

RESEARCH

Open Access



# Comparative assessment of genes driving cancer and somatic evolution in non-cancer tissues: an update of the Network of Cancer Genes (NCG) resource

Lisa Dressler<sup>1,2,3†</sup>, Michele Bortolomeazzi<sup>1,2†</sup>, Mohamed Reda Keddar<sup>1,2†</sup>, Hrvoje Misetić<sup>1,2†</sup>, Giulia Sartini<sup>1,2†</sup>, Amelia Acha-Sagredo<sup>1,2†</sup>, Lucia Montorsi<sup>1,2†</sup>, Neshika Wijewardhane<sup>1,2</sup>, Dimitra Repana<sup>1,2</sup>, Joel Nulsen<sup>1,2</sup>, Jacki Goldman<sup>4</sup>, Marc Pollitt<sup>4</sup>, Patrick Davis<sup>4</sup>, Amy Strange<sup>4</sup>, Karen Ambrose<sup>4</sup> and Francesca D. Ciccarelli<sup>1,2\*</sup> 

\* Correspondence: [francesca.ciccarelli@crick.ac.uk](mailto:francesca.ciccarelli@crick.ac.uk)

<sup>†</sup>Lisa Dressler, Michele Bortolomeazzi, Mohamed Reda Keddar, Hrvoje Misetić, Giulia Sartini, Amelia Acha-Sagredo and Lucia Montorsi contributed equally to this work.

<sup>1</sup>Cancer Systems Biology Laboratory, The Francis Crick Institute, London NW1 1AT, UK

<sup>2</sup>School of Cancer and Pharmaceutical Sciences, King's College London, London SE11UL, UK

Full list of author information is available at the end of the article

## Abstract

**Background:** Genetic alterations of somatic cells can drive non-malignant clone formation and promote cancer initiation. However, the link between these processes remains unclear and hampers our understanding of tissue homeostasis and cancer development.

**Results:** Here, we collect a literature-based repertoire of 3355 well-known or predicted drivers of cancer and non-cancer somatic evolution in 122 cancer types and 12 non-cancer tissues. Mapping the alterations of these genes in 7953 pan-cancer samples reveals that, despite the large size, the known compendium of drivers is still incomplete and biased towards frequently occurring coding mutations. High overlap exists between drivers of cancer and non-cancer somatic evolution, although significant differences emerge in their recurrence. We confirm and expand the unique properties of drivers and identify a core of evolutionarily conserved and essential genes whose germline variation is strongly counter-selected. Somatic alteration in even one of these genes is sufficient to drive clonal expansion but not malignant transformation.

**Conclusions:** Our study offers a comprehensive overview of our current understanding of the genetic events initiating clone expansion and cancer revealing significant gaps and biases that still need to be addressed. The compendium of cancer and non-cancer somatic drivers, their literature support, and properties are accessible in the Network of Cancer Genes and Healthy Drivers resource at <http://www.network-cancer-genes.org/>.

**Keywords:** Driver genes, Somatic evolution, Cancer initiation, Systems-level properties



## Background

Genetic alterations conferring selective advantages to cancer cells are the main drivers of cancer evolution and hunting for them has been at the core of international cancer genomic efforts [1–3]. Given the instability of the cancer genome, distinguishing driver alterations from the rest relies on analytical approaches that identify genes altered more frequently than expected or quantify the positive selection acting on them [4–6]. The results of these analyses have greatly expanded our understanding of the mechanisms driving cancer evolution, revealing high heterogeneity across and within cancers [7–9].

Recently, deep sequencing screens of non-cancer tissues have started to map positively selected genetic mutations in somatic cells that drive in situ formation of phenotypically normal clones [10, 11]. Many of these mutations hit cancer drivers, sometimes at a frequency higher than the corresponding cancer [12–16]. Yet, they do not drive malignant transformation. This conundrum poses fundamental questions on how genetic drivers of normal somatic evolution are related to and differ from those of cancer evolution. Addressing these questions will clarify the genetic relationship between tissue homeostasis and cancer initiation, with profound implications for cancer early detection.

To assess the extent of the current knowledge on cancer and non-cancer drivers, we undertook a systematic review of the literature and assembled a comprehensive repertoire of genes whose somatic alterations have been reported to drive cancer or non-cancer evolution. This allowed us to compare the current driver repertoire across and within cancer and non-cancer tissues and map their alterations in the large pancancer collection of samples from The Cancer Genome Atlas (TCGA). This revealed significant gaps and biases in our current knowledge of the driver landscape. We also computed an array of systems-level properties across driver groups, confirming the unique evolutionary path of driver genes and their central role in the cell.

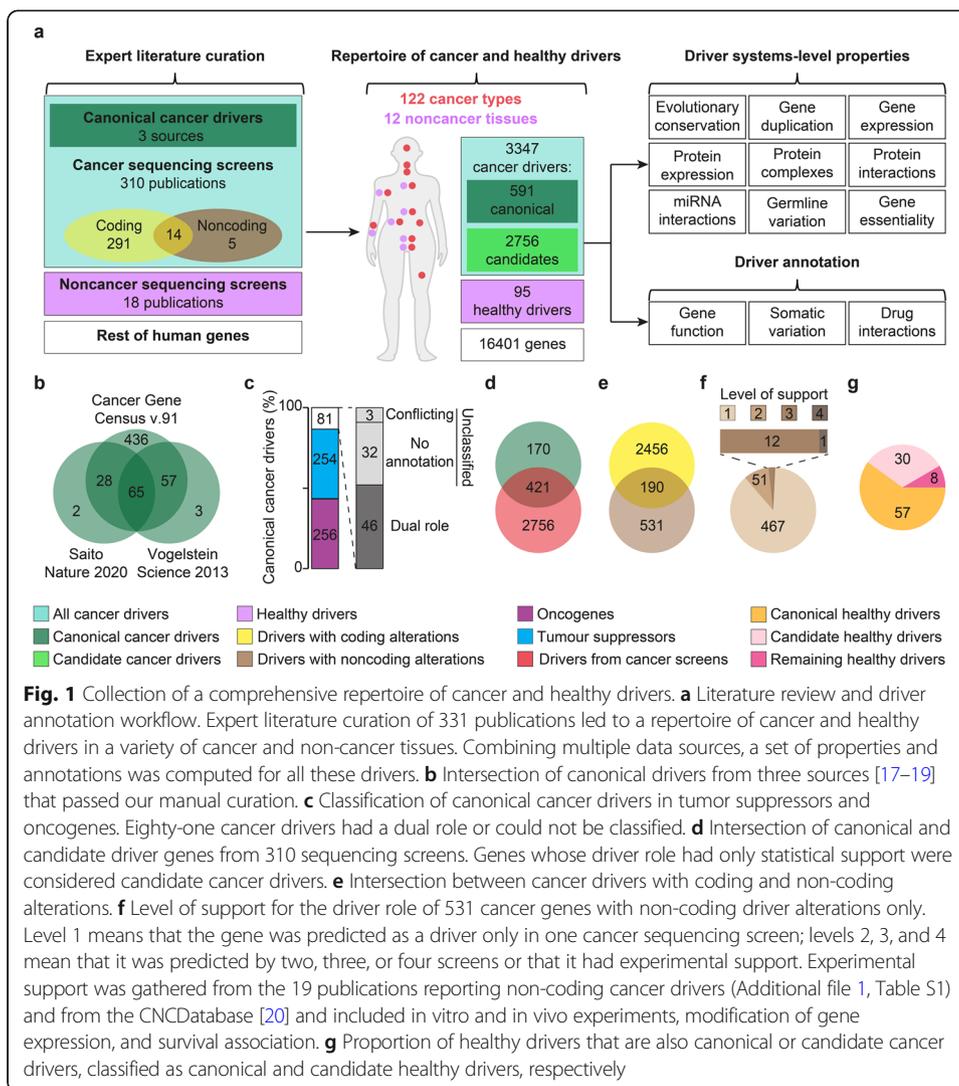
We collected all cancer and non-cancer driver genes, together with a large set of their properties, in the Network of Cancer Genes and Healthy Drivers (NCG<sup>HD</sup>) open-access resource.

## Results

### More than 3300 genes are canonical or candidate drivers of cancer and non-cancer somatic evolution

We conducted a census of currently known drivers through a comprehensive literature review of 331 scientific articles published between 2008 and 2020 describing somatically altered genes with a proven or predicted role in cancer or non-cancer somatic evolution (Fig. 1a). These publications included three sources of experimentally validated (canonical) cancer drivers, 311 sequencing screens of cancer (293) and non-cancer (18) tissues, and 17 pancancer studies (Additional file 1, Table S1). Each paper was assessed by at least two independent experts (Additional file 2, Fig. S1A–C) returning a total of 3355 drivers, 3347 in 122 cancer types and 95 in 12 non-cancer tissues, respectively (Fig. 1a). We further computed the systems-level properties of drivers and annotated their function, somatic variation, and drug interactions (Fig. 1a).

We reviewed the three sources of canonical cancer drivers [17–19] to exclude false positives (Additional file 3, Table S2) and fusion genes whose properties could not be



**Fig. 1** Collection of a comprehensive repertoire of cancer and healthy drivers. **a** Literature review and driver annotation workflow. Expert literature curation of 331 publications led to a repertoire of cancer and healthy drivers in a variety of cancer and non-cancer tissues. Combining multiple data sources, a set of properties and annotations was computed for all these drivers. **b** Intersection of canonical drivers from three sources [17–19] that passed our manual curation. **c** Classification of canonical cancer drivers in tumor suppressors and oncogenes. Eighty-one cancer drivers had a dual role or could not be classified. **d** Intersection of canonical and candidate driver genes from 310 sequencing screens. Genes whose driver role had only statistical support were considered candidate cancer drivers. **e** Intersection between cancer drivers with coding and non-coding alterations. **f** Level of support for the driver role of 531 cancer genes with non-coding driver alterations only. Level 1 means that the gene was predicted as a driver only in one cancer sequencing screen; levels 2, 3, and 4 mean that it was predicted by two, three, or four screens or that it had experimental support. Experimental support was gathered from the 19 publications reporting non-coding cancer drivers (Additional file 1, Table S1) and from the CNCCDatabase [20] and included in vitro and in vivo experiments, modification of gene expression, and survival association. **g** Proportion of healthy drivers that are also canonical or candidate cancer drivers, classified as canonical and candidate healthy drivers, respectively

mapped. Only 11% of the resulting 591 canonical drivers (Additional file 4, Table S3) were common to all three sources (Fig. 1b), indicating poor consensus even in well-known cancer genes. We further annotated the genetic mode of action for > 86% of canonical drivers, finding comparable proportions of oncogenes or tumor suppressors (Fig. 1c). The rest had a dual role or could not be univocally classified.

We extracted additional cancer drivers from the curation of 310 sequencing screens that applied a variety of statistical approaches (Additional file 2, Fig. S1 D) to identify cancer drivers among all altered genes. After removing possible false positives (Additional file 3, Table S2), the final list included 3177 cancer drivers, 2756 of which relied only on statistical support (candidate cancer drivers) and 421 were canonical drivers (Fig. 1d, Additional file 4, Table S3). Therefore, 170 canonical drivers have never been detected by any method, suggesting that they may elicit their role through non-mutational mechanisms or may fall below the detection limits of current approaches. Given the prevalence of cancer coding screens (Fig. 1a), only coding driver alterations have been reported for most genes (Fig. 1e) while 16% of them (531) were identified as drivers uniquely in non-coding screens. Since the prediction of drivers with non-coding

alterations remains challenging, we further investigated the type of support that these genes had for their driver activity. The overwhelming majority of them (467 genes, 87%) have been predicted as drivers in only one screen. The remaining 64 genes are canonical drivers, have been predicted as drivers in multiple screens, or have additional experimental support for their driver activity (Fig. 1f).

