

RESEARCH ARTICLE

Efficient representations of tumor diversity with paired DNA-RNA aberrations

Qian Ke¹, Wikum Dinalankara^{2,3}, Laurent Younes^{1*}, Donald Geman^{1*}, Luigi Marchionni^{2,3*}

1 Department of Applied Mathematics and Statistics, Johns Hopkins University, Baltimore, Maryland, United States of America, **2** Department of Oncology, Johns Hopkins University School of Medicine, Baltimore, Maryland, United States of America, **3** Department of Pathology and Laboratory Medicine, Weill Cornell Medicine, New York, New York, United States of America

☞ These authors contributed equally to this work.

* laurent.younes@jhu.edu (LY); geman@jhu.edu (DG); marchion@med.cornell.edu (LM)



OPEN ACCESS

Citation: Ke Q, Dinalankara W, Younes L, Geman D, Marchionni L (2021) Efficient representations of tumor diversity with paired DNA-RNA aberrations. *PLoS Comput Biol* 17(6): e1008944. <https://doi.org/10.1371/journal.pcbi.1008944>

Editor: Sushmita Roy, University of Wisconsin, Madison, UNITED STATES

Received: May 28, 2020

Accepted: April 7, 2021

Published: June 11, 2021

Copyright: © 2021 Ke et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: RNA-Seq data, somatic mutation data and copy number data for The Cancer Genome Atlas were obtained through the Xena Cancer Genome Browser database (<https://xenabrowser.net>) from individual cancer type cohorts. Computational functionality for the optimization procedure is provided at <https://github.com/wikum/lpcover> and the code for the analysis in the manuscript is provided at <https://github.com/wikum/CoveringAnalysis>. Processed data in the form of TAB delimited files, and selected tissue-level coverings (in excel format) are provided within the [Supporting information](#) files,

Abstract

Cancer cells display massive dysregulation of key regulatory pathways due to now well-catalogued mutations and other DNA-related aberrations. Moreover, enormous heterogeneity has been commonly observed in the identity, frequency and location of these aberrations across individuals with the same cancer type or subtype, and this variation naturally propagates to the transcriptome, resulting in myriad types of dysregulated gene expression programs. Many have argued that a more integrative and quantitative analysis of heterogeneity of DNA and RNA molecular profiles may be necessary for designing more systematic explorations of alternative therapies and improving predictive accuracy. We introduce a representation of multi-omics profiles which is sufficiently rich to account for observed heterogeneity and support the construction of quantitative, integrated, metrics of variation. Starting from the network of interactions existing in Reactome, we build a library of “paired DNA-RNA aberrations” that represent prototypical and recurrent patterns of dysregulation in cancer; each two-gene “Source-Target Pair” (STP) consists of a “source” regulatory gene and a “target” gene whose expression is plausibly “controlled” by the source gene. The STP is then “aberrant” in a joint DNA-RNA profile if the source gene is DNA-aberrant (e.g., mutated, deleted, or duplicated), and the downstream target gene is “RNA-aberrant”, meaning its expression level is outside the normal, baseline range. With M STPs, each sample profile has exactly one of the 2^M possible configurations. We concentrate on subsets of STPs, and the corresponding reduced configurations, by selecting tissue-dependent minimal coverings, defined as the smallest family of STPs with the property that every sample in the considered population displays at least one aberrant STP within that family. These minimal coverings can be computed with integer programming. Given such a covering, a natural measure of cross-sample diversity is the extent to which the particular aberrant STPs composing a covering vary from sample to sample; this variability is captured by the entropy of the distribution over configurations. We apply this program to data from TCGA for six distinct tumor types (breast, prostate, lung, colon, liver, and kidney cancer). This enables an efficient simplification of the complex landscape observed in cancer populations, resulting in the identification of novel signatures of molecular alterations which are not detected with

and are also available from the Marchionni laboratory website (<http://marchionnilab.org/signatures.html>).

Funding: All authors (QK, WD, LY, DG, and LM) were supported by National Institute of Health (NIH) National Cancer Institute (NCI) Grant R01CA200859. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

frequency-based criteria. Estimates of cancer heterogeneity across tumor phenotypes reveals a stable pattern: entropy increases with disease severity. This framework is then well-suited to accommodate the expanding complexity of cancer genomes and epigenomes emerging from large consortia projects.

Author summary

A large variety of genomic and transcriptomic aberrations are observed in cancer cells, and their identity, location, and frequency can be highly indicative of the particular subtype or molecular phenotype, and thereby inform treatment options. However, elucidating this association between sets of aberrations and subtypes of cancer is severely impeded by considerable diversity in the set of aberrations across samples from the same population. Most attempts at analyzing tumor heterogeneity have dealt with either the genome or transcriptome in isolation. Here we present a novel, multi-omics approach for quantifying heterogeneity by determining a small set of paired DNA-RNA aberrations that incorporates potential downstream effects on gene expression. We apply integer programming to identify a small set of paired aberrations such that at least one among them is present in every sample of a given cancer population. The resulting “coverings” are analyzed for six cancer cohorts from the Cancer Genome Atlas, and facilitate introducing an information-theoretic measure of heterogeneity. Our results identify many known facets of tumorigenesis as well as suggest potential novel genes and interactions of interest.

