

1 **Discovery of Cancer Driver Long Noncoding RNAs across 1112 Tumour Genomes: New Candidates**  
2 **and Distinguishing Features.**

3

4 **Authors**

5 Andrés Lanzós<sup>1,2,3</sup>, Joana Carlevaro-Fita<sup>1,2,3</sup>, Loris Mularoni<sup>4</sup>, Ferran Reverter<sup>1,2,3</sup>, Emilio Palumbo<sup>1,2,3</sup>,  
6 Roderic Guigó<sup>1,2,3</sup>, Rory Johnson<sup>1,2,3\*</sup>

7 1. Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Dr.  
8 Aiguader 88, Barcelona 08003, Spain

9 2. Universitat Pompeu Fabra (UPF), Barcelona, Spain.

10 3. Institut Hospital del Mar d'Investigacions Mèdiques (IMIM), 08003 Barcelona, Spain.

11 4. Research Unit on Biomedical Informatics, Department of Experimental and Health Sciences, Universitat  
12 Pompeu Fabra, Dr. Aiguader 88, Barcelona, Spain

13 \* Correspondence to [Rory.Johnson@crg.eu](mailto:Rory.Johnson@crg.eu)

14 **Keywords:** Cancer; cancer genome; somatic mutation; long non-coding RNA; driver mutation; The Cancer  
15 Genome Atlas; International Cancer Genome Consortium.

16

17

18

1 **Abstract**

2

3 Long noncoding RNAs (lncRNAs) represent a vast unexplored genetic space that may hold missing drivers  
4 of tumorigenesis, but few such “driver lncRNAs” are known. Until now, they have been discovered  
5 through changes in expression, leading to problems in distinguishing between causative roles and passenger  
6 effects. We here present a different approach for driver lncRNA discovery using mutational patterns in  
7 tumour DNA. Our pipeline, ExInAtor, identifies genes with excess load of somatic single nucleotide  
8 variants (SNVs) across panels of tumour genomes. Heterogeneity in mutational signatures between cancer  
9 types and individuals is accounted for using a simple local trinucleotide background model, which yields  
10 high precision and low computational demands. We use ExInAtor to predict drivers from the GENCODE  
11 annotation across 1112 entire genomes from 23 cancer types. Using a stratified approach, we identify 15  
12 high-confidence candidates: 9 novel and 6 known cancer-related genes, including *MALATI*, *NEATI* and  
13 *SAMMSON*. Both known and novel driver lncRNAs are distinguished by elevated gene length, evolutionary  
14 conservation and expression. We have presented a first catalogue of mutated lncRNA genes driving cancer,  
15 which will grow and improve with the application of ExInAtor to future tumour genome projects.

16

17

18

19

## 1 **Introduction**

2 Whole genome sequencing makes it possible to comprehensively discover the mutations, and the mutated  
3 genes, that are responsible for tumour formation. By sequencing pairs of normal and tumour genomes from  
4 large patient cohorts, projects such as the ICGC (International Cancer Genome Consortium) and TCGA  
5 (The Cancer Genome Atlas) aim to create definitive driver mutation catalogues for all common cancers  
6 (1,2). Focussing on entire genomes, rather than just captured exomes, these studies hope to identify driver  
7 elements amongst the ~98% DNA that does not encode protein. These noncoding regions contain a wealth  
8 of regulatory sequences and non-coding RNAs whose role in cancer has been neglected until now (3).

9 Amongst the most numerous, yet poorly understood of the latter are long noncoding RNAs (lncRNAs).  
10 These are long RNA transcripts that share many characteristics of mRNAs, with the key difference that  
11 they do not contain any recognizable Open Reading Frame (ORF), and thus are unlikely to encode protein  
12 (4). lncRNAs perform a diverse range of regulatory activities within both the nucleus and cytoplasm by  
13 interacting with protein complexes or other nucleic acids (5). While their expression tends to be lower than  
14 protein-coding mRNAs, lncRNAs are thought to be highly expressed in a subset of cells in a population  
15 (6). The number of lncRNA genes in the human genome is still uncertain, but probably lies in the range  
16 20,000 - 50,000 (7,8). This vast population of uncharacterized genes likely includes many with novel roles  
17 in cancer.

18 In recent years a small but growing number of lncRNA have been implicated in cancer progression through  
19 various mechanisms (9). lincRNA-P21, a tumour suppressor, acts downstream of P53 by recruiting the  
20 repressor hnRNP-K to target genes (10). Proto-oncogene lncRNAs include HOTAIR, upregulated in  
21 multiple cancers, which recruits the repressive PRC2 chromatin regulatory complex to hundreds of genes  
22 (11). Cancer-related lncRNA have features of functional genes, including sequence conservation,  
23 orthologues in other mammals, chromatin marks and regulated subcellular localisation (4). Moreover they  
24 display typical characteristics of cancer drivers, including influence on cellular phenotypes of proliferation  
25 and apoptosis, and in clinical features such as patient survival and altered expression across tumour  
26 collections (3,8,11).

27 The absence of whole-genome maps of somatic mutations has meant that searches for new cancer-related  
28 lncRNAs have relied on conventional transcriptomic approaches that reveal changes in their expression  
29 levels that accompany cancer. However such approaches are not capable of distinguishing passenger and  
30 driver effects, nor do they identify mutations in the mature lncRNA sequence that may drive tumorigenesis  
31 independent of upstream regulatory changes (8,12,13). Two recent studies clearly demonstrate that somatic  
32 mutations, in these cases amplifications of entire loci, can drive tumour formation (14,15). Nevertheless,  
33 we remain largely ignorant of the role that mutations in lncRNA genes play during the early stages of  
34 tumorigenesis.

35 The statistical analysis of somatic mutation patterns is a powerful means of identifying genes that drive  
36 early tumour formation. For protein-coding genes, tools have been developed to successfully identify  
37 mutations that result in gain or loss of function in translated peptide sequences (16). These tools take  
38 advantage of exome sequencing – the targeted capture and sequencing of approximately 2% of the genome  
39 encoding protein (17). A number of methods applied to this problem for protein-coding genes (16) take  
40 advantage of the fact that cancer-associated genes (both tumour-suppressors and proto-oncogenes) display

1 characteristic and non-random mutational patterns in their protein-coding sequence. They prioritise genes  
2 with mutations that are predicted to result from positive selection on the encoded protein.

3 As such, these methods are generally inapplicable to lncRNA, which do not encode protein, and for which  
4 we presently do not have maps of their functional sequence domains, nor an understanding of the molecular  
5 mechanism of such domains. While capable of discovering protein-coding driver genes, exome sequencing  
6 ignores mutations occurring in the multitude of noncoding regulatory elements known to exist in the human  
7 genome (18). Noncoding driver mutations can be comprehensively discovered for the first time by projects  
8 to sequence collections of entire cancer genomes (1). In the present study, we describe and characterise a  
9 tool, called ExInAtor, for the discovery of driver lncRNA genes. ExInAtor identifies genes with excess of  
10 exonic mutations, compared to the expected local neutral rate estimated from intronic and surrounding  
11 sequences. We present a comprehensive prediction of candidate lncRNAs across 1104 genomes from 23  
12 cancer types. These candidates have a series of features consistent with their being genuine drivers.

13

## 1 **Results**

### 2 **A method for discovering driver genes from cancer genomes**

3 Our aim was to develop a method to identify tumour driver long noncoding RNAs (lncRNAs) using  
4 short nucleotide variant (SNV) mutations from cancer genome sequencing projects. We define SNVs, from  
5 now on, as somatic substitutions or indels of length 1 nt. The majority of lncRNAs are spliced, and we  
6 assume throughout that their functional sequence resides in exonic regions that are incorporated into the  
7 mature transcript (19). Intronic sequence is removed during splicing and hence is not directly relevant to  
8 their function. Consequently, we hypothesised that driver lncRNAs will display an excess of somatic  
9 mutations in exons compared to the local background mutational rate, estimated by their introns and  
10 flanking genomic regions – henceforth referred to as “background regions”. This approach is conservative,  
11 given that background regions are likely to include functional regulatory elements that may themselves  
12 carry driver mutations.

13 We implemented this approach in a computational pipeline called ExInAtor (Fig. 1 and File S1).  
14 ExInAtor requires two principal inputs: an annotation of lncRNA genes and a catalogue of tumour  
15 mutations. At its heart, ExInAtor employs a parametric statistical test to identify genes that present a  
16 significantly elevated exonic mutation rate compared to local background regions. The latter are comprised  
17 of intronic and flanking genomic sequence. We took care to account for a key confounding factor: the  
18 unique mixture of mutational signatures that characterises every individual tumour, and every tumour type  
19 (20). Such signatures can be described as a probability for every nucleotide to mutate to every other,  
20 conditioned on the identity of flanking positions – summarised in a matrix of 96 trinucleotide substitution  
21 frequencies (20). In other words, mutation rates are dependent on nucleotide composition. The mutational  
22 signature must be taken into account when comparing mutational loads of exons to surrounding regions,  
23 because they tend to have marked differences in nucleotide composition – both for protein-coding genes  
24 and lncRNAs (21).

25 ExInAtor employs a subsampling approach to balance the trinucleotide content of exons and  
26 reference regions, thereby accounting for mutational signatures (Fig. 1A). Exonic regions of each gene are  
27 defined as the projection of all exons from the union of its transcripts. Next, the reference region is defined  
28 as all non-exonic nucleotides within the gene, in addition to upstream and downstream windows of defined  
29 length. Within these exonic and reference regions, the frequencies of trinucleotides are calculated. Then,  
30 nucleotides are randomly sampled (without replacement) from the reference region, until the maximum  
31 possible amount of sequence with identical trinucleotide composition has been collected. Now, the number  
32 of SNVs overlapping exons,  $M$ , and those overlapping remaining reference nucleotides,  $m$ , are compared  
33 using a contingency-table analysis and statistical significance is calculated according to hypergeometric  
34 distribution (Fig. 1B) (see Materials and Methods for more details).

