# 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)

PCAWG Working Group 2-5-9-14

Roderic Guigó (3,4,5)

Jakob Skou Pedersen (7)

Rory Johnson (1,2)


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**

1 **Abstract**

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

8     **CLC Version 1 comprises 122 lncRNAs implicated in 31 distinct cancers. LncRNAs are**
9 **included based on functional or genetic evidence for different causative roles in cancer**
10 **progression. All belong to the GENCODE reference annotation, to facilitate integration across**
11 **projects and datasets. For each entry, the evidence type, biological activity (oncogene or tumour**
12 **suppressor), source reference and cancer type are recorded. CLC genes are significantly**
13 **enriched amongst *de novo* predicted driver genes from PCAWG. CLC genes are distinguished**
14 **from other lncRNAs by a series of features consistent with biological function, including gene**
15 **length, expression and sequence conservation of both exons and promoters. We identify a trend**
16 **for CLC genes to be co-localised with known protein-coding cancer genes along the human**
17 **genome. Finally, by integrating data from transposon-mutagenesis functional screens, we show**
18 **that mouse orthologues of CLC genes tend also to be cancer driver genes.**

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

22

## Introduction

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

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

LncRNAs have been shown to play important roles in various cancers. For example, *MALAT1*, a potent oncogene across numerous cancers, is restricted to the nucleus and plays a housekeeping role in splicing (Gutschner & Diederichs 2012; Engreitz et al. 2014). *MALAT1* is overexpressed in a variety of cancer types, and its knockdown potently reduces not only proliferation but also metastasis *in vivo* (Gutschner et al. 2013).The *MALAT1* gene is subjected to elevated mutational rates in human tumours, suggesting that genetic alteration in its function is an important step in tumorigenesis (Lanzós et al. 2017). LncRNAs may also function as tumour suppressors. *LincRNA-p21* acts as a downstream effector of p53 regulation through recruitment of the repressor hnRNP-K (Huarte et al. 2010). These and other examples of lncRNAs linked to cancer, raise the question of how many more remain to be found amongst the ~99% of lncRNAs that are presently uncharacterised (Derrien et al. 2012; Quek et al. 2015; Iyer et al. 2015).

Recent tumour genome sequencing studies, in step with advanced bioinformatic driver-gene prediction methods, have yielded hundreds of new candidate protein-coding driver genes (Tamborero et al. 2013). For economic reasons, these studies initially restricted their attention to "exomes" or the ~2% of the genome covering protein-coding exons (Chang et al. 2013). Unfortunately such a strategy ignores mutations in the remaining ~98% of genomic sequence, home to the majority of lncRNAs (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

1    mutation patterns (Lawrence et al. 2014). It is critical to test the performance of such methods using

2    true-positive lists of known driver genes. For protein-coding genes, this role has been fulfilled by the

3    Cancer Gene Census (CGC) (Futreal et al. 2004), collected and regularly updated by manual

4    annotators based on experimentally-validated cases. Comparison of driver predictions to CGC genes

5    facilitates further method refinement and comparison between methods (Sjoblom et al. 2006; Redon

6    et al. 2006; Mularoni et al. 2016).

7    In addition to its benchmarking role, the CGC resource has also been useful in identifying unique

8    biological features of cancer genes (Furney et al. 2006; Furney et al. 2008). For example, CGC genes

9    tend to be more conserved and longer. Furthermore, they are enriched for genes with transcription

10    regulator activity and nucleic acid binding functions (Furney et al. 2006; Furney et al. 2008).

11    Until very recently, efforts to discover cancer lncRNAs have depended on classical functional

12    genomics approaches of differential expression using microarrays or RNA sequencing (Huarte et al.

13    2010; Iyer et al. 2015). While valuable, differential expression *per se* is not direct evidence for

14    causative roles in tumour evolution. To more directly identify lncRNAs that drive cancer progression,

15    methods have recently been developed to search for signals of positive selection using mutation maps

16    created using tumour genome sequences. OncodriveFML utilises nucleotide-level functional impact

17    scores inferred from predicted changes in RNA secondary structure (Sabarinathan et al. 2013)

18    together with an empirical significance estimate, to identify lncRNAs with an excess of high-impact

19    mutations (Mularoni et al. 2016). Another method, ExInAtor, identifies candidates with elevated

20    mutational load, using trinucleotide-adjusted local background (Lanzós et al. 2017). A clear

21    impediment in both cases has been the lack of true-positive set of known lncRNA driver genes,

22    analogous to CGC. Although there do exist databases of cancer lncRNAs, notably LncRNADisease

23    (Chen et al. 2013) and Lnc2Cancer (Ning et al. 2016), they mix unfiltered data from numerous

24    sources, resulting in inconsistent criteria for inclusion (including expression changes), and

25    inconsistent gene identifiers.

26    To facilitate the future discovery of cancer lncRNAs, and gain insights into their biology, we

27    have compiled a highly-curated set of cases with roles in cancer processes. Here we present the

28    *Cancer LncRNA Census* (CLC), the first compendium of lncRNAs with direct functional or genetic

29    evidence for cancer roles. We demonstrate the utility of CLC in assessing the performance of driver

30    lncRNA predictions. Through analysis of this geneset, we demonstrate that cancer lncRNAs have a

31    unique series of features that may in future be used to assist *de novo* predictions. Finally, we show

32    that CLC genes have conserved cancer roles across the approximately 80 million years of evolution

33    separating humans and rodents.

**Results**

**Definition of cancer related lncRNAs**

As part of recent efforts to identify driver lncRNAs within the Pancancer Analysis of Whole Genomes (PCAWG) project, we discovered the need for a high-confidence set of cancer-related lncRNA genes (CRL) (Lanzós et al. 2017), which we henceforth refer to as "cancer lncRNAs", and which we here present as 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 lncRNA 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 that of their 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 a total of 122 genes, or 0.76% of a total of 15,941 lncRNA gene loci annotated in GENCODE v24 (Derrien et al. 2012; Harrow et al. 2012) (Figure 1) (eight of the CLC genes are annotated as pseudogenes rather than lncRNAs by GENCODE). 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 CLC set is smaller than the *Lnc2Cancer* database (n= 667) (Ning et al. 2016), the *LncRNADisease* Database (n=121) (Chen et al. 2013) and *lncRNAdb* (n=191) (Quek et al. 2015). CLC covers between 17% and 31% of these databases (*Lnc2Cancer* and *LncRNADisease* respectively) but none of these resources contain the complete list of genes presented here (Figure 2). CLC is substantially larger than a previous "Cancer Related LncRNAs" set we produced (Lanzós et al. 2017)(Figure 2B). 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). 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 31(Figure 1, Supplementary Figure 1).

1    Altogether, CLC contains 337 unique lncRNA-cancer type relationships. Out of 122 genes, 77

2    (63.1%) were shown to function as oncogenes, 36 (29.5%) as tumour suppressors, and 9 (7.4%) with

3    evidence for both activities (Figure 1 and Supplementary Figure 1).

4    The most prolific lncRNAs, with >15 recorded cancer types, are *HOTAIR*, *MALAT1*, *MEG3* and

5    *H19* (Figure 1 and Supplementary Figure 1). It is not clear whether this reflects their unique pan-

6    cancer functionality, or is simply a result of their being amongst the most early-discovered and

7    widely-studied lncRNAs.