Applying a similar approach (Additional file 2, Fig. S1 A-C), we reviewed 18 sequencing screens of healthy or diseased (non-cancer) tissues. They collectively reported 95 genes whose somatic alterations could drive non-malignant clone formation (healthy drivers). Interestingly, only eight of them were not cancer drivers (Fig. 1g, Additional file 4, Table S3), suggesting a high overlap between genetic drivers of cancer and non-cancer evolution. However, since many non-cancer screens only re-sequenced cancer genes or applied methods developed for cancer genomics (Additional file 2, Fig. S1E), this overlap may be overestimated.

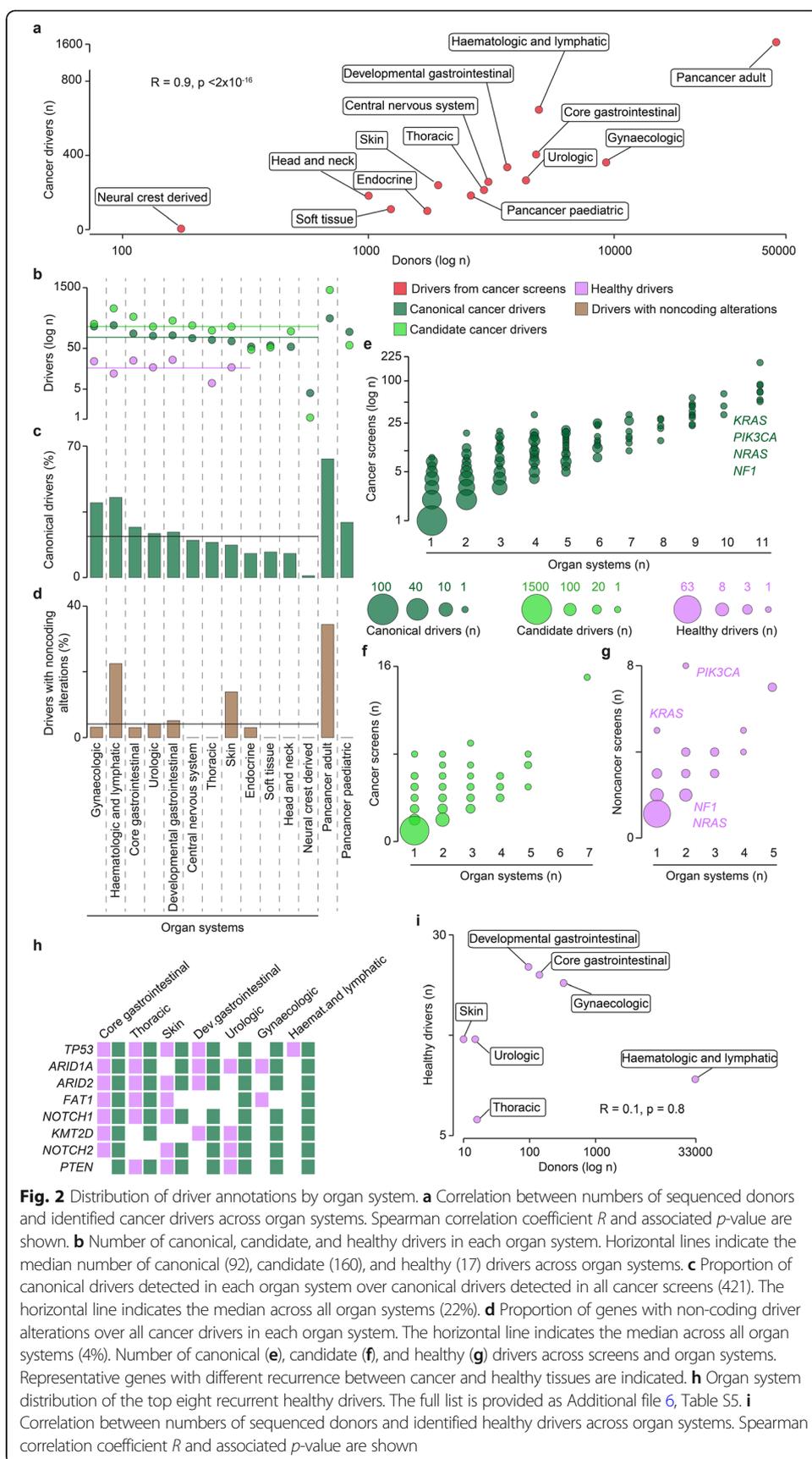
#### **The ability to capture cancer but not healthy driver heterogeneity increases with the donor sample size**

To compare cancer and healthy drivers across and within tissues, we grouped the 122 cancer types and 12 non-cancer tissues into 12 and seven organ systems, respectively (the “[Methods](#)” section).

Despite the high numbers of sequenced samples (Additional file 5, Table S4) and detected drivers (Fig. 1), several lines of evidence indicated that our knowledge of cancer drivers is still incomplete. First, we detected a strong positive correlation between cancer drivers and donors overall (Fig. 2a) and in individual organ systems (Additional file 2, Fig. S2). This suggests that the current ability to identify new drivers depends on the number of samples included in the analysis. Second, candidates outnumbered canonical drivers in all organ systems except those with a small sample size or low mutation rate such as pediatric cancers, where only the most recurrent canonical drivers could be identified (Fig. 2b). Third, large donor cohorts enabled the detection of a broader representation of canonical drivers than small cohorts (Fig. 2c). For example, pooling thousands of samples together led to >60% of canonical drivers being detected in adult pancancer re-analyses. Therefore, the size of the cohort influences the level of completeness and heterogeneity of the cancer driver repertoire. This is not surprising since all current approaches act at the cohort level, searching for positively selected genes altered more frequently than expected (Additional file 2, Fig. S1D).

Our analysis also showed that the contribution of non-coding driver alterations remains largely unappreciated and non-coding drivers have not yet been reported in several tumors, including all pediatric cancers (Fig. 2d). Owing to the re-analysis of large whole-genome collections [21–26], almost 40% of adult pancancer drivers were instead modified by non-coding alterations (Fig. 2d). Hematologic and skin tumors also had a high proportion of non-coding driver variants thanks to screens focused on non-coding mutations [27, 28]. Therefore, the re-analysis of already available whole-genome data and further sequencing screens of non-coding variants are needed to fully appreciate their driver contribution.

Compared to cancer, sequencing screens of non-cancer tissues are still in their infancy, as reflected by the lower numbers of screened tissues and detected drivers



**Fig. 2** Distribution of driver annotations by organ system. **a** Correlation between numbers of sequenced donors and identified cancer drivers across organ systems. Spearman correlation coefficient  $R$  and associated  $p$ -value are shown. **b** Number of canonical, candidate, and healthy drivers in each organ system. Horizontal lines indicate the median number of canonical (92), candidate (160), and healthy (17) drivers across organ systems. **c** Proportion of canonical drivers detected in each organ system over canonical drivers detected in all cancer screens (421). The horizontal line indicates the median across all organ systems (22%). **d** Proportion of genes with non-coding driver alterations over all cancer drivers in each organ system. The horizontal line indicates the median across all organ systems (4%). **e** Number of canonical (e), candidate (f), and healthy (g) drivers across screens and organ systems. Representative genes with different recurrence between cancer and healthy tissues are indicated. **h** Organ system distribution of the top eight recurrent healthy drivers. The full list is provided as Additional file 6, Table S5. **i** Correlation between numbers of sequenced donors and identified healthy drivers across organ systems. Spearman correlation coefficient  $R$  and associated  $p$ -value are shown

(Fig. 2b). Despite this, some similarities and differences with cancer drivers could already be observed. Like cancer drivers (Fig. 2e, f, Additional file 6, Table S5), also healthy drivers were mostly organ-specific (Fig. 2g) and the most recurrent healthy drivers were also cancer drivers in the same organ system (Fig. 2h, Additional file 6, Table S5). However, some recurrent cancer drivers (*KRAS*, *PI3KCA*, *NRAS*, *NFI*) were reported to drive non-cancer clonal expansion only in one or two organ systems (Fig. 2g). Therefore, differences start to emerge at the tissue level between drivers of cancer and non-cancer evolution. Moreover, unlike cancer drivers, no correlation existed between the numbers of drivers and donors (Fig. 2i). This is likely affected by the lower number of non-cancer sequencing studies available so far. If additional studies will confirm the absence of correlation, this may indicate that the healthy driver repertoire is easier to saturate since fewer drivers are needed to initiate and sustain non-cancer clonal expansion [10, 11].

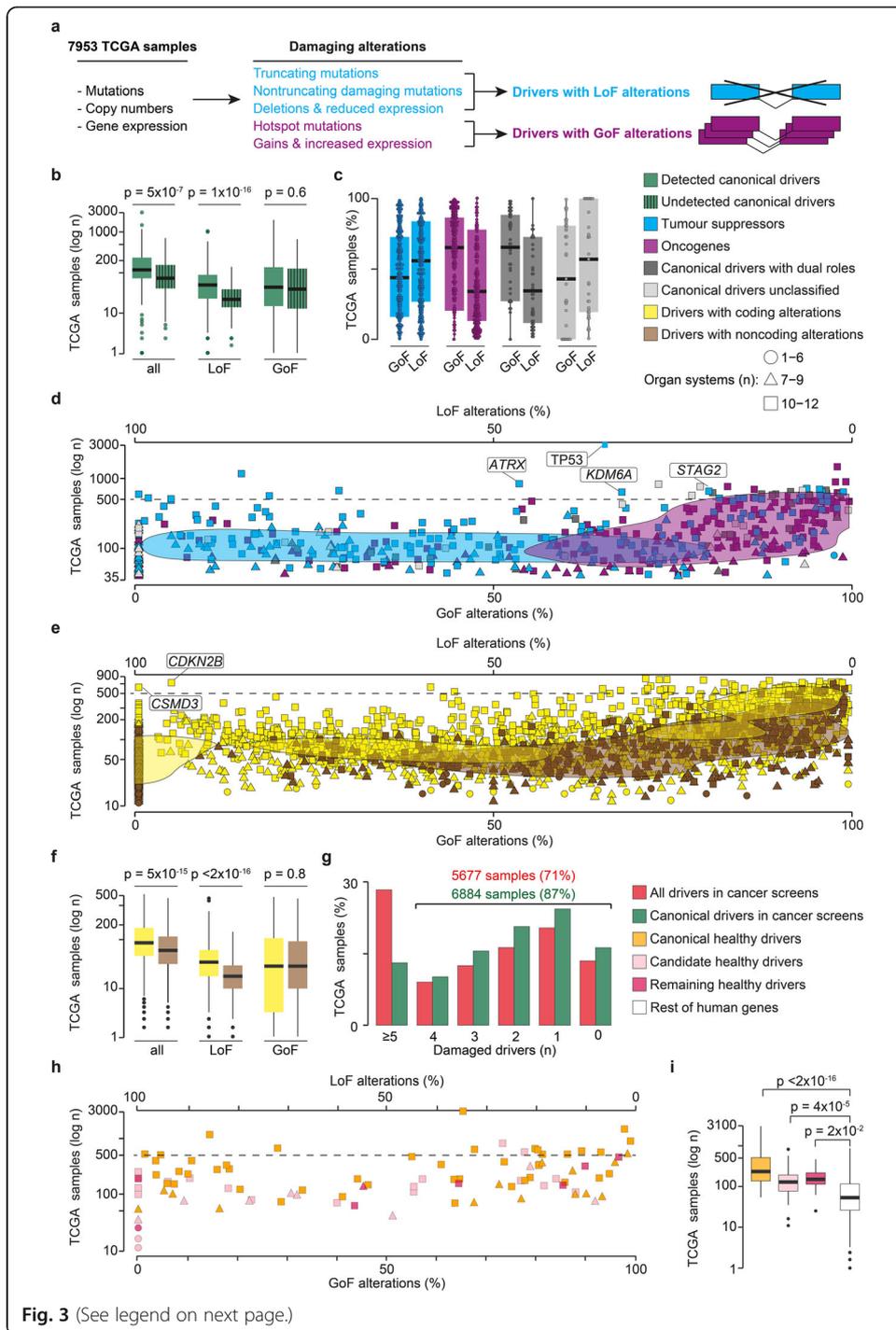
#### **Alteration pattern hints at driver mode of action and confirms the incompleteness of the driver repertoire**