35 We prepared a carefully-filtered lncRNA annotation, to avoid several potential sources of false  
36 positive predictions. We were particularly concerned by two potential confounding factors: first,  
37 misinterpretation of mutations that may affect protein-coding regions overlapping the same DNA as  
38 lncRNA exons; and second, the presence of mis-classified protein-coding transcripts among the GENCODE  
39 annotation (4). Thus, we removed genes of uncertain protein-coding potential, as judged by computational  
40 protein-coding potential classifiers (see Materials and Methods). We also removed any lncRNA genes, such  
41 as cis-antisense and intronic lncRNAs, that overlap annotated protein-coding genes. In this way we

1 narrowed the set of GENCODE v19 lncRNA genes from 13,870 to 5,887 intergenic, confidently-noncoding  
2 lncRNAs (Table 1). To this set we added back 27 cancer-related, GENCODE v19 lncRNAs from the  
3 literature (see below).

4 One advantage of ExInAator is its indifference to genes' biotype. This arises from its lack of reliance  
5 on measures of functional impact(22), meaning that it can equally be used on lncRNAs or protein-coding  
6 genes. Indeed, similar approaches have been used to discover coding driver genes in the past (23). We took  
7 advantage of this to assess its ability to discover known protein-coding driver genes from the Cancer Gene  
8 Census (24) amongst the GENCODE annotation. This provided us with a useful independent validation of  
9 ExInAator's precision, of particular value given the low number of known driver lncRNAs at present.

10

## 11 **Datasets of somatic mutations in cancer genomes**

12 To search for lncRNA driver genes, we took advantage of the two largest available sources of  
13 cancer genome mutations: one collected by the Cancer Genome Project at the Sanger Institute, hereafter  
14 named "Alexandrov" (20), and the other from The Cancer Genome Atlas (TCGA) (1) (Table 2). These data  
15 were aggressively filtered to remove potential artefacts arising from germline mutations (see Materials and  
16 Methods). The Alexandrov dataset comprises 9 cancers with between 15 and 119 individuals and 10,436  
17 and 2,796,863 mutations each. The TCGA dataset consists of 14 cancers with between 15 and 96 individuals  
18 and 21,113 to 4,680,653 mutations each. Of note is the large spread in sample sizes and mutation rates  
19 across tumour types. Taking all cancers together, we observed an excess of mutations in lncRNAs over  
20 protein-coding genes, and in background over exons, suggesting a general selective pressure against  
21 disruptive mutations in both gene classes (File S2).

22

## 23 **The landscape of driver lncRNAs across 23 tumour types**

24 To comprehensively discover candidate lncRNA drivers, ExInAator was run on the 23 tumour types  
25 described above. We adopted some analysis strategies to account for the relatively shallow nature of the  
26 data and our consequently weak statistical power to find driver genes. First, in order to discover both cancer-  
27 specific and ubiquitous driver genes, ExInAator was run on each dataset in distinct configurations: (1)  
28 grouping samples by tumour type ("Tumour Specific"), (2) pooling together the entire set of tumours within  
29 each of the two projects ("Pancancer") and (3) pooling data across both projects ("Superpancancer").

30 Second, we used sample stratification to boost sensitivity. This approach is commonly used when  
31 statistical power is reduced by multiple hypothesis testing (25,26). lncRNA genes were divided into two  
32 groups of different sizes, and each was treated independently during multiple hypothesis correction. This  
33 reduces the burden on resulting false discovery rate estimates. As a reference set, we curated 45  
34 experimentally-validated cancer-related lncRNAs from the scientific literature, henceforth "Cancer-Related  
35 lncRNAs" (CRLs) (File S3). All CRL genes belong to GENCODE v19 annotation. Remaining filtered  
36 lncRNAs are referred to as "Non-CRL" (File S4). Summary statistics of the gene sets used are shown in  
37 Table 1.

38 At a Q value (false discovery rate) cutoff of 0.1, we discovered a total of 15 lncRNAs (6 and 9 from  
39 CRL and non-CRL, respectively) (Fig. 2A) (Files S5 and S6) and 24 protein-coding genes (File S7).

1 Relaxing the cutoff to  $Q < 0.2$ , we discover 10 and 27 CRL and non-CRL lncRNAs, respectively. Henceforth  
2 we refer to these as driver genes, and a Q-value threshold of 0.1 is assumed unless stated otherwise.  
3 ExInAator predicted a total of five lncRNA driver genes in Alexandrov tumours, nine in TCGA and two in  
4 Superpancancer (one of them already detected in Pancancer TCGA). The greatest numbers of drivers  
5 predicted in individual tumours were three apiece in Breast and Kidney Chromophore (Fig. 2D).

6 Several findings suggest that false positive prediction rates are low. Reported P values closely  
7 follow the expected null distribution for the majority of genes (a full set of Quantile-quantile (QQ) plots  
8 can be found in File S8). Furthermore, while a number of tumour types display a small number of putative  
9 driver lncRNAs that strongly deviate from the null expectation (exemplified by Breast cancer sample in  
10 Fig. 2B), other samples yield no candidates at all (eg Liver cancer, Fig. 2C). In general, inspection of QQ  
11 plots shows a tendency for deflation of P values (File S8). To further test false discovery rates, we reran  
12 these analyses on tumour data that had been randomised using two different methods (see Materials and  
13 Methods for details). ExInAator predicted no lncRNA drivers in either dataset (grey dots in Figs. 2B&C and  
14 File S8). Together these data point to a rather conservative statistical model, which may discard some *bona*  
15 *fide* drivers. A comprehensive set of predictions across all analyses can be found in File S9.

16

## 17 **ExInAator identifies known and novel lncRNA driver genes**

18 ExInAator's sensitivity is demonstrated by its identification of altogether six CRL genes. These are:  
19 *MALAT1*, *NEATI*, *PCA3*, *BCAR4*, lncRNA-ATB (CTD-2314B22.3) and the recently-discovered  
20 melanoma driver *SAMMSON* (RP11-460N16.1) (Table 3). The latter was detected in stomach  
21 adenocarcinoma, and we found that it is also present in stomach RNAseq (File S10). The majority of  
22 candidates were found in tumour-specific analysis (Fig. 3A). Nevertheless, two CRL lncRNAs, *NEATI* and  
23 *MALAT1*, were identified in Pancancer analysis, consistent with a general role in tumorigenesis: both are  
24 long, unspliced and nuclear-retained lncRNAs with demonstrated roles across a range of cancer types (9).  
25 As shown in Fig. 3B, the *NEATI* exon region experiences an elevated mutation rate across cancers, when  
26 compared to its flanking background regions. *NEATI* was identified in a recent study of liver cancer  
27 genomes, and as the authors pointed out, it cannot be ruled out that it is identified through increased local  
28 mutation rate (27).

29 One important potential source of false positive signal in this study could be elevated mutational  
30 rates in DNA regulatory elements, such as enhancers, which happen to overlap the exon of a lncRNA  
31 annotation. Such cases would be expected to produce driver lncRNAs, where all mutations are concentrated  
32 in a single exon. This would be indistinguishable from *bona fide* driver lncRNAs that have an important  
33 functional domain located in a single exon. To investigate this further, we inspected the exon-level  
34 mutational density of all candidate lncRNAs (File S11). Intriguingly, we find at least two cases where  
35 mutations are elevated across multiple exons, but not intervening introns (Figs. 3C&D). Altogether of 13  
36 multi-exonic candidate lncRNAs, five have mutations in more than one exon. This supports the  
37 interpretation that, for these cases at least, mutations cause gain- or loss-of-function in mature lncRNA  
38 transcripts, and not through disruption of a DNA-encoded element.

39 Amongst the novel candidate driver genes were a number of intriguing cases with various  
40 characteristics of functionality, none of which have been described in the scientific literature. In Figure 3F  
41 we highlight one case, *RP11-820L6.1*, whose promoter is characterised by canonical histone modifications,

1 obvious evolutionary conservation and the recruitment of transcription factors. Most notably, the master  
2 tumour suppressor transcription factor and regulator of several cancer lncRNAs, P53, is bound within the  
3 first intron (28).

4 We further sought to establish the degree of overlap between ExInAator-predicted driver genes and  
5 candidates predicted by transcriptomic analyses. Two previous studies to identify cancer-related lncRNAs  
6 have searched for differentially-expressed transcripts in cancer transcriptome data from microarrays and  
7 RNA sequencing (8,12). From each study we extracted those transcripts that overlap the filtered geneset  
8 used here, retrieving a total of 80 and 186 genes from the Du et al and Iyer et al (“MiTranscriptome”)  
9 studies, respectively (Files S12 and S13)(8,12). Three genes are identified by both ExInAator and  
10 MiTranscriptome (*PCA3*, *NEAT1* and *MALAT1*) ( $P=0.0026$ , Chi-square with Yates' correction test) and  
11 another with Du et al (*PCA3*) ( $P=0.5$ , Chi-square with Yates' correction test) (File S14). It should be noted  
12 that all these genes belong to the CRL set. MiTranscriptome and Du share 11 genes ( $P\leq 0.0001$ , Chi-square  
13 with Yates' correction test). This surprising discordance of driver gene prediction between studies, in  
14 addition to their lack of overall intersection with the published CRL set, suggests that (1) these large-scale  
15 predictions have considerable false negative rates, and (2) that available catalogues of cancer-related  
16 lncRNAs, represented by the CRL set, are incomplete.

17 We searched for independent evidence of cancer roles for ExInAator-predicted candidates.  
18 Importantly, we separately considered (1) the entire set of candidates, including known CRL genes, and  
19 (2), the novel ExInAator candidates alone. This ensures that findings are not biased by the inclusion of  
20 experimentally-verified CRL drivers amongst candidate gene sets. We first tested the frequency with which  
21 candidates are affected by copy number variants (CNVs) across matched cancers (29). We found that all  
22 candidates, and novel candidates alone, both display a trend to have elevated rates of copy number variation  
23 (Fig. 3E). We also investigated whether candidates are more proximal to germline cancer mutations (29).  
24 Once more, we observe a trend for candidates to be more likely to be proximally located to such mutations  
25 than expected by chance. Although the small numbers involved do not generally reach statistical  
26 significance, these findings are additional evidence that ExInAator predictions, either including or excluding  
27 known cancer-related lncRNAs, are involved in tumour progression.