8    *In vitro* experiments were the most frequent evidence source, usually consisting of RNAi-

9    mediated knockdown in cultured cell lines, coupled to phenotypic assays such as proliferation or

10   migration (Supplementary Figure 1). Far fewer have been studied *in vivo*, or have cancer-associated

11   somatic or germline mutations. 19 lncRNAs had 3 or more independent evidence sources

12   (Supplementary Figure 1).

13

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

15   One of the primary motivations for CLC is as a true positive set for benchmarking and comparing

16   methods for identifying driver lncRNAs. In the domain of protein-coding driver gene predictions, the

17   Cancer Gene Census (CGC) has become the "gold standard" training set, against which new methods

18   are tested (Futreal et al. 2004). Typically, the predicted driver genes belonging to CGC are judged to

19   be true positives, and the fraction of these amongst predictions is used to estimate the Positive

20   Predictive Value (PPV), or precision. This measure can be calculated for increasing cutoff levels, to

21   assess the optimal cutoff.

22   First, we used CLC to examine the performance of the lncRNA driver predictor ExInAtor

23   (Lanzós et al. 2017) in recalling CLC genes using PCAWG tumour mutation data (PCAWG

24   Consortium, Manuscript In Preparation). A total of 2,687 GENCODE lncRNAs were included here,

25   of which 82 (3.1%) belong to CLC. Driver predictions on several cancers at the standard False

26   Discovery Rate ("$q$-value") cutoff of 0.1 are shown for selected cancers in Figure 3A. That panel

27   shows the CLC-defined precision ($y$-axis) as a function of predicted driver genes ranked by $q$-value

28   ($x$-axis). We observe rather heterogeneous performance across cancer cohorts. This may reflect a

29   combination of intrinsic biological differences and differences in cohort sizes, which differs widely

30   between the datasets shown. For the merged Pancancer dataset ("PANCAN"), ExInAtor predicted

31   three CLC genes amongst its top ten candidates ($q$-value < 0.1), a rate far in excess of the background

32   expectation. Similar enrichments are observed for other cancer types. These results support both the

33   predictive value of ExInAtor, and the usefulness of CLC in assessing lncRNA driver predictors.

1　Comprehensive CLC-based assessments of lncRNA driver discovery, across all methods and tumour

2　cohorts in PCAWG, may be found in the main PCAWG driver prediction publication (PCAWG

3　Consortium, Manuscript In Preparation).

4　　　Finally, we assessed the precision (i.e. positive predictive value) of PCAWG lncRNA and

5　protein-coding driver predictions across all cancers and all prediction methods (PCAWG Consortium,

6　Manuscript In Preparation). Using the same $q$-value cutoff of 0.1, we found that across all cancer

7　types and methods, a total of 8 (8.5%) of lncRNA predictions are in CLC (Figure 3B), while a total

8　of 139 (23.1%) of protein-coding predictions are CGC genes (Figure 3C). In terms of sensitivity,

9　9.8% and 25.1% of CLC and CGC genes are predicted as candidates, respectively. Despite the lower

10　detection of CLC genes in comparison to CGC genes, both sensitivity rates significantly exceed the

11　prediction rate of nonCLC and nonCGC genes, again highlighting the usefulness of the CLC geneset.

12

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

14　　　We recently found evidence, using a smaller geneset than presented here, that cancer lncRNAs

15　are distinguished by various genomic and expression features indicative of biological function

16　(Lanzós et al. 2017). We here extended these findings using a large series of potential gene features,

17　to search for those features distinguishing CLC from nonCLC lncRNAs (Figure 4A).

18　　　First, associations with expected cancer-related features were tested (Figure 4B). CLC genes are

19　significantly more likely to have their transcription start site (TSS) within 100 kb of cancer-associated

20　germline SNPs ("Cancer SNPs 100kb TSS"), and more likely to be either differentially-expressed or

21　epigenetically-silenced in tumours (Yan et al. 2015) (Figure 4B). Intriguingly, we observed a

22　tendency for CLC lncRNAs to be more likely to lie within 1 kb of known cancer protein-coding genes

23　("CGC 1kb TSS") – this is explored in more detail below.

24　　　We next investigated the properties of the genes themselves. As seen in Figure 4C, and consistent

25　with our previous findings (Lanzós et al. 2017), CLC genes ("Gene length") and their spliced products

26　("Exonic length") are significantly longer than average. No difference was observed in the ratio of

27　exonic to total length ("Exonic content"), nor repetitive sequence coverage ("Repeats coverage"), nor

28　GC content.

29　　　CLC genes also tend to have greater evidence of function, as inferred from evolutionary

30　conservation. Base-level conservation at various evolutionary depths was calculated for lncRNA

31　exons and promoters (Figure 4D). Across all measures tested, using either average base-level scores

32　or percent coverage by conserved elements, we found that CLC genes' exons are significantly more

33　conserved than other lncRNAs (Figure 4D). The same was observed for conservation of promoter

34　regions.

1    High levels of gene expression are known to correlate with lncRNA conservation, and are
2    hypothesized to be a reflection of functionality (Managadze et al. 2011). We found that CLC genes
3    have consistently higher steady-state expression levels, across both healthy organs and cultured cell
4    lines (Figure 4E). This difference tended to be more marked for cell lines compared to organ samples.
5    It is unclear whether this is a batch effect, since cell lines data comes from ENCODE (Djebali et al.
6    2012) and organ samples from Illumina Bodymap, or reflects a generalised upregulation of cancer
7    lncRNA expression in cultured cells. It should be noted that these mainly comprise immortalised
8    cancer cell lines but also includes lines such as IMR-90 that are neither immortal nor transformed.
9

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

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

16   Using TSS-TSS distances, we found that CLC genes on average tend to lie moderately closer to
17   protein-coding genes of all types, compared to nonCLC lncRNAs (median 29 kb vs 38 kb
18   respectively) (P=0.03, Wilcoxon test) (Supplementary Figure 2A). Interestingly, this effect is also
19   observed when assessing specifically for distance to CGC genes (median 1,122 kb vs 1,599 kb,
20   respectively) (P= 0.0004, Wilcoxon test) (Figure 5A).

21   It has been widely proposed that proximal lncRNA / protein-coding gene pairs are involved in
22   *cis*-regulatory relationships, which is reflected in expression correlation (Ponjavic et al. 2009). We
23   next asked whether proximal CLC-CGC pairs exhibit this behaviour. An important potential
24   confounding factor is the known positive correlation between nearby gene pairs (Marques et al. 2013),
25   and this must be controlled for. Using gene expression data across 11 human cell lines, we observed
26   a positive correlation between CLC-CGC gene pairs for each cell type (Figure 5B). To control for the
27   effect of proximity on correlation, we next randomly sampled a similar number of non-CLC lncRNAs
28   with matched distances (TSS-TSS) from the same CGC genes, and found that this correlation was
29   lost (Figure 5B, "nonCLC-CGC"). To further control for a possible correlation arising from the simple
30   fact that both CGC and CLC genes are involved in cancer, we randomly shuffled the CLC-CGC pairs
31   1000 times, again observing no correlation (Figure 5B, "Shuffled CLC-CGC"). Together these results
32   show that genomically-proximal protein-coding/non-coding gene pairs exhibit an expression
33   correlation that exceeds that expected by chance, even when controlling for genomic distance.

