

Unique genomic features and deeply-conserved functions of long non-coding RNAs in the Cancer LncRNA Census (CLC)

Joana Carlevaro-Fita* (1,2)

Andrés Lanzós* (3,4,5)

Lars Feuerbach (6)

Chen Hong (6)

David Mas-Ponte (3,4,5)

Jakob Skou Pedersen (7)

Rory Johnson (1,2)

On behalf of the PCAWG Drivers and Functional Interpretation Group and the ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Network

1. Department of Clinical Research, University of Bern, 3008 Bern, Switzerland
2. Department of Medical Oncology, Inselspital, University Hospital and University of Bern, 3010 Bern, Switzerland
3. Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Dr. Aiguader 88, Barcelona 08003, Spain
4. Universitat Pompeu Fabra (UPF), Barcelona, Spain.
5. Institut Hospital del Mar d'Investigacions Mèdiques (IMIM), Dr. Aiguader 88, 08003 Barcelona, Catalonia, Spain.
6. Applied Bioinformatics, Deutsches Krebsforschungszentrum, 69120 Heidelberg, Germany
7. Department for Molecular Medicine, Aarhus University Hospital, Palle Juul-Jensens Boulevard 99, 8200 Aarhus N, Denmark.

*** Equal contribution**

Correspondence to rory.johnson@dkf.unibe.ch

Keywords: long non-coding RNA; lncRNA; GENCODE; cancer driver; driver prediction; manual curation; ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Network; PCAWG.

Abstract

Long non-coding RNAs (lncRNAs) that drive tumorigenesis are a growing focus of cancer genomics studies. To facilitate further discovery, we have created the “Cancer LncRNA Census” (CLC), a manually-curated and strictly-defined compilation of lncRNAs with causative roles in cancer. CLC has two principle applications: first, as a resource for training and benchmarking *de novo* identification methods; and second, as a dataset for studying the fundamental properties of these genes.

CLC Version 1 comprises 122 lncRNAs implicated in 29 distinct cancers. LncRNAs are included based on functional or genetic evidence for causative roles in cancer progression. All belong to the GENCODE reference annotation, to enable integration across projects and datasets. For each entry, the evidence type, biological activity (oncogene or tumour suppressor), source reference and cancer type are recorded. Supporting its usefulness, CLC genes are significantly enriched amongst *de novo* predicted driver genes from PCAWG. CLC genes are distinguished from other lncRNAs by a series of features consistent with biological function, including gene length, high expression and sequence conservation of both exons and promoters. We identify a trend for CLC genes to be co-localised with known protein-coding cancer genes along the human genome. Finally, by integrating data from transposon-mutagenesis functional screens, we show that mouse orthologues of CLC genes tend also to be cancer genes.

Thus CLC represents a valuable resource for research into long non-coding RNAs in cancer. Their evolutionary and genomic properties have implications for understanding disease mechanisms and point to conserved functions across ~80 million years of evolution.

1 **Introduction**

2 Tumorigenesis is driven by a series of genetic mutations that promote cancer phenotypes and
3 consequently experience positive selection (Yates & Campbell 2012). The systematic discovery of
4 such driver mutations, and the genes whose functions they alter, has been made possible by tumour
5 genome sequencing. By collecting the entirety of such genes for every cancer type, we aim to develop
6 a comprehensive view of underlying processes and pathways, and thereby formulate effective,
7 targeted therapeutic strategies.

8 The cast of genetic elements implicated in tumorigenesis has recently grown as diverse new
9 classes of non-coding RNAs and regulatory features have been discovered. These include the long
10 non-coding RNAs (lncRNAs), of which tens of thousands have been catalogued (Guttman et al. 2009;
11 Jia et al. 2010; Cabili et al. 2011; Derrien et al. 2012). lncRNAs are >200 nt long transcripts with no
12 protein-coding capacity. Their evolutionary conservation and regulated expression, combined with a
13 number of well-characterised examples, have together led to the view that lncRNAs are *bona fide*
14 functional genes (Grote et al. 2013; Sauvageau et al. 2013; Ulitsky & Bartel 2013; Liu et al. 2017).
15 Current thinking holds that lncRNAs function by forming complexes with proteins and RNA both
16 inside and outside the nucleus (Guttman & Rinn 2012; Johnson & Guigó 2014).

17 lncRNAs have been shown to play important roles in various cancers. For example, *MALAT1*,
18 a potent oncogene across numerous cancers, is restricted to the nucleus and plays a housekeeping role
19 in splicing (Gutschner & Diederichs 2012; Engreitz et al. 2014). *MALAT1* is overexpressed in a
20 variety of cancer types, and its knockdown potently reduces not only proliferation but also metastasis
21 *in vivo* (Gutschner et al. 2013). *MALAT1* gene is subjected to elevated mutational rates in human
22 tumours, although it has not yet been established whether these mutations drive tumorigenesis
23 (Lanzós et al. 2017) (PCAWG Consortium, Manuscript in Preparation). On the other hand, lncRNAs
24 may also function as tumour suppressors. *LincRNA-p21* acts as a downstream effector of p53
25 regulation through recruitment of the repressor hnRNP-K (Huarte et al. 2010). These and other
26 examples of lncRNAs linked to cancer, raise the question of how many more remain to be found
27 amongst the ~99% of lncRNAs that are presently uncharacterised (Derrien et al. 2012; Quek et al.
28 2015; Iyer et al. 2015).

29 Recent tumour genome sequencing studies, in step with advanced bioinformatic driver-gene
30 prediction methods, have yielded hundreds of new candidate protein-coding driver genes (Tamborero
31 et al. 2013). For economic reasons, these studies initially restricted their attention to “exomes” or the
32 ~2% of the genome covering protein-coding exons (Chang et al. 2013). Unfortunately such a strategy
33 ignores mutations in the remaining ~98% of genomic sequence, home to the majority of lncRNAs
34 (Gutschner & Diederichs 2012; Derrien et al. 2012). Driver gene identification methods rely on

statistical models that make a series of assumptions about and simplifications of complex tumour mutation patterns (Lawrence et al. 2014). It is critical to test the performance of such methods using true-positive lists of known cancer driver genes. For protein-coding genes, this role has been fulfilled by the Cancer Gene Census (CGC) (Futreal et al. 2004), which is collected and regularly updated by manual annotators. Comparison of driver predictions to CGC genes facilitates further method refinement and comparison between methods (Sjoberg et al. 2006; Redon et al. 2006; Mularoni et al. 2016; Tokheim et al. 2016).

In addition to its benchmarking role, the CGC resource has also been useful in identifying unique biological features of cancer genes. For example, CGC genes tend to be more conserved and longer. Furthermore, they are enriched for genes with transcription regulator activity and nucleic acid binding functions (Furney et al. 2006; Furney et al. 2008).

Until very recently, efforts to discover cancer lncRNAs have depended on classical functional genomics approaches of differential expression using microarrays or RNA sequencing (Huarte et al. 2010; Iyer et al. 2015). While valuable, differential expression *per se* is not direct evidence for causative roles in tumour evolution. To more directly identify lncRNAs that drive cancer progression, a number of methods, including several within the PCAWG Network (PCAWG Consortium, Manuscript in Preparation), have recently been developed to search for signals of positive selection using mutation maps of tumour genomes. OncodriveFML utilises nucleotide-level functional impact scores inferred from predicted changes in RNA secondary structure (Sabarinathan et al. 2013) together with an empirical significance estimate, to identify lncRNAs with an excess of high-impact mutations (Mularoni et al. 2016). Another method, ExInAator, identifies candidates with elevated mutational load, using trinucleotide-adjusted local background (Lanzós et al. 2017). A clear impediment in both cases has been the lack of true-positive set of known lncRNA driver genes, analogous to CGC. Although there do exist databases of cancer lncRNAs, notably lncRNADisease (Chen et al. 2013) and lnc2Cancer (Ning et al. 2016), they mix unfiltered data from numerous sources, resulting in inconsistent criteria for inclusion (including expression changes), and inconsistent gene identifiers.

To facilitate the future discovery of cancer lncRNAs, and gain insights into their biology, we have compiled a highly-curated set of cases with roles in cancer processes. Here we present the *Cancer lncRNA Census* (CLC), the first compendium of lncRNAs with direct functional or genetic evidence for cancer roles. We demonstrate the utility of CLC in assessing the performance of driver lncRNA predictions. Through analysis of this geneset, we demonstrate that cancer lncRNAs have a unique series of features that may in future be used to assist *de novo* predictions. Finally, we show

- 1 that CLC genes have conserved cancer roles across the approximately 80 million years of evolution
- 2 separating humans and rodents.

Results

Definition of cancer related lncRNAs

As part of recent efforts to identify driver lncRNAs by the Drivers and Functional Interpretation Group (PCAWG-2-5-9-14) within the ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Network (henceforth PCAWG), we discovered the need for a high-confidence reference set of cancer-related lncRNA genes, which we henceforth refer to as “cancer lncRNAs”. We here present Version 1 of the *Cancer LncRNA Census* (CLC).

Cancer lncRNAs were identified from the literature using defined and consistent criteria, being direct experimental or genetic evidence for roles in cancer progression or phenotypes (see Materials and Methods). Alterations in expression alone were not considered sufficient evidence. Importantly, only lncRNAs with GENCODE identifiers were included. For every cancer lncRNA, one or more associated cancer types were collected.

Attesting to the value of this approach, we identified several cases in semi-automatically annotated cancer lncRNA databases of lncRNAs that were misassigned GENCODE identifiers, usually with an overlapping protein-coding gene (Chen et al. 2013). We also excluded a number of published lncRNAs for which we could not find evidence to meet our criteria, for example *CONCR*, *SRA1* and *KCNQ1OT1* (Marchese et al. 2016; Lanz et al. 1999; Higashimoto et al. 2006).

Version 1 of CLC contains 122 lncRNA genes, however, eight of them are annotated as pseudogenes rather than lncRNAs by GENCODE. The remaining 114 CLC genes correspond to 0.72% of a total of 15,941 lncRNA gene loci annotated in GENCODE v24 (Derrien et al. 2012; Harrow et al. 2012) (Figure 1). For comparison, the Cancer Gene Census (CGC) (COSMIC v78, downloaded Oct, 3, 2016) lists 561 or 2.8% of protein-coding genes (Futreal et al. 2004). The entire remaining set of 15,827 lncRNA loci is henceforth referred to as “nonCLC” (Figure 1). The full CLC dataset is found in Supplementary Table 1.

The cancer classification terminology used amongst the source literature for CLC was not uniform. Therefore, using the International Classification of Diseases for Oncology (World Health Organization 2013), we reassigned the cancer types described in the original research articles to a reduced set of 29 (Figure 1 and Supplementary Figure 1).

Altogether, CLC contains 333 unique lncRNA-cancer type relationships. Out of 122 genes, 77 (63.1%) were shown to function as oncogenes, 36 (29.5%) as tumour suppressors, and 9 (7.4%) with evidence for both activities (Figure 1 and Supplementary Figure 1).

The most prolific lncRNAs, with ≥ 16 recorded cancer types, are *HOTAIR*, *MALAT1*, *MEG3* and *H19* (Figure 1 and Supplementary Figure 1). It is not clear whether this reflects their unique pan-

1 cancer functionality, or is simply a result of their being amongst the most early-discovered and
2 widely-studied lncRNAs.

3 *In vitro* experiments were the most frequent evidence source, usually consisting of RNAi-
4 mediated knockdown in cultured cell lines, coupled to phenotypic assays such as proliferation or
5 migration (Supplementary Figure 1). Far fewer have been studied *in vivo*, or have cancer-associated
6 somatic or germline mutations. 19 lncRNAs had 3 or more independent evidence sources
7 (Supplementary Figure 1).

8

9 **CLC and other databases**

10 There are a number of relevant lncRNA databases presently available: the *Lnc2Cancer* database
11 (n= 654) (Ning et al. 2016), the *LncRNADisease* Database (n=121) (Chen et al. 2013), *lncRNAdb*
12 (n=191) (Quek et al. 2015) and the “Cancer Related lncRNAs” set we recently produced (n=45)
13 (Lanzós et al. 2017). CLC covers between 17% and 31% of these databases (*Lnc2Cancer* and
14 *LncRNADisease* respectively) but none of these resources contain the complete list of genes presented
15 here (Figure 2A). We sought to use recent unbiased proliferation screen data to independently
16 compare cancer lncRNA databases (Zhu et al. 2016; Liu et al. 2017). Using only GENCODE-
17 annotated genes, CLC is the resource that overall has the highest fraction of independently-identified
18 proliferation lncRNAs, although the sparse nature of the data means that this conclusion is not
19 definitive (Figure 2B).

