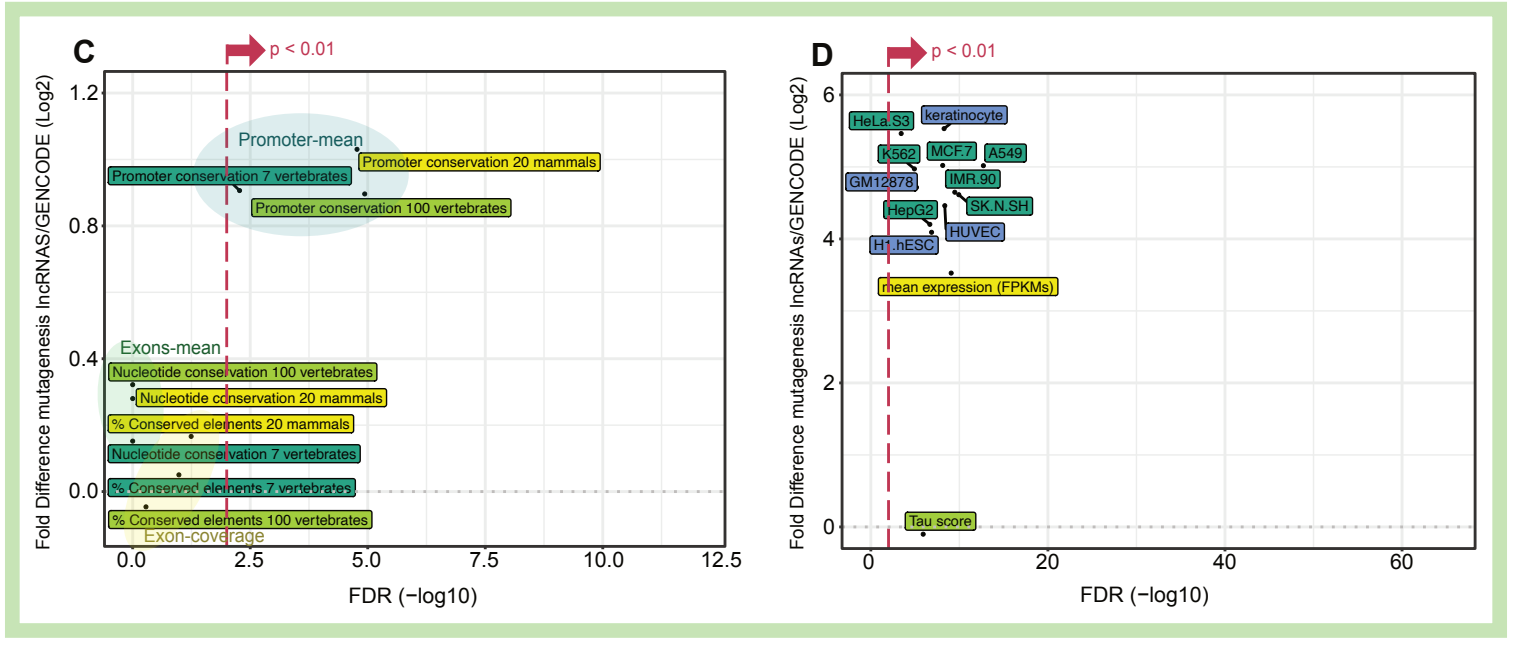


Figure 4

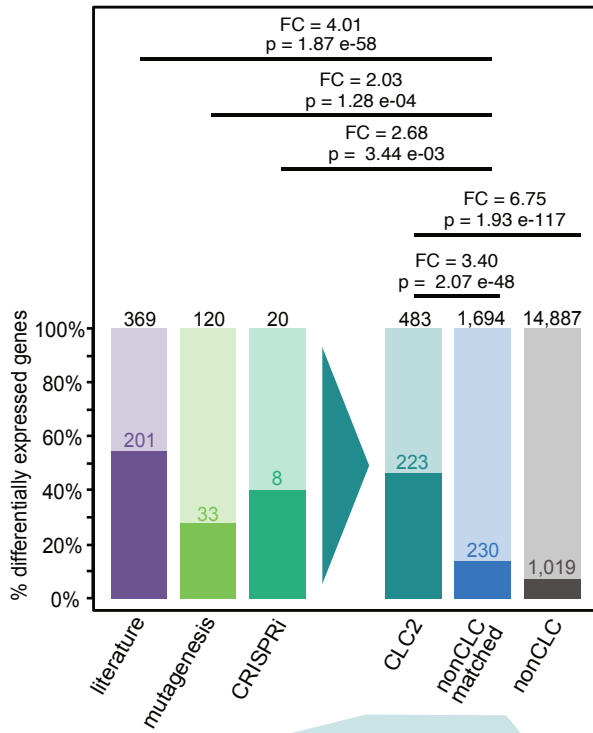
literature lncRNAs
mutagenesis lncRNAs
CRISPRi lncRNAs