1    These results prompted us to further explore the genomic localization of CLC genes relative to
2    their proximal protein-coding gene and the nature of their neighbouring genes. First, considering gene
3    pairs as proximal when both TSS lie within 10 kb, we find, in agreement with our previous results,
4    that CLC genes are significantly more likely to lie near to a CGC gene, compared to other lncRNAs
5    (6.6% vs 2.1%, respectively) (P=0.004, Fisher's exact test) (Supplementary Figure 2B). Next, we
6    observed an unexpected difference in the genomic organisation of CLC genes: when classified by
7    orientation with respect to nearest protein-coding gene (Derrien et al. 2012), we found a significant
8    enrichment of CLC genes immediately downstream and on the same strand as protein-coding genes
9    ("Samestrand, pc up", Figure 5C). Moreover, CLC genes are approximately twice as likely to lie in
10   an upstream, divergent orientation to a protein-coding gene ("Divergent", Figure 5C). Of these CLC
11   genes, 20% are divergent to a CGC gene, compared to 5% for nonCLC genes (P=0.018, Fisher's
12   exact test) (Figure 5D), and several are divergent to protein-coding genes that have also been linked
13   or defined to be involved in cancer, despite not being classified as CGCs (Supplementary Table 2).

14   Given this noteworthy enrichment of CGC genes in a divergent configuration to protein-coding
15   genes, we next inspected the latters' function annotation. Examining their Gene Ontology (GO) terms,
16   molecular pathways and other gene function related terms, we found this group of genes to be
17   enriched in GO terms for "sequence-specific DNA binding", "DNA binding", "tube development"
18   and "transcriptional misregulation in cancer" (Figure 5E). These results were confirmed by another,
19   independent GO-analysis suite (see Materials and Methods). Interestingly, three out of the top four
20   functional groups were observed previously in a study of protein-coding genes divergent to long
21   upstream antisense transcripts (LUAT) in primary mouse tissues (Lepoivre et al. 2013).

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

24

## Evidence for anciently conserved cancer roles of lncRNAs

26   Numerous studies have employed unbiased forward genetic screens to identify genes that either
27   inhibit or promote tumorigenesis in mice (Copeland & Jenkins 2010). These studies use engineered,
28   randomly-integrating transposons, carrying bidirectional polyadenylation sites as well as strong
29   promoters. Insertions, or clusters of insertions, called "common insertion sites" (CIS) that are
30   identified in sequenced tumour DNA, implicate the overlapping or neighbouring gene locus as either
31   an oncogene or tumour-suppressor gene. Although these studies have traditionally been focused on
32   identifying protein-coding genes, they can in principle also identify non-coding driver loci.

33   We reasoned that comparison of mouse CISs to orthologous human regions could yield
34   independent evidence for the functionality of human cancer lncRNAs (Figure 6A). To test this, we

1    collected a comprehensive set of CISs in mouse (Abbott et al. 2015), consisting of 2,906 loci from 7

2    distinct cancer types (Supplementary Table 4). These sites were then mapped to in orthologous

3    regions in the human genome, resulting in 1,309 human CISs, or hCISs. 7.3% of these CISs lie outside

4    of protein-coding gene boundaries, which were used for the following analyses.

5      Mapping hCISs to lncRNA annotations, we discovered altogether eight CLC genes (6.6%)

6    carrying an insertion within their gene span: *DLEU2*, *GAS5*, *MONC*, *NEAT1*, *PINT*, *PVT1*, *SLNCR1*,

7    *XIST* (Table 1). In contrast, just 61 (0.4%) nonCLC genes contained hCISs (Figure 6B). A good

8    example is *SLINCR1*, shown in Figure 6C, which drives invasiveness of human melanoma cells

9    (Schmidt et al. 2016), and whose mouse orthologue contains a CIS discovered in pancreatic cancer.

10      This analysis would suggest that CLC genes are enriched for hCISs; however, there remains the

11    possibility that this is confounded by their greater length. To account for this, we performed two

12    separate validations. First, sets of nonCLC genes with CLC-matched length were randomly sampled,

13    and the number of intersecting hCISs per unit gene length (Mb) was counted (Figure 6D). Second,

14    CLC genes were randomly relocated in the genome, and the number of genes intersecting at least one

15    hCIS was counted (Figure 6E). Both analyses showed that the number of intersecting hCISs per Mb

16    of CLC gene span is far greater than expected. In contrast, nonCLC genes show a depletion for hCIS

17    sites (Figure 6F).

18      Together these analyses demonstrate that CLC genes are orthologous to mouse cancer-causing

19    genomic loci at a rate greater than expected by random chance. These identified cases, and possibly

20    other CLC genes, display cancer functions that have been conserved over tens of millions of years

21    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 any 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 ExInAtor pipeline (Lanzós et al. 2017). Supporting the usefulness of this approach, we found a strong 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%). This difference may be due to technical or biological factors. For example, CLC genes that are not identified as drivers by PCAWG may be false negatives, whose identification simply awaits greater statistical power afforded by more genomes; or else many CLC genes may contribute to cancer phenotypes, but are not targeted

1    by somatic mutations during tumorigenesis. Addressing these questions will be a key objective in

2    future studies.

3          Analysis of CLC geneset has broadened our understanding of the unique features of cancer

4    lncRNAs, and generally supports the notion that lncRNAs have intrinsic biological functionality.

5    Cancer lncRNAs are distinguished by a series of features that are consistent with both (a) roles in

6    cancer (eg tumour expression changes), and (b) general biological functionality (eg high expression,

7    evolutionary conservation). Elevated evolutionary conservation in the exons of CLC genes would

8    appear to support their functionality as a mature RNA transcript, in contrast to the act of their

9    transcription alone (Latos et al. 2012). Another intriguing observation has been the colocalisation of

10    cancer lncRNAs with known protein-coding cancer genes: these are genomically proximal and

11    exhibit elevated expression correlation, even after normalising for proximity. This points to a

12    regulatory link between cancer lncRNAs and protein-coding genes, perhaps through chromatin

13    looping, as described in previous reports for *CCAT1* and *MYC*, for example (Xiang et al. 2014).

14    One important caveat for all features discussed here is ascertainment bias: almost all lncRNAs

15    discussed have been curated from published, single-gene studies. It is entirely possible that selection

16    of genes for initial studies was highly non-random, and influenced by a number of factors – including

17    high expression, evolutionary conservation and proximity to known cancer genes – that could bias

18    our inference of lncRNA features. This may be the explanation for the observed excess of cancer

19    lncRNAs in divergent configuration to protein-coding genes. However, the general validity of some

20    of the CLC-specific features described here – including high expression and evolutionary

21    conservation – were also observed recent unbiased genome-wide screens (Lanzós et al. 2017; Liu et

22    al. 2017), suggesting that they are genuine.