Introduction

Cancer cells evade the normal mechanisms controlling cellular growth and tissue homeostasis through the disruption of key regulatory pathways controlling these processes. Such dysregulation results from genetic and epigenetic aberrations, encompassing mutations, copy number alterations, and changes in chromatin states, which affect the genes participating in such regulatory networks.

Over the past several decades, the list of known genetic and genomic aberrations in cancer has greatly expanded, thanks to large-scale projects such as the The Cancer Genome Atlas (TCGA, [1]), the Catalogue Of Somatic Mutations In Cancer (COSMIC, [2]), the MSK/IMPACT study [3], and recent efforts from the ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium [4].

Whereas the number of aberrations which suffice for progression to an advanced cancer is thought to be rather small, at least for solid tumors [5, 6] and at the pathway level [7], the number of ways (combinations of aberrations) for which this can be actualized is very large. In particular, the landscape collectively emerging from these studies exhibits a high degree of variation in the identity, frequency, and location of these aberrations, as well as tissue- and expression-dependency [8, 9]. These differences—collectively referred to as *tumor heterogeneity*—are “context-specific”, differing among tissue types and epigenetic conditions [8], across different cells within a lesion (*intra-tumor heterogeneity*), between tumor lesions within the same individual (*inter-tumor heterogeneity*), and across distinct individuals with the same cancer type or sub-type (*across-sample or population-level heterogeneity*).

In addition, such DNA defects, in order to be “functional” (*i.e.*, manifest themselves) and ultimately alter the cellular phenotype, must propagate through the signaling and regulatory

network and alter the downstream gene expression programs [10, 11]. These downstream transcriptional changes are in fact also context-specific, varying within and among cancers, local environments, and individuals. Most importantly, it has been speculated that transcriptionally heterogeneous tumors may be more adaptable to changes in the tumor microenvironment and therefore more likely to acquire new properties such as metastatic potential and resistance to treatments, leading to dismal patient outcomes; in addition, for predicting the response to targeted therapies, gene expression profiles may be more discriminating than mutational status [12]. The analysis of heterogeneity of molecular profiles, both DNA and RNA, is therefore of paramount importance. Consequently, a deeper, integrative and quantitative analysis of tumor heterogeneity is necessary for achieving a better understanding of the underlying biology, for designing more systematic explorations of candidate therapies, and for improving the accuracy of prognosis and treatment response predictions.

Unsurprisingly, even *representing* such high-dimensional variability poses great challenges, especially if a major goal is to find suitable metrics to quantify the level of tumor heterogeneity. We assume that large-scale projects (see above) and studies (*e.g.*, [13]) have already provided reasonably comprehensive lists of the most important recurrent molecular alterations driving cancer initiation and progression. But merely counting or cataloging aberrations will not suffice to precisely measure heterogeneity in a tumor population, and to quantify how this differs across diverse contexts (*e.g.*, between cancer arising in distinct organs, or between tumor subtypes). In order to identify functional aberrations potentially exploitable as biomarkers and therapeutic targets, it is necessary to go well beyond frequency estimates to more powerful representations rooted in biological mechanism and accounting for statistical dependency among aberrations.

We introduce a representation of *omics* profiles which is sufficiently rich to account for observed heterogeneity and to support the construction of quantitative, integrated metrics. Our framework is centered on the joint analysis of “paired DNA-RNA aberrations” that represent prototypical and recurrent patterns of dysregulation in cancer. Specifically, we represent the space of gene alterations that result in network perturbations and downstream changes of gene expression levels as a catalogue of mechanistic, two-gene “Source-Target Pairs” (STPs), each consisting of a “source” gene (important driver) and a “target” gene for which the mRNA expression is controlled either directly by the source gene or indirectly by a close descendant of the source.

We extend STPs from a network property to a sample property (like the existence of individual aberrations) by declaring an STP to be “aberrant” in a joint DNA-RNA profile if the source gene is DNA-aberrant (*e.g.*, mutated, deleted, or duplicated), and the target gene is RNA-aberrant, meaning its expression level is “divergent” (*i.e.*, outside the normal, baseline range [14]). This defines one binary random variable per STP, of which there are typically hundreds of thousands, most of which have a very small probability to be realized in a sample.