To gain further insights into their mode of action, we mapped the type of alterations acquired by cancer and healthy drivers in 34 cancer types from TCGA. After predicting the damaging alterations in 7953 TCGA samples with matched mutation, copy number, and gene expression data (the “Methods” section), we identified the drivers with loss-of-function (LoF) and gain-of-function (GoF) alterations in these samples, respectively (Fig. 3a).

The comparison between canonical cancer drivers detected and undetected in sequencing screens (Fig. 1d) revealed that the latter were damaged in a significantly lower number of samples, due to fewer LoF alterations (Fig. 3b, Additional file 2, Fig. S3A). GoF alterations were instead comparable between the two groups, suggesting that current driver detection methods fail to identify drivers that undergo copy number gains but are rarely mutated.

We confirmed that the driver alteration patterns reflected their mode of action, with canonical tumor suppressors and oncogenes showing a prevalence of LoF and GoF alterations, respectively (Fig. 3c). Canonical drivers with a dual role resembled the alteration pattern of oncogenes while those still unclassified had a prevalence of LoF alterations, suggesting a putative tumor suppressor role (Fig. 3c). While all frequently altered (> 500 samples) oncogenes were overwhelmingly modified by GoF alterations (Additional file 7, Table S6), 16 of the 22 most frequently altered tumor suppressors had a prevalence of GoF alterations (Fig. 3d). In the majority of cases, this was due to different alteration patterns across organ systems (Additional file 2, Fig. S3B), and a possible oncogenic role has been documented for some others [29–38].

Since candidate drivers had no annotation of their mode of action, we reasoned that their alteration pattern could hint at their role as tumor suppressors or oncogenes. According to their prevalent pancancer alterations, 1318 candidates could be classified as putative tumor suppressors and 1405 as putative oncogenes (Additional file 7, Table S6). Interestingly, while candidates with predicted coding driver alterations showed similar distributions of LoF and GoF alterations (Fig. 3e), those with only non-coding driver alterations had a significantly lower occurrence of LoF alterations (Fig. 3f,



**Fig. 3** (See legend on next page.)

(See figure on previous page.)

**Fig. 3** Damaging alteration pattern of drivers in TCGA. **a** Identification of damaged drivers in 7953 TCGA samples. Mutations, gene deletions, and amplifications were annotated according to their predicted damaging effect. This allowed to distinguish drivers acquiring loss-of-function (LoF) or gain-of-function (GoF) alterations. **b** Number of TCGA samples with damaging alterations (all, LoF, GoF) in canonical drivers that were detected (421) or undetected (170) by cancer driver detection methods. **c** Proportion of TCGA samples with GoF and LoF alterations in tumor suppressors, oncogenes, and canonical drivers with a dual or unclassified role. Proportion of TCGA samples with GoF and LoF alterations in **(d)** canonical drivers and **(e)** candidate drivers. Genes mentioned in the text are highlighted. The two-dimensional Gaussian kernel density estimations were calculated for each driver group using the R density function. **f** Number of TCGA samples with damaging alterations (all, LoF, GoF) in drivers previously reported in coding and non-coding sequences. **g** Proportion of samples with variable numbers of all damaged drivers or only canonical drivers. **h** Proportion of TCGA samples with GoF and LoF alterations in healthy drivers. Canonical and candidate healthy drivers correspond to genes with a known or predicted cancer driver role. **i** Number of TCGA samples with damaged canonical, candidate, and remaining healthy drivers and the rest of human genes. All distributions were compared using a two-sided Wilcoxon rank-sum test

Additional file 2, Fig. S3C). This may suggest an activating role for their non-coding alterations too. Almost all candidates damaged in  $\geq 500$  samples (111/115) were putative oncogenes (Fig. 3e, Additional file 7, Table S6). Of the four putative tumor suppressors, *CSMD3* has a disputed cancer role [39–41] and a likely inflated mutation rate [42], while *CDKN2B* cooperates with its paralog *CDKN2A* to inhibit cell cycle [43], supporting its tumor suppressor role.

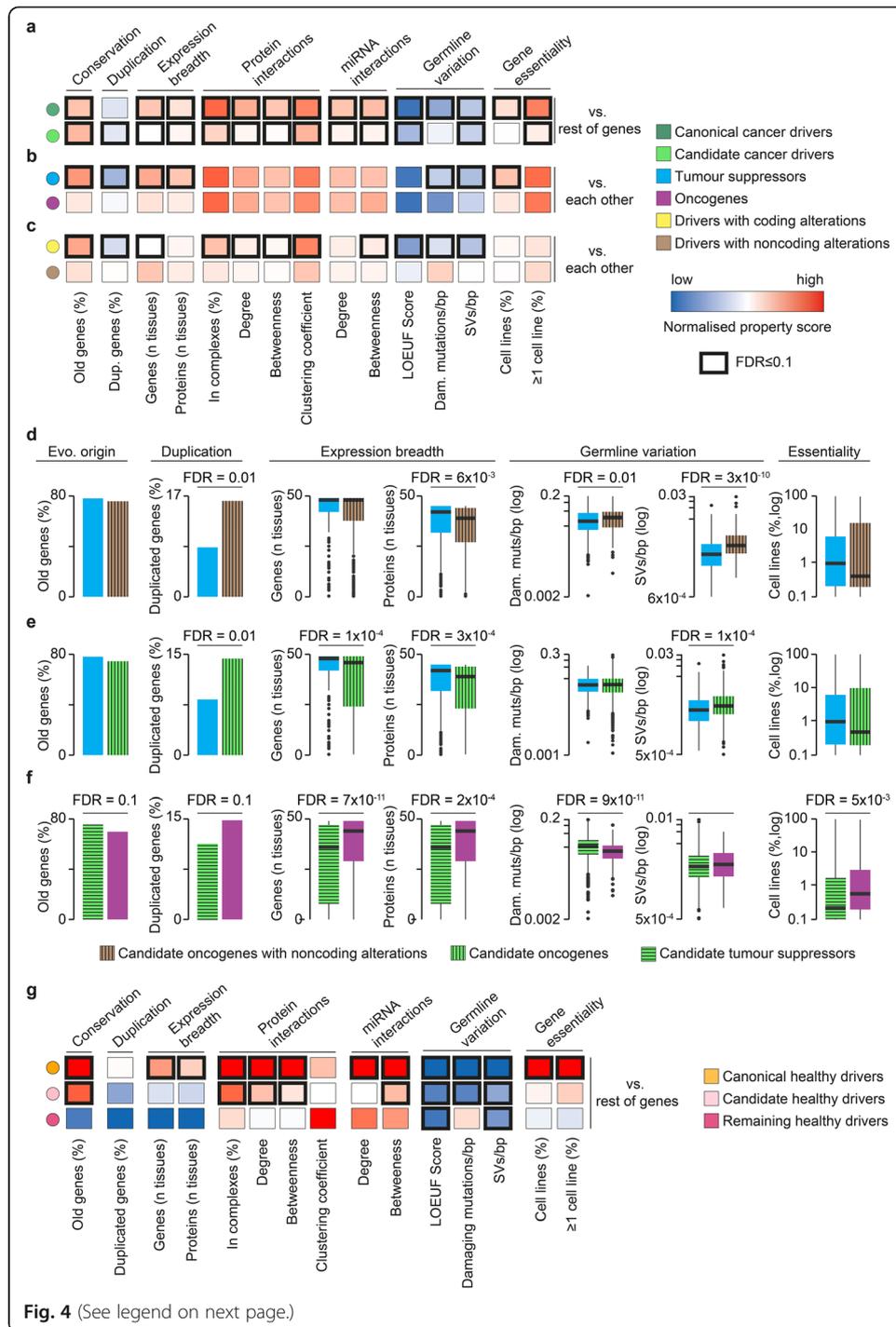
The number of damaged cancer drivers in individual TCGA samples confirmed that, despite all efforts, the current driver repertoire is still largely incomplete. The large majority of samples (71% and 87%, considering all drivers or only canonical drivers, respectively) had less than five damaged drivers, and  $\sim 15\%$  of them had no damaged driver (Fig. 3g).

Given their high overlap with cancer drivers, most healthy drivers were recurrently damaged in cancer samples with no prevalence of GoF or LoF alterations (Fig. 3h, Additional file 7, Table S6). Interestingly, all healthy drivers, even the eight with no cancer involvement, were damaged in significantly more cancer samples than the rest of human genes (Fig. 3i). Moreover, 57% of TCGA samples had at least two altered drivers, one of which was a healthy driver, further supporting the hypothesis that more than one driver may be needed to promote the transformation of non-malignant clones into cancer [10, 11].

### Properties of cancer and healthy drivers support their central role in the cell

A substantial body of work including our own [44–53] has shown that cancer drivers differ from the rest of the genes for an array of systems-level properties (Fig. 1a) that are a consequence of their unique evolutionary path and role in the cell. Using our granular annotation of drivers, we set out to check for similarities and differences across the driver groups.

We confirmed that cancer drivers, and in particular canonical drivers, were more conserved throughout evolution and less likely to retain gene duplicates than other human genes (Fig. 4a, Additional file 8, Table S7). They also showed broader tissue expression, engaged in a larger number of protein complexes, and occupied more central and highly connected positions in the protein-protein and miRNA-gene networks (Fig. 4a). We reported substantial differences between tumor suppressors and



**Fig. 4** (See legend on next page.)

(See figure on previous page.)