28

## 29 **ExInAator identifies known protein-coding cancer genes**

30 Although ExInAator was designed with lncRNAs in mind, it makes no use of functional impact  
31 predictions and hence is agnostic to the protein-coding potential of the genes it analyses. We took advantage  
32 of this versatility to further test ExInAator's precision, by comparing predictions to the gold-standard  
33 catalogue of the Cancer Gene Census (CGC) (24). CGC is a manually-curated and regularly-updated  
34 annotation of genes whose somatic mutations have been associated with cancer. CGC genes represent a  
35 subset of 545 genes (File S15) (2.7%) of the entire GENCODE set of 20,314 studied here (File S16) (Table  
36 1).

37 We ran ExInAator using protein-coding gene annotations, without stratification. Altogether, a total  
38 of 24 protein-coding drivers were identified at a false discovery rate cutoff of  $Q<0.1$ . Of these, 38% are  
39 CGC genes (indicated in red, Fig. 4A). This represents enrichment of 14-fold over random expectation  
40 ( $P\leq 0.0001$ , Chi-square with Yates' correction test). The most significantly enriched gene in this analysis  
41 is *TP53*, the most frequently mutated across cancers and identified in previous exome sequencing projects

1 (16). *TP53* exons display an obvious and consistent enrichment of somatic mutations in both datasets,  
2 clustered in exons 4 and 7-11 (Fig. 4D). This *TP53* signal is observed in both Pancancer and multiple  
3 individual cancer types.

4 Several of the 15 non-CGC genes identified have evidence for cancer roles: *ANKRD18A* in lung  
5 cancer (30), *DDX3X* and *PBRM1* in various cancers (31), *HPSE2* in thyroid carcinoma (32), *MYO5B* in  
6 gastric cancer (33). These findings suggest that ExInAator precision may be better than implied by the  
7 analysis of CGC genes alone.

8 We examined the performance of ExInAator, in terms of the percent of predicted genes that belong  
9 to CGC, at a series of Q value thresholds (Fig. 4B) (File S17). Shown are separate analyses for all cancer  
10 types (expressed as mean prediction per cancer), and various pancancer combinations. These show that,  
11 although the number of predicted genes are low, they tend to have far higher rate than that 2.7% expected  
12 by random chance, even at a Q value threshold of 0.1.

13 In summary, ExInAator performs well in identifying known cancer related genes at high precision  
14 from a protein-coding training set ~10 times larger than CRL lncRNAs.

15

## 16 **ExInAator is competitive with tools designed for protein-coding genes**

17 Next we compared ExInAator to a series of well-known pipelines for identification of protein-coding  
18 drivers: MutSig (17), OncodriveFM (22) and OncodriveClust (34). In side-by-side analyses on identical  
19 Alexandrov Pancancer data, we found that ExInAator has low sensitivity (ie makes few predictions), but has  
20 excellent precision. In fact, its predictions contain a higher percentage of CGC genes than the other methods  
21 (Fig. 4C and File S18). For example, at a cutoff of  $Q < 0.1$ , ExInAator predicts 3 genes (of which 2 are known  
22 drivers), compared to 4 known drivers out of 39 for MutSig, 11 known drivers out of 104 for  
23 OncoDriveClust and 59 known drivers out of 589 for OncodriveFM (Fig. 4C). Furthermore, comparing  
24 the top 30 candidates detected at several cutoffs (File S19), the majority of genes detected by ExInAator are  
25 also detected by at least one other method.

26 We also compared the four programs' performance on real and simulated Pancancer data, displayed  
27 as Q-Q plots in Files S8 and S20. Again, ExInAator performs relatively well: its predictions on true data  
28 mirror the expected distribution quite well, and true P values are smaller than for simulated data. ExInAator  
29 predictions appear to be conservative, having a tendency for moderately deflated P values. In contrast, other  
30 methods tend to perform worse, being either strongly deflated (MutSig), inflated (OncodriveFM) or  
31 predicting less in true than randomised data (OncodriveClust). In summary, despite not employing any  
32 information from protein-coding sequence to inform its predictions, ExInAator is surprisingly competitive  
33 with existing methods in the identification of coding driver genes. In particular, its predictions have low  
34 sensitivity (possibly many false negatives) but high precision (a high fraction of true positives). This lends  
35 weight to the accuracy of ExInAator's lncRNA predictions.

36

## 37 **lncRNAs are predicted as drivers at higher rates compared to coding genes**

38 We were interested in the overall rates of prediction of lncRNAs and protein-coding genes, as well  
39 as their apparent tumour-specificity. Known driver genes are highly variable with respect to their tumour-

1 type specificity. *TP53* mutations are found across a wide range of cancers, while other drivers are only  
2 mutated in single tumour types (16,31). In this analysis, we detected no lncRNAs in more than one tumour  
3 (File S21). In contrast, two coding genes were discovered in two independent cancer types, while *TP53* was  
4 identified in no less than 9. Interestingly, a higher fraction of lncRNAs was predicted as driver genes than  
5 protein coding: 0.25% and 0.11%, respectively. These figures are likely to be strongly influenced by both  
6 the low sensitivity of ExInAator discussed above and by the sparse data. In future, many more genes are  
7 likely to be identified in multiple cancers when deeper data is available. Nevertheless these findings suggest  
8 that lncRNA are mutated in cancer at a rate similar to, or higher than protein-coding genes.

9

## 10 **Novel and known driver lncRNAs share distinctive features of functionality**

11 Returning to the driver lncRNAs identified by ExInAator, we next asked whether any features  
12 distinguish these from other lncRNAs. Previous studies of lncRNA have used features such as evolutionary  
13 conservation and expression as proxies for functionality (35,36). Furthermore, previous research on protein-  
14 coding cancer genes showed that their genes and their processed transcripts tend to be longer than average  
15 (37).

16 We compiled a series of features and, for each one, asked to what extent it differs between the CRL  
17 genes and all other lncRNAs. The full set of results, plotted by magnitude of difference and statistical  
18 significance, are shown in Figure 5A. It is clear that CRL genes are distinguished by a diverse range of  
19 features. They are transcribed from longer genes, and have longer mature transcripts (“exonic length”).  
20 They are under stronger evolutionary constraint: their promoters and exons are more evolutionarily  
21 conserved across a range of evolutionary distances. Their steady state RNA levels are higher and more  
22 variable across human tissues. Finally, they are also more likely to be proximal to a binding site of the P53  
23 tumour suppressor. In contrast, there is no difference in genic or exonic GC content between CRLs and  
24 other genes.

25 Having established a series of cancer lncRNA-specific features, we asked whether these features  
26 are also present in ExInAator candidate genes. We were particularly interested in whether novel candidates  
27 (ie non-CRL) share these characteristics, since this would represent an independent test for the value of  
28 ExInAator predictions. Therefore we compared the features of three gene groups: CRL lncRNAs, all  
29 ExInAator candidate genes, and novel ExInAator candidates alone. These groups were compared to the null  
30 set of genes, represented by the entire set of remaining Gencode lincRNAs (“All other genes”).

31 In Figure 5B are shown the results across seven selected features. ExInAator candidates, in common  
32 with CRL, have longer genes and transcripts than lincRNAs in general ( $P=4E-8$ ,  $P=6E-4$ , respectively,  
33 Wilcoxon test). Surprisingly, and in contrast to CRL genes, ExInAator candidates have significantly lower  
34 GC content ( $P=7E-3$ ), and higher repetitive sequence content ( $P=0.03$ ). Finally, for features of evolutionary  
35 conservation of both promoter and exon, in addition to steady-state RNA levels, we find that novel  
36 candidates display a similar trend as CRL genes, although these do not reach statistical significance  
37 ( $P>0.05$ ). In summary, and pending future replication with larger gene sets, it appears that novel ExInAator  
38 predicted cancer genes share a number of distinguishing features with known cancer lncRNAs, consistent  
39 with being *bona fide* driver genes.

## 1 **Materials and Methods**

### 2 **Gene annotation and filtering**

3 The GENCODE v19 lncRNA catalogue was downloaded in GTF format from  
4 ([www.genecodegenes.org](http://www.genecodegenes.org)) (4,38), and comprises 13,870 genes. A number of filtering steps were applied to  
5 this list. First, only intergenic genes (having no transcripts overlapping protein-coding genes on the opposite  
6 strand, or within 10 kb at their closest point on the same strand) were retained (6,308). Second, any lncRNA  
7 gene with transcripts of uncertain protein-coding potential were removed, leaving 5,887 genes (see below  
8 for details). Third, we included several cancer-related lncRNAs from the scientific literature, resulting in a  
9 final set of 5,914 lncRNA genes (Table 1 and Files S3 and S4). Note that literature genes may violate the  
10 two filters above, but must have a GENCODE identifier. This set of filtered lncRNAs is used throughout.

11 The protein-coding gene catalogue was also obtained in GTF format from GENCODE v19 (38).  
12 From this annotation, all genes with biotype “protein-coding” were selected, resulting in 20,314 genes and  
13 145,518 transcripts. Finally, all transcripts not having biotype “protein-coding” were removed, reducing  
14 the transcripts to 81,702. (Table 1 and File S16).

15

### 16 **Somatic mutation data curation**

17 Whole-genome cancer somatic mutations were obtained in BED format from two sources: 10  
18 cancers described in Alexandrov et al (20), and 14 cancers from TCGA (1). In addition, we included an  
19 additional dataset of 100 stomach adenocarcinoma (STAD) with the Alexandrov dataset (39), resulting in  
20 an original set of 22,877,059 mutations. Only single nucleotide somatic mutations and indels of length 1  
21 were retained (97.7% of the total somatic variants), hereafter referred to as “mutations”. AML and ALL  
22 cancers from the Alexandrov dataset were removed due to their low number of genomes and mutations.  
23 Statistics on the remaining cancers can be found in Table 2. Both mutation datasets were prefiltered in order  
24 to remove possible misannotated germline SNPs. First, any mutations identical to an entry in dbSNP 146  
25 “common” (>1% frequency) were removed, leaving 22,128,594 mutations (96.7%). Second, any recurrent  
26 mutations, having the same nucleotide change observed in the same location more than once, were collapsed  
27 and treated as a single event, resulting in a final set of 20,837,263 mutations (91.1%).