CLC2 lncRNAs

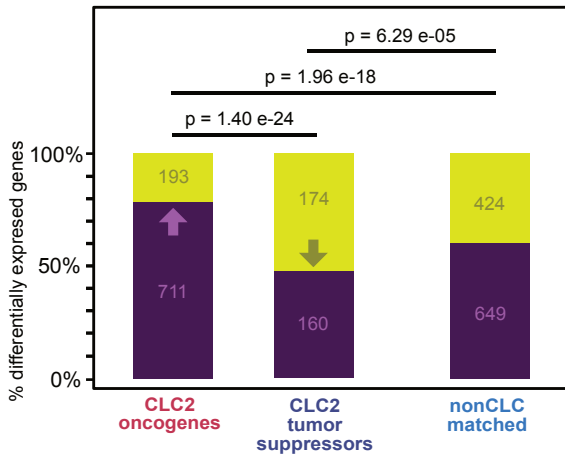


nonCLC matched lncRNAs
expression-matched lncRNA genes
nonCLC lncRNAs
lncRNA genes

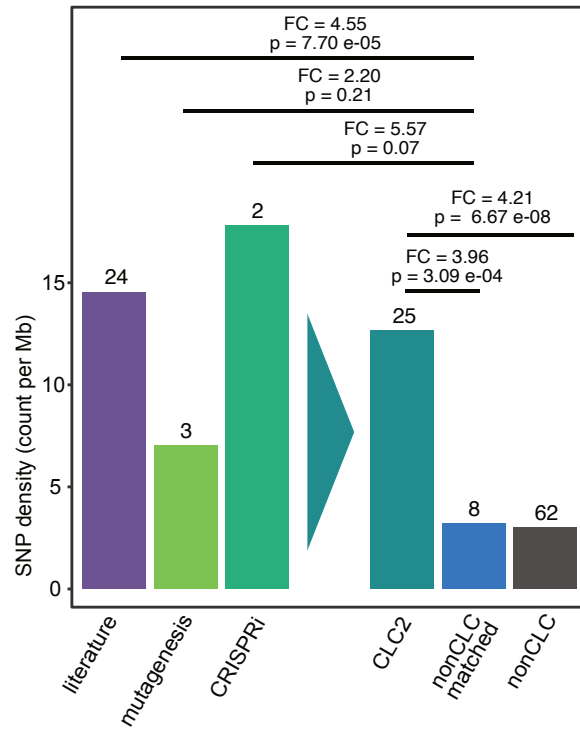
A Differentially expressed genes



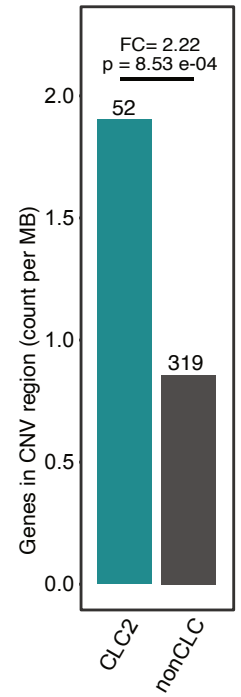
↑ upregulated expression
↓ downregulated expression