**Fig. 4** Systems-level properties of cancer and healthy drivers. Comparisons of systems-level properties between (a) canonical or candidate cancer drivers and the rest of human genes, (b) tumor suppressors and oncogenes, and (c) cancer genes with coding driver alterations and cancer genes with non-coding driver alterations. The normalized property score was calculated as the normalized difference between the median (continuous properties) or proportion (categorical properties) values in each driver group and the rest of human genes (the “Methods” section). Comparisons of systems-level properties between (d) candidate oncogenes with non-coding driver alterations (324) and canonical tumor suppressors, (e) candidate oncogenes (1405) and canonical tumor suppressors, and (f) candidate tumor suppressors (1318) and canonical oncogenes. g. Comparisons of systems-level properties between canonical healthy, candidate healthy, and remaining healthy drivers and the rest of human genes. Proportions of old (pre-metazoan), duplicated, essential genes, and proteins involved in the complexes were compared using a two-sided Fisher’s exact test. Distributions of gene and protein expression, protein-protein, miRNA-gene interactions, and germline variation were compared using a two-sided Wilcoxon rank-sum test. False discovery rate (FDR) was corrected for using Benjamini-Hochberg

oncogenes, with the former enriched in old and single-copy genes showing broader tissue expression (Fig. 4b, Additional file 8, Table S7).

We further expanded the systems-level properties of cancer drivers by exploring their tolerance towards germline variation, because this may indicate their essentiality. Using germline data from healthy individuals [54], we compared the loss-of-function observed/expected upper bound fraction (LOEUF) score, which quantifies selection towards LoF variation [54] as well as the number of damaging mutations and structural variants (SVs) per coding base pairs (bp) between drivers and the rest of genes (the “Methods” section). Cancer drivers, and in particular canonical drivers, had a significantly lower LOEUF score and retained fewer damaging germline mutations and SVs than the rest of the genes (Fig. 4a). This indicates that they are indispensable for cell survival in the germline. Selection against harmful variation was stronger in tumor suppressors than oncogenes (Fig. 4b). This was supported by a significantly higher proportion of cell lines where cancer drivers, and in particular tumor suppressors, were essential (Fig. 4a, b), as gathered from the integration of nine genome-wide essentiality screens [55–63] (the “Methods” section).

Genes with non-coding driver alterations had weaker systems-level properties than those with coding alterations (Fig. 4c, Additional file 8, Table S7) and the subset of them with > 50% GoF alterations resembled the property profile of oncogenes when compared to tumor suppressors (Fig. 4d, Additional file 8, Table S7). In general, all candidate drivers with a prevalence of GoF were similar to oncogenes, showing a higher proportion of duplicated genes, narrower tissue expression, and higher tolerance to germline variation than tumor suppressors (Fig. 4e, Additional file 8, Table S7). Conversely, candidate drivers with a prevalence of LoF were older, less duplicated, and less tolerant to germline variation than oncogenes (Fig. 4f, Additional file 8, Table S7).

Systems-level properties of healthy drivers varied according to the overlap with cancer drivers (Fig. 4g, Additional file 8, Table S7). Intriguingly, canonical healthy drivers showed stronger systems-level properties than any other group of drivers. In particular, they were enriched in evolutionarily conserved and broadly expressed genes encoding highly inter-connected proteins are regulated by many miRNAs. Moreover, these genes showed a strong selection against germline variation and high enrichment in essential genes (Fig. 4g). They therefore represent a core of genes with a very central role in the cell, whose modifications are not tolerated in the germline but are selected for in

somatic cells because they confer selective growth advantages. Candidate healthy drivers and those not involved in cancer had a substantially different property profile (Fig. 4g). Although numbers are too low for any robust conclusion, it is tempting to speculate that genes able to initiate non-cancer clonal expansion but not tumorigenesis may follow a different evolutionary path.

### **The Network of Cancer Genes: an open-access repository of annotated drivers**

We collected the whole repertoire of 3347 cancer and 95 healthy drivers, their literature support, and properties in the seventh release of the Network of Cancer Genes and Healthy Drivers (NCG<sup>HD</sup>) database. NCG<sup>HD</sup> is accessible through an open-access portal that enables interactive queries of drivers (Fig. 5a) as well as the bulk download of the database content.

In addition to the known or predicted mode of action and systems-level properties of cancer and healthy drivers, NCG<sup>HD</sup> 7.0 also annotates their function, alteration pattern, and gene expression profile in TCGA and cancer cell lines, reported interactions with antineoplastic drugs, and potential role as treatment biomarkers (Fig. 5b). Altogether, this constitutes an extensive compendium of annotation of driver genes, including information relevant for planning experiments involving them.

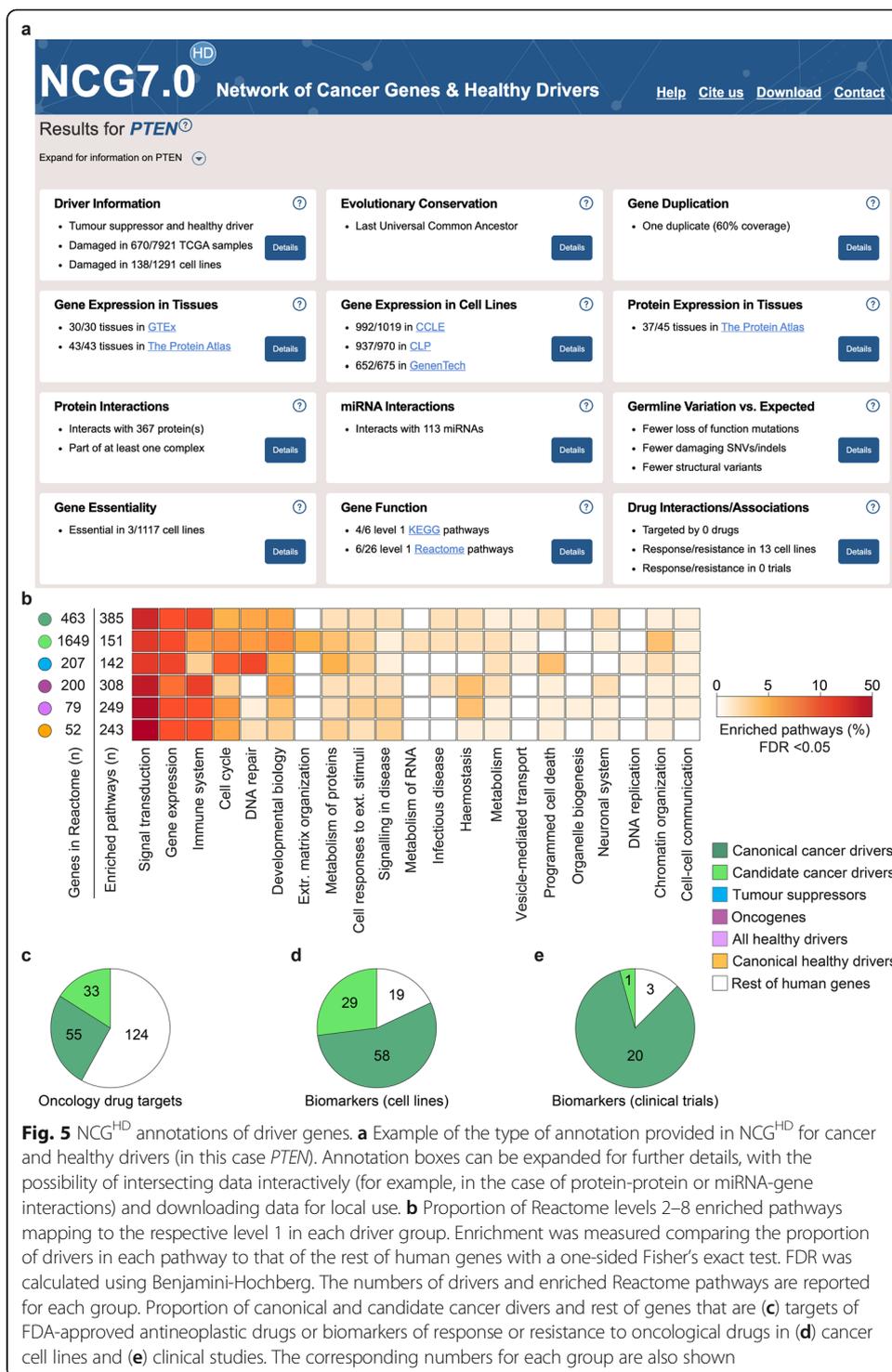
Functional gene set enrichment analysis showed that at least 60% of enriched pathways (FDR < 0.05) in any driver group converge to five broad functional processes (signal transduction, gene expression, immune system, cell cycle, and DNA repair, Fig. 5b, Additional file 9, Table S8). Within these, tumor suppressors showed a prevalence in cell cycle and DNA repair pathways, while oncogenes were enriched in the gene expression and immune system-related pathways (Additional file 9, Table S8). Healthy drivers closely resembled the functional profile of cancer drivers, given the high overlap (Fig. 5b). Because of the low number, it was not possible to assess the functional enrichment of healthy drivers not involved in cancer.

More than 9% of canonical cancer drivers are targets of anti-cancer drugs and cancer drivers constitute around 40% of their targets (Fig. 5c). Moreover, most of the genes used as biomarkers of resistance or response to treatment in cell lines (Fig. 5d) or clinical trials (Fig. 5e) are cancer drivers, with an overwhelming prevalence of canonical cancer drivers.

### **Discussion**

The wealth of cancer genomic data and the availability of increasingly sophisticated analytical approaches for their interpretation have substantially improved the understanding of how cancer starts and develops. However, our in-depth analysis of the vast repertoire of drivers that have been collected so far shows clear limits in the current knowledge of the driver landscape.

The identification of drivers as genes under positive selection or with a higher than expected mutation frequency within a cohort of patients has biased the current cancer driver repertoire towards genes whose coding point mutations or small indels frequently recur across patients. This strongly impairs the ability to map the full extent of driver heterogeneity leading to an underappreciation of the driver contribution of rarely altered genes and those modified through non-coding or gene copy number alterations,



particularly amplifications. It also results in a sizeable fraction of samples with very few or no cancer drivers. This gap can be solved by complementing cohort-level approaches with methods that account for all types of alterations and predict drivers in individual samples, for example identifying their network deregulations [64–66] or applying machine learning to identify driver alterations [67]. Alternatively, we have shown that

systems-level properties capture the main features of cancer drivers, justifying their use for patient-level driver detection [68, 69].

Our comprehensive study has also shown that cancer sequencing screens have so far mostly focused on resequencing and analyzing the protein-coding portion of cancer genomes, leaving the contribution of non-coding drivers mostly uncovered. This bias may be addressed by performing additional cancer whole genome sequencing screens and improving analytical methods for the prediction of non-coding driver alterations.

Biases are starting to emerge also in the knowledge of healthy drivers. Many non-cancer sequencing screens only targeted cancer genes and healthy driver detection methods used so far were originally developed for cancer genomics. Both these factors may contribute at least in part to explain the high overlap between drivers of cancer and non-cancer evolution. An unbiased investigation of altered genes able to promote clonal expansion but not tumorigenesis could confirm whether their properties are indeed different from cancer drivers as suggested by our initial analysis on the few of them that have been identified so far. Additionally, the investigation of somatically mutated clones in non-cancer tissues has just started and new screens are continuously published. The integrated analysis of these new studies will broaden our understanding of non-cancer clonal expansion and further clarify its relationship with cancer transformation.

Our literature review did not cover driver genes deriving from chromosomal rearrangements or epigenetic changes because of their scattered annotations in the literature and difficulty in mapping their properties. Adding these genes to the repertoire when their knowledge will be mature will help close the gaps in the knowledge of the genetic drivers of tumorigenesis.