28

### 29 **Assessing the protein-coding potential of lncRNA**

30 All GENCODE v19 lncRNA transcripts were tested for protein-coding potential with CPAT (40)  
31 at default settings. Any gene having one or more transcripts predicted to be protein-coding (coding potential  
32  $\geq 0.364$ ) was removed from further analysis.

33

### 34 **ExInAator design**

35 ExInAator requires eight mandatory inputs: (1) a gene annotation in GTF format containing  
36 information on genes and exons to analyse (transcript information is ignored), (2) a catalogue of mutations  
37 in BED format, (3) the number of individual genomes or samples represented by the BED file, (4) the output  
38 folder destination, (5) a file with two columns showing the name of each chromosome and its nucleotide

1 length, (6) a gene annotation in GTF format containing information on genes and exons of the whole  
2 genome (transcript information is ignored), (7) FASTA file of the whole genome and (8) a file containing  
3 all the possible trinucleotides. Optional inputs are: (1) a minimum number of exonic and/or (2) background  
4 mutations that each gene must have to be analysed, (3) the number of CPU cores to use in the analysis and  
5 (4) the extension length of the background region that includes all introns

6 The ExInAator workflow can be divided into the following steps: exon and background definition,  
7 mutations mapping, sub-sampling of background region, gene filtering by mutation counts and statistical  
8 analysis (Fig. 1).

9 Exon / Background definition (File S1-A-B-C): The full set of exons from all transcripts belonging  
10 to a gene are merged. The remaining genic space is then defined as background, which is extended to both  
11 sides of the gene according to the window length parameter. In the present study, this value was set at 10  
12 kb throughout. Regions overlapping exons from any other gene are removed from this background region.  
13 The coordinates of non-overlapping exons and background regions are saved in BED format. The total  
14 exonic and background nucleotide length is calculated.

15 Mutations mapping (File S1-D): Mutations are mapped to exons and background regions, then  
16 counted. Despite in this study we collapsed the recurrent mutations in only one, if two or more mutations  
17 fall in the same position they are counted separately.

18 Sub-sampling of background region (File S1-E-F): The trinucleotide content of the exonic and  
19 background regions are calculated. Then, regions of identical size and trinucleotide composition to the  
20 exonic region are sampled from the background region. This is performed sequentially, without  
21 replacement, until it is impossible to continue. At every iteration, the sampled positions are added to a new  
22 background region, along with their associated mutations. In this way, a new background region of maximal  
23 size and identical composition to the exonic region is assembled for every gene.

24 Gene filtering by mutation counts: Mutation data are sporadic and of low density, potentially  
25 resulting in inflated P values. To avoid this, ExInAator accepts a user-defined minimum number of exonic  
26 and background mutations, below which lncRNAs will not be considered. These cutoffs may be defined by  
27 the user, with the default filter (used in the present study) discarding genes with less than 1 exonic mutations  
28 or 1 background mutations.

29 Statistical analysis: Statistical enrichment of exonic mutations is determined using the  
30 hypergeometric test (Fig. 1B). The following contingency table is compiled for each gene, with the total  
31 exonic and background lengths, N and n respectively:

32  $M$  = number of exonic positions mutated

33  $N - M$  = number of exonic positions not mutated

34  $m$  = number of background positions mutated

35  $n - m$  = number of background positions not mutated

36 This is the starting point for calculations of statistical significance of enrichment of exonic  
37 mutations using the hypergeometric distribution, which describes the probability of obtaining a given  
38 number of successes in a given number of draws without replacement from a finite population of a specific

1 size. It is important to note that the positions corresponding to each genome are counted independently,  
2 meaning that the total gene length  $N$  is defined as gene length multiplied by the number of genomes.  $n$  is  
3 treated similarly. Statistical significance is estimated for a gene to have that many or more exonic mutations,  
4 then are corrected for multiple hypothesis testing using the Benjamini-Hochberg procedure, which controls  
5 the False Discovery Rate (FDR), here indicated by “ $Q$ ”.

6 ExInAtor returns the input gene list with mutation counts and associated exonic enrichment  $Q$ -  
7 values. The latest ExInAtor version is freely available for download here:  
8 <https://github.com/alanzos/ExInAtor/>

9

## 10 **Creation of a simulated mutation dataset**

11 Two distinct methods were used to create trinucleotide-aware simulations of tumour mutations. In  
12 the first method (“Fixed window reassignment”), the genome was divided into fixed partitions of 50 kb.  
13 Mutations were randomly assigned to another genomic location with the same reference trinucleotide and  
14 surrounding nucleotides for substitutions and indels, respectively. In the second method (“Sliding window  
15 reassignment”), a 50 kb window is centred on each individual mutation. The mutation is then reassigned to  
16 another position with identical reference trinucleotide within its window. These simulations, while  
17 maintaining approximately the same number of single nucleotide substitutions and indels of the original  
18 Alexandrov dataset as well as the same mutation trinucleotide signature, constitute neutral datasets that are  
19 not expected to be enriched in cancer related lncRNA.

20

## 21 **Visual inspection and validation of candidates’ mutations**

22 To verify the quality of the mutation calling, we visually validated 12 single somatic mutations  
23 from 4 candidates. First, we downloaded a SAM file of the surrounding regions of each mutation (+/- 2kb)  
24 with the BAM Slicer tool from CGHUB (<https://cghub.ucsc.edu/>); then we opened those files with IGV to  
25 check the reads supporting the mutations (File S22) (41).

26

## 27 **Comparison of cancer features**

28 We obtained the list of lncRNA genes proximal to cancer-related germline SNPs from Table S5 of  
29 (29). The enrichment of indicated lncRNA genesets with respect to these genes were assessed by  
30 contingency table analysis using Fisher’s Exact test. For the analysis of CNVs, the set of regions were  
31 obtained from Table S3 of the same paper, and statistical enrichments were calculated similarly. Only data  
32 from cancers corresponding to those in the present study were considered.

33

## 34 **Comparison of lncRNA features**

35 To assess which features may distinguish cancer lncRNAs, we collected different genomic and  
36 expression data for all the genes and divided them into the four groups of interest:

- 1 1) non-CRL non-candidates (non-CRL gene list excluding ExInAtor discoveries at a  $Q \leq 0.2$ , “All other
- 2 lincRNAs”)
- 3 2) CRL genes
- 4 3) all ExInAtor candidates (discovered at a  $Q \leq 0.2$ )
- 5 4) novel ExInAtor candidates (candidate genes that do not appear in the CRL list).

6 For each feature we compared all groups to non-CRL non-candidates. Statistical tests were performed using  
7 R. Features were compiled from the following sources:

8 Gene Sequence: Gene sequence features were calculated based on Gencode v19 annotations. Exonic  
9 regions of each gene were defined as the projection of all exons from the union of its transcripts. Promoter  
10 regions of each gene were defined as a window of +/- 100 nucleotides from the reference transcription start  
11 site (TSS).

12 Conservation: PhastCons scores from vertebrate and primate species alignments and PhastCons Elements  
13 from vertebrate, mammals and primate species alignments were downloaded from UCSC Genome Browser.  
14 Two separate analyses were performed, using either base-level scores, or conserved element regions. We  
15 separately computed the average exonic base-level conservation score of each gene for primates and  
16 vertebrates PhastCons scores. We merged conserved elements annotations from primate, mammal and  
17 vertebrate species alignments and intersected these regions with promoters and exonic regions. We then  
18 computed the percent of nucleotides (from promoters or exonic regions) covered by conserved elements for  
19 each gene.

20 Repeat Elements: We downloaded the 2013 version of RepeatMasker human genomic repetitive element  
21 annotations and converted it to BED format. These annotations were intersected with exonic regions of  
22 lincRNAs. For each gene we calculated the percent of exonic nucleotides overlapping repetitive elements.

23 Tissue Expression Analysis: We extracted tissue expression values for 16 human tissues from Human Body  
24 Map (HBM) RNAseq data, downloaded from ArrayExpress under accession number E-MTAB-513. These  
25 data were used to quantify Gencode v19 genes using the GRAPE pipeline (42). Considering only genes that  
26 are expressed (RPKM>0) at least in one tissue we described the mean, the maximum and the variance of  
27 RPKM expression values across tissues. The percent of expressed genes for a given group represents the  
28 total number of genes that are expressed at least in one tissue compared to the total number of genes of the  
29 given group.

30 P53 analysis: We obtained CHIP data for p53 binding sites from (28). Binding maps from the two available  
31 timepoints were merged. We attempted to assess a possible link between cancer driver lincRNAs and p53  
32 binding site regions in two different ways. We first analysed whether the position of CRL genes in the  
33 genome tend to be closer to p53 binding site regions compared to non-CRL genes. To this aim, we  
34 calculated the nucleotide distance from the promoter of the gene (defined as explained before) to the closest  
35 p53 binding site region for all CRL and non-CRL genes. As an alternative, we compared the probability of  
36 finding a p53 binding site close to a TSS for CRL and non-CRL genes: for each we counted how many  
37 genes out of the total contain at least one predicted p53 binding site region within a window of 100kb,  
38 centred on the TSS.

## 1 Discussion

2 Here we have presented ExInAtor, to our knowledge the first method specifically designed to  
3 identify cancer driver lncRNAs from tumour genome cohorts. ExInAtor aims to address the unique  
4 opportunity of comprehensively discovering cancer driver lncRNAs within and across tumour types using  
5 mutation data generated by projects such as TCGA and ICGC.