Samples are then characterized by their entire set of paired DNA-RNA aberrations, or aberrant STPs. Therefore, given there are M STPs, exactly one of the 2^M possible configurations is assigned to each sample. The extent to which these subsets vary from sample to sample is then a natural measure of heterogeneity in the population from which the samples are drawn.

Due to the difficulty of estimating rare events with the modest sample sizes available in cancer genomics today, any multivariate property of the probability distribution over the 2^M STP configurations (for example, its entropy) cannot be accurately approximated without a substantial further reduction of complexity. Such a reduction is provided by the concept of *minimal coverings* of a population (previously employed for modeling networks [15]). Here, we focus on smallest collections C of paired aberrations with the property that (nearly) every tumor sample has at least one aberrant STP in C . Indeed, since nearly all tumor samples exhibit

multiple aberrant STPs, a *minimal covering* necessarily exists (perhaps not unique), which can be found using well-known algorithms for formulating “optimal set covering” as the solution of an integer-programming problem (see [Methods](#)).

Our main contribution is then a method for integrating DNA and RNA data which yields novel insights about regulatory mechanisms in cancer, and consists of three parts:

1. A representation of network dysregulation based on matched pairs of genes, one gene with aberrant DNA and the other gene downstream, with aberrant RNA expression.
2. An algorithm for finding the minimal covering of a cancer (sub)population by aberrant genes or gene pairs.
3. An information-theoretic characterization of inter-sample heterogeneity as the entropy of the distribution of covering states.

Our methods are described in more detail in the next sections, followed by a presentation of our results. We conclude this paper with a discussion and provide additional results in supplementary material (see [S1 Text](#)).

1 Methods

1.1 Overall strategy

Identifying and quantifying the cross-sample heterogeneity of *omics* datasets with large numbers of random variables requires making simplifying assumptions and approximations on the joint distribution of the considered variables to make it feasible. We performed our analyses using matched DNA mutations, copy number alterations, and RNA expression data, pre-processed with the method previously described in [14]. In the present study we specifically focused on six distinct tumor types (TCGA code in parenthesis): breast invasive carcinoma (BRCA), prostate adenocarcinoma (PRAD), lung adenocarcinoma (LUAD), liver hepatocellular carcinoma (LIHC), kidney renal clear cell carcinoma (KIRC), and colon adenocarcinoma (COAD). For simplicity, hereafter, we will refer to these tumor types according to the tissue of origin (breast, prostate, lung, kidney, liver, and colon).

Our definition of aberrant expression of RNA [14] requires expression data from a baseline population, taken here as corresponding normal tissue (see 1.2.1). Consequently, our selection of cancer types was constrained by having enough normal samples in TCGA to estimate the “normal expression range” of the RNA-Seq data. In addition, we also consider a variety of clinical scenarios across different patient populations. Our approach is depicted in the schematic of [Fig 1](#).

Cancer phenotypes. We also focused on specific patient subgroups defined based on standard clinical and pathological variables (which are ordinal) routinely used for patient risk stratification. Tumor stage (from I to IV) indicates extension of a cancer and whether it has spread beyond the site of origin. The lymph node status (positive versus negative) indicates the presence of lymph node metastases. Tumor T status (from T1 to T4) indicates the size of the primary tumor. Tumor histologic grade (from G1 to G3 or G4, depending on the tumor type) captures the progressive departure from the the normal tissue and cellular architecture observed under a microscope. The Gleason grading system [16] is specific to prostate cancer and it accounts for 5 grades. The Gleason sum results from the two predominant grade patterns observed (*i.e.*, “primary” and “secondary”), with a sum of 6 (3+3) corresponding to indolent tumors, and sums from 7 to 10 associated with increasingly aggressive phenotypes. Finally, the PAM50 breast cancer subtypes [17] and the colorectal cancer CRIS classes [18]