20

21 **CLC for benchmarking lncRNA driver prediction methods**

22 One of the primary motivations for CLC is to develop a true positive set for benchmarking and
23 comparing methods for identifying driver lncRNAs. In the domain of protein-coding driver gene
24 predictions, the Cancer Gene Census (CGC) has become such a “gold standard” training set (Futreal
25 et al. 2004). Typically, the predicted driver genes belonging to CGC are judged to be true positives,
26 and the fraction of these amongst predictions is used to estimate the Positive Predictive Value (PPV),
27 or precision. This measure can be calculated for increasing cutoff levels, to assess the optimal cutoff.

28 First, we used CLC to examine the performance of the lncRNA driver predictor ExInAtor
29 (Lanzós et al. 2017) in recalling CLC genes using PCAWG tumour mutation data (PCAWG
30 Consortium, Manuscript In Preparation). A total of 2,687 GENCODE lncRNAs were tested here, of
31 which 82 (3.1%) belong to CLC. Driver predictions on several cancers at the standard False Discovery
32 Rate (“*q*-value”) cutoff of 0.1 are shown for selected cancers in Figure 3A. That panel shows the
33 CLC-defined precision (*y*-axis) as a function of predicted driver genes ranked by *q*-value (*x*-axis). We
34 observe rather heterogeneous performance across cancer cohorts. This may reflect a combination of

intrinsic biological differences and differences in cohort sizes, which differs widely between the datasets shown. For the merged pan-cancer dataset, ExInAator predicted three CLC genes amongst its top ten candidates (q -value < 0.1), a rate far in excess of the background expectation (“Baseline”, being the fraction all lncRNAs being in CLC). Similar enrichments are observed for other cancer types. These results support both the predictive value of ExInAator, and the usefulness of CLC in assessing lncRNA driver predictors. Comprehensive CLC-based assessments of lncRNA driver discovery, across all methods and tumour cohorts in PCAWG, may be found in the main PCAWG driver prediction publication (PCAWG Consortium, Manuscript In Preparation).

Finally, we assessed the precision (i.e. positive predictive value) of PCAWG lncRNA and protein-coding driver predictions across all cancers and all prediction methods (PCAWG Consortium, Manuscript In Preparation). Using the same q -value cutoff of 0.1, we found that across all cancer types and methods, a total of 8 (8.5%) of lncRNA predictions belong to CLC (Figure 3B), while a total of 139 (23.1%) of protein-coding predictions belong to CGC (Figure 3C). In terms of sensitivity, 9.8% and 25.1% of CLC and CGC genes are predicted as candidates, respectively. Despite the lower detection of CLC genes in comparison to CGC genes, both sensitivity rates significantly exceed the prediction rate of nonCLC and nonCGC genes ($P=0.007$ and $P<0.001$ Fisher’s exact tests, respectively), again highlighting the usefulness of the CLC geneset (Figure 3C).

CLC genes are distinguished by function- and disease-related features

We recently found evidence, using a smaller set of Cancer Related lncRNAs (CRLs), that cancer lncRNAs are distinguished by various genomic and expression features indicative of biological function (Lanzós et al. 2017). We here extended these findings using a large series of potential gene features, to search for those features distinguishing CLC from nonCLC lncRNAs (Figure 4A).

First, associations with expected cancer-related features were tested (Figure 4B). CLC genes are significantly more likely to have their transcription start site (TSS) within 100 kb of cancer-associated germline SNPs (“Cancer SNPs 100kb TSS”), and more likely to be either differentially-expressed or epigenetically-silenced in tumours (Yan et al. 2015) (Figure 4B). Intriguingly, we observed a tendency for CLC lncRNAs to be more likely to lie within 1 kb of known cancer protein-coding genes (“CGC 1kb TSS”) – this is explored in more detail below. Furthermore, we found that CLC genes are also significantly closer to non-cancer, phenotype-associated germline SNPs (“NonCancer SNPs 100kb TSS”) in comparison to nonCLC genes (Figure 4B), supporting the biological functionality of CLC genes.

We next investigated the properties of the genes themselves. As seen in Figure 4C, and consistent with our previous findings (Lanzós et al. 2017), CLC genes (“Gene length”) and their spliced products

1 (“Exonic length”) are significantly longer than average. No difference was observed in the ratio of
2 exonic to total length (“Exonic content”), nor overall exon repetitive sequence coverage (“Repeats
3 coverage”), nor GC content.

4 CLC genes also tend to have greater evidence of function, as inferred from evolutionary
5 conservation. Base-level conservation at various evolutionary depths was calculated for lncRNA
6 exons and promoters (Figure 4D). Across all measures tested, using either average base-level scores
7 or percent coverage by conserved elements, we found that CLC genes’ exons are significantly more
8 conserved than other lncRNAs (Figure 4D). The same was observed for conservation of promoter
9 regions.

10 High levels of gene expression in normal tissues are known to correlate with lncRNA
11 conservation, and are hypothesized to be a reflection of functionality (Managadze et al. 2011).
12 Additionally, genes with oncogenic roles tend to be highly expressed in cancer samples (Furney et al.
13 2006). We found that CLC has consistently higher steady-state expression levels across PCAWG
14 tumours (Figure 4E), as well as healthy organs and cultured cell lines (Supplementary Figure 2).

15 Finally, we investigated whether CLC transcripts might be initiated by any types of Transposable
16 Elements (TEs) (see Materials and methods). We found that CLC TSSs are enriched for one category,
17 “Simple repeats” (Supplementary Figure 3).

18

19 **Evidence for genomic clustering of non-coding and protein-coding cancer genes**

20 In light of recent evidence for colocalisation and coexpression of disease-related lncRNAs and
21 protein-coding genes (Tan et al. 2017), we were curious whether such an effect holds for cancer-
22 related lncRNAs and protein-coding genes. We asked, more specifically, whether CLC genes tend to
23 be closer to CGC genes than expected by chance, and whether this is manifested in a more co-
24 regulated expression.

25 To this aim, we computed TSS-TSS distances from lncRNAs to protein-coding genes and we
26 found that CLC genes on average tend to lie moderately closer to protein-coding genes of all types,
27 compared to nonCLC lncRNAs (Supplementary Figure 4A, B). Since CLC genes are enriched for
28 functional features (i.e. expression and conservation), we couldn’t rule out the possibility that
29 proximity to protein-coding genes is a feature of functional lncRNAs rather than cancer lncRNA
30 genes. In order to further investigate this possibility, we repeated the analysis dividing the nonCLC
31 set into potentially functional nonCLC genes (PF-nonCLC) (nonCLC genes sampled to match CLC
32 expression and conservation, N=149, Supplementary Figure 5) and “other nonCLC” (the rest of
33 nonCLC). Interestingly, when comparing distances to any type of protein-coding genes, both CLC
34 and PF-nonCLC are significantly closer than the rest of lncRNA (Wilcoxon test, $P=0.03$, 0.007 ,

1 respectively), being the PF-nonCLC genes the closest ones (median 21.9 kb, 29 kb and 37.8 kb, for
2 PF-nonCLC, CLC, and other nonCLC, respectively) (Supplementary Figure 4C). However, when
3 assessing specifically for distance to CGC genes, only CLC set is significantly closer than the rest of
4 lncRNAs (Wilcoxon test, $P=0.0008$) and it represents the group with the lowest distance (median
5 1,122 kb, 1,330kb and 1,607 kb for CLC, PF-nonCLC, and other nonCLC, respectively) (Figure 5A).
6 Thus, although proximity to protein-coding genes seems to be a feature of potentially functional
7 lncRNAs, CLC genes are closer to cancer genes compared to other lncRNAs with similar function-
8 like properties.

9 It has been widely proposed that proximal lncRNA / protein-coding gene pairs are involved in
10 *cis*-regulatory relationships, which is reflected in expression correlation (Ponjavic et al. 2009). We
11 next asked whether proximal CLC-CGC pairs exhibit this behaviour. An important potential
12 confounding factor, is the known positive correlation between nearby gene pairs (Marques et al.
13 2013), and this must be controlled for. Using gene expression data across 11 human cell lines, we
14 observed a positive correlation between CLC-CGC gene pairs for each cell type (Figure 5B). To
15 control for the effect of proximity on correlation, we next randomly sampled a similar number of non-
16 CLC lncRNAs with matched distances (TSS-TSS) from the same CGC genes, and found that this
17 correlation was lost (Figure 5B, “nonCLC-CGC”). To further control for a possible correlation arising
18 from the simple fact that both CGC and CLC genes are involved in cancer, and CLC genes are in
19 general enriched for conservation and expression, we next randomly shuffled the CLC-CGC pairs
20 1000 times, again observing no correlation (Figure 5B, “Shuffled CLC-CGC”). Together these results
21 show that genomically-proximal protein-coding/non-coding gene pairs exhibit an expression
22 correlation that exceeds that expected by chance, even when controlling for genomic distance.

23 These results prompted us to further explore the genomic localization of CLC genes relative to
24 their proximal protein-coding gene and the nature of their neighbouring genes. Next, we observed an
25 unexpected difference in the genomic organisation of CLC genes: when classified by orientation with
26 respect to nearest protein-coding gene (Derrien et al. 2012), we found a significant enrichment of
27 CLC genes immediately downstream and on the same strand as protein-coding genes (“Samestrand,
28 pc up”, Figure 5C). Moreover, CLC genes are approximately twice as likely to lie in an upstream,
29 divergent orientation to a protein-coding gene (“Divergent”, Figure 5C). Of these CLC genes, 20%
30 are divergent to a CGC gene, compared to 5% for nonCLC genes ($P=0.018$, Fisher’s exact test)
31 (Figure 5D), and several are divergent to protein-coding genes that have also been linked or defined
32 to be involved in cancer, despite not being classified as CGCs (Supplementary Table 2).

33 Given this noteworthy enrichment of CGC genes in a divergent configuration to protein-coding
34 genes, we next inspected the latters’ function annotation. Examining their Gene Ontology (GO) terms,

1 molecular pathways and other gene function related terms, we found this group of genes to be
2 enriched in GO terms for “sequence-specific DNA binding”, “DNA binding”, “tube development”
3 and “transcriptional misregulation in cancer” (Figure 5E). These results were confirmed by another,
4 independent GO-analysis suite (see Materials and Methods). Interestingly, three out of the top four
5 functional groups were observed previously in a study of protein-coding genes divergent to long
6 upstream antisense transcripts in primary mouse tissues (Lepoivre et al. 2013).

7 Thus, CLC genes appear to be non-randomly distributed with respect to protein-coding genes,
8 and particularly their CGC subset.

9

10 **Evidence for anciently conserved cancer roles of lncRNAs**

11 In mouse, numerous studies have employed unbiased forward genetic screens to identify genes
12 that either inhibit or promote tumorigenesis (Copeland & Jenkins 2010). These studies use
13 engineered, randomly-integrating transposons carrying bidirectional polyadenylation sites as well as
14 strong promoters. Insertions, or clusters of insertions, called “common insertion sites” (CIS) that are
15 identified in sequenced tumour DNA, implicate the overlapping or neighbouring gene locus as either
16 an oncogene or tumour-suppressor gene. Although these studies have traditionally been focused on
17 identifying protein-coding genes, they can in principle also identify non-coding RNA driver loci.

18 We thus reasoned that comparison of mouse CISs to orthologous human regions could yield
19 independent evidence for the functionality of human cancer lncRNAs (Figure 6A). To test this, we
20 collected a comprehensive set of CISs in mouse (Abbott et al. 2015), consisting of 2,906 loci from 7
21 distinct cancer types (Supplementary Table 4). These sites were then mapped to orthologous regions
22 in the human genome, resulting in 1,309 human CISs, or hCISs. 7.3% of these CISs lie outside of
23 protein-coding gene boundaries, and were used for the following analyses.

24 Mapping hCISs to lncRNA annotations, we discovered altogether eight CLC genes (6.6%)
25 carrying at least one insertion within their gene span: *DLEU2*, *GAS5*, *MONC*, *NEAT1*, *PINT*, *PVT1*,
26 *SLNCRI*, *XIST* (Table 1). In contrast, just 61 (0.4%) nonCLC genes contained hCISs (Figure 6B). A
27 good example is *SLNCRI*, shown in Figure 6C, which drives invasiveness of human melanoma cells
28 (Schmidt et al. 2016), and whose mouse orthologue contains a CIS discovered in pancreatic cancer.
29 We examined the possibility that hCIS insertions in these CLC genes could in fact be caused by
30 nearby, protein-coding cancer genes. However, none of these eight CLC genes are within 100 kb of
31 a CGC gene, with the exception of *CCAT1* lncRNA, lying 58 kb from *c-MYC* oncogene.

32 This analysis would suggest that CLC genes are enriched for hCISs; however, there remains the
33 possibility that this is confounded by their greater length. To account for this, we performed two
34 separate validations. First, sets of nonCLC genes with CLC-matched length were randomly sampled,

1 and the number of intersecting hCISs per unit gene length (Mb) was counted (Supplementary Figure
2 6A). Second, CLC genes were randomly relocated in the genome, and the number of genes
3 intersecting at least one hCIS was counted (Supplementary Figure 6B). Both analyses showed that
4 the number of intersecting hCISs per Mb of CLC gene span is far greater than expected. In contrast,
5 nonCLC genes show a depletion for hCIS sites (Supplementary Figure 6C).