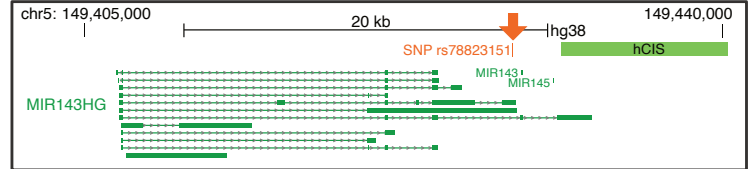
C Exonic cancer SNPs



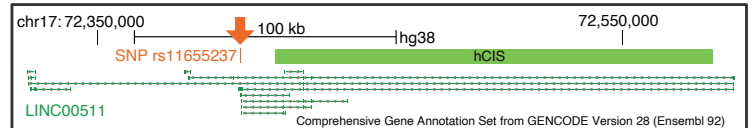
E CNV genes



D MIR143HG/CARMN



LINC00511



LINC01488

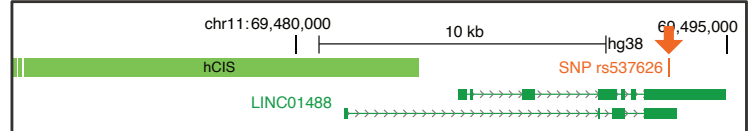
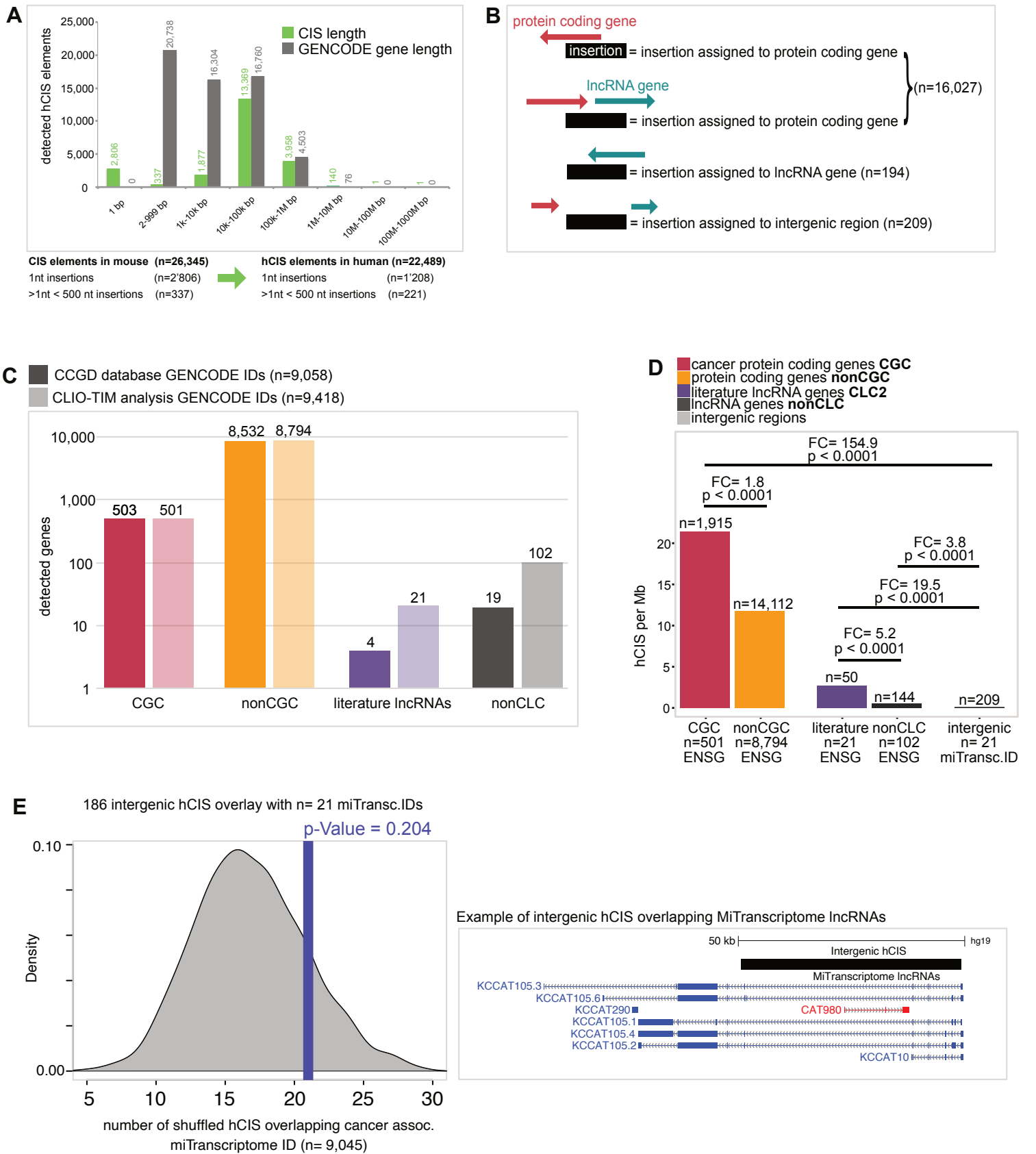
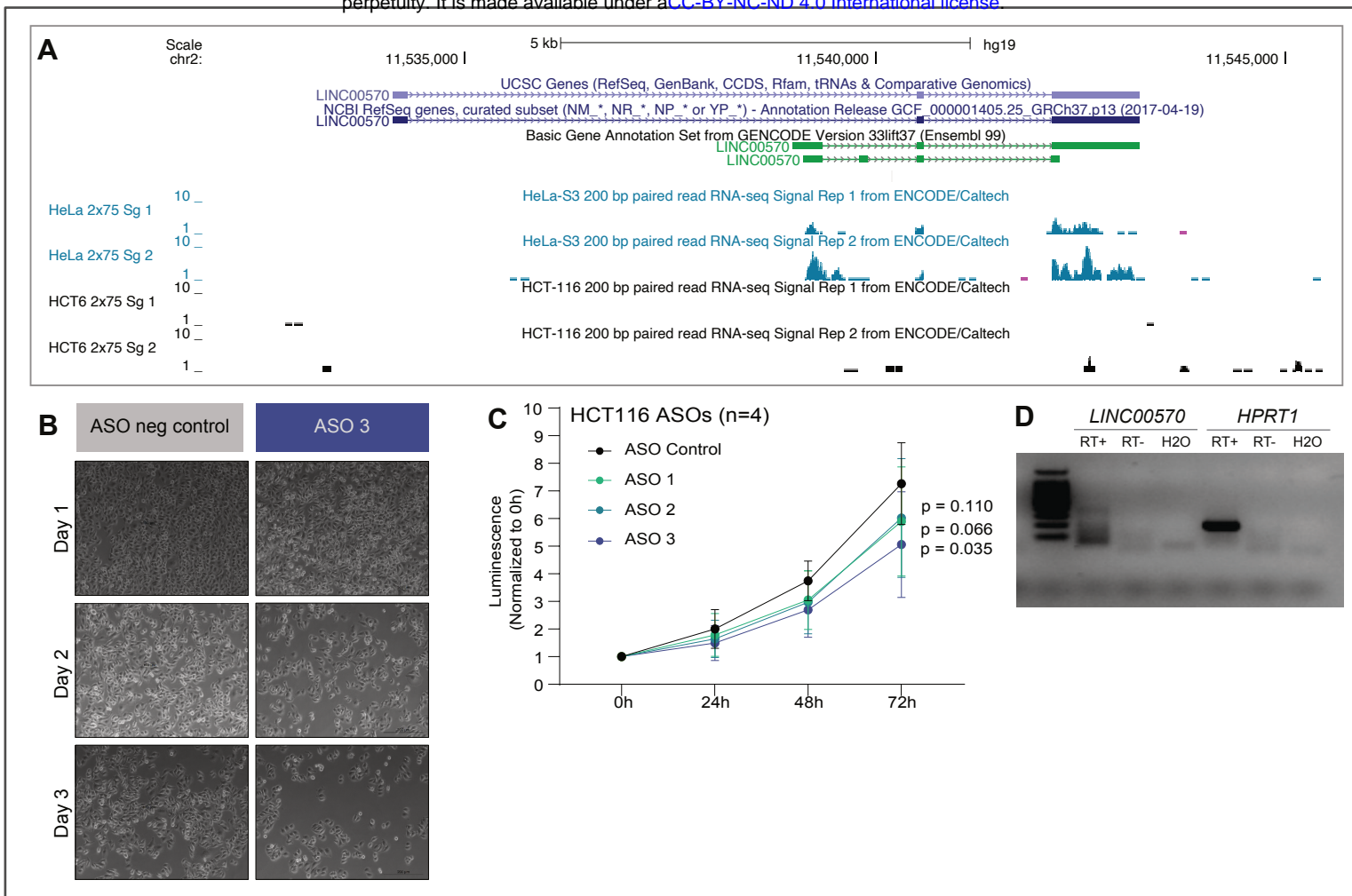