23    In principle, lncRNAs may play two distinct roles in cancer: first, as "driver genes" whose

24    mutations are early and positively-selected events in tumorigenesis; or second, as "downstream

25    genes", which do make a genuine contribution to cancer phenotypes, but through non-genetic

26    alterations in cellular networks resulting from changes in expression, localisation or molecular

27    interactions. These downstream genes may not be expected to display positively-selected mutational

28    patterns, but would be expected to display cancer-specific alterations in expression. A key question

29    for the future is how lncRNAs break down between these two categories, and whether this ratio is

30    distinct from protein-coding genes. Despite the care taken to annotate CLC, it might be expected to

31    contain both types of cancer lncRNA.

32    Despite the relatively low concordance of CLC genes with PCAWG driver predictions, the

33    results of this study strongly support the value and key cancer role of identified lncRNAs in cancer.

34    Most notably, the existence of a core set of eight lncRNAs with independently-identified mouse

1    orthologues with similar cancer functions, is powerful evidence that these genes are *bona fide* cancer

2    genes, whose overexpression or silencing can drive tumour formation. To our knowledge this is the

3    most direct demonstration to date of anciently-conserved functions and disease roles for lncRNAs. It

4    will be intriguing to investigate in future whether more human-mouse orthologous lncRNAs have

5    been identified in such screens.

6

**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 criterion was allowed to link an lncRNA to a cancer type, only if 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).

14

1    CLC set at this stage relies on GENCODE v24 annotation, and therefore all CLC genes have a

2    GENCODE v24 ID assigned. However, data relative to GENCODE v24 was not available for all

3    types of data used in this study (ie. all data relative to PCAWG is based on GENCODE v19). Thus,

4    for some analysis only genes also present in GENCODE v19 could be used (specified in the

5    corresponding methods section) and the total number of genes analysed in these cases is slightly lower

6    (107 instead of 122 CLC genes and 13503 instead of 15827 nonCLC).

7

**8    LncRNA and protein-coding driver prediction analysis**

9    LncRNA and protein-coding predictions for ExInAtor and the rest of PCAWG methods, as well

10    as the combined list of drivers, were extracted from the consortium database (PCAWG Consortium,

11    Manuscript In Preparation). Parameters and details about each individual methods and the combined

12    list of drivers can be found on the main PCAWG driver publication (PCAWG Consortium,

13    Manuscript In Preparation) and false discovery rate correction was applied on each individual cancer

14    type for each individual method in order to define candidates ($q$-value cutoff of 0.1). To calculate

15    sensitivity (percentage of true positives that are predicted as candidates) and precision (percentage of

16    predicted candidates that are true positives) for lncRNA and protein-coding predictions we used the

17    CLC and CGC (COSMIC v78, downloaded Oct, 3, 2016) sets, respectively. To assess the statistical

18    significance of sensitivity rates, we used Fisher's exact test.

19

**20    Feature Identification**

21    We compiled several quantitative and qualitative traits of GENCODE lncRNAs and used them

22    to compare CLC genes to the rest of lncRNAs (referred to as "nonCLC"). Analysis of quantitative

23    traits were performed using Wilcoxon test while qualitative traits were tested using Fisher' exact test.

24    These methods principally refer to Figure 4.

25    <u>Sequence / gene properties:</u> Exonic positions of each gene were defined as the projection of the

26    union of exons from all its transcripts. Introns were defined as all remaining non-exonic nucleotides

27    within the gene span. Repeats coverage refers to the percent of exonic nucleotides of a given gene

28    overlapping repeats and low complexity DNA sequence regions obtained from RepeatMasker data

29    housed in the UCSC Genome Browser (Tyner et al. 2017). Exonic content refers to the fraction of

30    total gene span covered by exons. For this section we used GENCODE v19.

31    <u>Evolutionary conservation:</u> Two types of PhastCons conservation data were used: base-level

32    scores and conserved elements. These data for different multispecies alignments (GRCh38/hg38)

33    were downloaded from UCSC genome browser (Tyner et al. 2017). Mean scores and percent overlap

15

1    by elements were calculated for exons and promoter regions (GENCODE v24). Promoters were

2    defined as the 200nt region centred on the annotated gene start.

3    Expression: We used polyA+ RNA-seq data from 10 human cell lines produced by ENCODE

4    (Djebali et al. 2012; ENCODE Project Consortium et al. 2012) and from various human tissues by

5    the Illumina Human Body Map Project (HBM) (www.illumina.com; ArrayExpress ID: E-MTAB-

6    513) (GENCODE v19).

7    Cancer SNPs: On October, 4, 2016, we collected all 2,192 SNPs related to "cancer", "tumour"

8    and "tumor" terms in the NHGRI-EBI Catalog of published genome-wide association studies

9    (Hindorff et al. 2009; Welter et al. 2014) (https://www.ebi.ac.uk/gwas/home). Then we calculated the

10    closest SNP to each lncRNA TSS using *closest* function from Bedtools v2.19 (GENCODE v24).

11    Epigenetically-silenced lncRNAs: We obtained a published list of 203 cancer-associated

12    epigenetically-silenced lncRNA genes present in GENCODE v24 (Yan et al. 2015). These candidates

13    were identified due to DNA methylation alterations in their promoter regions affecting their

14    expression in several cancer types.

15    Differentially expressed in cancer: We collected a list of 3,533 differentially-expressed lncRNAs

16    in cancer (Yan et al. 2015) (GENCODE v24).

17    Distance to protein-coding genes and CGC genes: For each lncRNA we calculated the TSS to

18    TSS distance to the closest protein-coding gene (GENCODE v24) or CGC gene (downloaded on

19    October, 3, 2016 from Cosmic database)(Futreal et al. 2004) using *closest* function from Bedtools

20    v2.19.

21    Coexpression with closest CGC gene: We took CLC-CGC gene pairs whose TSS-TSS distance

22    was <200kb. RNAseq data from 11 human cell lines from ENCODE was used to assess expression

23    levels (Djebali et al. 2012; ENCODE Project Consortium et al. 2012). ENCODE RNAseq data were

24    obtained from ENCODE Data Coordination Centre (DCC) in September 2016,

25    https://www.encodeproject.org/matrix/?type=Experiment. All data is relative to GENCODE v24. We

26    calculated the expression correlation of gene pairs within each of the 11 cell lines, using the Pearson

27    measure. To control for the effect of proximity, we randomly sampled a subset of nonCLC-CGC pairs

28    matching the same TSS-TSS distance distribution as above, and performed the same expression

29    correlation analysis ("nonCLC-CGC"). Finally, to further control for the fact that CLC and CGC are

30    both cancer genes, which may influence their expression correlation, we shuffled CLC-CGC pairs

31    1000 times, and tested expression correlation for each set ("Shuffled CLC-CGC"). Finally, we tested

32    if correlation coefficients of the three groups were significantly different (Kolmogorov Smirnov test).

33    Genomic classification: We used an in-house script to classify lncRNA transcripts into different

34    genomic categories based on their orientation and proximity to the closest protein-coding gene

16

1   (GENCODE v24): a 10 kb distance was used to distinguish "genic" from "intergenic" lncRNAs.

2   When transcripts belonging to the same gene had different classifications, we used the category

3   represented by the largest number of transcripts.