6 We further compared the enrichment of hCIS in protein-coding genes, lncRNA genes and other
7 intergenic space. Compared to the genomic space they occupy, there is a clear enrichment of hCIS
8 elements in both protein-coding CGC genes, as well as CLC lncRNAs (Figure 6D). Expressed as
9 insertion rate per megabase of gene span, it is clear that CLC genes are targeted more frequently than
10 background intergenic DNA and non-cancer-related protein-coding genes. Of note are the non-
11 background insertion rates for non-cancer-related protein-coding and lncRNA genes, suggesting that
12 there remain substantial numbers of undiscovered cancer genes in both groups.

13 Together these analyses demonstrate that CLC genes are orthologous to mouse cancer-causing
14 genomic loci at a rate greater than expected by random chance. These identified cases, and possibly
15 other CLC genes, display cancer functions that have been conserved over tens of millions of years
16 since human-rodent divergence.

Discussion

We have presented the Cancer lncRNA Census, the first controlled set of GENCODE-annotated lncRNAs with demonstrated roles in tumorigenesis or cancer phenotypes.

The present state of knowledge of lncRNAs in cancer, and indeed lncRNAs generally, remains highly incomplete. Consequently, our aim was to create a geneset with the greatest possible confidence, by eliminating the relatively large number of published “cancer lncRNAs” with as-yet unproven causative roles in disease processes. Thus, we used a rather strict definition of cancer lncRNA, being those having direct experimental or genetic evidence for a causative role in cancer phenotypes. By this measure, gene expression changes alone do not suffice. By introducing these well-defined inclusion criteria, we hope to ensure that CLC contains the highest possible proportion of *bona fide* cancer genes, giving it maximum utility for *de novo* predictor benchmarking. In addition, its basis in GENCODE ensures portability across datasets and projects. Inevitably some well-known lncRNAs did not meet these criteria (including *SRA1*, *CONCR*, *KCNQ1OT1*) (Marchese et al. 2016; Lanz et al. 1999; Higashimoto et al. 2006); these may be included in future when more validation data becomes available. We believe that CLC will complement the established lncRNA databases such as *lncRNAdb*, *lncRNADisease* and *lnc2Cancer*, which are more comprehensive, but are likely to have a higher false-positive rate due to their more relaxed inclusion criteria (Chen et al. 2013; Quek et al. 2015; Ning et al. 2016).

De novo lncRNA driver gene discovery is likely to become increasingly important as the number of sequenced tumours grow. The creation and refinement of statistical methods for driver gene discovery will depend on the available of high-quality true positive genesets such as CLC. It will be important to continue to maintain and improve the CLC in step with anticipated growth in publications on validated cancer lncRNAs. Very recently, CRISPR-based screens (Zhu et al. 2016; Liu et al. 2017) have catalogued large numbers of lncRNAs contributing to proliferation in cancer cell lines, which will be incorporated in future versions.

We used CLC to estimate the performance of *de novo* driver lncRNA predictions from the PCAWG project, made using the ExInAor pipeline (Lanzós et al. 2017). Supporting the usefulness of this approach, we found an enrichment for CLC genes amongst the top-ranked driver predictions. Extending this to the full set of PCAWG driver predictors, approximately ten percent of CLC genes (9.8%) are called as drivers by at least one method (PCAWG Consortium, Manuscript In Preparation), which is lower to the rate of CGC genes identified (25.1%).

The low rate of concordance between *de novo* predictions and CLC genes may be due to technical or biological factors. Indeed, it is important to state that we do not yet know whether CLC holds “cancer driver” lncRNAs, and indeed, how many such genes exist. In principle, lncRNAs may

play two distinct roles in cancer: first, as driver genes, defined as those whose mutations are early and positively-selected events in tumorigenesis; or second, as “downstream genes”, which do make a genuine contribution to cancer phenotypes, but through non-genetic alterations in cellular networks resulting from changes in expression, localisation or molecular interactions. These downstream genes may not display positively-selected mutational patterns, but would be expected to display cancer-specific alterations in expression. A key question for the future is how lncRNAs break down between these two categories, and the utility of CLC in benchmarking *de novo* driver predictions will depend on this. However, the identification of lncRNAs whose silencing or overexpression is sufficient for tumour formation in mouse, would seem to suggest that they are true “driver genes”.

Analysis of the CLC geneset has broadened our understanding of the unique features of cancer lncRNAs, and generally supports the notion that lncRNAs have intrinsic biological functionality. Cancer lncRNAs are distinguished by a series of features that are consistent with both (a) roles in cancer (eg tumour expression changes), and (b) general biological functionality (eg high expression, evolutionary conservation). Elevated evolutionary conservation in the exons of CLC genes would appear to support their functionality as a mature RNA transcript, in contrast to the act of their transcription alone (Latos et al. 2012). Another intriguing observation has been the colocalisation of cancer lncRNAs with known protein-coding cancer genes: these are genomically proximal and exhibit elevated expression correlation. This points to a regulatory link between cancer lncRNAs and protein-coding genes, perhaps through chromatin looping, as described in previous reports for *CCAT1* and *MYC*, for example (Xiang et al. 2014).

One important caveat for all features discussed here is ascertainment bias: almost all lncRNAs discussed have been curated from published, single-gene studies. It is entirely possible that selection of genes for initial studies was highly non-random, and influenced by a number of factors – including high expression, evolutionary conservation and proximity to known cancer genes – that could bias our inference of lncRNA features. This may be the explanation for the observed excess of cancer lncRNAs in divergent configuration to protein-coding genes. However, the general validity of some of the CLC-specific features described here – including high expression and evolutionary conservation – were also observed recent unbiased genome-wide screens (Lanzós et al. 2017; Liu et al. 2017), suggesting that they are genuine.

Despite the relatively low concordance of CLC genes with PCAWG driver predictions, the results of this study strongly support the value and key cancer role of identified lncRNAs in cancer. Most notably, the existence of a core set of eight lncRNAs with independently-identified mouse orthologues with similar cancer functions, is a powerful evidence that these genes are *bona fide* cancer genes, whose overexpression or silencing can drive tumour formation. To our knowledge this is the

- 1 most direct demonstration to date of anciently-conserved functions and disease roles for lncRNAs. It
- 2 will be intriguing to investigate in future whether more human-mouse orthologous lncRNAs have
- 3 been identified in such screens.
- 4

Materials and Methods

Manual Curation

All lncRNAs in lncRNADB and those listed in Schmitt and Chang's recent review article were collected (Quek et al. 2015; Schmitt et al. 2016). To these were added all cases from *lncRNADisease* and *lnc2Cancer* databases (Chen et al. 2013; Ning et al. 2016). This primary list formed the basis for a manual literature search: all available publications for each gene were identified by keyword search in Pubmed. If publications were found conforming to at least one of the inclusion criteria (below) and the gene has a GENCODE ID, then it was added to CLC, with appropriate information on the associated cancer, biological activity. For the numerous cases where no GENCODE ID was supplied in the original publication, any available ID, or primer or siRNA sequence was used to identify the gene using the UCSC Genome Browser Blat tool (Kent et al. 2002).

Inclusion criteria sufficient to define a cancer lncRNA and link it to a cancer type were:

1. Class t: *In vitro* demonstration that their knockdown and/or overexpression in cultured cancer cells results in changes to cancer-associated phenotypes. These typically include proliferation rates, migration, sensitivity to apoptosis, or anchorage-independent growth.

2. Class v: *In vivo* demonstration that their knockdown and/or overexpression in cancer cells alters their tumorigenicity when injected into animal models.

3. Class g: Germline mutations or variants that predispose humans to cancer.

4. Class s: Somatic mutations that show evidence for positive selection during tumour formation.

An additional criteria was allowed to link an lncRNA to a cancer type, only if at least one of the above criteria was already met for another cancer:

5. Class p: Prognosis: The lncRNAs expression is statistically linked to disease progression or response to treatment.

If an lncRNA was found to promote tumorigenesis or cancer phenotype, it was defined as "oncogene" (og). Conversely those found to inhibit such phenotypes were defined as "tumour suppressor" (tsg). Several lncRNAs were found to have both activities recorded in different cancer types, and were given both labels (og/tsg). For every lncRNA-cancer association, a single representative publication is recorded. Finally, it is important to note that no lncRNAs were included based on evidence from previous driver gene discovery studies of the types represented by OncodriveFML, ExInAtoR, ncdDetect or others described in PCAWG (Mularoni et al. 2016; Lanzós et al. 2017; Juul et al. 2017) (PCAWG Consortium, Manuscript In Preparation).

CLC set at this stage relies on GENCODE v24 annotation, and therefore all CLC genes have a GENCODE v24 ID assigned. However, data relative to GENCODE v24 was not available for all types of data and analysis used in this study (ie all data relative to PCAWG is based on GENCODE v19). Thus, for some analysis only genes also present in GENCODE v19 could be used (specified in the corresponding methods section) and the total number of genes analysed in these cases is slightly lower (107 instead of 122 CLC genes and 13,503 instead of 15,827 nonCLC).

LncRNA and protein-coding driver prediction analysis

LncRNA and protein-coding predictions for ExInAator and the rest of PCAWG methods, as well as the combined list of drivers, were extracted from the consortium database (PCAWG Consortium, Manuscript In Preparation). Parameters and details about each individual methods and the combined list of drivers can be found on the main PCAWG driver publication (PCAWG Consortium, Manuscript In Preparation) and false discovery rate correction was applied on each individual cancer type for each individual method in order to define candidates (q -value cutoff of 0.1). This way, we combined the predicted candidates of each individual method in each individual cancer type (including pan-cancer). To calculate sensitivity (percentage of true positives that are predicted as candidates) and precision (percentage of predicted candidates that are true positives) for lncRNA and protein-coding predictions we used the CLC and CGC (COSMIC v78, downloaded Oct, 3, 2016) sets, respectively. To assess the statistical significance of sensitivity rates, we used Fisher's exact test.

Feature Identification

We compiled several quantitative and qualitative traits of GENCODE lncRNAs and used them to compare CLC genes to the rest of lncRNAs (referred to as "nonCLC"). Analysis of quantitative traits were performed using Wilcoxon test while qualitative traits were tested using Fisher's exact test. These methods principally refer to Figure 4 and 5 as well as Supplementary Figures 2, 3, 4 and 5.

Cancer SNPs: On October, 4, 2016, we collected all 2,192 SNPs related to "cancer", "tumour" and "tumor" terms in the NHGRI-EBI Catalog of published genome-wide association studies (Hindorff et al. 2009; Welter et al. 2014) (<https://www.ebi.ac.uk/gwas/home>). Then we calculated the closest SNP to each lncRNA TSS using *closest* function from Bedtools v2.19 (Quinlan & Hall 2010) (GENCODE v24).

NonCancer SNPs: On July, 31, 2017, we collected all 29,813 SNPs not related to "cancer", "tumour" and "tumor" terms in the NHGRI-EBI Catalog of published genome-wide association studies (Hindorff et al. 2009; Welter et al. 2014) (<https://www.ebi.ac.uk/gwas/home>). Then we

1 calculated the closest SNP to each lncRNA TSS using *closest* function from Bedtools v2.19 (Quinlan
2 & Hall 2010)(GENCODE v24).

3 Epigenetically-silenced lncRNAs: We obtained a published list of 203 cancer-associated
4 epigenetically-silenced lncRNA genes present in GENCODE v24 (Yan et al. 2015). These candidates
5 were identified due to DNA methylation alterations in their promoter regions affecting their
6 expression in several cancer types.

7 Differentially expressed in cancer: We collected a list of 3,533 differentially-expressed lncRNAs
8 in cancer compared to normal samples (Yan et al. 2015) (GENCODE v24).

9 Sequence / gene properties: Exonic positions of each gene were defined as the projection of the
10 union of exons from all its transcripts. Introns were defined as all remaining non-exonic nucleotides
11 within the gene span. Repeats coverage refers to the percent of exonic nucleotides of a given gene
12 overlapping repeats and low complexity DNA sequence regions obtained from RepeatMasker data
13 housed in the UCSC Genome Browser (Tyner et al. 2017). Exonic content refers to the fraction of
14 total gene span covered by exons. For this section we used GENCODE v19.

15 Evolutionary conservation: Two types of PhastCons conservation data were used: base-level
16 scores and conserved elements. These data for different multispecies alignments (GRCh38/hg38)
17 were downloaded from UCSC genome browser (Tyner et al. 2017). Mean scores and percent overlap
18 by elements were calculated for exons and promoter regions (GENCODE v24). Promoters were
19 defined as the 200nt region centred on the annotated gene start.