## Conclusions

Our comprehensive analysis of cancer sequencing screens showed that the current repertoire of cancer driver genes is still incomplete and biased towards frequent mutations altering the gene coding sequence. This calls for the need for additional screens and methods to identify further coding and non-coding cancer drivers at single patient resolution. We confirmed the central role of cancer drivers within the cell, which is reflected in their evolutionary path and is shared by the majority of known healthy drivers. Further sequencing screens of healthy tissues are needed to clarify whether this is a feature of all genes whose mutations can driver non-cancer clonal expansion or there is a group of healthy drivers that underwent a different evolutionary path.

## Methods

### Literature curation

A literature search was carried out in PubMed, TCGA (<https://www.cancer.gov/tcga>) and ICGC (<https://dcc.icgc.org/>) to retrieve cancer screens published between 2018 and 2020 (Additional file 2, Fig. S1A). This resulted in 135 coding and 154 non-coding cancer screens. Of these, only 80 and 37 were retained after examining abstracts and full text, respectively. Criteria for removal were the absence of driver genes or driver detection methods and the impossibility to map non-coding driver alterations to genes. The 37 new cancer screens were added to 273 publications previously curated by our team

[70], totaling 310 publications (Additional file 1, Table S1). A similar literature search retrieved 24 sequencing screens of non-cancer tissues publications, 18 of which were retained after the abstract and full-text examination (Additional file 2, Fig. S1A; Additional file 1, Table S1). Each paper was reviewed independently by two experts and further discussed if annotations differed to extract the list of driver genes, the number of donors, the type of screen (whole-genome, whole-exome, target gene re-sequencing), the cancer or non-cancer tissues, and the driver detection method (Additional file 2, Fig. S1B).

Canonical cancer drivers were extracted from two publications [17, 18] and the Cancer Gene Census [71] v.91. In the latter case, all tiers 1 and 2 genes were retained, except those from genomic rearrangements leading to gene fusion (Additional file 2, Fig. S1B). Collected genes were further classified as tumor suppressor, oncogene, or having a dual role according to the annotation in the majority of sources. Genes with conflicting or unavailable annotation were left unclassified.

Drivers from cancer screens and canonical sources underwent further filtering (Additional file 2, Fig. S1C). First, they were intersected with a list of 148 possible false positives [18, 42]. After a manual check of the supporting evidence, two drivers were retained as canonical, five were considered as candidates, and 41 were removed (Additional file 3, Table S2). The three resulting lists (canonical drivers, drivers from cancer screens, and healthy drivers) were intersected to annotate canonical drivers in cancer screens, remaining drivers in cancer screens (candidate cancer drivers), canonical healthy drivers, candidate healthy drivers, and remaining healthy drivers (Additional file 2, Fig. S1C; Additional file 4, Table S3).

Cancer types and non-cancer tissues were mapped to organ systems using previous classification [72]. Cancer types not included in this classification were mapped based on their histopathology (retinoblastoma to central nervous system, vascular and peripheral nervous system cancers to soft tissue, penile tumors to urologic system).

### **Pancancer TCGA data**

A dataset of 7953 TCGA samples with quality-controlled mutation (SNVs and indels), copy number, and gene expression data in 34 cancer types was assembled from the Genomic Data Commons portal I [73] (<https://portal.gdc.cancer.gov/>). Mutations were annotated with ANNOVAR [74] (April 2018) and dbNSFP [75] v3.0 and only those identified as exonic or splicing were retained. Damaging mutations included (1) truncating (stopgain, stoploss, frameshift) mutations, (2) missense mutations predicted by at least seven out of 10 predictors (SIFT [76], PolyPhen-2 HDIV [77], PolyPhen-2 HVAR, MutationTaster [78], MutationAssessor [79], LRT [80], FATHMM [81], PhyloP [82], GERP++RS [83], and SiPhy [84]), (3) splicing mutations predicted by at least one of two splicing-specific methods (ADA [75] and RF [75]), and (4) hotspot mutations identified with OncodriveCLUST [85] v1.0.0.

Copy number variant (CNV) segments, sample ploidy, and sample purity values were obtained from TCGA SNP arrays using ASCAT [86] v.2.5.2. Segments were intersected with the exonic coordinates of 19,756 human genes in hg19 and genes were considered to have CNV if at least 25% of their transcribed length was covered by a CNV segment. RNA-Seq data were used to filter out false-positive CNVs. Putative gene gains were

defined as copy number (CN) > 2 times sample ploidy and the levels of expression were compared between samples with and without each gene gain using a two-sided Wilcoxon rank-sum test and corrected for multiple testing using Benjamini-Hochberg. Only gene gains with a false discovery rate (FDR) < 0.05 were retained. Homozygous gene losses had CN = 0 and fragments per kilobase per million (FPKM) values < 1 over sample purity. Heterozygous gene losses had CN = 1 or CN = 0 but FPKM values > 1 over sample purity. This resulted in 2,192,832 redundant genes damaged in 7921 TCGA samples.

In total, 518,115 genes were considered to acquire LoF alterations because they underwent homozygous deletion or had truncating, missense damaging, splicing mutations, or double hits (CN = 1 and LoF damaging mutation), while 1,674,717 genes were considered to acquire GoF alterations because they had a hotspot mutation or underwent gene gain with increased expression (Fig. 3a).

### Systems-level properties

Protein sequences from RefSeq [87] v.99 were aligned to hg38 using BLAT [88]. Unique genomic loci were identified for 19,756 genes based on gene coverage, span, score, and identity [89]. Genes sharing at least 60% of their protein sequence were considered as duplicates [46].

Evolutionary conservation was assessed for 18,922 human genes using their orthologs in EggNOG [90] v.5.0. Genes were considered to have a pre-metazoan origin (and therefore conserved in evolution) if they had orthologs in prokaryotes, eukaryotes, or opisthokonts [53].

Gene expression for 19,231 genes in 49 healthy tissues was derived from the union of Protein Atlas [91] v.19.3 and GTEx [92] v.8. Genes were considered to be expressed in a tissue if their expression value was  $\geq 1$  transcript per million (TPM). Protein expression for 13,229 proteins in 45 healthy tissues was derived from Protein Atlas [91] v.19.3 retaining the highest value when multiple expression values were available.

A total of 542,397 non-redundant binary interactions between 17,883 proteins were gathered from the integration of five sources (BioGRID [93] v.3.5.185, IntAct [94] v.4.2.14, DIP [95] (February 2018), HPRD [96] v.9 and Bioplex [97] v.3.0). Data on 9476 protein complexes involving 8504 proteins were derived from CORUM [98] v.3.0, HPRD [96] v.9 and Reactome [99] v.72. Experimentally supported interactions between 14,747 genes and 1758 miRNAs were acquired from miRTarBase [100] v.8.0 and miRecords [101] v.4.0. Degree, betweenness, and clustering coefficient were calculated for protein and miRNA networks using the igraph R package [102] v.1.2.6.

The loss-of-function observed/expected upper bound fraction (LOEUF) score for 18,392 genes was obtained from gnomAD [54] v.2.1.1. Germline mutations (SNVs and indels) were obtained from the union of 2504 samples from the 1000 Genomes Project Phase 3 [103] v.5a and 125,748 samples from gnomAD [54] v.2.1.1. Mutations were annotated with ANNOVAR [74] (October 2019), and 18,812 genes were considered as damaged using the same definitions as for TCGA samples. A total of 32,558 germline SVs for 14,158 genes were derived using 15,708 samples from gnomAD [54] v.2.1.1. The numbers of damaging mutations and SVs per base pairs (bp) were calculated for each gene.

Essentiality data for 19,013 genes in 1122 cell lines were obtained integrating three RNAi knockdown and six CRISPR Cas9 knockout screens [55–63]. Genes with CERES [57] or DEMETER [63] scores  $< -1$  or Bayes score [104]  $> 5$  were considered as essential.

Proportions of duplicated, pre-metazoan, essential genes, and proteins engaging in complexes were compared between the gene groups using two-sided Fisher's exact test. Distributions of tissues where genes or proteins were expressed, protein and miRNA network properties, LOEUF scores, damaging mutations, and SVs per bp were compared between the gene groups using a two-sided Wilcoxon test. Multiple comparisons within each property were corrected using Benjamini-Hochberg. For each systems-level property in each driver group ( $d$ ), a normalized property score was calculated as:

$$\text{Normalised property score} = \text{sgn}(\Delta_d) \times \frac{|\Delta_d| - \min_t |\Delta_t|}{\max_t |\Delta_t| - \min_t |\Delta_t|}$$

where  $t$  represents 11 gene groups (canonical drivers, candidate drivers, tumor suppressors, oncogenes, drivers with coding alterations, drivers with non-coding alterations, canonical healthy drivers, candidate healthy drivers, remaining healthy drivers, and the rest of human genes);  $\text{sgn}(\Delta_d)$  is the sign of the difference; and  $\Delta_d$  indicates the difference of medians (continuous properties) or proportions (categorical properties) between each driver group and the rest of human genes. Minima and maxima were taken over all 11 gene groups for each property.

#### Pancancer cell line data

Mutation, CNV and gene expression data for 1291 cell lines were obtained from DepMap [56, 105] v. 20Q3. Mutations were functionally annotated using ANNOVAR [74] and LoF mutations were identified as described for TCGA samples. Hotspot mutations were detected using hotspot positions derived from TCGA. Homozygous gene deletions were defined as  $\text{CN} < 0.25$  times cell line ploidy and expression  $< 1$  TPM; heterozygous gene deletions were defined as  $0.25 < \text{CN} < 0.75$  times cell line ploidy; gene gains were defined as  $\text{CN} > 2$  times cell line ploidy and significantly higher expression relative to cell lines with no gene gains. Genes with LoF or GoF alterations were defined as for TCGA samples. To map cell lines to organ systems, they were first associated with the TCGA cancer types and then the same classification as for TCGA was used [72].

#### Driver functional annotation

Gene functions were collected for 11,778 proteins from Reactome [99] v.72 and KEGG [106] v.94.1 (levels 1 and 2). Driver enrichment in Reactome pathways (levels 2–8) compared to the rest of human genes was assessed using a one-sided Fisher's exact test and corrected for multiple testing with Benjamini-Hochberg. Enriched pathways were then mapped to the corresponding Reactome level 1.

#### Drug interactions

A total of 247 FDA-approved, antineoplastic, and immunomodulating drugs targeting 212 human genes were downloaded from DrugBank [107] v.5.1.8. Genetic biomarkers of response and resistance to drugs in cancer cell lines were obtained from Genomics

of Drug Sensitivity in Cancer (GDSC) [108] v.8.2. Of those, only 467 associations with  $FDR \leq 0.25$  involving 129 drugs and 106 genes were retained. Genetic biomarkers of response and resistance in clinical studies were obtained from the Variant Interpretation for Cancer Consortium Meta-Knowledgebase [109] v.1. A total of 868 associations between drugs and genomic features involving 64 anti-cancer drugs and drug combinations and 24 human genes were retained [109].

### Database and website implementation

All annotations of driver genes were entered into a relational database based on MySQL [110] v.8.0.21 connected to a web interface enabling interactive retrieval of information through gene identifiers. The frontend was developed with PHP [111] v.7.4.15. The interactive displays of miRNA-gene and protein-protein interactions were implemented with the R packages Shiny [112] v.1.6.0 and igraph [102] v.1.2.6 and ran on Shiny Server v1.5.16.958.

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-022-02607-z>.