4       Functional enrichment analysis: The list of protein-coding genes (GENCODE v24) that are

5   divergent and closer than 10 kb to CLC genes (or nonCLC) was used for a functional enrichment

6   analysis (20 unique genes in the case of CLC analysis and 1202 in the case of nonCLC analysis). We

7   show data obtained using g:Profiler web server (Reimand et al. 2016), g:GOSt, with default

8   parameters for functional enrichment analysis of protein-coding genes divergent to CLC and using

9   Bonferroni correction for protein-coding gene divergent to nonCLC. For CLC analysis we performed

10  the same test with independent methods: Metascape (http://metascape.org) (Tripathi et al. 2015) and

11  GeneOntoloy (Panther classification system)(Mi et al. 2013; Mi et al. 2017). In both cases similar

12  results were found.

13

14  **Mouse mutagenesis screen analysis**

15      We extracted the genomic coordinates of transposon common insertion sites (CISs) in Mouse

16  (GRCm38/mm10) http://ccgd-starrlab.oit.umn.edu/about.php (Abbott et al. 2015). This database

17  contains target sites identified by transposon-based forward genetic screens in mice. LiftOver (Kent

18  et al. 2002) was used at default settings to obtain aligned human genome coordinates (hCISs)

19  (GRCh38/hg38). We discarded hCIS regions longer than 1000 nucleotides and intersecting protein

20  coding genes, and intersected the remainders with the genomic coordinates CLC and nonCLC genes.

21      To correctly assess the statistical enrichment of CLC in hCIS regions, we performed two control

22  analyses:

23      Randomly repositioning of CLC and nonCLC genes: We randomly relocated CLC/nonCLC

24  genes 10,000 times within the whole genome using the tool *shuffle* from BedTools v19. In each

25  iteration, we calculated the number of genes that intersected at least one hCIS, and created the

26  distribution of these simulated values. Finally, we calculated an empirical *p*-value by counting how

27  many of the simulated values were higher or equal than the real values. This analysis was performed

28  separately for CLC and nonCLC genes.

29      Length-matched sampling: To calculate if the enrichment of hCIS intersecting genes in CLC set

30  is higher and statistically different from nonCLC set, while controlling by gene length, we created

31  1000 samples of nonCLC genes with the same gene length distribution as CLC genes. Each sample

32  was intersected with hCIS, and the number of intersecting hCISs per Mb of gene length was

33  calculated. To create the nonCLC samples we used the *matchDistribution* script:

34  https://github.com/julienlag/matchDistribution.

1 **<u>Acknowledgements</u>**

9

## **Contributions**

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

**Figure Legends**

**Figure 1: Overview of the Cancer LncRNA Census.** Rows represent the 122 CLC genes, columns represent 31 cancer types. 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) Shown are the total number of unique human lncRNAs contained in each of the three indicated databases (note that for LncRNADisease, numbers refer only to cancer-related genes) and in the CRL set (cancer-related lncRNA genes defined by Lanzós et al. 2017). Colours represent the number of genes that: are present in CLC (green); belong to GENCODE v24 annotation but not CLC (pink); do not belong to GENCODE annotation, and therefore wouldn't meet the criteria for CLC (red). (B) Barplot shows the inverse overlap, being the proportion of CLC genes present in the other databases.

**Figure 3: CLC as benchmark for cancer driver predictions.** (A) CLC benchmarking of ExInAtor 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 genes 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 hypothetic feature analysis example to illustrate the content of the following figures. All panels in figure 4 display features (dots), plotted by their log fold difference (Odds Ratio in case of panel "B") between CLC / nonCLC genesets ($y$-axis) and statistical significance ($x$-axis). In all plots dark and light green dashed lines indicate 0.05 and 0.01 significance thresholds, respectively. (B) Cancer-related data from indicated sources: $y$-axis shows the log2 of the Odds Ratio obtained by comparing CLC to nonCLC

20

1    by Fisher's exact test; *x*-axis displays the estimated *p*-value from the same test. "CGC 1 kb TSS"

2    refers to the fraction of genes that have a nearby known CGC cancer protein-coding gene. This is

3    explored in more detail in the next Figure. (C) Sequence and gene properties: *y*-axis shows the log2

4    fold-difference of CLC / nonCLC means; *x*-axis represents the *p*-value obtained. (D) Evolutionary

5    conservation: "Phastc mean" indicates average base-level PhastCons score; "Elements" indicates

6    percent coverage by PhastCons conserved elements (see Materials and Methods). Colours distinguish

7    exons (blue) and promoters (purple). (E) Gene expression: cell lines (blue) and tissues (purple). For

8    B-D, statistical significance was calculated using Wilcoxon test.

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

10    (A) Cumulative distribution of the genomic distance of lncRNA transcription start site (TSS) to the

11    closest Cancer Gene Census (CGC) (protein-coding) gene TSS. *p*-value was calculated using

12    Wilcoxon test. (B) Boxplot shows the distribution of the gene expression correlation between CLC

13    and their closest CGC genes in 11 human cell lines, including two control analyses (distance-matched

14    nonCLC-CGC pairs, and shuffled CLC-CGC pairs). Correlation was calculated for gene pairs within

15    each cell type, using Pearson method. *p*-value for Kolmogorov-Smirnov test is shown. (C) Genomic

16    classification of lncRNAs. Genes are classified according to distance and orientation to the closest

17    protein-coding gene, and these are grouped into three categories: genes closer than 10kb to closest

18    protein-coding gene, genes overlapping a protein-coding gene and intergenic genes (>10kb from

19    closest protein-coding gene). *p*-values for Fisher's exact tests are shown. (D) The percentage of

20    divergent CLC (left bar) and nonCLC (right bar) genes divergent to a cancer protein-coding gene

21    (CGC). Numbers represent numbers of genes with which the percentage is calculated. *p*-value for

22    Fisher's exact test is shown. (E) Functional annotations of the 20 protein-coding genes divergent to

23    CLC genes from Panel C. Bars indicate the –log10 (corrected) *p*-value (see Materials and Methods)

24    and are coloured based on the "enrichment": the number of genes that contain the functional term

25    divided by the total number of queried genes. Numbers at the end of the bars correspond to the number

26    of genes that fall into the category.

27    **Figure 6: Evidence for ancient conserved cancer roles of lncRNAs.** (A) Functional

28    conservation of human CLC genes was inferred by the presence of Common Insertion Sites (CIS),

29    identified in transposon mutagenesis screens, at orthologous regions in the mouse genome. Orthology

30    was inferred from Chain alignments and identified using LiftOver utility. (B) Number of CLC and

31    nonCLC genes that contain human orthologous common insertion sites (hCIS) (see Table 1).

32    Significance was calculated using Fisher's exact test. (C) UCSC browser screenshot of a CLC gene

33    (*SLNCR1*, ENSG00000227036) intersecting a CIS (yellow arrow). (D) Distribution of the number of

34    intersecting hCIS per Megabase (Mb) of total gene length, for 1000 subsets of nonCLC genes with

1    same length distribution as CLC genes (grey). Vertical blue line represents the overall value for CLC

2    geneset: 1.42 hCIS sites per Mb of gene span. (E) Distribution of the number of genes overlapping a

3    hCIS after 10,000 genomic randomizations of CLC genes (blue). Vertical black line represents the

4    observed number of CLC genes (8) that intersect a hCIS. (F) Distribution of the number of

5    intersecting genes with a hCIS after 10,000 genomic randomizations of nonCLC genes (grey).