20 Expression: We used polyA+ RNA-seq data from 10 human cell lines produced by ENCODE
21 (Djebali et al. 2012; ENCODE Project Consortium et al. 2012), from various human tissues by the
22 Illumina Human Body Map Project (HBM) (www.illumina.com; ArrayExpress ID: E-MTAB-513),
23 and from cancer samples from PCAWG RNAseq expression data (PCAWG Consortium, Manuscript
24 In Preparation). In this last case, for each cancer type we computed the expression mean of genes
25 across all RNAseq samples belonging to that cancer type (GENCODE v19).

26 Transposable elements: We downloaded 5,520,016 transposable elements from the UCSC table
27 browser (Karolchik et al. 2004) on August, 3, 2017. We separated them by element types and counted
28 how many of them intersected or not with the transcription start sites of CLC and nonCLC genes, in
29 order to detect any association with the Fisher' exact test.

30 Distance to protein-coding genes and CGC genes: For each lncRNA we calculated the TSS to
31 TSS distance to the closest protein-coding gene (GENCODE v24) or CGC gene (downloaded on
32 October, 3, 2016 from Cosmic database) (Futreal et al. 2004) using *closest* function from Bedtools
33 v2.19 (Quinlan & Hall 2010). In order to divide nonCLC genes into potentially functional nonCLC
34 (PF-nonCLC) and others, we sampled the list of all nonCLC genes to get a subsample that has a

1 matched distribution to CLC genes in conservation (% of conserved elements, from Vertebrate Multiz
2 Alignment 100 Species from UCSC genome browser data, in exonic regions). Then we sampled again
3 the resulting subset to get a final subset that also matches CLC genes in terms of expression (median
4 of expression across 16 human tissues, data from Illumina Human Body Map Project (HBM)). To
5 create the nonCLC samples we used the *matchDistribution* script:
6 <https://github.com/julienlag/matchDistribution>.

7 Coexpression with closest CGC gene: We took CLC-CGC gene pairs whose TSS-TSS distance
8 was <200kb. RNAseq data from 11 human cell lines from ENCODE was used to assess expression
9 levels (Djebali et al. 2012; ENCODE Project Consortium et al. 2012). ENCODE RNAseq data were
10 obtained from ENCODE Data Coordination Centre (DCC) in September 2016,
11 <https://www.encodeproject.org/matrix/?type=Experiment>. All data is relative to GENCODE v24. We
12 calculated the expression correlation of gene pairs within each of the 11 cell lines, using the Pearson
13 measure. To control for the effect of proximity, we randomly sampled a subset of nonCLC-CGC pairs
14 matching the same TSS-TSS distance distribution as above, and performed the same expression
15 correlation analysis (“nonCLC-CGC”). Finally, to further control for the fact that CLC and CGC are
16 both cancer genes, which may influence their expression correlation, we shuffled CLC-CGC pairs
17 1000 times, and tested expression correlation for each set (“Shuffled CLC-CGC”).

18 Genomic classification: We used an in-house script to classify lncRNA transcripts into different
19 genomic categories based on their orientation and proximity to the closest protein-coding gene
20 (GENCODE v24): a 10 kb distance was used to distinguish “genic” from “intergenic” lncRNAs.
21 When transcripts belonging to the same gene had different classifications, we used the category
22 represented by the largest number of transcripts.

23 Functional enrichment analysis: The list of protein-coding genes (GENCODE v24) that are
24 divergent and closer than 10 kb to CLC genes (or nonCLC) was used for a functional enrichment
25 analysis (20 unique genes in the case of CLC analysis and 1202 in the case of nonCLC analysis). We
26 show data obtained using g:Profiler web server (Reimand et al. 2016), g:GOST, with default
27 parameters for functional enrichment analysis of protein-coding genes divergent to CLC and using
28 Bonferroni correction for protein-coding gene divergent to nonCLC. For CLC analysis we performed
29 the same test with independent methods: Metascape (<http://metascape.org>) (Tripathi et al. 2015) and
30 GeneOntology (Panther classification system)(Mi et al. 2013; Mi et al. 2017). In both cases similar
31 results were found.

32
33
34

1 **Mouse mutagenesis screen analysis**

2 We extracted the genomic coordinates of transposon common insertion sites (CISs) in Mouse
3 (GRCm38/mm10) <http://ccgd-starrlab.oit.umn.edu/about.php> (Abbott et al. 2015). This database
4 contains target sites identified by transposon-based forward genetic screens in mice. LiftOver (Kent
5 et al. 2002) was used at default settings to obtain aligned human genome coordinates (hCISs)
6 (GRCh38/hg38). We discarded hCIS regions longer than 1000 nucleotides and those that overlap
7 protein-coding genes, and intersected the remainders with the genomic coordinates CLC and nonCLC
8 genes.

9 To correctly assess the statistical enrichment of CLC in hCIS regions, we performed two control
10 analyses:

11 Randomly repositioning of CLC and nonCLC genes: We randomly relocated CLC/nonCLC
12 genes 10,000 times within the whole genome using the tool *shuffle* from BedTools v19 (Quinlan &
13 Hall 2010). In each iteration, we calculated the number of genes that intersected at least one hCIS,
14 and created the distribution of these simulated values. Finally, we calculated an empirical *p*-value by
15 counting how many of the simulated values were higher or equal than the real values. This analysis
16 was performed separately for CLC and nonCLC genes.

17 Length-matched sampling: To calculate if the enrichment of hCIS intersecting genes in CLC set
18 is higher and statistically different from nonCLC set, while controlling by gene length, we created
19 1000 samples of nonCLC genes with the same gene length distribution as CLC genes. Each sample
20 was intersected with hCIS, and the number of intersecting hCISs per Mb of gene length was
21 calculated. To create the nonCLC samples we used the *matchDistribution* script:
22 <https://github.com/julienlag/matchDistribution>.

Acknowledgements

We wish to thanks Julien Lagarde (CRG) for help and advice in bioinformatic analysis. We acknowledge Romina Garrido (CRG), Deborah Re (DKF), Silvia Roesselet (DKF) and Marianne Zahn (Inselspital) for administrative support. We thank Ivo Buchhalter (DKFZ) and Sandra Koser (DKFZ) for preprocessing the SNV and expression data for the integrated analysis. Iñigo Martincorena (Sanger Institute) kindly provided the script for analysing driver prediction sensitivity. A.L. is supported by pre-doctoral fellowship FPU14/03371. This research was supported by the NCCR RNA & Disease funded by the Swiss National Science Foundation.

9

Contributions

RJ conceived the project, performed manual annotation of CLC, and supervised with advice and suggestions of JS-P, LF and CH. JCF and AL performed the feature analysis and evolutionary analysis. AL performed mutation analysis. RJ, AL and JCF drafted the manuscript and prepared the figures and supplementary material. All authors read and approved the final draft.

6

Figure Legends

Figure 1: Overview of the Cancer lncRNA Census. Rows represent the 122 CLC genes, columns represent 29 cancer types. Asterisk next to gene names indicate that they are predicted as drivers by PCAWG, based either on gene or promoter evidence (see Supplementary Table 1).. Blue cells indicate evidence for the involvement of a given lncRNA in that cancer type. Left column indicates functional classification: tumour suppressor (TSG), oncogene (OG) or both (OG/TSG). Above and to the right, barplots indicate the count totals of each column / row. The piechart shows the fraction that CLC within GENCODE v24 lncRNAs. Note that 8 CLC genes are classified as “pseudogenes” by GENCODE. “nonCLC” refers to all other GENCODE-annotated lncRNAs, which are used as background in comparative analyses.

Figure 2: Intersection of CLC with public databases. (A) Proportional Venn diagrams displaying the overlap between CLC set and the three indicated databases. Shown are the total number of unique human lncRNAs contained in each intersection (note that for LncRNADisease, numbers refer only to cancer-related genes). Databases are divided into genes that belong to GENCODE v24 annotation and others. (B) Barplot shows the percent of GENCODE v24 lncRNAs of each database that is present in the final list of cancer lncRNA candidates of two CRISPR cancer screenings (Liu et al. 2016 and Zhu et al. 2016). N represents the number of GENCODE v24 lncRNAs that could be used for the analysis. Names of the genes that overlap between the databases and the screenings are shown in each bar.

Figure 3: CLC as benchmark for cancer driver predictions. (A) CLC benchmarking of ExInAto driver lncRNA predictions using PCAWG whole genome tumours at q -value (false discovery rate) cutoff of 0.1. Genes sorted increasingly by q -value are ranked on x axis. Percentage of CLC genes amongst cumulative set of predicted candidates at each step of the ranking (precision), are shown on the y axis. Black line shows the baseline, being the percentage of CLC genes in the whole list of genes tested. Coloured dots represent the number of candidates predicted under the q -value cutoff of 0.1. “n” in the legend shows the number of CLC and total candidates for each cancer type. (B) Rate of driver gene predictions amongst CLC and nonCLC genesets (q -value cutoff of 0.1) by all the individual methods and the combined list of drivers developed in PCAWG. p -value is calculated using Fisher’s exact test for the difference between CLC and nonCLC genesets. (C) Rate of driver gene predictions amongst CGC and nonCGC genesets (q -value cutoff of 0.1) by all the individual methods and the combined list of drivers developed in PCAWG. p -value is calculated using Fisher’s exact test for the difference between CGC and nonCGC genesets.

Figure 4: Distinguishing features of CLC genes. (A) Panel showing a hypothetical feature analysis example to illustrate the content of the following figures. All panels in figure 4 display

1 features (dots), plotted by their log fold difference (Odds Ratio in case of panel “B”) between CLC /
2 nonCLC genesets (y-axis) and statistical significance (x-axis). In all plots dark and light green dashed
3 lines indicate 0.05 and 0.01 significance thresholds, respectively. (B) Cancer and noncancer disease
4 related data from indicated sources: y-axis shows the log2 of the Odds Ratio obtained by comparing
5 CLC to nonCLC by Fisher’s exact test; x-axis displays the estimated *p*-value from the same test.
6 “CGC 1 kb TSS” refers to the fraction of genes that have a nearby known CGC cancer protein-coding
7 gene. This is explored in more detail in the next Figure. “Noncancer SNPs” refers to GWAS SNPs
8 associated with diseases/traits other than cancer (C) Sequence and gene properties: y-axis shows the
9 log2 fold-difference of CLC / nonCLC means; x-axis represents the *p*-value obtained. (D)
10 Evolutionary conservation: “Phastc mean” indicates average base-level PhastCons score; “Elements”
11 indicates percent coverage by PhastCons conserved elements (see Materials and Methods). Colours
12 distinguish exons (blue) and promoters (purple). (E) Tumour RNAseq: expression levels of lncRNA
13 genes in different cancer tissues obtained from RNAseq expression data from PCAWG. For B-D,
14 statistical significance was calculated using Wilcoxon test.

15 **Figure 5: Evidence for genomic clustering of non-coding and protein-coding cancer genes.**

16 (A) Cumulative distribution of the genomic distance of lncRNA transcription start site (TSS) to the
17 closest Cancer Gene Census (CGC) (protein-coding) gene TSS. LncRNAs are divided into CLC
18 (n=122), potentially functional nonCLC genes (PF-nonCLC) (n=149), and other nonCLC genes
19 (n=15678) (B) Boxplot shows the distribution of the gene expression correlation between CLC and
20 their closest CGC genes in 11 human cell lines, including two control analyses (distance-matched
21 nonCLC-CGC pairs, and shuffled CLC-CGC pairs). Correlation was calculated for gene pairs within
22 each cell type, using Pearson method. *p*-value for Kolmogorov-Smirnov test is shown. (C) Genomic
23 classification of lncRNAs. Genes are classified according to distance and orientation to the closest
24 protein-coding gene, and these are grouped into three categories: genes closer than 10kb to closest
25 protein-coding gene, genes overlapping a protein-coding gene and intergenic genes (>10kb from
26 closest protein-coding gene). *p*-values for Fisher’s exact tests are shown. (D) The percentage of
27 divergent CLC (left bar) and nonCLC (right bar) genes divergent to a cancer protein-coding gene
28 (CGC). Numbers represent numbers of genes with which the percentage is calculated. *p*-value for
29 Fisher’s exact test is shown. (E) Functional annotations of the 20 protein-coding genes divergent to
30 CLC genes from Panel C. Bars indicate the $-\log_{10}$ (corrected) *p*-value (see Materials and Methods)
31 and are coloured based on the “enrichment”: the number of genes that contain the functional term
32 divided by the total number of queried genes. Numbers at the end of the bars correspond to the number
33 of genes that fall into the category.