6 We have presented the results of scans across the two most substantial tumour genome sequencing  
7 cohorts presently available, the Alexandrov and TCGA datasets, altogether comprising more than 1000  
8 genomes from 23 cancer types. In addition to successfully retrieving at nine known protein coding drivers  
9 (38% of total predictions) and six published cancer-related lncRNAs (40% of predictions), we identify for  
10 the first time a total of nine novel lncRNA driver genes at low false positive rates (0.1 FDR). These novel  
11 candidates share with known cancer lncRNAs a series of features including evolutionary conservation,  
12 normal tissue expression and gene length. They also tend to be proximal to germline cancer SNPs and have  
13 increased probability of lying in CNV regions, lending weight to their association with tumourigenesis.  
14 Together these observations lend weight to the idea that ExInAtor predicts *bona fide* driver lncRNAs. The  
15 true test of these predictions must await experimental validation in cell lines and animal models.

16 The distinguishing features of cancer-related lncRNAs are reminiscent of similar findings for  
17 protein coding genes (37). Evolutionary conservation and high steady-state RNA levels are generally  
18 interpreted in this context as evidence for functionality of lncRNAs (35,36). The significance of other  
19 features is less clear, and we should be careful to consider possible non-biological factors. In the case of  
20 gene length, it is likely that ExInAtor has greater statistical power for longer genes, possibly explaining the  
21 significantly elevated lengths of known and novel candidates. Furthermore, it is likely that the annotated  
22 length of lncRNAs is correlated with their expression, since higher expressed genes have more supporting  
23 ESTs and cDNAs, and hence are more complete.

24 Other observations were unexpected: the exons of novel candidate drivers have elevated repetitive  
25 content and reduced GC content. Furthermore, and in contrast to the above, these features are not shared  
26 with known CRL driver genes. It is unclear whether this reflects technical artefacts of the analysis, or a  
27 genuine biological insight. We can think of no bias in ExInAtor, or the cancer mutation datasets, that may  
28 favour gene models with these properties, although it is entirely feasible. On the other hand, transposable  
29 elements have been linked to both cancer (43,44) and lncRNA functionality (45). It is attractive to  
30 hypothesise that repeat-rich lncRNAs play roles in tumourigenesis and are preferentially mutated during  
31 this process. Further study will be required to establish the significance of these findings.

32 At present, our understanding of how lncRNA function is encoded in sequence motifs and  
33 structures is limited (19). Consequently, advanced approaches for scoring the functional effect of mutations,  
34 such as those used for protein sequences, are unavailable. Nevertheless, future improvements to ExInAtor  
35 may include information on RNA structures, protein binding sites, post-transcriptional processing and  
36 evolutionary conservation to weight mutations based on their likely impact on lncRNA function.  
37 Furthermore, more sensitive statistical methods employing information on mutation clustering and cancer-  
38 specific mutational signatures will likely improve predictions.

39 We expect that future studies will yield many more candidate lncRNAs than produced here.  
40 Although the datasets used represent a large proportion of all presently available tumour genomes, future  
41 projects will likely be larger and produce mutation calls of better quality. For example, the upcoming

1 PCAWG project will likely produce several fold more genomes than used here, and with more sophisticated  
2 mutation calling (1,46).

3         The increasing scale of cancer genome projects will place a growing emphasis on computational  
4 efficiency. One of the benefits of ExInAator is its ability to handle data with complex trinucleotide biases  
5 uses a simple subsampling algorithm, and without any functional impact predictions. This simplicity has  
6 the unintended benefit that ExInAator is capable of identifying protein-coding drivers with precision  
7 comparable to the best methods. Another outcome is that ExInAator makes very low computational  
8 demands: analyses for this paper were executed on a workstation running Intel Core i7 processors. 25  
9 minutes were required to analyse protein coding genes in Superpancancer (the largest dataset tested here)  
10 using a single core and 2,050 MB of RAM. It required just three minutes to analyse Pilocytic astrocytoma  
11 with six cores and 648 MB of RAM. Together, these features make ExInAator suited to future, large-scale  
12 cancer genome sequencing projects.

13

14

15

16

17

1 **Tables**

2 **Table 1: Filtered gene sets.**

3

<b>Element</b>	<b>CRL</b>	<b>LncRNA</b>		<b>CGC</b>	<b>Protein coding</b>	
		<b>Non-CRL</b>	<b>Total</b>		<b>Not CGC</b>	<b>Total</b>
Genes	45	5,869	5,914	545	19,769	20,314
Transcripts	297	9,086	9,383	3,239	78,463	81,702
Exons	1,259	27,025	28,284	35,902	702,974	738,876
Merged Exons	267	19,153	19,420	9,326	218,186	227,512

4

5

6

1 **Table 2: Cancer datasets used in this study.**

<b>Dataset</b>	<b>Cancer</b>	<b>Mutations</b>	<b>Genomes</b>
Alexandrov	Breast	655,823	119
Alexandrov	CLL	51,377	28
Alexandrov	Liver	867,080	88
Alexandrov	Lung_adeno	1,520,078	24
Alexandrov	Lymphoma_B-cell	126,581	24
Alexandrov	Medulloblastoma	123,642	100
Alexandrov	Pancreas	110,944	15
Alexandrov	Pilocytic_astrocytoma	10,436	101
Alexandrov	Stad	2,796,863	100
Alexandrov	Pancancer	6,259,996	607
TCGA	BLCA	385,128	21
TCGA	BRCA	620,238	96
TCGA	CRC	4,680,653	42
TCGA	GBM	180,896	27
TCGA	HNSC	295,709	27
TCGA	KICH	24,508	15
TCGA	KIRC	131,828	29
TCGA	LGG	35,474	18
TCGA	LUAD	1,237,722	46
TCGA	LUSC	1,626,973	45
TCGA	PRAD	21,113	20
TCGA	SKCM	3,538,750	38
TCGA	THCA	37,882	34
TCGA	UCEC	2,268,210	47
TCGA	Pancancer	14,841,279	505
Both	Superpancancer	20,837,263	1112

2

1 **Table 3: List of predicted lncRNA drivers at Q<0.1.** Known cancer genes from CRL are marked in  
 2 bold face. “Pc” - PanCancer.

3

Cancer	Gene Name	Gene ID	Ex mut	Ex len	Intr mut	Intr len	Pval	Qval	Ex mut rate	Intr mut rate	Ratio
Breast	AP000469.2	ENSG00000224832	9	238824	2	1176908	4.38E-06	5.05E-03	3.77E-02	1.70E-03	22.2
Stad	<b>PCA3</b>	ENSG00000225937	10	392190	20	2702680	2.86E-03	5.49E-02	2.55E-02	7.40E-03	3.4
KICH	RP11-308N19.1	ENSG00000234323	1	26384	1	2437289	2.13E-02	6.39E-02	3.79E-02	4.10E-04	92.4
Stad	<b>SAMMSON</b>	ENSG00000240405	5	208795	1	697899	3.14E-03	5.49E-02	2.39E-02	1.43E-03	16.7
GBM	<b>lncRNA-ATB</b>	ENSG00000244306	4	314168	4	1347296	4.68E-02	9.36E-02	1.27E-02	2.97E-03	4.3
Super_Pc	<b>NEAT1</b>	ENSG00000245532	163	25316741	16	5559984	4.98E-04	2.19E-02	6.44E-03	2.88E-03	2.2
Pc_TCGA	<b>NEAT1</b>	ENSG00000245532	96	11497239	7	2524993	9.28E-04	4.08E-02	8.35E-03	2.77E-03	3.0
PRAD	RP11-455B3.1	ENSG00000248202	7	16573	1	380379	1.70E-09	1.19E-08	4.22E-01	2.63E-03	160.7
KICH	RP11-332J15.1	ENSG00000249734	3	8277	15	314985	1.03E-02	6.39E-02	3.62E-01	4.76E-02	7.6
Breast	RP11-707A18.1	ENSG00000250125	6	485395	9	6908893	2.38E-04	9.16E-02	1.24E-02	1.30E-03	9.5
KICH	RP11-6C14.1	ENSG00000250488	1	5369	1	687929	1.54E-02	6.39E-02	1.86E-01	1.45E-03	128.1
Super_Pc	<b>MALAT1</b>	ENSG00000251562	83	9683213	40	7768392	4.45E-03	9.79E-02	8.57E-03	5.15E-03	1.7
Breast	RP11-1101K5.1	ENSG00000253434	6	165285	19	4878981	1.28E-04	7.38E-02	3.63E-02	3.89E-03	9.3
HNSC	RP11-354A14.1	ENSG00000254689	3	12579	4	697784	1.84E-04	8.14E-02	2.38E-01	5.73E-03	41.6
BRCA	RP11-189E14.4	ENSG00000261623	9	352599	1	1252991	9.53E-06	9.25E-03	2.55E-02	7.98E-04	32.0
SKCM	<b>BCAR4</b>	ENSG00000262117	6	50458	8	533056	6.81E-04	2.11E-02	1.19E-01	1.50E-02	7.9

4

5

1 **Acknowledgements**

2 The authors wish to thank Núria López-Bigas for her support throughout the project. We also thank  
3 Marta Melé (Harvard University) for insightful discussions, and Maite Huarte (CIMA) for insightful  
4 discussions and providing P53 binding maps, and Darek Kedra (CNAG) for the help with the visual  
5 inspection of the somatic mutations. We especially thank Erik Larsson (University of Gothenburg) for  
6 sharing TCGA mutation calls ahead of publication. We thank Romina Garrido (CRG) for administrative  
7 support. We acknowledge support of the Spanish Ministry of Economy and Competitiveness, ‘Centro de  
8 Excelencia Severo Ochoa 2013-2017’, SEV-2012-0208. R.J. is supported by Ramón y Cajal RYC-2011-  
9 08851 and Plan Nacional BIO2011-27220. A.L. is supported by pre-doctoral fellowship FPU14/03371.