**Additional file 1: Table S1.** Publications describing driver genes.

**Additional file 2: Figure S1.** Literature search, review and annotation workflow; **Figure S2.** Correlation between numbers of donors and cancer drivers in individual organ systems; **Figure S3.** Patterns of driver damaging alterations in TCGA samples.

**Additional file 3: Table S2.** Putative false positive cancer drivers.

**Additional file 4: Table S3.** Canonical cancer drivers.

**Additional file 5: Table S4.** Donors in cancer and noncancer sequencing screens.

**Additional file 6: Table S5.** Drivers reported in cancer and non-cancer screens.

**Additional file 7: Table S6.** Cancer and non-cancer drivers damaged in TCGA.

**Additional file 8: Table S7.** Systems-level properties of driver genes.

**Additional file 9: Table S8.** Proportion of enriched pathways across driver groups.

**Additional file 10.** Review history.

### Acknowledgements

We thank Steve Hindmarsh and Stefan Boeing for their contribution to the development of the NCG database and website.

### Peer review information

Anahita Bishop and Barbara Cheifet were the primary editors of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

### Review history

The review history is available as Additional file 10.

### Authors' contributions

LD analyzed the protein-protein interactions, protein complex, gene essentiality, and cancer cell line data. MB analyzed the gene conservation. MRK analyzed the gene duplicability. HM analyzed the TCGA data. MB and HM analyzed the miRNA-target interactions. GS analyzed the gene function, RNA and protein expression, and drug interactions. AAS, LM, NW, and DR curated the literature. JN analyzed the germline variation. GS, MB, JG, and KA developed the database. MRK, HM, LD, MB, MP, PD, and AS developed the website. LD, MB, MRK, HM, GS, AAS, and FDC analyzed the data. FDC conceived and supervised the study. MB, AAS, GS, and FDC wrote the manuscript with contributions from LD and HM. The authors reviewed and approved the final manuscript.

### Funding

This work was supported by the Cancer Research UK [C43634/A25487], the Cancer Research UK King's Health Partners Centre at King's College London [C604/A25135], the Cancer Research UK City of London Centre [C7893/A26233], the innovation programme under the Marie Skłodowska-Curie grant agreement [CONTRA-766030], the EPSRC Centre for Doctoral Training in Cross-Disciplinary Approaches to Non-Equilibrium Systems (CANES, EP/L015854/1), the Health Education England Genomics Education Programme, and the Francis Crick Institute, which receives its core funding from Cancer Research UK (FC001002), the UK Medical Research Council (FC001002), and the Wellcome Trust (FC001002). For

the purpose of Open Access, the authors have applied a CC BY public copyright license to any Author Accepted Manuscript version arising from this submission.

#### Availability of data and materials

The whole content of NCG<sup>HD</sup> can be freely downloaded from the website (<http://network-cancer-genes.org/>). No license is required.

Original data were obtained from the following online sources:

1000 Genomes Project Phase 3 [103] v.5a: <https://www.internationalgenome.org/category/phase-3/>  
 BioGRID [93] v.3.5.185: <https://thebiogrid.org/>  
 Bioplex [97] v.3.0: <https://bioplex.hms.harvard.edu/interactions.php>  
 CORUM [98] v.3.0: <http://mips.helmholtz-muenchen.de/corum/>  
 Depmap [59, 60] v20Q3: <https://depmap.org/portal/>  
 DIP [95] (February 2018): <https://dip.doe-mbi.ucla.edu/dip/Main.cgi>  
 DrugBank [107] v.5.1.8: <https://go.drugbank.com/>  
 EggNog [90] v.5: <http://eggnog5.embl.de/#/app/home>  
 GDSC [108] v.8.2: <https://www.cancerxgene.org/>  
 GnomAD [54] v.2.1.1: <https://gnomad.broadinstitute.org/>  
 GTEx [92] v.8: <https://gtexportal.org/home/>  
 HPRD [96] v.9: <https://www.hprd.org/>  
 IntAct [94] v.4.2.14: <https://www.ebi.ac.uk/intact/home>  
 KEGG [106] v.94.1: <https://www.genome.jp/kegg/>  
 Meta-KB [109] v.1: <https://cancervariants.org/>  
 MiRecords [101] v.8.0: <http://c1.accurascience.com/miRecords/>  
 MiTarBase [100] v.4.0: [https://mirtarbase.cuhk.edu.cn/~miRTarBase/miRTarBase\\_2022/php/index.php](https://mirtarbase.cuhk.edu.cn/~miRTarBase/miRTarBase_2022/php/index.php)  
 NCI Genomics Data Commons Portal [73]: <https://gdc.cancer.gov/>  
 PICKLES [61]: <https://pickles.hart-lab.org/>  
 Protein Atlas [91] v.19.3: <https://www.proteinatlas.org/>  
 Reactome [99] v.72: <https://reactome.org/>  
 RefSeq [87] v.99: <https://www.ncbi.nlm.nih.gov/refseq/>

#### Declarations

##### Ethics approval and consent to participate

Not applicable.

##### Consent for publication

The results shown here are in part based upon data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>.

##### Competing interests

The authors declare that they have no competing interests.

##### Author details

<sup>1</sup>Cancer Systems Biology Laboratory, The Francis Crick Institute, London NW1 1AT, UK. <sup>2</sup>School of Cancer and Pharmaceutical Sciences, King's College London, London SE11UL, UK. <sup>3</sup>Department of Medical and Molecular Genetics, King's College London, London SE1 9RT, UK. <sup>4</sup>Scientific Computing, The Francis Crick Institute, London NW1 1AT, UK.

Received: 30 August 2021 Accepted: 10 January 2022

Published online: 26 January 2022

#### References

1. Network CGAR. Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature*. 2008;455(7216):1061–8. <https://doi.org/10.1038/nature07385>.
2. International Cancer Genome C, Hudson TJ, Anderson W, Artez A, Barker AD, Bell C, et al. International network of cancer genome projects. *Nature*. 2010;464:993–8.
3. Hutter C, Zenklusen JC. The Cancer Genome Atlas: creating lasting value beyond its data. *Cell*. 2018;173(2):283–5. <https://doi.org/10.1016/j.cell.2018.03.042>.
4. Pon JR, Marra MA. Driver and passenger mutations in cancer. *Annu Rev Pathol*. 2015;10(1):25–50. <https://doi.org/10.1146/annurev-pathol-012414-040312>.
5. Porta-Pardo E, Kamburov A, Tamborero D, Pons T, Grases D, Valencia A, et al. Comparison of algorithms for the detection of cancer drivers at subgene resolution. *Nat Methods*. 2017;14(8):782–8. <https://doi.org/10.1038/nmeth.4364>.
6. Martínez-Jiménez F, Muiños F, Sentís I, Deu-Pons J, Reyes-Salazar I, Arnedo-Pac C, et al. A compendium of mutational cancer driver genes. *Nat Rev Cancer*. 2020;20(10):555–72. <https://doi.org/10.1038/s41568-020-0290-x>.
7. Bailey MH, Tokheim C, Porta-Pardo E, Sengupta S, Bertrand D, Weerasinghe A, et al. Comprehensive characterization of cancer driver genes and mutations. *Cell*. 2018;173(2):371–85 e18. <https://doi.org/10.1016/j.cell.2018.02.060>.
8. Consortium ITP-CAoWG. Pan-cancer analysis of whole genomes. *Nature*. 2020;578(7793):82–93. <https://doi.org/10.1038/s41586-020-1969-6>.
9. Kandath C, McLellan MD, Vandin F, Ye K, Niu B, Lu C, et al. Mutational landscape and significance across 12 major cancer types. *Nature*. 2013;502(7471):333–9. <https://doi.org/10.1038/nature12634>.
10. Wijewardhane N, Dressler L, Ciccarelli FD. Normal somatic mutations in cancer transformation. *Cancer Cell*. 2021;39(2):125–9. <https://doi.org/10.1016/j.ccell.2020.11.002>.