Figure 6: Evidence for ancient conserved cancer roles of lncRNAs. (A) Functional conservation of human CLC genes was inferred by the presence of Common Insertion Sites (CIS), identified in transposon mutagenesis screens, at orthologous regions in the mouse genome. Orthology was inferred from Chain alignments and identified using LiftOver utility. (B) Number of CLC and nonCLC genes that contain human orthologous common insertion sites (hCIS) (see Table 1). Significance was calculated using Fisher's exact test. (C) UCSC browser screenshot of a CLC gene (*SLNCRI*, ENSG00000227036) intersecting a CIS (yellow arrow). (D) Number of basepairs and number of overlapping hCIS for cancer driver protein-coding genes (CGC), non cancer driver protein-coding genes (nonCGC), cancer related lncRNAs (CLC), rest of GENCODE lncRNAs (nonCLC) and the rest of the genome that do not overlap any of the previous element types (intergenic). Arrows indicate the number of hCIS and the percentage for each element type. (E) Number of overlapping hCIS per megabase of genomic span for each gene class.

Supplementary Figure Legends

Supplementary Figure 1: CLC summary statistics. (A) Barplot showing the non-redundant number of genes in CLC broken down by supporting evidence types. p: prognostic; t: in vitro; v: in vivo; g: germline mutations; s: somatic mutations. (B) Similar as previous, but with (redundant) number of genes per individual evidence type. (C) Histogram of genes broken down by their number of associated cancer types. (D) Histogram of cancer types, by their (redundant) number of associated lncRNAs.

Supplementary Figure 2: CLC genes are highly expressed. Panels display feature analysis results similar to Figure 4 using other datasets. (A) Panel displaying the log fold difference between CLC and nonCLC genesets (y-axis) and statistical significance by Wilcoxon test (x-axis) when comparing RNAseq expression levels in human tissues (each dot represents a different tissue) from Human Body Map data. (B) Same than previous panel for expression data in human cell lines instead of tissues, from ENCODE RNAseq data.

Supplementary Figure 3: CLC TSSs association with Transposable Elements. Figure shows the comparison of the intersection of each category of Transposable elements with transcription start sites (TSS) of CLC and nonCLC genes. y-axis shows the log2 of the Odds Ratio obtained by comparing CLC to nonCLC by Fisher's exact test; x-axis displays the estimated p-value from the same test.

Supplementary Figure 4: CLC lncRNAs tend to be closer to protein-coding genes. (A) Cumulative distribution of the genomic distance from CLC and nonCLC genes, to the closest protein-coding gene (NB this may be a CGC gene or not). Distances are defined as the distance of the annotated transcription start site (TSS) of each gene in the pair. p-value for Wilcoxon test is shown. (B) Same as (A) for genomic distance to closest CGC genes. (C) Same than A dividing nonCLC genes into two groups: potentially functional nonCLC (PF-nonCLC) (those nonCLC genes that are in the same range of expression and conservation than CLC genes) and other nonCLC (the rest of nonCLC genes).

Supplementary Figure 5: sampling nonCLC genes. (A) Density plot comparing the percentage of PhastCons conserved elements in lncRNA exons of CLC genes and a subset of nonCLC sampled to match CLC conservation distribution. (B) Same than A but comparing the median of RNAseq expression values across 16 human tissues. NonCLC subset here is sampled from the subset obtained after matching conservation distribution.

Supplementary Figure 6: hCIS enrichment corrected by gene length. (A) Distribution of the number of intersecting hCIS per Megabase (Mb) of total gene length, for 1000 subsets of nonCLC genes with same length distribution as CLC genes (grey). Vertical blue line represents the overall

- 1 value for CLC geneset: 1.42 hCIS sites per Mb of gene span. (B) Distribution of the number of genes
- 2 overlapping a hCIS after 10,000 genomic randomizations of CLC genes (blue). Vertical black line
- 3 represents the observed number of CLC genes (8) that intersect a hCIS. (C) Distribution of the number
- 4 of intersecting genes with a hCIS after 10,000 genomic randomizations of nonCLC genes (grey).
- 5 Vertical black line represents the observed number of nonCLC genes that intersect a hCIS (64).
- 6

1 **Supplementary Tables:**

2 **Supplementary Table 1: full CLC set.**

3 **Supplementary Table 2: CLC – protein-coding pairs.**

4 **Supplementary Table 3: GO analysis for protein-coding genes divergent to nonCLC genes.**

5 **Supplementary Table 4: Counts of mouse CIS per cancer type.**

References

- Abbott, K.L. et al., 2015. The candidate cancer gene database: A database of cancer driver genes from forward genetic screens in mice. *Nucleic Acids Research*, 43(D1), pp.D844–D848.
- Cabili, M.N. et al., 2011. Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses. *Genes & development*, 25(18), pp.1915–27.
Available at: <http://www.ncbi.nlm.nih.gov/pubmed/21890647> [Accessed March 7, 2017].
- Chang, K. et al., 2013. The Cancer Genome Atlas Pan-Cancer analysis project. *Nature genetics*, 45(10), pp.1113–1120. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/24071849>.
- Chen, G. et al., 2013. LncRNADisease: a database for long-non-coding RNA-associated diseases. *Nucleic acids research*, 41(Database issue), pp.D983-6. Available at:
<http://www.ncbi.nlm.nih.gov/pubmed/23175614> [Accessed March 2, 2017].
- Copeland, N.G. & Jenkins, N.A., 2010. Harnessing transposons for cancer gene discovery. *Nature reviews. Cancer*, 10(10), pp.696–706. Available at:
<http://www.nature.com/doifinder/10.1038/nrc2916> [Accessed March 7, 2017].
- Derrien, T. et al., 2012. The GENCODE v7 catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and expression. *Genome research*, 22(9), pp.1775–89.
Available at: <http://genome.cshlp.org/content/22/9/1775.long> [Accessed May 23, 2014].
- Djebali, S. et al., 2012. Landscape of transcription in human cells. *Nature*, 489(7414), pp.101–108.
- ENCODE Project Consortium, T. et al., 2012. An integrated encyclopedia of DNA elements in the human genome. *Nature*, 488.
- Engreitz, J.M. et al., 2014. RNA-RNA interactions enable specific targeting of noncoding RNAs to nascent Pre-mRNAs and chromatin sites. *Cell*, 159(1), pp.188–99. Available at:
<http://www.ncbi.nlm.nih.gov/pubmed/25259926> [Accessed March 2, 2017].
- Furney, S. et al., 2006. Structural and functional properties of genes involved in human cancer. *BMC Genomics*, 7(1), p.3. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/16405732>
[Accessed February 24, 2017].
- Furney, S.J. et al., 2008. Distinct patterns in the regulation and evolution of human cancer genes. *In Silico Biol*, 8(December 2007), pp.33–46.
- Futreal, P. et al., 2004. A census of human cancer genes. *Nat Rev Cancer*, 4(3), pp.177–183.
- Grote, P. et al., 2013. *The Tissue-Specific lncRNA Fendrr Is an Essential Regulator of Heart and Body Wall Development in the Mouse*, Available at:
<http://www.sciencedirect.com/science/article/pii/S1534580712005862> [Accessed April 25, 2017].

- 1 Gutschner, T. et al., 2013. The noncoding RNA MALAT1 is a critical regulator of the metastasis
2 phenotype of lung cancer cells. *Cancer research*, 73(3), pp.1180–9. Available at:
3 <http://www.ncbi.nlm.nih.gov/pubmed/23243023> [Accessed March 2, 2017].
- 4 Gutschner, T. & Diederichs, S., 2012. The hallmarks of cancer: a long non-coding RNA point of
5 view. *RNA biology*, 9(6), pp.703–19. Available at:
6 <http://www.ncbi.nlm.nih.gov/pubmed/22664915> [Accessed June 28, 2016].
- 7 Guttman, M. et al., 2009. Chromatin signature reveals over a thousand highly conserved large non-
8 coding RNAs in mammals. *Nature*, 458(7235), pp.223–7. Available at:
9 <http://www.ncbi.nlm.nih.gov/pubmed/19182780> [Accessed March 7, 2017].
- 10 Guttman, M. & Rinn, J.L., 2012. Modular regulatory principles of large non-coding RNAs. *Nature*,
11 482(7385), pp.339–46. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/22337053>
12 [Accessed March 2, 2017].
- 13 Harrow, J. et al., 2012. GENCODE: The reference human genome annotation for The ENCODE
14 Project. *Genome Research*, 22(9), pp.1760–1774. Available at:
15 <http://www.ncbi.nlm.nih.gov/pubmed/22955987> [Accessed April 25, 2017].
- 16 Higashimoto, K. et al., 2006. Imprinting disruption of the CDKN1C/KCNQ1OT1 domain: the
17 molecular mechanisms causing Beckwith-Wiedemann syndrome and cancer. *Cytogenetic and*
18 *genome research*, 113(1–4), pp.306–12. Available at:
19 <http://www.ncbi.nlm.nih.gov/pubmed/16575194> [Accessed April 27, 2017].
- 20 Hindorff, L. a et al., 2009. Potential etiologic and functional implications of genome-wide
21 association loci for human diseases and traits. *Proceedings of the National Academy of*
22 *Sciences of the United States of America*, 106(23), pp.9362–7. Available at:
23 <http://www.ncbi.nlm.nih.gov/pubmed/19474294>.
- 24 Huarte, M. et al., 2010. A large intergenic noncoding RNA induced by p53 mediates global gene
25 repression in the p53 response. *Cell*, 142(3), pp.409–19. Available at:
26 <http://www.ncbi.nlm.nih.gov/pubmed/20673990> [Accessed February 25, 2017].
- 27 Iyer, M.K. et al., 2015. The landscape of long noncoding RNAs in the human transcriptome. *Nature*
28 *GeNetics*, 47(3). Available at: <https://www.nature.com/ng/journal/v47/n3/pdf/ng.3192.pdf>
29 [Accessed April 19, 2017].
- 30 Jia, H. et al., 2010. Genome-wide computational identification and manual annotation of human
31 long noncoding RNA genes. *RNA (New York, N.Y.)*, 16(8), pp.1478–87. Available at:
32 <http://rnajournal.cshlp.org/cgi/doi/10.1261/rna.1951310> [Accessed March 7, 2017].
- 33 Johnson, R. & Guigó, R., 2014. The RIDL hypothesis: transposable elements as functional domains
34 of long noncoding RNAs. *RNA (New York, N.Y.)*, 20(7), pp.959–76. Available at:

1 <http://www.ncbi.nlm.nih.gov/pubmed/24850885> [Accessed March 2, 2017].

2 Juul, M. et al., 2017. Non-coding cancer driver candidates identified with a sample- and position-
3 specific model of the somatic mutation rate. *eLife*, 6. Available at:
4 <http://www.ncbi.nlm.nih.gov/pubmed/28362259> [Accessed April 27, 2017].

5 Karolchik, D. et al., 2004. The UCSC Table Browser data retrieval tool. *Nucleic acids research*,
6 32(Database issue), pp.D493-6. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/14681465>
7 [Accessed August 15, 2017].

8 Kent, W.J. et al., 2002. The human genome browser at UCSC. *Genome research*, 12(6), pp.996–
9 1006. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/12045153> [Accessed April 20,
10 2017].

11 Lanz, R.B. et al., 1999. A steroid receptor coactivator, SRA, functions as an RNA and is present in
12 an SRC-1 complex. *Cell*, 97(1), pp.17–27. Available at:
13 <http://www.ncbi.nlm.nih.gov/pubmed/10199399> [Accessed April 25, 2017].

14 Lanzós, A. et al., 2017. Discovery of Cancer Driver Long Noncoding RNAs across 1112 Tumour
15 Genomes: New Candidates and Distinguishing Features. *Scientific Reports*, 7, p.41544.
16 Available at: <http://www.ncbi.nlm.nih.gov/pubmed/28128360> [Accessed February 24, 2017].

17 Latos, P.A. et al., 2012. Airn Transcriptional Overlap, But Not Its lncRNA Products, Induces
18 Imprinted Igf2r Silencing. *Science*, 338(6113), pp.1469–1472. Available at:
19 <http://www.sciencemag.org/cgi/doi/10.1126/science.1228110> [Accessed April 25, 2017].

20 Lawrence, M.S. et al., 2014. Discovery and saturation analysis of cancer genes across 21 tumour
21 types. *Nature*, 505(7484), pp.495–501. Available at:
22 <http://www.nature.com/doi/10.1038/nature12912> [Accessed January 13, 2017].

23 Lepoivre, C. et al., 2013. Divergent transcription is associated with promoters of transcriptional
24 regulators. *BMC Genomics*, 14(1), p.914. Available at:
25 <http://www.ncbi.nlm.nih.gov/pubmed/24365181> [Accessed April 25, 2017].

26 Liu, S.J. et al., 2017. CRISPRi-based genome-scale identification of functional long noncoding
27 RNA loci in human cells. *Science*, 355(6320), p.eaah7111. Available at:
28 <http://www.sciencemag.org/lookup/doi/10.1126/science.aah7111> [Accessed April 25, 2017].