10

11 **Author contributions**

12 R.J. conceived the project, and supervised with advice and suggestions of R.G.. Primary  
13 development of the tool was carried out by A.L., as well as the creation of one of the simulated datasets.  
14 Statistical and technical assistance for running analysis were provided by F.R and E.P., respectively. A.L.  
15 and J.C. performed the secondary analysis. L.M. created one of the simulated datasets, executed the analysis  
16 of MutSig, OncodriveFM and OncodriveClust on protein coding genes. R.J., A.L. and J.C. drafted the  
17 manuscript and prepared the figures and supplementary material. All authors read and approved the final  
18 draft.

19

20 **Competing interests**

21 The authors declare that they have no competing interests.

22

23

24

25

26

27

## 1 **Figure Legends**

2 **Fig. 1: Outline of the ExInAator method.** (A) The steps of gene definition, subsampling and  
3 analysis performed to quantify exonic and background mutations. Sampling is performed in such a way  
4 that, at the end, the trinucleotide frequency of the background region is identical to the exonic region. (B)  
5 The number of mutations in background and exonic regions is compared by a contingency table analysis.

6 **Fig. 2: The landscape of driver lncRNAs across 23 tumour types.** (A) The numbers and  
7 proportion of literature-reported cancer-related long noncoding RNAs (CRL) and non-CRL candidates  
8 identified in this study. Q value is equivalent to false discovery rate (FDR). (B) A Quantile-Quantile (QQ)  
9 plot showing the performance of ExInAator on Breast cancer mutations. Each point represents one gene.  
10 Note the deviation of real P values from the theoretical distribution in a small tail of cases. Simulated data  
11 was created by randomising mutations while maintaining trinucleotide context. (C) As for B, for Liver  
12 cancer. Note the lack of candidates in this dataset. (D) The number of driver genes discovered at a Q cutoff  
13 of 0.1 across the Alexandrov and TCGA collections. Cancer Gene Census (CGC) are true positive, known  
14 protein-coding cancer driver genes.

15 **Fig. 3: LncRNA cancer driver genes predicted by ExInAator across cancer genomes.** (A) All  
16 driver lncRNAs ( $Q \leq 0.1$ ) and the tumour type in which they are identified. Gene names in blue indicate  
17 those belonging to CRL. (B) A mutation density plot for NEAT1 in all cancers, plotting the SNVs per  
18 kilobase as a function of gene regions. Grey represent background regions, while colours represent the  
19 mutational contribution of each cancer type to the single exon. The x-axis represents position, in bp, with  
20 respect to the start of the background region, defined here to be at 10 kb upstream of the gene's annotated  
21 TSS. (C) The Breast mutation profile of RP11-1101K5.1, a gene with mutations in four exons. Rectangles  
22 depict mutational density of exons (blue) and introns (grey). The gene structure is indicated below, where  
23 wider portions represent exons, separated by narrower introns. (D) The Breast mutation frequency in  
24 *BCAR4*. (E) Percentage of genes and candidates in CNV regions and proximal to cancer-related germline  
25 SNPs. Numbers above bars indicate the absolute numbers of genes represented by each percentage.  
26 Statistical significance in each case was estimated using Fisher's Exact test. (F) An example of an ExInAator-  
27 predicted novel candidate gene, *RP11-820L6.1*. Note the presence of promoter-like histone marks (red,  
28 ChromHMM track), evolutionary conservation (PhastCons Primate conservation), and cancer SNVs around  
29 the gene TSS, as well as a proximal P53 binding site ("P53\_merged").

30 **Fig. 4: ExInAator discovers known protein-coding drivers at high precision.** (A) The  $-\log_{10} Q$ -  
31 values of all candidates at  $Q \leq 0.1$  cutoff in the Alexandrov and TCGA datasets. Gene names in red indicate  
32 known drivers belonging to the Cancer Gene Census (CGC). (B) The precision of ExInAator predictions  
33 was estimated as the percent of predicted driver genes that also belong to CGC. Bars are coloured by the  
34 Q-value cutoff used, and the fraction of all known genes belonging to CGC is shown in blue as a reference.  
35 "Mean" displays the average overlap across all individual cancer types. The numbers above each bar  
36 indicate the total number of predicted driver genes at that cutoff. For example, in "Superpancer", a total of  
37 three candidates are identified at a cutoff of 0.1, of which two (66%) belong to CGC. (C) Comparison of  
38 the performance of ExInAator to other methods for protein-coding driver gene discovery, using the  
39 Alexandrov Pancancer dataset. Plot description as for Panel B. (D) The mutational profile of the TP53  
40 tumour suppressor gene across all cancers.  $x$  axis indicates the position within the gene,  $y$  axis shows the  
41 mutation frequency.

1           **Fig. 5: Features distinguishing cancer lncRNAs are also found in ExInAator candidates.** A)  
2 Identification of cancer lncRNA features by comparing literature-curated Cancer Related LncRNA (CRL)  
3 genes to others. Dots represent 16 features that were compared in CRL and non-CRL genes. *y* axis shows  
4 the log<sub>2</sub> fold difference of CRL vs non-CRL means for the values of the given feature. *x* axis represents the  
5 P value obtained from the statistical test applied when comparing CRL and non-CRL. Features are coloured  
6 depending on whether they are discrete features, analysed by Fisher Exact test, or continuous features  
7 analysed by Wilcoxon test. B) Cancer lncRNA features in ExInAator-predicted driver genes. Shown are  
8 cumulative distributions for seven selected features. Dashed vertical lines indicate the mean value of each  
9 group. Genes are grouped by: literature-reported cancer CRL lncRNAs, all ExInAator candidates (both CRL  
10 and not) (“All candidates”), only novel ExInAator candidates that are not included in CRL (“novel  
11 candidates”), and non-CRL non-candidates, being all other GENCODE lincRNAs. Candidates here were  
12 defined at a  $Q \leq 0.2$ . Groups significantly different from the latter at a threshold of  $P=0.05$  (Wilcoxon test)  
13 are represented by a thick line.  
14  
15  
16  
17

## References

1. Weinstein, J.N. *et al.* (2013) The Cancer Genome Atlas Pan-Cancer analysis project. *Nature genetics*, **45**, 1113-1120.
2. Hudson, T.J. *et al.* (2010) International network of cancer genome projects. *Nature*, **464**, 993-998.
3. Gutschner, T. and Diederichs, S. (2012) The hallmarks of cancer: a long non-coding RNA point of view. *RNA biology*, **9**, 703-719.
4. 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**, 1775-1789.
5. Ulitsky, I. and Bartel, D.P. (2013) lincRNAs: genomics, evolution, and mechanisms. *Cell*, **154**, 26-46.
6. Liu, S.J. *et al.* (2016) Single-cell analysis of long non-coding RNAs in the developing human neocortex. *Genome biology*, **17**, 67.
7. Hangauer, M.J., Vaughn, I.W. and McManus, M.T. (2013) Pervasive Transcription of the Human Genome Produces Thousands of Previously Unidentified Long Intergenic Noncoding RNAs. *PLoS genetics*, **9**, e1003569.
8. Iyer, M.K. *et al.* (2015) The landscape of long noncoding RNAs in the human transcriptome. *Nature genetics*.
9. Chen, G. *et al.* (2013) LncRNADisease: a database for long-non-coding RNA-associated diseases. *Nucleic acids research*, **41**, D983-986.
10. Huarte, M. *et al.* (2010) A large intergenic noncoding RNA induced by p53 mediates global gene repression in the p53 response. *Cell*, **142**, 409-419.
11. Gupta, R.A. *et al.* (2010) Long non-coding RNA HOTAIR reprograms chromatin state to promote cancer metastasis. *Nature*, **464**, 1071-1076.
12. Du, Z. *et al.* (2013) Integrative genomic analyses reveal clinically relevant long noncoding RNAs in human cancer. *Nature structural & molecular biology*, **20**, 908-913.
13. Zheng, J. *et al.* (2016) Pancreatic cancer risk variant in LINC00673 creates a miR-1231 binding site and interferes with PTPN11 degradation. *Nature genetics*, **48**, 747-757.
14. Akrami, R. *et al.* (2013) Comprehensive analysis of long non-coding RNAs in ovarian cancer reveals global patterns and targeted DNA amplification. *PloS one*, **8**, e80306.
15. Leucci, E. *et al.* (2016) Melanoma addiction to the long non-coding RNA SAMMSON. *Nature*, **531**, 518-522.
16. Tamborero, D. *et al.* (2013) Comprehensive identification of mutational cancer driver genes across 12 tumor types. *Scientific reports*, **3**, 2650.
17. Lawrence, M.S. *et al.* (2013) Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature*, **499**, 214-218.
18. The ENCODE Project Consortium. (2012) An integrated encyclopedia of DNA elements in the human genome. *Nature*, **489**, 57-74.
19. Guttman, M. and Rinn, J.L. (2012) Modular regulatory principles of large non-coding RNAs. *Nature*, **482**, 339-346.
20. Alexandrov, L.B. *et al.* (2013) Signatures of mutational processes in human cancer. *Nature*, **500**, 415-421.
21. Haerty, W. and Ponting, C.P. (2015) Unexpected selection to retain high GC content and splicing enhancers within exons of multiexonic lincRNA loci. *RNA*, **21**, 333-346.
22. Gonzalez-Perez, A. and Lopez-Bigas, N. (2012) Functional impact bias reveals cancer drivers. *Nucleic acids research*, **40**, e169.
23. Hodis, E. *et al.* (2012) A landscape of driver mutations in melanoma. *Cell*, **150**, 251-263.
24. Futreal, P.A. *et al.* (2004) A census of human cancer genes. *Nature reviews. Cancer*, **4**, 177-183.
25. Xu, C., Ciampi, A. and Greenwood, C.M. (2014) Exploring the potential benefits of stratified false discovery rates for region-based testing of association with rare genetic variation. *Frontiers in genetics*, **5**, 11.
26. Sun, L., Craiu, R.V., Paterson, A.D. and Bull, S.B. (2006) Stratified false discovery control for large-scale hypothesis testing with application to genome-wide association studies. *Genetic epidemiology*, **30**, 519-530.