11. Kakiuchi N, Ogawa S. Clonal expansion in non-cancer tissues. *Nat Rev Cancer*. 2021;21(4):239–56. <https://doi.org/10.1038/s41568-021-00335-3>.
12. Martincorena I, Roshan A, Gerstung M, Ellis P, Van Loo P, McLaren S, et al. Tumor evolution. High burden and pervasive positive selection of somatic mutations in normal human skin. *Science*. 2015;348(6237):880–6. <https://doi.org/10.1126/science.aaa6806>.
13. Tang J, Fewings E, Chang D, Zeng H, Liu S, Jorapur A, et al. The genomic landscapes of individual melanocytes from human skin. *Nature*. 2020;586(7830):600–5. <https://doi.org/10.1038/s41586-020-2785-8>.
14. Yokoyama A, Kakiuchi N, Yoshizato T, Nannya Y, Suzuki H, Takeuchi Y, et al. Age-related remodelling of oesophageal epithelia by mutated cancer drivers. *Nature*. 2019;565(7739):312–7. <https://doi.org/10.1038/s41586-018-0811-x>.
15. Martincorena I, Fowler JC, Wabik A, Lawson ARJ, Abascal F, Hall MWJ, et al. Somatic mutant clones colonize the human esophagus with age. *Science*. 2018;362(6417):911–7. <https://doi.org/10.1126/science.aau3879>.
16. Suda K, Nakaoka H, Yoshihara K, Ishiguro T, Tamura R, Mori Y, et al. Clonal expansion and diversification of cancer-associated mutations in endometriosis and normal endometrium. *Cell Rep*. 2018;24(7):1777–89. <https://doi.org/10.1016/j.celrep.2018.07.037>.
17. Vogelstein B, Papadopoulos N, Velculescu VE, Zhou S, Diaz LA, Kinzler KW. Cancer genome landscapes. *Science*. 2013;339(6127):1546–58. <https://doi.org/10.1126/science.1235122>.
18. Saito Y, Koya J, Araki M, Kogure Y, Shingaki S, Tabata M, et al. Landscape and function of multiple mutations within individual oncogenes. *Nature*. 2020;582(7810):95–9. <https://doi.org/10.1038/s41586-020-2175-2>.
19. Sondka Z, Bamford S, Cole CG, Ward SA, Dunham I, Forbes SA. The COSMIC Cancer Gene Census: describing genetic dysfunction across all human cancers. *Nat Rev Cancer*. 2018;18(11):696–705. <https://doi.org/10.1038/s41568-018-0060-1>.
20. Liu EM, Martinez-Fundichely A, Bollapragada R, Spiewack M, Khurana E. CNCDatabase: a database of non-coding cancer drivers. *Nucleic Acids Res*. 2021;49(D1):D1094–D101. <https://doi.org/10.1093/nar/gkaa915>.
21. Campbell PJ, Getz G, Korbel JO, Stuart JM, Jennings JL, Stein LD, et al. Pan-cancer analysis of whole genomes. *Nature*. 2020;578(7793):82–93. <https://doi.org/10.1038/s41586-020-1969-6>.
22. Hornshoj H, Nielsen MM, Sinnott-Armstrong NA, Switnicki MP, Juul M, Madsen T, et al. Pan-cancer screen for mutations in non-coding elements with conservation and cancer specificity reveals correlations with expression and survival. *NPJ Genom Med*. 2018;3(1):1. <https://doi.org/10.1038/s41525-017-0040-5>.
23. Juul M, Bertl J, Guo Q, Nielsen MM, Switnicki M, Hornshoj H, et al. Non-coding cancer driver candidates identified with a sample- and position-specific model of the somatic mutation rate. *Elife*. 2017;6. <https://doi.org/10.7554/eLife.21778>.
24. Zhu H, Uuskula-Reimand L, Isaev K, Wadi L, Alizada A, Shuai S, et al. Candidate cancer driver mutations in distal regulatory elements and long-range chromatin interaction networks. *Mol Cell*. 2020;77(6):1307–21 e10. <https://doi.org/10.1016/j.molcel.2019.12.027>.
25. Lanzos A, Carlevaro-Fita J, Mularoni L, Reverter F, Palumbo E, Guigo R, et al. Discovery of cancer driver long noncoding RNAs across 1112 tumour genomes: new candidates and distinguishing features. *Sci Rep*. 2017;7(1):41544. <https://doi.org/10.1038/srep41544>.
26. Mularoni L, Sabarinathan R, Deu-Pons J, Gonzalez-Perez A, Lopez-Bigas N. OncodriveFML: a general framework to identify coding and non-coding regions with cancer driver mutations. *Genome Biol*. 2016;17(1):128. <https://doi.org/10.1186/s13059-016-0994-0>.
27. Cornish AJ, Hoang PH, Dobbins SE, Law PJ, Chubb D, Orlando G, et al. Identification of recurrent noncoding mutations in B-cell lymphoma using capture Hi-C. *Blood Adv*. 2019;3(1):21–32. <https://doi.org/10.1182/bloodadvances.2018026419>.
28. Hayward NK, Wilmott JS, Waddell N, Johansson PA, Field MA, Nones K, et al. Whole-genome landscapes of major melanoma subtypes. *Nature*. 2017;545(7653):175–80. <https://doi.org/10.1038/nature22071>.
29. Botlagunta M, Vesuna F, Mironchik Y, Raman A, Lisok A, Winnard P Jr, et al. Oncogenic role of DDX3 in breast cancer biology. *Oncogene*. 2008;27(28):3912–22. <https://doi.org/10.1038/onc.2008.33>.
30. Pu J, Wang J, Qin Z, Wang A, Zhang Y, Wu X, et al. IGF2BP2 promotes liver cancer growth through an m6A-FEN1-dependent mechanism. *Front Oncol*. 2020;10:578816. <https://doi.org/10.3389/fonc.2020.578816>.
31. Sun X, Jia M, Sun W, Feng L, Gu C, Wu T. Functional role of RBM10 in lung adenocarcinoma proliferation. *Int J Oncol*. 2019;54(2):467–78. <https://doi.org/10.3892/ijo.2018.4643>.
32. Soussi T, Wiman KG. TP53: an oncogene in disguise. *Cell Death Differ*. 2015;22(8):1239–49. <https://doi.org/10.1038/cdd.2015.53>.
33. Yang MH, Chang SY, Chiou SH, Liu CJ, Chi CW, Chen PM, et al. Overexpression of NBS1 induces epithelial-mesenchymal transition and co-expression of NBS1 and Snail predicts metastasis of head and neck cancer. *Oncogene*. 2007;26(10):1459–67. <https://doi.org/10.1038/sj.onc.1209929>.
34. Manandhar S, Kim CG, Lee SH, Kang SH, Basnet N, Lee YM. Exostosin 1 regulates cancer cell stemness in doxorubicin-resistant breast cancer cells. *Oncotarget*. 2017;8(41):70521–37. <https://doi.org/10.18632/oncotarget.19737>.
35. Li A, Zhu X, Wang C, Yang S, Qiao Y, Qiao R, et al. Upregulation of NDRG1 predicts poor outcome and facilitates disease progression by influencing the EMT process in bladder cancer. *Sci Rep*. 2019;9(1):5166. <https://doi.org/10.1038/s41598-019-41660-w>.
36. Meacham CE, Lawton LN, Soto-Feliciano YM, Pritchard JR, Joughin BA, Ehrenberger T, et al. A genome-scale in vivo loss-of-function screen identifies Phf6 as a lineage-specific regulator of leukemia cell growth. *Genes Dev*. 2015;29(5):483–8. <https://doi.org/10.1101/gad.254151.114>.
37. Sesen J, Casaos J, Scotland SJ, Seva C, Eisinger-Mathason TS, Skuli N. The bad, the good and eIF3e/INT6. *Front Biosci (Landmark Ed)*. 2017;22:1–20.
38. Shi J, Zhang L, Zhou D, Zhang J, Lin Q, Guan W, et al. Biological function of ribosomal protein L10 on cell behavior in human epithelial ovarian cancer. *J Cancer*. 2018;9(4):745–56. <https://doi.org/10.7150/jca.21614>.
39. Liu P, Morrison C, Wang L, Xiong D, Vedell P, Cui P, et al. Identification of somatic mutations in non-small cell lung carcinomas using whole-exome sequencing. *Carcinogenesis*. 2012;33(7):1270–6. <https://doi.org/10.1093/carcin/bgs148>.
40. Lai MW, Liang KH, Lin WR, Huang YH, Huang SF, Chen TC, et al. Hepatocarcinogenesis in transgenic mice carrying hepatitis B virus pre-S/S gene with the sW172\* mutation. *Oncogenesis*. 2016;5(12):e273. <https://doi.org/10.1038/oncsis.2016.77>.

41. Cai C, Cooper GF, Lu KN, Ma X, Xu S, Zhao Z, et al. Systematic discovery of the functional impact of somatic genome alterations in individual tumors through tumor-specific causal inference. *PLoS Comput Biol*. 2019;15(7):e1007088. <https://doi.org/10.1371/journal.pcbi.1007088>.
42. Lawrence MS, Stojanov P, Polak P, Kryukov GV, Cibulskis K, Sivachenko A, et al. Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature*. 2013;499(7457):214–8. <https://doi.org/10.1038/nature12213>.
43. Hannon GJ, Beach D. p15INK4B is a potential effector of TGF- $\beta$ -induced cell cycle arrest. *Nature*. 1994;371:257–61.
44. Syed AS, D'Antonio M, Ciccarelli FD. Network of Cancer Genes: a web resource to analyze duplicability, orthology and network properties of cancer genes. *Nucleic Acids Res*. 2010;38(suppl\_1):D670–D75. <https://doi.org/10.1093/nar/gkp957>.
45. Trigos AS, Pearson RB, Papenfuss AT, Goode DL. Somatic mutations in early metazoan genes disrupt regulatory links between unicellular and multicellular genes in cancer. *eLife*. 2019;8:e40947. <https://doi.org/10.7554/eLife.40947>.
46. Rambaldi D, Giorgi FM, Capuani F, Ciliberto A, Ciccarelli FD. Low duplicability and network fragility of cancer genes. *Trends Genet*. 2008;24(9):427–30. <https://doi.org/10.1016/j.tig.2008.06.003>.
47. Domazet-Lošo T, Tautz D. Phylostratigraphic tracking of cancer genes suggests a link to the emergence of multicellularity in metazoa. *BMC Biol*. 2010;8(1):66. <https://doi.org/10.1186/1741-7007-8-66>.
48. D'Antonio M, Ciccarelli FD. Integrated analysis of recurrent properties of cancer genes to identify novel drivers. *Genome Biol*. 2013;14(5):R52. <https://doi.org/10.1186/gb-2013-14-5-r52>.
49. Ostrow SL, Barshir R, DeGregori J, Yeager-Lotem E, Hershsberg R. Cancer evolution is associated with pervasive positive selection on globally expressed genes. *PLoS Genet*. 2014;10(3):e1004239. <https://doi.org/10.1371/journal.pgen.1004239>.
50. An O, Dall'Olio GM, Mourikis TP, Ciccarelli FD. NCG 5.0: updates of a manually curated repository of cancer genes and associated properties from cancer mutational screenings. *Nucleic Acids Res*. 2016;44:D992–9.
51. Jonsson PF, Bates PA. Global topological features of cancer proteins in the human interactome. *Bioinformatics*. 2006;22(18):2291–7. <https://doi.org/10.1093/bioinformatics/btl390>.
52. Xia J, Sun J, Jia P, Zhao Z. Do cancer proteins really interact strongly in the human protein-protein interaction network? *Comput Biol Chem*. 2011;35(3):121–5. <https://doi.org/10.1016/j.compbiolchem.2011.04.005>.
53. D'Antonio M, Ciccarelli FD. Modification of gene duplicability during the evolution of protein interaction network. *PLoS Comput Biol*. 2011;7(4):e1002029. <https://doi.org/10.1371/journal.pcbi.1002029>.
54. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*. 2020;581(7809):434–43. <https://doi.org/10.1038/s41586-020-2308-7>.
55. Dempster JM, Rossen J, Kazachkova M, Pan J, Kugener G, Root DE, et al. Extracting biological insights from the project Achilles genome-scale CRISPR screens in cancer cell lines. *BioRxiv*. 2019;720243. <https://doi.org/10.1101/720243>.
56. Broad D. DepMap 20Q3 Public, figshare. Dataset. 2020. <https://doi.org/10.6084/m9.figshare.12931238.v1>.
57. Meyers RM, Bryan JG, McFarland JM, Weir BA, Sizemore AE, Xu H, et al. Computational correction of copy number effect improves specificity of CRISPR-Cas9 essentiality screens in cancer cells. *Nat Genet*. 2017;49(12):1779–84. <https://doi.org/10.1038/ng.3984>.
58. Behan FM, Iorio F, Picco G, Goncalves E, Beaver CM, Migliardi G, et al. Prioritization of cancer therapeutic targets using CRISPR-Cas9 screens. *Nature*. 2019;568(7753):511–6. <https://doi.org/10.1038/s41586-019-1103-9>.
59. DepMap Broad. Project SCORE processed with CERES. figshare. Dataset. 2019. <https://doi.org/10.6084/m9.figshare.9116732>.
60. DepMap Broad. DepMap GeCKO 19Q1. figshare. Fileset. 2019. <https://doi.org/10.6084/m9.figshare.7668407>.
61. Lenoir WF, Lim TL, Hart T. PICKLES: the database of pooled in-vitro CRISPR knockout library essentiality screens. *Nucleic Acids Res*. 2018;46(D1):D776–D80. <https://doi.org/10.1093/nar/gkx993>.
62. McFarland JM, Ho ZV, Kugener G, Dempster JM, Montgomery PG, Bryan JG, et al. Improved estimation of cancer dependencies from large-scale RNAi screens using model-based normalization and data integration. *Nat Commun*. 2018;9(1):4610. <https://doi.org/10.1038/s41467-018-06916-5>.
63. Tsherniak A, Vazquez F, Montgomery PG, Weir BA, Kryukov G, Cowley GS, et al. Defining a cancer dependency map. *Cell*. 2017;170(3):564–76 e16. <https://doi.org/10.1016/j.cell.2017.06.010>.
64. Bertrand D, Chng KR, Sherbat FG, Kiesel A, Chia BK, Sia YY, et al. Patient-specific driver gene prediction and risk assessment through integrated network analysis of cancer omics profiles. *Nucleic Acids Res*. 2015;43(7):e44. <https://doi.org/10.1093/nar/gku1393>.
65. Bashashati A, Haffari G, Ding J, Ha G, Lui K, Rosner J, et al. DriverNet: uncovering the impact of somatic driver mutations on transcriptional networks in cancer. *Genome Biol*. 2012;13(12):R124. <https://doi.org/10.1186/gb-2012-13-12-r124>.
66. Hou JP, Ma J. DawnRank: discovering personalized driver genes in cancer. *Genome Med*. 2014;6(7):56. <https://doi.org/10.1186/s13073-014-0056-8>.
67. Dong C, Guo Y, Yang H, He Z, Liu X, Wang K. iCAGES: integrated CAncer GEnome Score for comprehensively prioritizing driver genes in personal cancer genomes. *Genome Med*. 2016;8(1):135. <https://doi.org/10.1186/s13073-016-0390-0>.
68. Nulsen J, Missetic H, Yau C, Ciccarelli FD. Pan-cancer detection of driver genes at the single-patient resolution. *Genome Med*. 2021;13(1):12. <https://doi.org/10.1186/s13073-021-00830-0>.
69. Mourikis TP, Benedetti L, Foxall E, Temelkovski D, Nulsen J, Perner J, et al. Patient-specific cancer genes contribute to recurrently perturbed pathways and establish therapeutic vulnerabilities in esophageal adenocarcinoma. *Nat Commun*. 2019;10(1):3101. <https://doi.org/10.1038/s41467-019-10898-3>.
70. Repana D, Nulsen J, Dressler L, Bortolomeazzi M, Venkata SK, Tourna A, et al. The Network of Cancer Genes (NCG): a comprehensive catalogue of known and candidate cancer genes from cancer sequencing screens. *Genome Biol*. 2019;20(1):1. <https://doi.org/10.1186/s13059-018-1612-0>.
71. Tate JG, Bamford S, Jubb HC, Sondka Z, Beare DM, Bindal N, et al. COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res*. 2019;47(D1):D941–D47. <https://doi.org/10.1093/nar/gky1015>.
72. Hoadley KA, Yau C, Hinoue T, Wolf DM, Lazar AJ, Drill E, et al. Cell-of-origin patterns dominate the molecular classification of 10,000 tumors from 33 types of cancer. *Cell*. 2018;173:291–304.e6.
73. Grossman RL, Heath AP, Ferretti V, Varmus HE, Lowy DR, Kibbe WA, et al. Toward a shared vision for cancer genomic data. *N Engl J Med*. 2016;375(12):1109–12. <https://doi.org/10.1056/NEJMp1607591>.
74. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res*. 2010;38(16):e164. <https://doi.org/10.1093/nar/gkq603>.