29 Managadze, D. et al., 2011. Negative Correlation between Expression Level and Evolutionary Rate
30 of Long Intergenic Noncoding RNAs. *Genome Biology and Evolution*, 3(0), pp.1390–1404.
31 Available at: <http://www.ncbi.nlm.nih.gov/pubmed/22071789> [Accessed March 26, 2017].

32 Marchese, F.P. et al., 2016. A Long Noncoding RNA Regulates Sister Chromatid Cohesion.
33 *Molecular Cell*, 63(3), pp.397–407. Available at:
34 <http://www.ncbi.nlm.nih.gov/pubmed/27477908> [Accessed April 25, 2017].

- 1 Marques, A.C. et al., 2013. Chromatin signatures at transcriptional start sites separate two equally
2 populated yet distinct classes of intergenic long noncoding RNAs. *Genome biology*, 14(11),
3 p.R131. Available at: [http://genomebiology.biomedcentral.com/articles/10.1186/gb-2013-14-](http://genomebiology.biomedcentral.com/articles/10.1186/gb-2013-14-11-r131)
4 11-r131 [Accessed March 20, 2017].
- 5 Mi, H. et al., 2013. Large-scale gene function analysis with the PANTHER classification system.
6 *Nature Protocols*, 8(8), pp.1551–1566. Available at:
7 <http://www.ncbi.nlm.nih.gov/pubmed/23868073> [Accessed April 23, 2017].
- 8 Mi, H. et al., 2017. PANTHER version 11: expanded annotation data from Gene Ontology and
9 Reactome pathways, and data analysis tool enhancements. *Nucleic Acids Research*, 45(D1),
10 pp.D183–D189. Available at: [https://academic.oup.com/nar/article-](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkw1138)
11 lookup/doi/10.1093/nar/gkw1138 [Accessed April 23, 2017].
- 12 Mularoni, L. et al., 2016. OncodriveFML: a general framework to identify coding and non-coding
13 regions with cancer driver mutations. *Genome biology*, 17(1), p.128. Available at:
14 <http://www.ncbi.nlm.nih.gov/pubmed/27311963> [Accessed June 28, 2016].
- 15 Ning, S. et al., 2016. Lnc2Cancer: a manually curated database of experimentally supported
16 lncRNAs associated with various human cancers. *Nucleic acids research*, 44(D1), pp.D980-5.
17 Available at: <http://www.ncbi.nlm.nih.gov/pubmed/26481356> [Accessed March 2, 2017].
- 18 Ponjavic, J. et al., 2009. Genomic and transcriptional co-localization of protein-coding and long
19 non-coding RNA pairs in the developing brain. Y. Hayashizaki, ed. *PLoS genetics*, 5(8),
20 p.e1000617. Available at: <http://dx.plos.org/10.1371/journal.pgen.1000617> [Accessed March
21 20, 2017].
- 22 Quek, X.C. et al., 2015. lncRNADB v2.0: expanding the reference database for functional long
23 noncoding RNAs. *Nucleic acids research*, 43(Database issue), pp.D168-73. Available at:
24 <http://www.ncbi.nlm.nih.gov/pubmed/25332394> [Accessed March 2, 2017].
- 25 Quinlan, A.R. & Hall, I.M., 2010. BEDTools: a flexible suite of utilities for comparing genomic
26 features. *Bioinformatics*, 26(6), pp.841–842. Available at:
27 <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btq033>
28 [Accessed August 16, 2017].
- 29 Redon, R. et al., 2006. Global variation in copy number in the human genome. *Nature*, 444(7118),
30 pp.444–54. Available at: <http://dx.doi.org/10.1038/nature05329>.
- 31 Reimand, U. et al., 2016. g:Profiler—a web server for functional interpretation of gene lists (2016
32 update). *Nucleic Acids Research*. Available at:
33 http://biit.cs.ut.ee/gprofiler/doc/papers/gprofiler_nar_2016.pdf [Accessed April 20, 2017].
- 34 Sabarinathan, R. et al., 2013. RNAsnp: Efficient Detection of Local RNA Secondary Structure

1 Changes Induced by SNPs. *Human Mutation*, 34(4), p.n/a-n/a. Available at:
2 <http://doi.wiley.com/10.1002/humu.22273> [Accessed April 19, 2017].
3 Sauvageau, M. et al., 2013. Multiple knockout mouse models reveal lincRNAs are required for life
4 and brain development. *eLife*, 2, p.e01749. Available at:
5 <http://www.ncbi.nlm.nih.gov/pubmed/24381249> [Accessed April 25, 2017].
6 Schmidt, K. et al., 2016. The lincRNA SLNCR1 Mediates Melanoma Invasion through a Conserved
7 SRA1-like Region. *Cell reports*, 15(9), pp.2025–37. Available at:
8 <http://linkinghub.elsevier.com/retrieve/pii/S2211124716304314> [Accessed March 20, 2017].
9 Schmitt, A.M. et al., 2016. Long Noncoding RNAs in Cancer Pathways. *Cancer Cell*, 29(4),
10 pp.452–463. Available at: <http://www.ncbi.nlm.nih.gov/pubmed/27070700> [Accessed April
11 20, 2017].
12 Sjoblom, T. et al., 2006. The Consensus Coding Sequences of Human Breast and Colorectal
13 Cancers. *Science*, 314(5797), pp.268–274. Available at:
14 <http://www.ncbi.nlm.nih.gov/pubmed/16959974> [Accessed February 24, 2017].
15 Tamborero, D. et al., 2013. Comprehensive identification of mutational cancer driver genes across
16 12 tumor types. *Scientific reports*, 3, p.2650. Available at:
17 <http://www.ncbi.nlm.nih.gov/pubmed/24084849> [Accessed June 28, 2016].
18 Tan, J.Y. et al., 2017. cis -Acting Complex-Trait-Associated lincRNA Expression Correlates with
19 Modulation of Chromosomal Architecture. *Cell Reports*, 18(9), pp.2280–2288. Available at:
20 <http://www.ncbi.nlm.nih.gov/pubmed/28249171> [Accessed April 25, 2017].
21 Tokheim, C.J. et al., 2016. Evaluating the evaluation of cancer driver genes. *Proceedings of the*
22 *National Academy of Sciences of the United States of America*, 113(50), pp.14330–14335.
23 Available at: <http://www.ncbi.nlm.nih.gov/pubmed/27911828> [Accessed August 16, 2017].
24 Tripathi, S. et al., 2015. Meta- and Orthogonal Integration of Influenza "OMICs" Data
25 Defines a Role for UBR4 in Virus Budding. *Cell host & microbe*, 18(6), pp.723–35. Available
26 at: <http://linkinghub.elsevier.com/retrieve/pii/S1931312815004564> [Accessed April 23, 2017].
27 Tyner, C. et al., 2017. The UCSC Genome Browser database: 2017 update. *Nucleic Acids Research*,
28 45. Available at: [https://oup.silverchair-](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/nar/45/D1/10.1093_nar_gkw1134/3/gkw1134.pdf?Expires=1493043886&Signature=VkpqKibgiT9kty13UqDbQabocXaBUaciIshI~KU4D010dRoa-9n7qLF0kT3eT2HZRiqrI7W74jyxPg1eyKhuPrIozHFwCJWxa3tO3wh95deEodxSyEwvjA64rFAZFbmV0EtdUoWRqL5nhsJqLSPCZmPXukDpSxBH7SrrmYX33UqRcVo6jq-ICvde4XmPNDgacH4BwRTU2K0~D-OeV~kzq6s-zshWmPjUIJM-)
29 [cdn.com/oup/backfile/Content_public/Journal/nar/45/D1/10.1093_nar_gkw1134/3/gkw1134.p](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/nar/45/D1/10.1093_nar_gkw1134/3/gkw1134.pdf?Expires=1493043886&Signature=VkpqKibgiT9kty13UqDbQabocXaBUaciIshI~KU4D010dRoa-9n7qLF0kT3eT2HZRiqrI7W74jyxPg1eyKhuPrIozHFwCJWxa3tO3wh95deEodxSyEwvjA64rFAZFbmV0EtdUoWRqL5nhsJqLSPCZmPXukDpSxBH7SrrmYX33UqRcVo6jq-ICvde4XmPNDgacH4BwRTU2K0~D-OeV~kzq6s-zshWmPjUIJM-)
30 [df?Expires=1493043886&Signature=VkpqKibgiT9kty13UqDbQabocXaBUaciIshI~KU4D01](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/nar/45/D1/10.1093_nar_gkw1134/3/gkw1134.pdf?Expires=1493043886&Signature=VkpqKibgiT9kty13UqDbQabocXaBUaciIshI~KU4D010dRoa-9n7qLF0kT3eT2HZRiqrI7W74jyxPg1eyKhuPrIozHFwCJWxa3tO3wh95deEodxSyEwvjA64rFAZFbmV0EtdUoWRqL5nhsJqLSPCZmPXukDpSxBH7SrrmYX33UqRcVo6jq-ICvde4XmPNDgacH4BwRTU2K0~D-OeV~kzq6s-zshWmPjUIJM-)
31 [0dRoa-](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/nar/45/D1/10.1093_nar_gkw1134/3/gkw1134.pdf?Expires=1493043886&Signature=VkpqKibgiT9kty13UqDbQabocXaBUaciIshI~KU4D010dRoa-9n7qLF0kT3eT2HZRiqrI7W74jyxPg1eyKhuPrIozHFwCJWxa3tO3wh95deEodxSyEwvjA64rFAZFbmV0EtdUoWRqL5nhsJqLSPCZmPXukDpSxBH7SrrmYX33UqRcVo6jq-ICvde4XmPNDgacH4BwRTU2K0~D-OeV~kzq6s-zshWmPjUIJM-)
32 [9n7qLF0kT3eT2HZRiqrI7W74jyxPg1eyKhuPrIozHFwCJWxa3tO3wh95deEodxSyEwvjA64](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/nar/45/D1/10.1093_nar_gkw1134/3/gkw1134.pdf?Expires=1493043886&Signature=VkpqKibgiT9kty13UqDbQabocXaBUaciIshI~KU4D010dRoa-9n7qLF0kT3eT2HZRiqrI7W74jyxPg1eyKhuPrIozHFwCJWxa3tO3wh95deEodxSyEwvjA64rFAZFbmV0EtdUoWRqL5nhsJqLSPCZmPXukDpSxBH7SrrmYX33UqRcVo6jq-ICvde4XmPNDgacH4BwRTU2K0~D-OeV~kzq6s-zshWmPjUIJM-)
33 [rFAZFbmV0EtdUoWRqL5nhsJqLSPCZmPXukDpSxBH7SrrmYX33UqRcVo6jq-](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/nar/45/D1/10.1093_nar_gkw1134/3/gkw1134.pdf?Expires=1493043886&Signature=VkpqKibgiT9kty13UqDbQabocXaBUaciIshI~KU4D010dRoa-9n7qLF0kT3eT2HZRiqrI7W74jyxPg1eyKhuPrIozHFwCJWxa3tO3wh95deEodxSyEwvjA64rFAZFbmV0EtdUoWRqL5nhsJqLSPCZmPXukDpSxBH7SrrmYX33UqRcVo6jq-ICvde4XmPNDgacH4BwRTU2K0~D-OeV~kzq6s-zshWmPjUIJM-)
34 [ICvde4XmPNDgacH4BwRTU2K0~D-OeV~kzq6s-zshWmPjUIJM-](https://oup.silverchair-cdn.com/oup/backfile/Content_public/Journal/nar/45/D1/10.1093_nar_gkw1134/3/gkw1134.pdf?Expires=1493043886&Signature=VkpqKibgiT9kty13UqDbQabocXaBUaciIshI~KU4D010dRoa-9n7qLF0kT3eT2HZRiqrI7W74jyxPg1eyKhuPrIozHFwCJWxa3tO3wh95deEodxSyEwvjA64rFAZFbmV0EtdUoWRqL5nhsJqLSPCZmPXukDpSxBH7SrrmYX33UqRcVo6jq-ICvde4XmPNDgacH4BwRTU2K0~D-OeV~kzq6s-zshWmPjUIJM-)

1 XCB7Mpx2kd5JVIN7IPVCt0vs8gK~BDHHdryJEkWLf9L4ZbFxzvLvtulvQUZnrNSNZFLR
2 PPDHFnJ5C7YZM8-U5g27ma2eGbwIQyjZ1qqbcxOg6QLkLw__&Key-Pair-
3 Id=APKAIUCZBIA4LVPAVW3Q [Accessed April 20, 2017].
4 Ulitsky, I. & Bartel, D.P., 2013. lincRNAs: Genomics, Evolution, and Mechanisms. *Cell*, 154(1),
5 pp.26–46. Available at: <http://linkinghub.elsevier.com/retrieve/pii/S0092867413007599>
6 [Accessed November 20, 2016].
7 Welter, D. et al., 2014. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations.
8 *Nucleic Acids Research*, 42(D1), pp.1001–1006.
9 World Health Organization, 2013. *International Classification of Diseases for Oncology (ICD-O).*
10 *Third Edition. First Revision.*,
11 Xiang, J.-F. et al., 2014. Human colorectal cancer-specific CCAT1-L lncRNA regulates long-range
12 chromatin interactions at the MYC locus. *Cell Research*, 24(5), pp.513–531. Available at:
13 <http://www.nature.com/doi/10.1038/cr.2014.35> [Accessed April 25, 2017].
14 Yan, X. et al., 2015. Comprehensive Genomic Characterization of Long Non-coding RNAs across
15 Human Cancers. *Cancer Cell*.
16 Yates, L.R. & Campbell, P.J., 2012. Evolution of the cancer genome. *Nature Reviews Genetics*,
17 13(11), pp.795–806. Available at: <http://www.nature.com/doi/10.1038/nrg3317>
18 [Accessed April 25, 2017].
19 Zhu, S. et al., 2016. Genome-scale deletion screening of human long non-coding RNAs using a
20 paired-guide RNA CRISPR–Cas9 library. *Nature Biotechnology*, 34(12), pp.1279–1286.
21 Available at: <http://www.nature.com/doi/10.1038/nbt.3715> [Accessed December 21,
22 2016].