6    Vertical black line represents the observed number of nonCLC genes that intersect a hCIS (64).

## 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 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) Number of CLC and nonCLC genes that are proximal to the nearest CGC protein-coding gene. Proximity is defined as a TSS-TSS distance of <10 kb. Significance was calculated using Fisher's exact test.

1    **<u>Supplementary Table Legends:</u>**

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.**

1 **References**

2

3 Abbott, K.L. et al., 2015. The candidate cancer gene database: A database of cancer driver genes

4     from forward genetic screens in mice. *Nucleic Acids Research*, 43(D1), pp.D844–D848.

5 Cabili, M.N. et al., 2011. Integrative annotation of human large intergenic noncoding RNAs reveals

6     global properties and specific subclasses. *Genes & development*, 25(18), pp.1915–27.

7     Available at: http://www.ncbi.nlm.nih.gov/pubmed/21890647 [Accessed March 7, 2017].

8 Chang, K. et al., 2013. The Cancer Genome Atlas Pan-Cancer analysis project. *Nature genetics*,

9     45(10), pp.1113–1120. Available at: http://www.ncbi.nlm.nih.gov/pubmed/24071849.

10 Chen, G. et al., 2013. LncRNADisease: a database for long-non-coding RNA-associated diseases.

11     *Nucleic acids research*, 41(Database issue), pp.D983-6. Available at:

12     http://www.ncbi.nlm.nih.gov/pubmed/23175614 [Accessed March 2, 2017].

13 Copeland, N.G. & Jenkins, N.A., 2010. Harnessing transposons for cancer gene discovery. *Nature*

14     *reviews. Cancer*, 10(10), pp.696–706. Available at:

15     http://www.nature.com/doifinder/10.1038/nrc2916 [Accessed March 7, 2017].

16 Derrien, T. et al., 2012. The GENCODE v7 catalog of human long noncoding RNAs: analysis of

17     their gene structure, evolution, and expression. *Genome research*, 22(9), pp.1775–89.

18     Available at: http://genome.cshlp.org/content/22/9/1775.long [Accessed May 23, 2014].

19 Djebali, S. et al., 2012. Landscape of transcription in human cells. *Nature*, 489(7414), pp.101–108.

20 ENCODE Project Consortium, T. et al., 2012. An integrated encyclopedia of DNA elements in the

21     human genome. *Nature*, 488.

22 Engreitz, J.M. et al., 2014. RNA-RNA interactions enable specific targeting of noncoding RNAs to

23     nascent Pre-mRNAs and chromatin sites. *Cell*, 159(1), pp.188–99. Available at:

24     http://www.ncbi.nlm.nih.gov/pubmed/25259926 [Accessed March 2, 2017].

25 Furney, S. et al., 2006. Structural and functional properties of genes involved in human cancer.

26     *BMC Genomics*, 7(1), p.3. Available at: http://www.ncbi.nlm.nih.gov/pubmed/16405732

27     [Accessed February 24, 2017].

28 Furney, S.J. et al., 2008. Distinct patterns in the regulation and evolution of human cancer genes. *In*

29     *Silico Biol*, 8(December 2007), pp.33–46.

30 Futreal, P. et al., 2004. A census of human cancer genes. *Nat Rev Cancer*, 4(3), pp.177–183.

31 Grote, P. et al., 2013. *The Tissue-Specific lncRNA Fendrr Is an Essential Regulator of Heart and*

32     *Body Wall Development in the Mouse*, Available at:

33     http://www.sciencedirect.com/science/article/pii/S1534580712005862 [Accessed April 25,

34     2017].

Gutschner, T. et al., 2013. The noncoding RNA MALAT1 is a critical regulator of the metastasis phenotype of lung cancer cells. *Cancer research*, 73(3), pp.1180–9. Available at: http://www.ncbi.nlm.nih.gov/pubmed/23243023 [Accessed March 2, 2017].

Gutschner, T. & Diederichs, S., 2012. The hallmarks of cancer: a long non-coding RNA point of view. *RNA biology*, 9(6), pp.703–19. Available at: http://www.ncbi.nlm.nih.gov/pubmed/22664915 [Accessed June 28, 2016].

Guttman, M. et al., 2009. Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in mammals. *Nature*, 458(7235), pp.223–7. Available at: http://www.ncbi.nlm.nih.gov/pubmed/19182780 [Accessed March 7, 2017].

Guttman, M. & Rinn, J.L., 2012. Modular regulatory principles of large non-coding RNAs. *Nature*, 482(7385), pp.339–46. Available at: http://www.ncbi.nlm.nih.gov/pubmed/22337053 [Accessed March 2, 2017].

Harrow, J. et al., 2012. GENCODE: The reference human genome annotation for The ENCODE Project. *Genome Research*, 22(9), pp.1760–1774. Available at: http://www.ncbi.nlm.nih.gov/pubmed/22955987 [Accessed April 25, 2017].

Higashimoto, K. et al., 2006. Imprinting disruption of the CDKN1C/KCNQ1OT1 domain: the molecular mechanisms causing Beckwith-Wiedemann syndrome and cancer. *Cytogenetic and genome research*, 113(1–4), pp.306–12. Available at: http://www.ncbi.nlm.nih.gov/pubmed/16575194 [Accessed April 27, 2017].

Hindorff, L. a et al., 2009. Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proceedings of the National Academy of Sciences of the United States of America*, 106(23), pp.9362–7. Available at: http://www.ncbi.nlm.nih.gov/pubmed/19474294.

Huarte, M. et al., 2010. A large intergenic noncoding RNA induced by p53 mediates global gene repression in the p53 response. *Cell*, 142(3), pp.409–19. Available at: http://www.ncbi.nlm.nih.gov/pubmed/20673990 [Accessed February 25, 2017].

Iyer, M.K. et al., 2015. The landscape of long noncoding RNAs in the human transcriptome. *Nature GeNetics*, 47(3). Available at: https://www.nature.com/ng/journal/v47/n3/pdf/ng.3192.pdf [Accessed April 19, 2017].

Jia, H. et al., 2010. Genome-wide computational identification and manual annotation of human long noncoding RNA genes. *RNA (New York, N.Y.)*, 16(8), pp.1478–87. Available at: http://rnajournal.cshlp.org/cgi/doi/10.1261/rna.1951310 [Accessed March 7, 2017].

Johnson, R. & Guigó, R., 2014. The RIDL hypothesis: transposable elements as functional domains 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   Kent, W.J. et al., 2002. The human genome browser at UCSC. *Genome research*, 12(6), pp.996–

6       1006. Available at: http://www.ncbi.nlm.nih.gov/pubmed/12045153 [Accessed April 20,

7       2017].

8   Lanz, R.B. et al., 1999. A steroid receptor coactivator, SRA, functions as an RNA and is present in

9       an SRC-1 complex. *Cell*, 97(1), pp.17–27. Available at:

10      http://www.ncbi.nlm.nih.gov/pubmed/10199399 [Accessed April 25, 2017].