75. Liu X, Wu C, Li C, Boerwinkle E. dbNSFP v3.0: a one-stop database of functional predictions and annotations for human nonsynonymous and splice-site SNVs. *Hum Mutat.* 2016;37(3):235–41. <https://doi.org/10.1002/humu.22932>.
76. Ng PC, Henikoff S. SIFT: predicting amino acid changes that affect protein function. *Nucleic Acids Res.* 2003;31(13):3812–4. <https://doi.org/10.1093/nar/gkg509>.
77. Adzhubei I, Jordan DM, Sunyaev SR. Predicting functional effect of human missense mutations using PolyPhen-2. *Curr Protoc Hum Genet.* 2013;Chapter 7:Unit7.20.
78. Schwarz JM, Rödelsperger C, Schuelke M, Seelow D. MutationTaster evaluates disease-causing potential of sequence alterations. *Nat Methods.* 2010;7(8):575–6. <https://doi.org/10.1038/nmeth0810-575>.
79. Reva B, Antipin Y, Sander C. Predicting the functional impact of protein mutations: application to cancer genomics. *Nucleic Acids Res.* 2011;39(17):e118. <https://doi.org/10.1093/nar/gkr407>.
80. Chun S, Fay JC. Identification of deleterious mutations within three human genomes. *Genome Res.* 2009;19(9):1553–61. <https://doi.org/10.1101/gr.092619.109>.
81. Shihab HA, Gough J, Cooper DN, Stenson PD, Barker GL, Edwards KJ, et al. Predicting the functional, molecular, and phenotypic consequences of amino acid substitutions using hidden Markov models. *Hum Mutat.* 2013;34(1):57–65. <https://doi.org/10.1002/humu.22225>.
82. Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res.* 2010;20(1):110–21. <https://doi.org/10.1101/gr.097857.109>.
83. Davydov EV, Goode DL, Sirota M, Cooper GM, Sidow A, Batzoglou S. Identifying a high fraction of the human genome to be under selective constraint using GERP++. *PLoS Comput Biol.* 2010;6(12):e1001025. <https://doi.org/10.1371/journal.pcbi.1001025>.
84. Garber M, Guttman M, Clamp M, Zody MC, Friedman N, Xie X. Identifying novel constrained elements by exploiting biased substitution patterns. *Bioinformatics.* 2009;25(12):i54–62. <https://doi.org/10.1093/bioinformatics/btp190>.
85. Tamborero D, Gonzalez-Perez A, Lopez-Bigas N. OncodriveCLUST: exploiting the positional clustering of somatic mutations to identify cancer genes. *Bioinformatics.* 2013;29(18):2238–44. <https://doi.org/10.1093/bioinformatics/btt395>.
86. Van Loo P, Nordgard SH, Lingjaerde OC, Russnes HG, Rye IH, Sun W, et al. Allele-specific copy number analysis of tumors. *Proc Natl Acad Sci U S A.* 2010;107(39):16910–5. <https://doi.org/10.1073/pnas.1009843107>.
87. O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, et al. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* 2016;44(D1):D733–D45. <https://doi.org/10.1093/nar/gkv1189>.
88. Kent WJ. BLAT—the BLAST-like alignment tool. *Genome Res.* 2002;12(4):656–64.
89. Bhagwat M, Young L, Robison RR. Using BLAT to find sequence similarity in closely related genomes. *Curr Protoc Bioinformatics.* 2012;37(1):10.8.1–10.8.24. <https://doi.org/10.1002/0471250953.bi1008s37>.
90. Huerta-Cepas J, Szklarczyk D, Heller D, Hernández-Plaza A, Forslund SK, Cook H, et al. eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Res.* 2019;47(D1):D309–D14. <https://doi.org/10.1093/nar/gky1085>.
91. Uhlén M, Fagerberg L, Hallström BM, Lindskog C, Oksvold P, Mardinoglu A, et al. Proteomics. Tissue-based map of the human proteome. *Science.* 2015;347(6220):1260419. <https://doi.org/10.1126/science.1260419>.
92. Consortium G. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science.* 2020;369(6509):1318–30. <https://doi.org/10.1126/science.aaz1776>.
93. Oughtred R, Stark C, Breitkreutz BJ, Rust J, Boucher L, Chang C, et al. The BioGRID interaction database: 2019 update. *Nucleic Acids Res.* 2019;47(D1):D529–D41. <https://doi.org/10.1093/nar/gky1079>.
94. Orchard S, Ammari M, Aranda B, Breuza L, Briganti L, Broackes-Carter F, et al. The MIntAct project—IntAct as a common curation platform for 11 molecular interaction databases. *Nucleic Acids Res.* 2014;42(D1):D358–63. <https://doi.org/10.1093/nar/gkt1115>.
95. Salwinski L, Miller CS, Smith AJ, Pettit FK, Bowie JU, Eisenberg D. The database of interacting proteins: 2004 update. *Nucleic Acids Res.* 2004;32(90001):D449–51. <https://doi.org/10.1093/nar/gkh086>.
96. Keshava Prasad TS, Goel R, Kandasamy K, Keerthikumar S, Kumar S, Mathivanan S, et al. Human Protein Reference Database—2009 update. *Nucleic Acids Res.* 2009;37(Database):D767–72. <https://doi.org/10.1093/nar/gkn892>.
97. Huttlin EL, Bruckner RJ, Navarrete-Perea J, Cannon JR, Baltier K, Gebreab F, et al. Dual proteome-scale networks reveal cell-specific remodeling of the human interactome. *Cell.* 2021;184(11):3022–40 e28. <https://doi.org/10.1016/j.cell.2021.04.011>.
98. Giurgiu M, Reinhard J, Brauner B, Dunger-Kaltenbach I, Fobo G, Frishman G, et al. CORUM: the comprehensive resource of mammalian protein complexes-2019. *Nucleic Acids Res.* 2019;47(D1):D559–D63. <https://doi.org/10.1093/nar/gky973>.
99. Jassal B, Matthews L, Viteri G, Gong C, Lorente P, Fabregat A, et al. The reactome pathway knowledgebase. *Nucleic Acids Res.* 2020;48(D1):D498–503. <https://doi.org/10.1093/nar/gkz1031>.
100. Huang H-Y, Lin Y-C-D, Li J, Huang K-Y, Shrestha S, Hong H-C, et al. miRTarBase 2020: updates to the experimentally validated microRNA–target interaction database. *Nucleic Acids Res.* 2020;48:D148–D54.
101. Xiao F, Zuo Z, Cai G, Kang S, Gao X, Li T. miRecords: an integrated resource for microRNA–target interactions. *Nucleic Acids Res.* 2009;37(Database):D105–D10. <https://doi.org/10.1093/nar/gkn851>.
102. Csardi G, Nepusz T. The igraph software package for complex network research. *InterJ Complex Syst.* 2006;1695:1–9.
103. Auton A, Abecasis GR, Altshuler DM, Durbin RM, Abecasis GR, Bentley DR, et al. A global reference for human genetic variation. *Nature.* 2015;526(7571):68–74. <https://doi.org/10.1038/nature15393>.
104. Hart T, Moffat J. BAGEL: a computational framework for identifying essential genes from pooled library screens. *BMC Bioinformatics.* 2016;17(1):164. <https://doi.org/10.1186/s12859-016-1015-8>.
105. Ghandi M, Huang FW, Jane-Valbuena J, Kryukov GV, Lo CC, McDonald ER 3rd, et al. Next-generation characterization of the Cancer Cell Line Encyclopedia. *Nature.* 2019;569(7757):503–8. <https://doi.org/10.1038/s41586-019-1186-3>.
106. Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K. KEGG: new perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* 2016;45(D1):D353–D61. <https://doi.org/10.1093/nar/gkw1092>.
107. Wishart DS, Feunang YD, Guo AC, Lo EJ, Marcu A, Grant JR, et al. DrugBank 5.0: a major update to the DrugBank database for 2018. *Nucleic Acids Res.* 2018;46(D1):D1074–d82. <https://doi.org/10.1093/nar/gkx1037>.

108. Iorio F, Knijnenburg TA, Vis DJ, Bignell GR, Menden MP, Schubert M, et al. A Landscape of pharmacogenomic interactions in cancer. *Cell*. 2016;166(3):740–54. <https://doi.org/10.1016/j.cell.2016.06.017>.
109. Wagner AH, Walsh B, Mayfield G, Tamborero D, Sonkin D, Krysiak K, et al. A harmonized meta-knowledgebase of clinical interpretations of somatic genomic variants in cancer. *Nat Genet*. 2020;52(4):448–57. <https://doi.org/10.1038/s41588-020-0603-8>.
110. MySQL 8.0 reference manual. <https://dev.mysql.com/doc/refman/8.0/en/>.
111. Bakken S, Suraski Z, Schmid E. PHP manual; 2020.
112. Chang W, Cheng J, Allaire J, Sievert C, Schloerke B, Xie Y, et al. shiny: web application framework for R. 2021.

### **Publisher's Note**

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