- 1 27. Fujimoto, A. *et al.* (2016) Whole-genome mutational landscape and characterization of noncoding  
2 and structural mutations in liver cancer. *Nature genetics*, **48**, 500-509.
- 3 28. Sanchez, Y. *et al.* (2014) Genome-wide analysis of the human p53 transcriptional network unveils  
4 a lncRNA tumour suppressor signature. *Nature communications*, **5**, 5812.
- 5 29. Yan, X. *et al.* (2015) Comprehensive Genomic Characterization of Long Non-coding RNAs across  
6 Human Cancers. *Cancer cell*, **28**, 529-540.
- 7 30. Liu, W.B. *et al.* (2012) ANKRD18A as a novel epigenetic regulation gene in lung cancer.  
8 *Biochemical and biophysical research communications*, **429**, 180-185.
- 9 31. Gonzalez-Perez, A. *et al.* (2013) IntOGen-mutations identifies cancer drivers across tumor types.  
10 *Nature methods*, **10**, 1081-1082.
- 11 32. Matos, L.L. *et al.* (2015) The Profile of Heparanase Expression Distinguishes Differentiated Thyroid  
12 Carcinoma from Benign Neoplasms. *PloS one*, **10**, e0141139.
- 13 33. Dong, W. *et al.* (2012) Inactivation of MYO5B promotes invasion and motility in gastric cancer cells.  
14 *Digestive diseases and sciences*, **57**, 1247-1252.
- 15 34. Tamborero, D., Gonzalez-Perez, A. and Lopez-Bigas, N. (2013) OncodriveCLUST: exploiting the  
16 positional clustering of somatic mutations to identify cancer genes. *Bioinformatics*, **29**, 2238-2244.
- 17 35. Marques, A.C. *et al.* (2013) Chromatin signatures at transcriptional start sites separate two equally  
18 populated yet distinct classes of intergenic long noncoding RNAs. *Genome biology*, **14**, R131.
- 19 36. Hezroni, H. *et al.* (2015) Principles of long noncoding RNA evolution derived from direct comparison  
20 of transcriptomes in 17 species. *Cell reports*, **11**, 1110-1122.
- 21 37. Furney, S.J., Higgins, D.G., Ouzounis, C.A. and Lopez-Bigas, N. (2006) Structural and functional  
22 properties of genes involved in human cancer. *BMC genomics*, **7**, 3.
- 23 38. Harrow, J. *et al.* (2012) GENCODE: the reference human genome annotation for The ENCODE  
24 Project. *Genome research*, **22**, 1760-1774.
- 25 39. Wang, K. *et al.* (2014) Whole-genome sequencing and comprehensive molecular profiling identify  
26 new driver mutations in gastric cancer. *Nature genetics*, **46**, 573-582.
- 27 40. Wang, L. *et al.* (2013) CPAT: Coding-Potential Assessment Tool using an alignment-free logistic  
28 regression model. *Nucleic acids research*, **41**, e74.
- 29 41. Robinson, J.T. *et al.* (2011) Integrative genomics viewer. *Nature biotechnology*, **29**, 24-26.
- 30 42. Knowles, D.G., Roder, M., Merkel, A. and Guigo, R. (2013) Grape RNA-Seq analysis pipeline  
31 environment. *Bioinformatics*, **29**, 614-621.
- 32 43. Criscione, S.W., Zhang, Y., Thompson, W., Sedivy, J.M. and Neretti, N. (2014) Transcriptional  
33 landscape of repetitive elements in normal and cancer human cells. *BMC genomics*, **15**, 583.
- 34 44. Ferreira, P.G. *et al.* (2014) Transcriptome characterization by RNA sequencing identifies a major  
35 molecular and clinical subdivision in chronic lymphocytic leukemia. *Genome research*, **24**, 212-  
36 226.
- 37 45. Johnson, R. and Guigo, R. (2014) The RIDL hypothesis: transposable elements as functional  
38 domains of long noncoding RNAs. *RNA*, **20**, 959-976.
- 39 46. Alioto, T.S. *et al.* (2015) A comprehensive assessment of somatic mutation detection in cancer  
40 using whole-genome sequencing. *Nature communications*, **6**, 10001.

41

42

# Figure 1

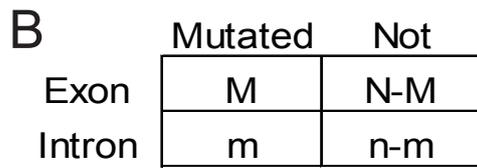
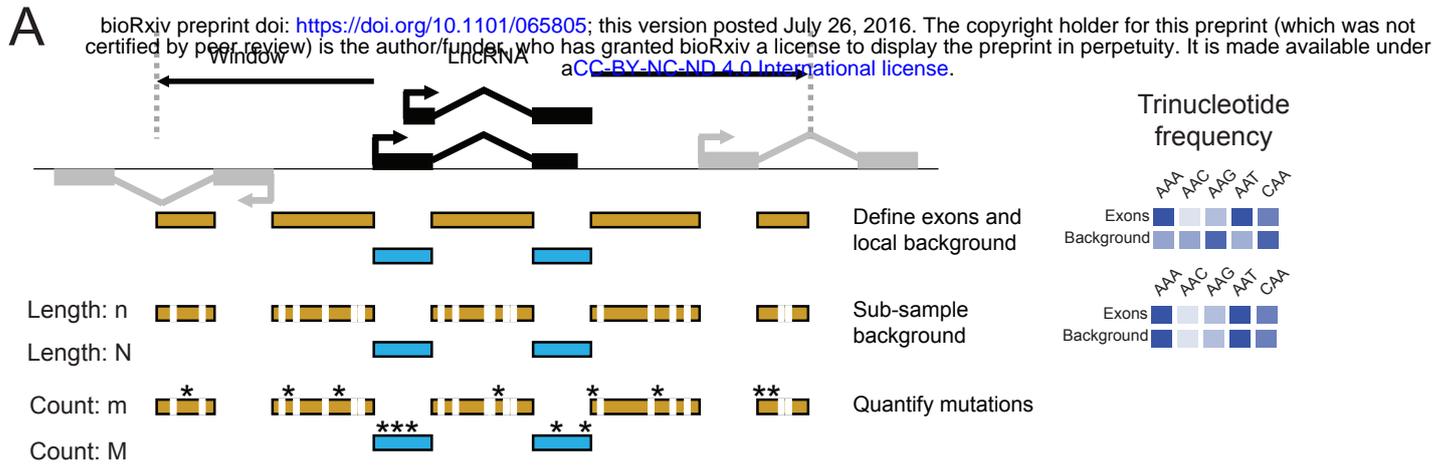
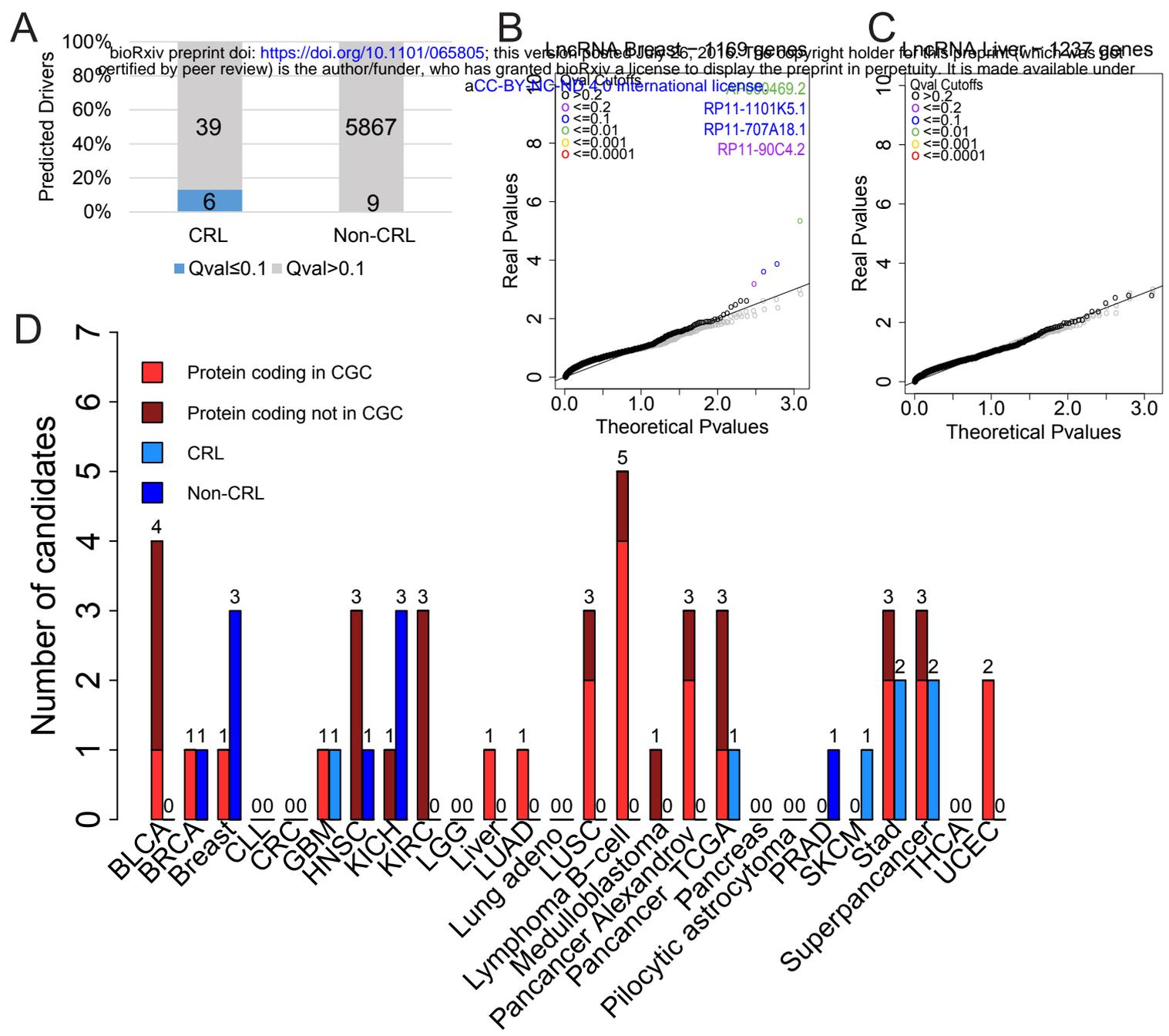
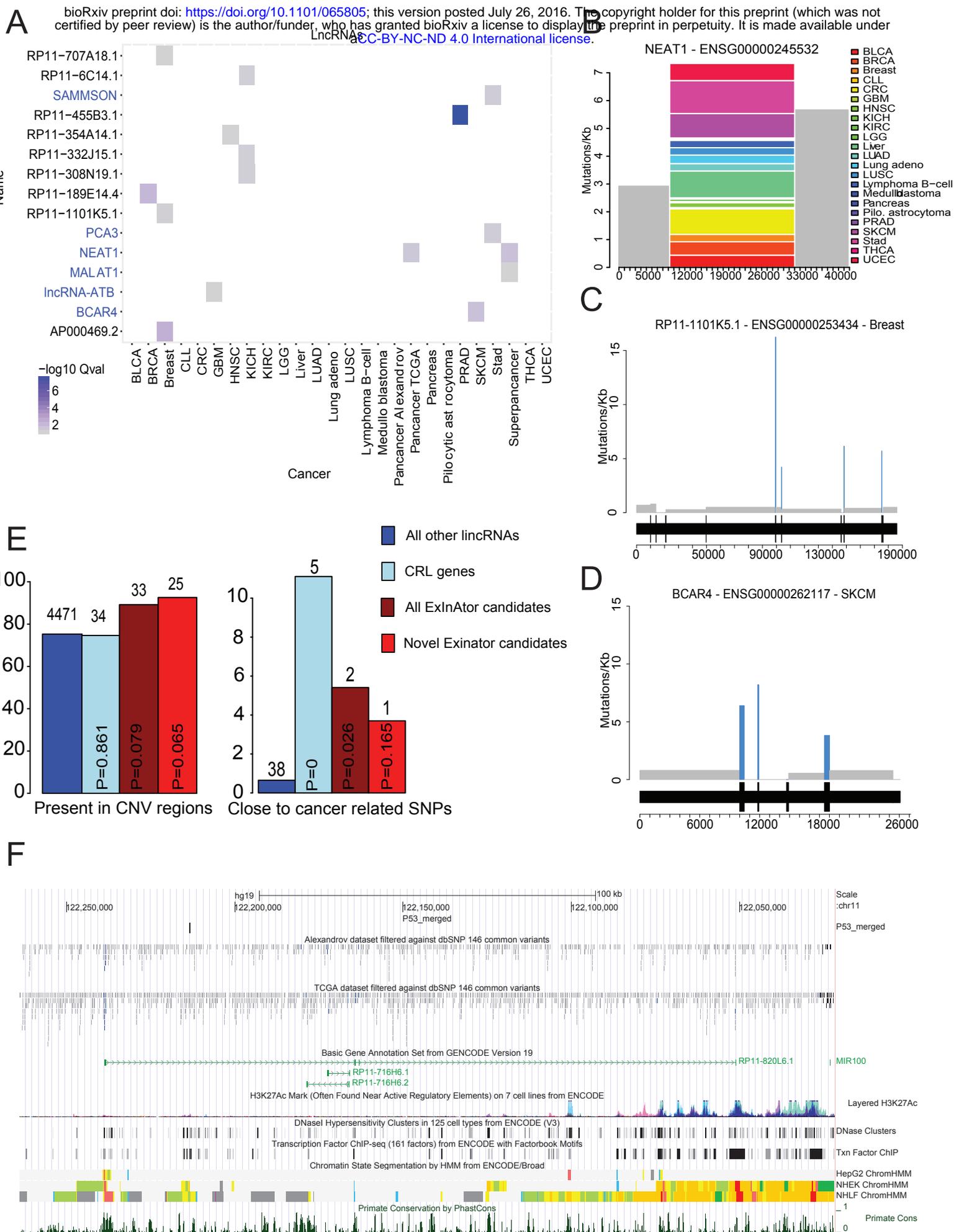