11  Lanzós, A. et al., 2017. Discovery of Cancer Driver Long Noncoding RNAs across 1112 Tumour

12      Genomes: New Candidates and Distinguishing Features. *Scientific Reports*, 7, p.41544.

13      Available at: http://www.ncbi.nlm.nih.gov/pubmed/28128360 [Accessed February 24, 2017].

14  Latos, P.A. et al., 2012. Airn Transcriptional Overlap, But Not Its lncRNA Products, Induces

15      Imprinted Igf2r Silencing. *Science*, 338(6113), pp.1469–1472. Available at:

16      http://www.sciencemag.org/cgi/doi/10.1126/science.1228110 [Accessed April 25, 2017].

17  Lawrence, M.S. et al., 2014. Discovery and saturation analysis of cancer genes across 21 tumour

18      types. *Nature*, 505(7484), pp.495–501. Available at:

19      http://www.nature.com/doifinder/10.1038/nature12912 [Accessed January 13, 2017].

20  Lepoivre, C. et al., 2013. Divergent transcription is associated with promoters of transcriptional

21      regulators. *BMC Genomics*, 14(1), p.914. Available at:

22      http://www.ncbi.nlm.nih.gov/pubmed/24365181 [Accessed April 25, 2017].

23  Liu, S.J. et al., 2017. CRISPRi-based genome-scale identification of functional long noncoding

24      RNA loci in human cells. *Science*, 355(6320), p.eaah7111. Available at:

25      http://www.sciencemag.org/lookup/doi/10.1126/science.aah7111 [Accessed April 25, 2017].

26  Managadze, D. et al., 2011. Negative Correlation between Expression Level and Evolutionary Rate

27      of Long Intergenic Noncoding RNAs. *Genome Biology and Evolution*, 3(0), pp.1390–1404.

28      Available at: http://www.ncbi.nlm.nih.gov/pubmed/22071789 [Accessed March 26, 2017].

29  Marchese, F.P. et al., 2016. A Long Noncoding RNA Regulates Sister Chromatid Cohesion.

30      *Molecular Cell*, 63(3), pp.397–407. Available at:

31      http://www.ncbi.nlm.nih.gov/pubmed/27477908 [Accessed April 25, 2017].

32  Marques, A.C. et al., 2013. Chromatin signatures at transcriptional start sites separate two equally

33      populated yet distinct classes of intergenic long noncoding RNAs. *Genome biology*, 14(11),

34      p.R131. Available at: http://genomebiology.biomedcentral.com/articles/10.1186/gb-2013-14-

11-r131 [Accessed March 20, 2017].

Mi, H. et al., 2013. Large-scale gene function analysis with the PANTHER classification system. *Nature Protocols*, 8(8), pp.1551–1566. Available at: http://www.ncbi.nlm.nih.gov/pubmed/23868073 [Accessed April 23, 2017].

Mi, H. et al., 2017. PANTHER version 11: expanded annotation data from Gene Ontology and Reactome pathways, and data analysis tool enhancements. *Nucleic Acids Research*, 45(D1), pp.D183–D189. Available at: https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkw1138 [Accessed April 23, 2017].

Mularoni, L. et al., 2016. OncodriveFML: a general framework to identify coding and non-coding regions with cancer driver mutations. *Genome biology*, 17(1), p.128. Available at: http://www.ncbi.nlm.nih.gov/pubmed/27311963 [Accessed June 28, 2016].

Ning, S. et al., 2016. Lnc2Cancer: a manually curated database of experimentally supported lncRNAs associated with various human cancers. *Nucleic acids research*, 44(D1), pp.D980-5. Available at: http://www.ncbi.nlm.nih.gov/pubmed/26481356 [Accessed March 2, 2017].

Ponjavic, J. et al., 2009. Genomic and transcriptional co-localization of protein-coding and long non-coding RNA pairs in the developing brain. Y. Hayashizaki, ed. *PLoS genetics*, 5(8), p.e1000617. Available at: http://dx.plos.org/10.1371/journal.pgen.1000617 [Accessed March 20, 2017].

Quek, X.C. et al., 2015. lncRNAdb v2.0: expanding the reference database for functional long noncoding RNAs. *Nucleic acids research*, 43(Database issue), pp.D168-73. Available at: http://www.ncbi.nlm.nih.gov/pubmed/25332394 [Accessed March 2, 2017].

Redon, R. et al., 2006. Global variation in copy number in the human genome. *Nature*, 444(7118), pp.444–54. Available at: http://dx.doi.org/10.1038/nature05329.

Reimand, U. et al., 2016. g:Profiler—a web server for functional interpretation of gene lists (2016 update). *Nucleic Acids Research*. Available at: http://biit.cs.ut.ee/gprofiler/doc/papers/gprofiler_nar_2016.pdf [Accessed April 20, 2017].

Sabarinathan, R. et al., 2013. RNAsnp: Efficient Detection of Local RNA Secondary Structure Changes Induced by SNPs. *Human Mutation*, 34(4), p.n/a-n/a. Available at: http://doi.wiley.com/10.1002/humu.22273 [Accessed April 19, 2017].

Sauvageau, M. et al., 2013. Multiple knockout mouse models reveal lincRNAs are required for life and brain development. *eLife*, 2, p.e01749. Available at: http://www.ncbi.nlm.nih.gov/pubmed/24381249 [Accessed April 25, 2017].

Schmidt, K. et al., 2016. The lncRNA SLNCR1 Mediates Melanoma Invasion through a Conserved SRA1-like Region. *Cell reports*, 15(9), pp.2025–37. Available at:

1    http://linkinghub.elsevier.com/retrieve/pii/S2211124716304314 [Accessed March 20, 2017].

2    Schmitt, A.M. et al., 2016. Long Noncoding RNAs in Cancer Pathways. *Cancer Cell*, 29(4),

3        pp.452–463. Available at: http://www.ncbi.nlm.nih.gov/pubmed/27070700 [Accessed April

4        20, 2017].

5    Sjoblom, T. et al., 2006. The Consensus Coding Sequences of Human Breast and Colorectal

6        Cancers. *Science*, 314(5797), pp.268–274. Available at:

7        http://www.ncbi.nlm.nih.gov/pubmed/16959974 [Accessed February 24, 2017].

8    Tamborero, D. et al., 2013. Comprehensive identification of mutational cancer driver genes across

9        12 tumor types. *Scientific reports*, 3, p.2650. Available at:

10       http://www.ncbi.nlm.nih.gov/pubmed/24084849 [Accessed June 28, 2016].

11   Tan, J.Y. et al., 2017. cis -Acting Complex-Trait-Associated lincRNA Expression Correlates with

12       Modulation of Chromosomal Architecture. *Cell Reports*, 18(9), pp.2280–2288. Available at:

13       http://www.ncbi.nlm.nih.gov/pubmed/28249171 [Accessed April 25, 2017].