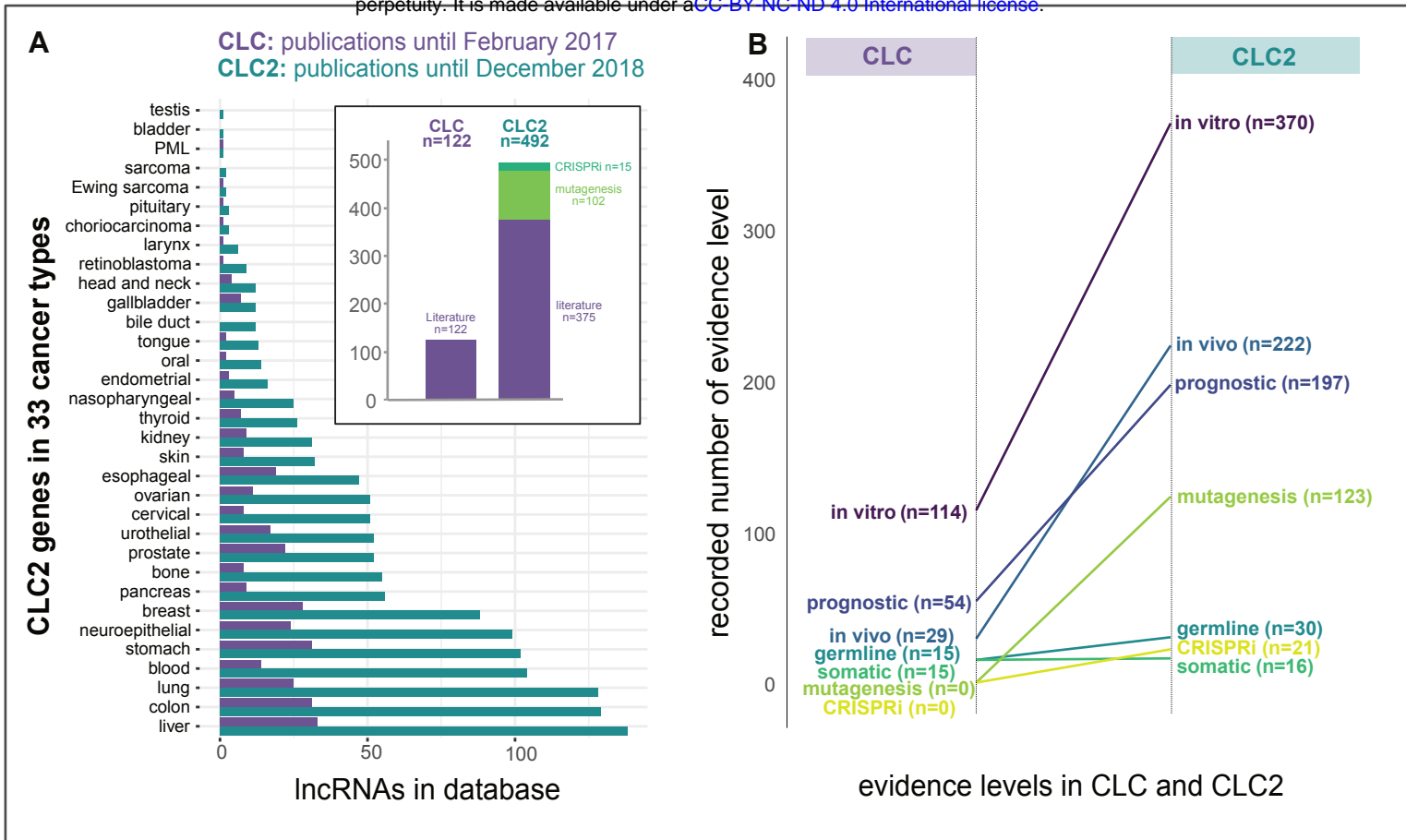
Figure 5



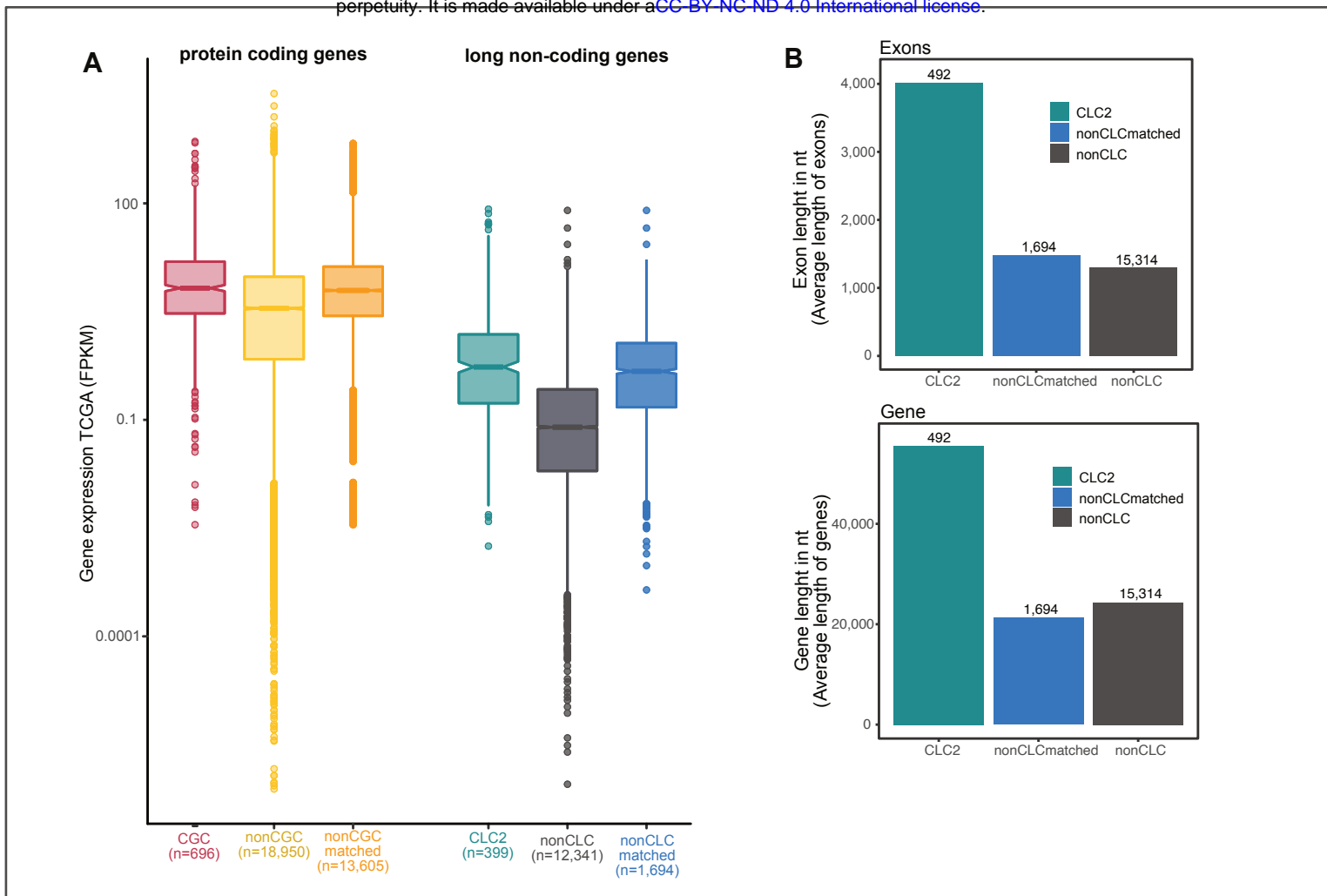
SUPP. Figure 1



SUPP. Figure 2



SUPP. Figure 3

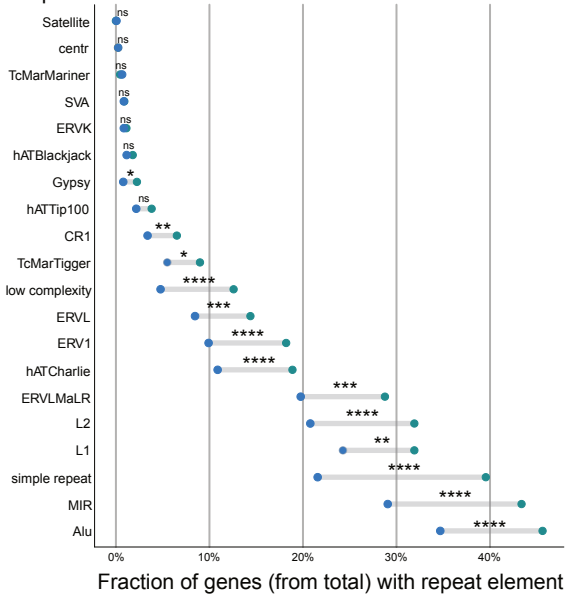


SUPP. Figure 4

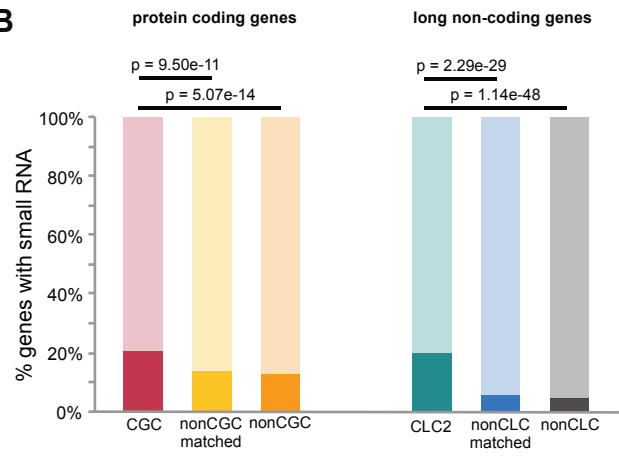
Gene type

- cancer protein coding genes **CGC**
- highly expressed protein coding genes **nonCGC matched**
- protein coding genes **nonCGC**
- cancer lncRNA genes **CLC2**
- highly expressed lncRNA genes **nonCLC matched**
- lncRNA genes **nonCLC**