Figure 2

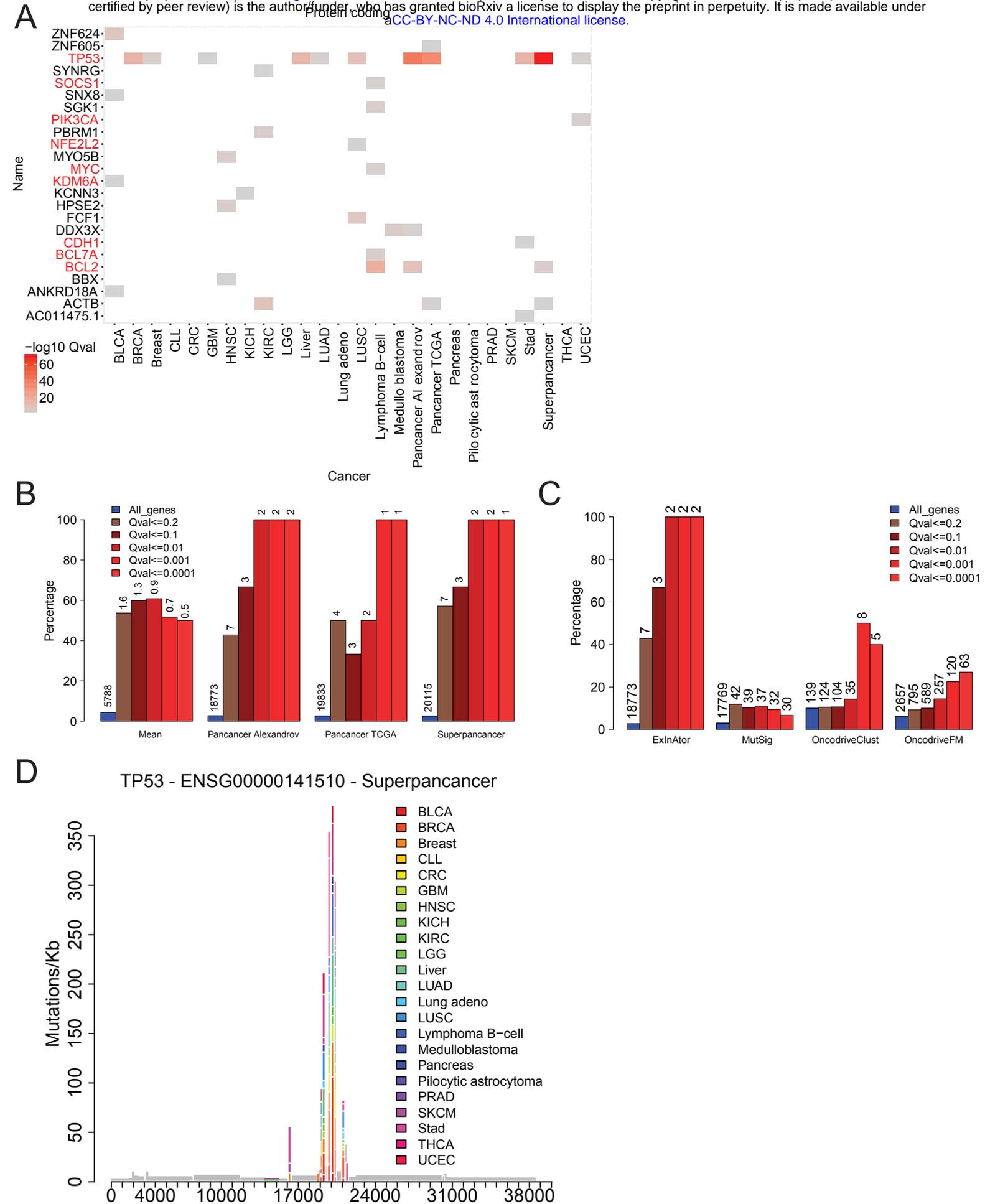


# Figure 3



# Figure 4

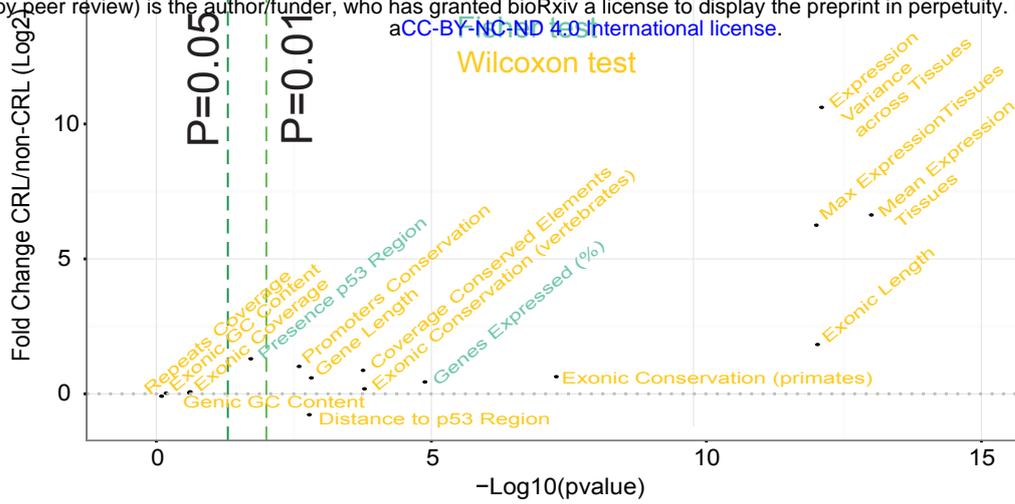
bioRxiv preprint doi: <https://doi.org/10.1101/065805>; this version posted July 26, 2016. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



# Figure 5

**A**

bioRxiv preprint doi: <https://doi.org/10.1101/065805>; this version posted July 26, 2016. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



**B**