23
24
25

Tables

2

Table 1: List of intergenic CIS human (GRCh38) / mouse (GRCm38) gene pairs.

4

Human CLC Name	Human CLC ID	Chr Human	Start Human	End Human	Chr Mouse	Start Mouse	End Mouse	PubMed ID	Cancer Type Mouse
DLEU2	ENSG00000231607	chr13	50,048,971	50,049,063	chr14	61,631,880	61,631,972	24316982	Liver
DLEU2	ENSG00000231607	chr13	50,049,117	50,049,206	chr14	61,632,026	61,632,110	24316982	Liver
GAS5	ENSG00000234741	chr1	173,864,370	173,864,435	chr1	161,038,091	161,038,156	25961939	Sarcoma
MONC	ENSG00000215386	chr21	16,539,096	16,539,161	chr16	77,598,935	77,599,000	23685747	Nervous System
MONC	ENSG00000215386	chr21	16,561,654	16,561,655	chr16	77,616,439	77,616,440	24316982	Liver
NEAT1	ENSG00000245532	chr11	65,444,511	65,444,512	chr19	5,825,497	5,825,498	24316982	Liver
PINT	ENSG00000231721	chr7	131,049,455	131,049,456	chr6	31,179,149	31,179,150	22699621	Pancreatic
PVT1	ENSG00000249859	chr8	128,007,970	128,007,971	chr15	62,186,646	62,186,647	22699621	Pancreatic
SLNCR1	ENSG00000227036	chr17	72,507,275	72,507,276	chr11	113,137,613	113,137,614	22699621	Pancreatic
XIST	ENSG00000229807	chrX	73,841,539	73,841,540	chrX	103,473,862	103,473,863	24316982	Liver
XIST	ENSG00000229807	chrX	73,841,539	73,841,540	chrX	103,473,862	103,473,863	24316982	Liver

5

Figure 1

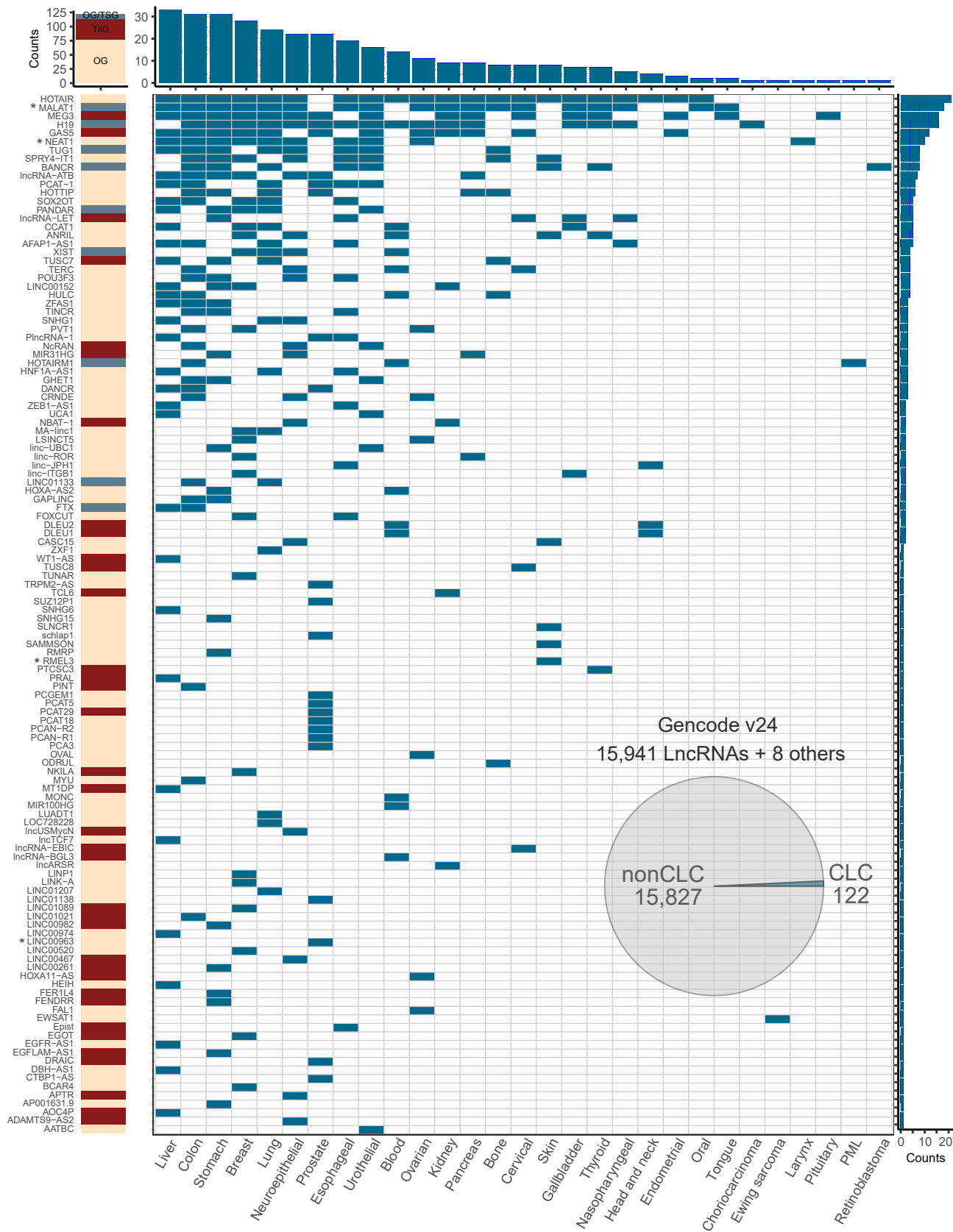


Figure 2

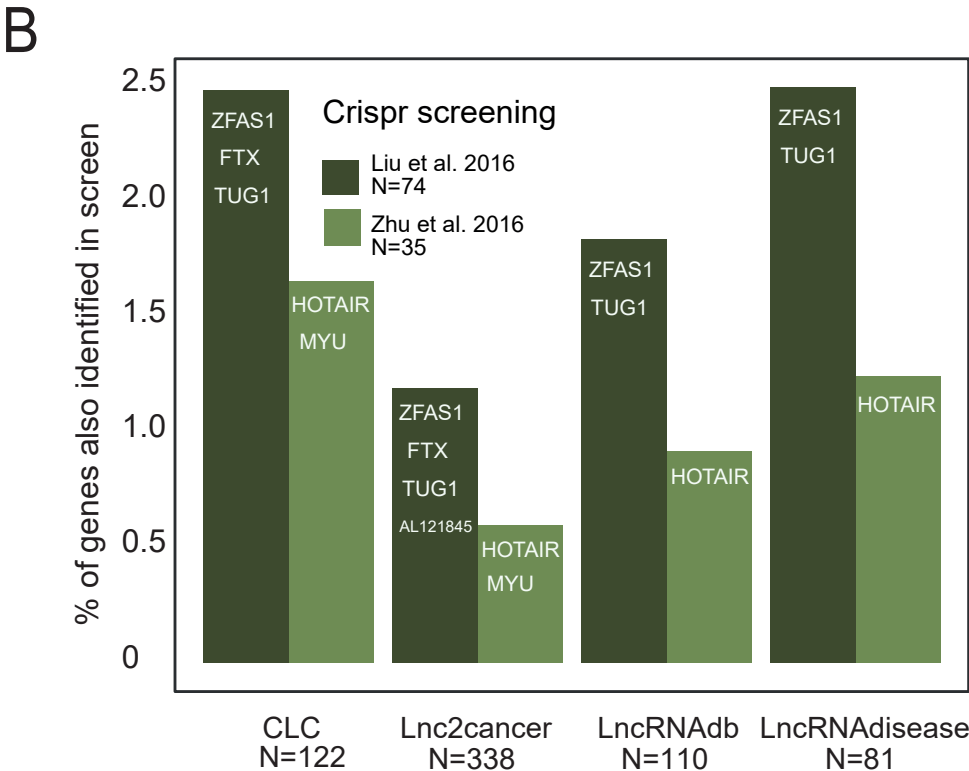
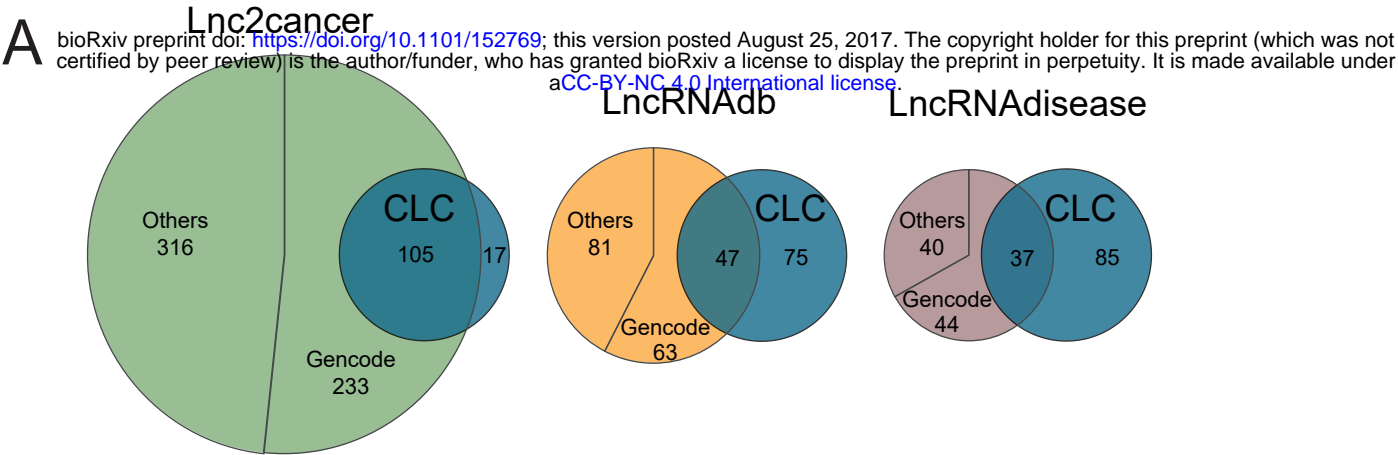