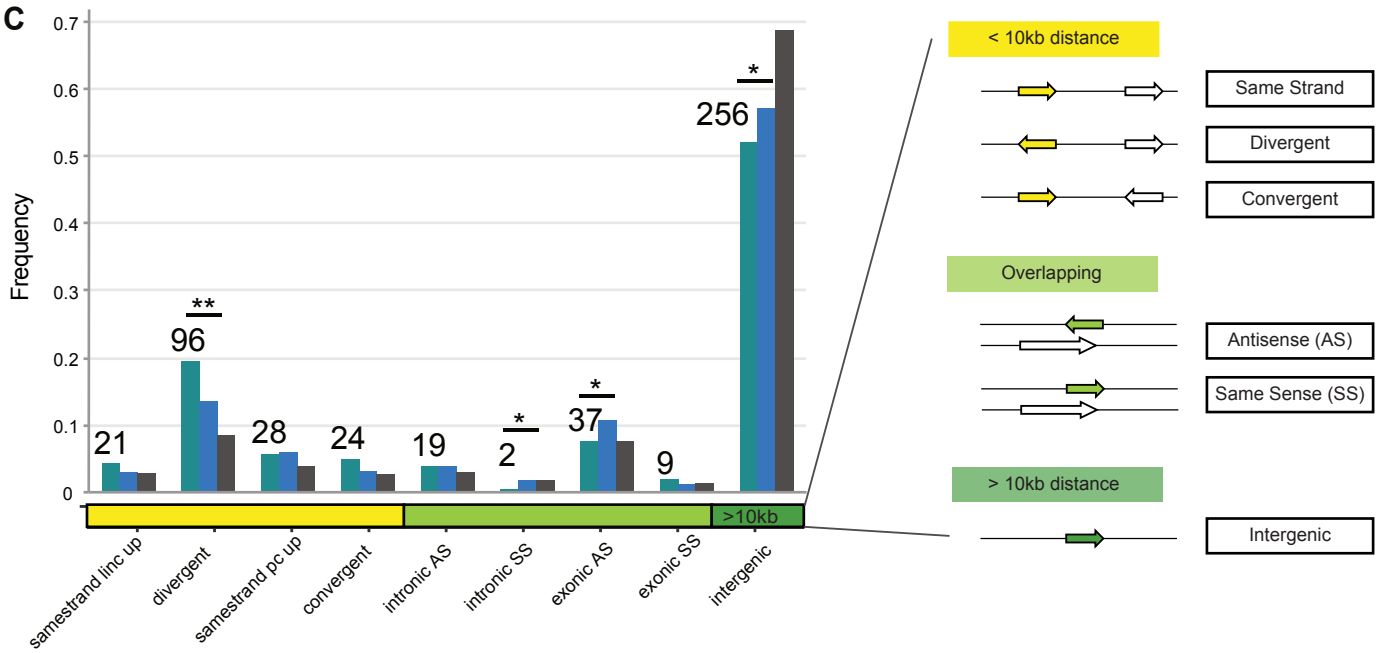
A with repeat:



B



C



SUPP. Figure 5

A

CLC2 cancer type	TCGA cancer type
bile duct	CHOL
blood	DLBC
blood	LAML
breast	BRCA
cervical	CESC
colon	COAD
endometrial	UCEC
endometrial	UCS
esophageal	ESCA
head and neck	HNSC
kidney	KIRP
kidney	KIRC
kidney	KICH
liver	LIHC
lung	LUAD
lung	LUSC
neuroepithelial	LGG
neuroepithelial	PCPG
neuroepithelial	GBM
ovarian	OV
pancreas	PAAD
prostate	PRAD
skin	UVM
skin	SKCM
stomach	STAD
testis	TGCT
thyroid	THCA
urothelial	BLCA
sarcoma	SARC