14   Tripathi, S. et al., 2015. Meta- and Orthogonal Integration of Influenza &quot;OMICs&quot; Data

15       Defines a Role for UBR4 in Virus Budding. *Cell host & microbe*, 18(6), pp.723–35. Available

16       at: http://linkinghub.elsevier.com/retrieve/pii/S1931312815004564 [Accessed April 23, 2017].

17   Tyner, C. et al., 2017. The UCSC Genome Browser database: 2017 update. *Nucleic Acids Research*,

18       45. Available at: https://oup.silverchair-

19       cdn.com/oup/backfile/Content_public/Journal/nar/45/D1/10.1093_nar_gkw1134/3/gkw1134.p

20       df?Expires=1493043886&Signature=VkpgKibgiT9kty13UqDbQabocXaBUaciIshI~KU4D01

21       0dRoa-

22       9n7qLF0kT3eT2HZRiqrI7W74jyxPg1eyKhuPrlOzHFwCJWxa3tO3wh95deEodxSyEwvjA64

23       rFAZFbmv0EtdUoWRqL5nhsJqLSPCZmPXukDpSxBH7SrrmYX33UqRcVo6jq-

24       ICvde4XmPNDgacH4BwRTU2K0~D-OeV~kzq6s-zshWmPjUIJM-

25       XCB7Mpx2kd5JVIN7lPVCt0vs8gK~BDHHdryJEkWLf9L4ZbFxzvLvtulvQUZnrNSNZFLR

26       PPDHFnJ5C7YZM8-U5g27ma2eGbwIQyjZ1qqbcxOg6QLkLw__&Key-Pair-

27       Id=APKAIUCZBIA4LVPAVW3Q [Accessed April 20, 2017].

28   Ulitsky, I. & Bartel, D.P., 2013. lincRNAs: Genomics, Evolution, and Mechanisms. *Cell*, 154(1),

29       pp.26–46. Available at: http://linkinghub.elsevier.com/retrieve/pii/S0092867413007599

30       [Accessed November 20, 2016].

31   Welter, D. et al., 2014. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations.

32       *Nucleic Acids Research*, 42(D1), pp.1001–1006.

33   World Health Organization, 2013. *International Classification of Diseases for Oncology (ICD-O).*

34       *Third Edition. First Revision.*,

1  Xiang, J.-F. et al., 2014. Human colorectal cancer-specific CCAT1-L lncRNA regulates long-range

2  chromatin interactions at the MYC locus. *Cell Research*, 24(5), pp.513–531. Available at:

3  http://www.nature.com/doifinder/10.1038/cr.2014.35 [Accessed April 25, 2017].

4  Yan, X. et al., 2015. Comprehensive Genomic Characterization of Long Non-coding RNAs across

5  Human Cancers. *Cancer Cell*.

6  Yates, L.R. & Campbell, P.J., 2012. Evolution of the cancer genome. *Nature Reviews Genetics*,

7  13(11), pp.795–806. Available at: http://www.nature.com/doifinder/10.1038/nrg3317

8  [Accessed April 25, 2017].

9  Zhu, S. et al., 2016. Genome-scale deletion screening of human long non-coding RNAs using a

10  paired-guide RNA CRISPR–Cas9 library. *Nature Biotechnology*, 34(12), pp.1279–1286.

11  Available at: http://www.nature.com/doifinder/10.1038/nbt.3715 [Accessed December 21,

12  2016].

13

14

15

1   **Tables**

2

3   **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

## A

### Genes from lncRNA databases present in CLC

Legend:
- ■ Not present in GENCODE v24
- ■ Not present in CLC
- ■ Present in CLC

## B

### CLC genes present in other lncRNA databases



Covered  ■ NO  ■ YES

# Figure 3



**A** — ExInAtor $q$-value < 0.1
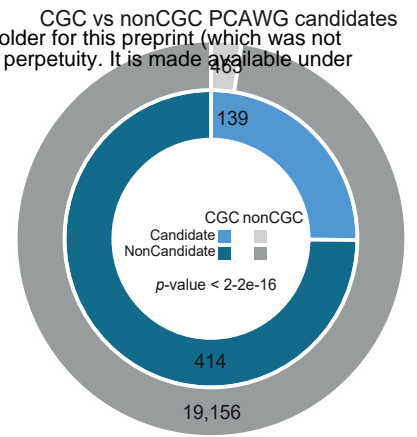
ColoRect AdenoCA (n=0/3)
Lymph tumors (n=2/5)
PANCAN (n=3/10)
Skin Melanoma (n=1/2)
Uterus AdenoCA (n=0/2)
Baseline

Percentage of genes in CLC list

Genes ranked by $q$−value

**B** — CLC vs nonCLC PCAWG candidates

CLC nonCLC
Candidate
NonCandidate

$p$-value = 0.007

86
8
74
2,519

**C** — CGC vs nonCGC PCAWG candidates

CGC nonCGC
Candidate
NonCandidate

$p$-value < 2-2e-16

463
139
414
19,156

# Figure 4

**A**

Figure 4

# Figure 5

A

B Genes correlation within cells

C Genomic Classification

D

E 20 pc-genes divergent to CLC

# Figure 6



**A**

LncRNA

Chain alignment

Insertional mutagenesis hotspot

Cancer phenotype

**B** CLC and hCIS loci

64

8

CLC  nonCLC
Overlap
Not Overlap

*p*-value = 7.39e-08

114

15,763

**C**

SLNCR1

Human hg38

Scale  200 kb  hg38
chr17:  72,150,000  72,200,000  72,250,000  72,300,000  72,350,000  72,400,000  72,450,000  72,500,000  72,550,000  72,600,000  72,650,000
Basic Gene Annotation Set from GENCODE Version 24 (Ensembl 83)

SOX9-AS1  SOX9-AS1  LINC00511  LINC00511  SLC39A11
SOX9-AS1  SOX9-AS1  LINC00511  RP11-1124B17.1  LINC00511  SLC39A11
SOX9-AS1  LINC00511  LINC00511  LINC00511  CTD-3010D24.3
SOX9-AS1  LINC00511  RN7SKP180
AC007461.2  SOX9-AS1  RP11-57A1.1
SOX9-AS1
SOX9-AS1  SOX9

Mouse mm10

Scale  200 kb  mm10
chr11:  112,800,000  112,850,000  112,900,000  112,950,000  113,000,000  113,050,000  113,100,000  113,150,000  113,200,000  113,250,000
Basic Gene Annotation Set from ENCODE/GENCODE Version M11 (Ensembl 86)

Sox9  2610035D17Rik  Slc39a11
Gm11681  2610035D17Rik  Slc39a11
4933434M16Rik  4732490B19Rik  Slc39a11

hCIS clusters  CCGD clusters

**D**

Sampled length-matched nonCLC vs CLC

Number of hCIS intersecting a CLC gene / CLC genes length (Mb) = 1.42

Counts

600

400

200

0

0  1  2  3  4

#hCIS intersecting a nonCLC gene / nonCLC genes length (Mb)

**E**

Random reposition of CLC

Observed = 8
*p*-value = 0.0022

Counts

2000

1000

0

0  3  6

#CLC genes intersecting hCIS

**F**

Random reposition of nonCLC

Observed = 64
*p*-value = 0

Counts

400

300

200

100

0

60  90  120  150

#nonCLC genes intersecting hCIS