Figure 3

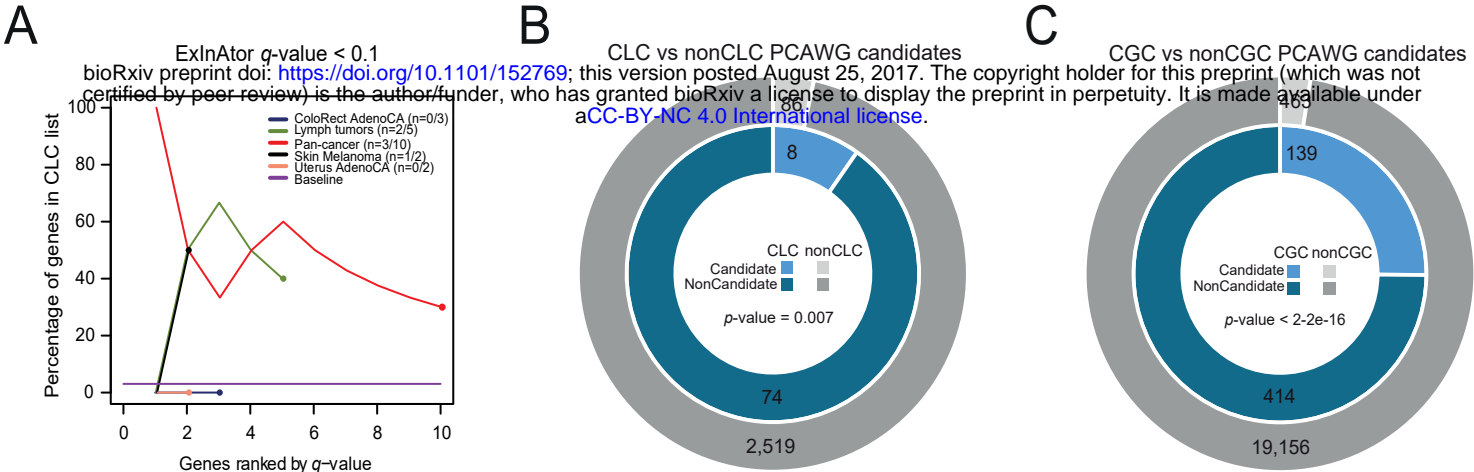
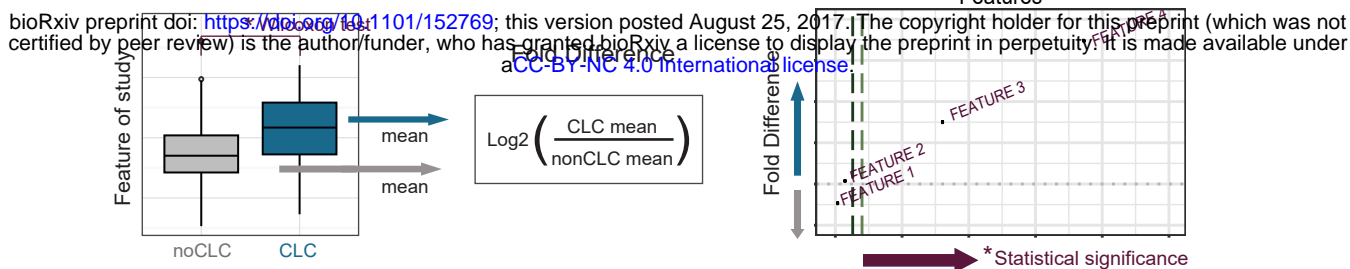
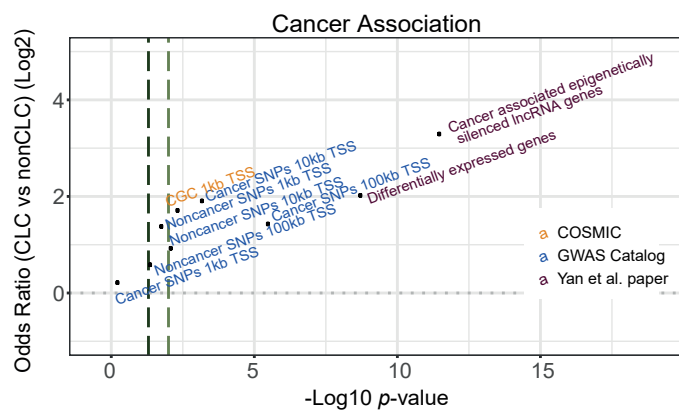


Figure 4

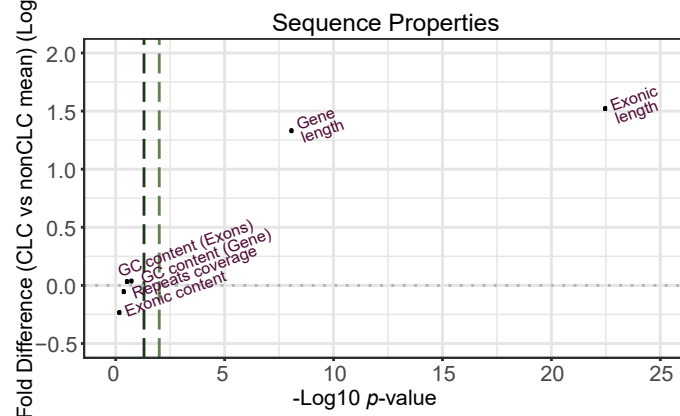
A



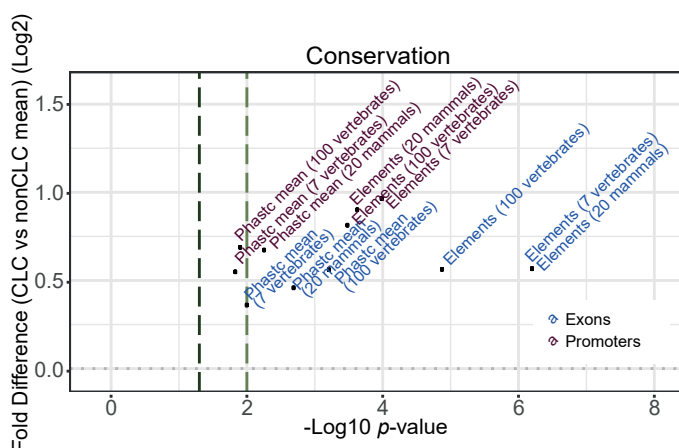
B



C



D



E

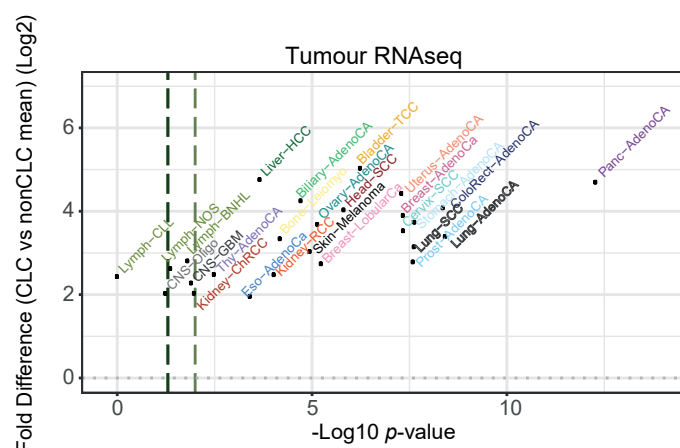
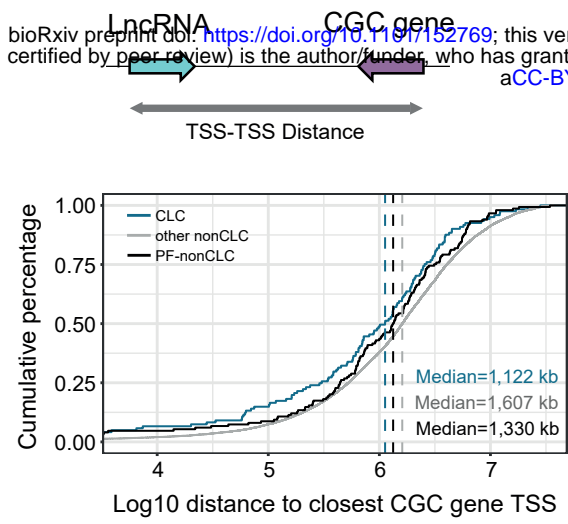
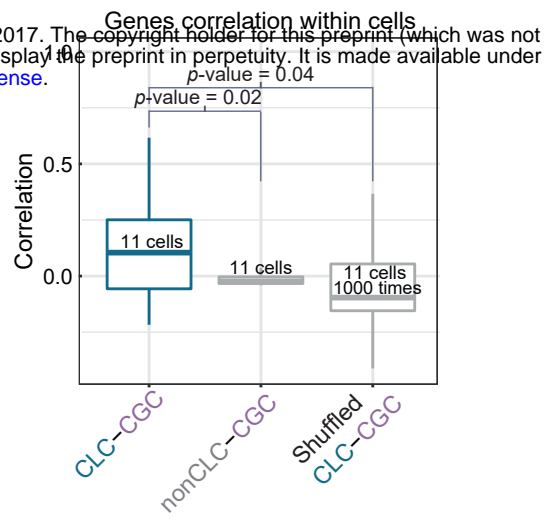


Figure 5

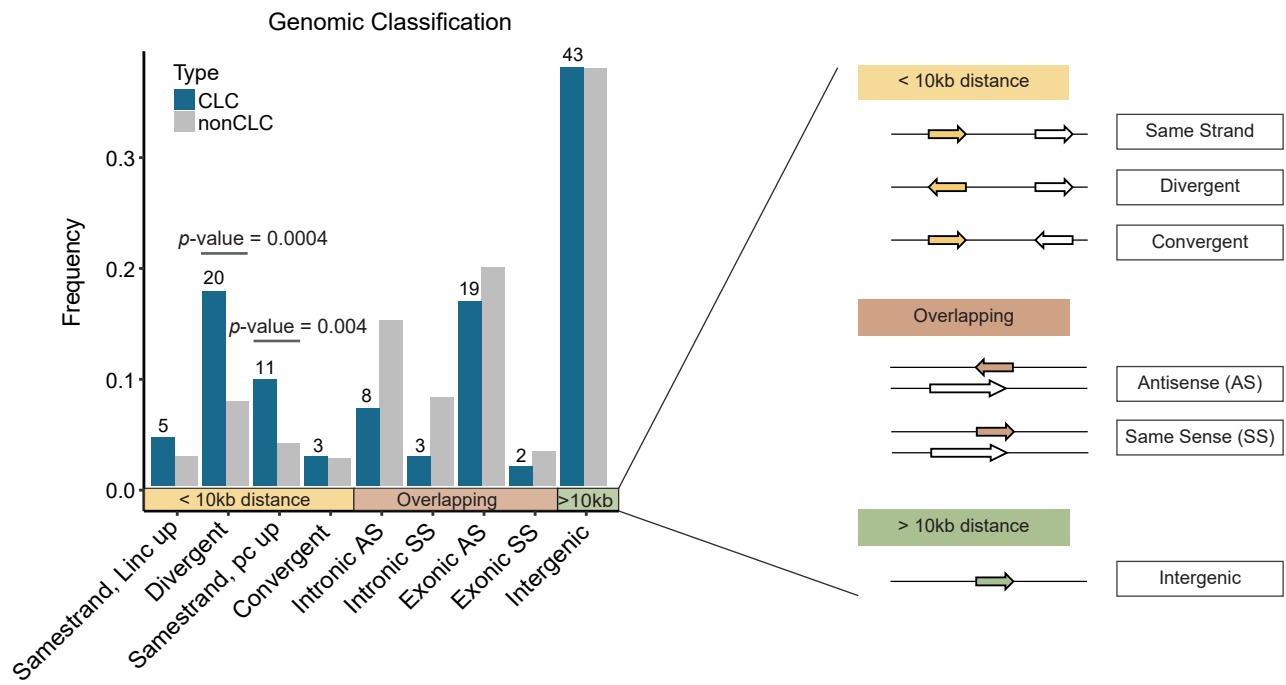
A



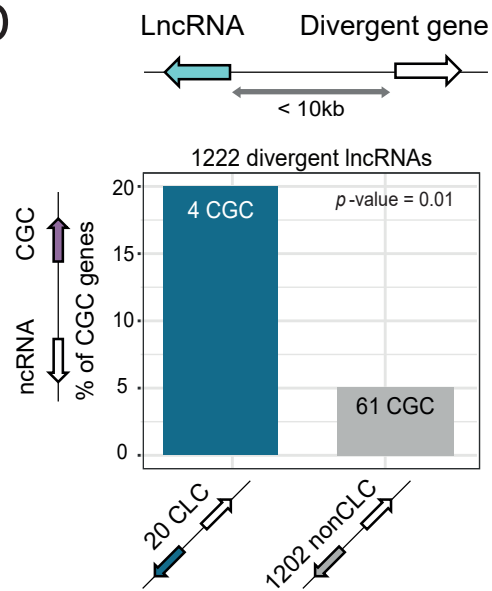
B



C



D



E

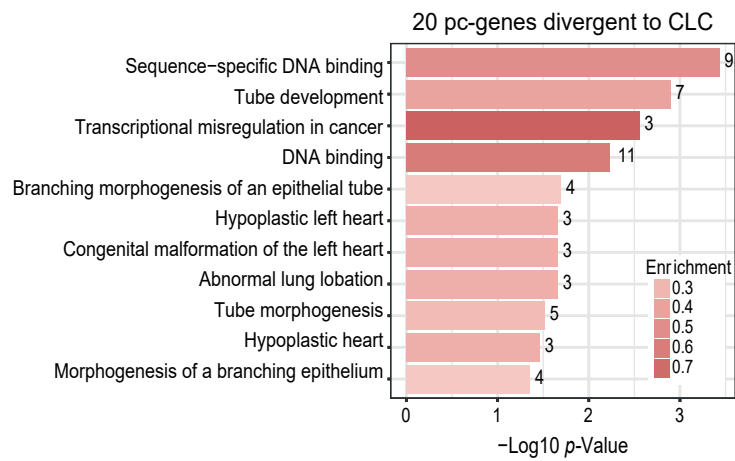


Figure 6