B

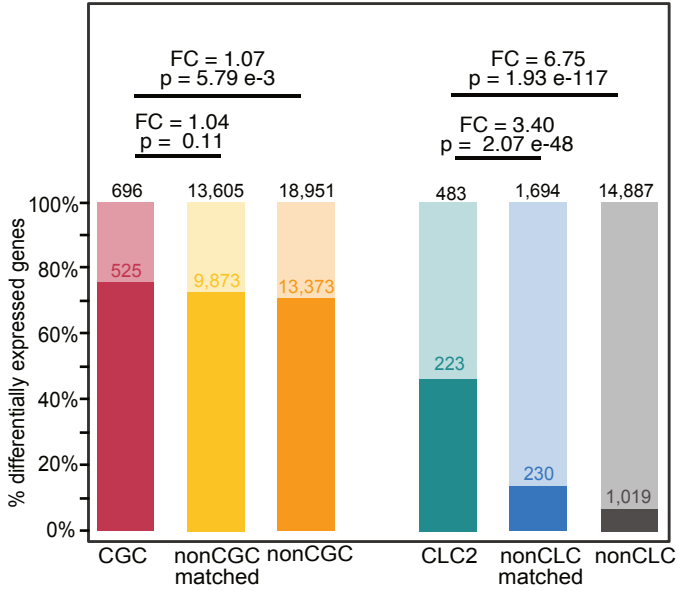
TCGA cancer type	tumor (01)	healthy (11)	total samples
TCGA-BLCA	19	19	38
TCGA-BRCA	112	112	224
TCGA-CESC	3	3	6
TCGA-CHOL	9	9	18
TCGA-COAD	41	41	82
TCGA-ESCA	8	8	16
TCGA-HNSC	43	43	86
TCGA-KICH	23	23	46
TCGA-KIRC	72	72	144
TCGA-KIRP	31	31	62
TCGA-LIHC	50	50	100
TCGA-LUAD	57	57	114
TCGA-LUSC	49	49	98
TCGA-PAAD	4	4	8
TCGA-PCPG	3	3	6
TCGA-PRAD	52	52	104
TCGA-SARC	2	2	4
TCGA-STAD	27	27	54
TCGA-THCA	58	58	116
TCGA-UCEC	23	23	46
TOTAL			1372

SUPP. Figure 6

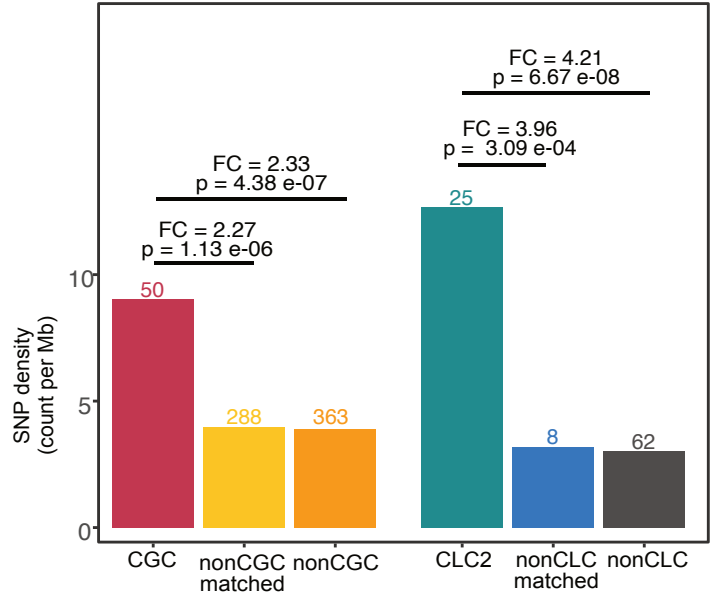
Gene type

- cancer protein coding genes **CGC**
- highly expressed protein coding genes **nonCGC matched**
- protein coding genes **nonCGC**
- cancer lncRNA genes **CLC2**
- highly expressed lncRNA genes **nonCLC matched**
- lncRNA genes **nonCLC**

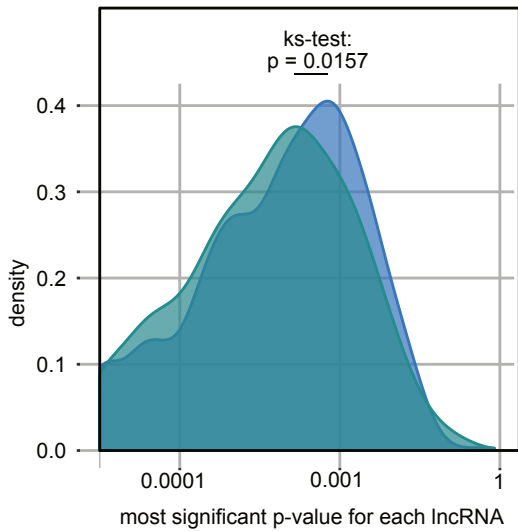
A Differentially expressed genes



B Exonic cancer SNPs



C Survival analysis



SUPP. Figure 7