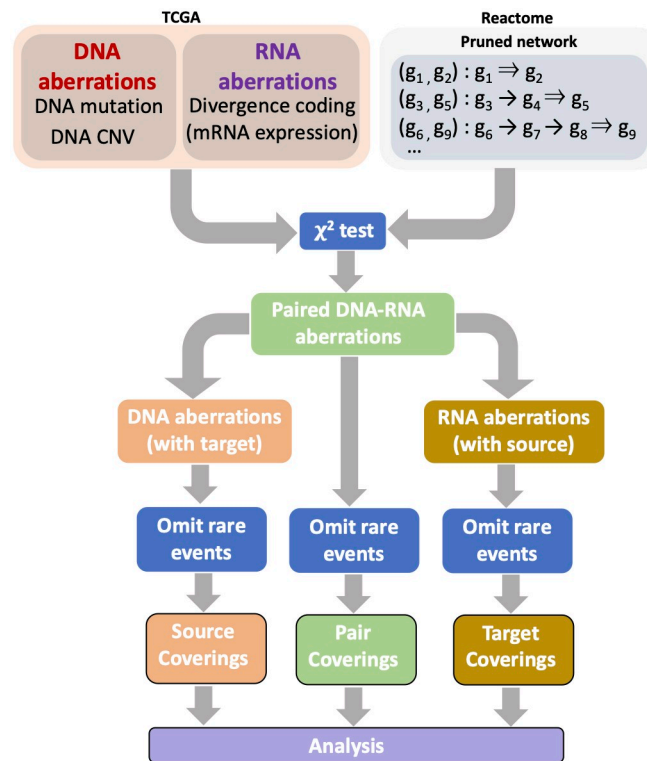


Fig 1. Overall analytical workflow. Source-target pairs (STPs) are constructed using the links available in Reactome [19]. In the TCGA cancer cohorts, the mutation and copy number variation data are used to construct binary DNA aberration profiles; the presence of either a mutation or high/low copy number variation at a given gene is treated as an aberration for the given gene for that sample *omics* profile. The gene expression data are used to construct binary RNA aberration profiles based on falling outside the “normal” expression range (in quantiles) for each gene based on TCGA normal tissue expression data, as previously described [14]. The binary profiles are combined to produce paired DNA-RNA aberrations, following which filtering is performed by selecting pairs that are determined to be significant (two-sided χ^2 test). The selected STPs then give rise to individual source (DNA) and target (RNA) aberrations, providing binary *omics* profiles at the level of source, target, and pairs. STPs that are present in less than 2% of samples for a given tissue are omitted. Then coverings are computed at the pair, source and target levels and subtype analysis and heterogeneity analysis carried out.

<https://doi.org/10.1371/journal.pcbi.1008944.g001>

are patient subgroups with distinct prognosis defined based on specific gene expression signatures.

Aberration detection. We reduce the data to binary variables indicating deviations from normal behavior, and the resulting indicators are furthermore filtered using an STP-based analysis requiring plausible mechanisms leading to the aberrations.

Covering estimation. In order to reduce the number of variables under consideration, we estimate subsets of “important variables,” called coverings, defined as minimal sets of variables from which, with high probability, cancer samples have at least one aberrant observation (see Section 1.3 on their computation). Because such optimal coverings are generally not unique, we include the consideration of variables that are present in at least one covering (union), or the restriction to variables that appear in all of them (intersection), that we refer to as “core” variables, or the use of a single covering, for example the one maximizing the sum of frequencies of aberrations among its variables.

Entropy estimation. We assess the heterogeneity of a population of samples by computing the entropy over a limited family of configurations determined by a covering of this population. This computation is not straightforward; even though reduced profiles involve a relatively

small number m of binary variables (typically a few dozen) indicating aberration of STPs, the observed sample size remains insufficient to allow for the estimation of the probabilities of the 2^m joint configurations of these variables. Some approximations are necessary and are described in Section 1.4.

Code-based reduction. Using a tree-based decomposition, we decompose tumor samples resulting from a given cancer type into cells, or bins, associated with a small number of conjunctions and disjunctions of aberrations. It is then possible to visualize and compare the resulting histograms in sub-populations defined by specific subtypes or phenotypes. This is described in Section 1.5.

1.2 Aberrations

1.2.1 Univariate deviation from normality in omics data. We transform the original data into sparse binary vectors indicating whether each variable deviates from a reference state or normal range when observed on a given sample.

Our pre-processing of DNA data, which already provides deviations from wild type, is quite simple. We consider that a gene g is aberrant at the DNA level if it includes a mutation that differs from the wild type, or if its copy number corresponds to a homozygous deletion (which would entail a complete gene inactivation), or a gain of 2 copies (to increase the plausibility of aberrant over-expression). We exclude heterozygous deletions and single copy gains since their impact is more difficult to interpret biologically. We will write X_g^{dna} for the corresponding binary random variable, so $X_g^{dna} = 1$ when g is DNA-aberrant.

The binarization of RNA data is more involved, and is based on the notion of “divergence” we previously developed [14]. Briefly, following a rank transformation, the range of RNA expression is estimated for normal samples for genes of interest. Then for each tumor sample and each gene, there is a binary variable with values 1 or 0 depending on whether the expression of the gene is outside or inside the expected normal region. Thus a gene is declared as RNA-aberrant if its ranking among other genes in the same sample falls outside of its normal range estimated from baseline data. Let X_g^{rna} be the corresponding binary random variable. This dichotomization requires a training step, solely based on normal tissue data, in order to estimate these normal intervals of variation. This being done, the decision for a gene to be RNA-divergent in a tumor sample only involves the RNA profile of this sample and is in particular independent of other tumor observations in the dataset.

1.2.2 Building source-target pairs. These binary *omics* variables are filtered by requiring that the deviations they represent have a plausible causal explanation as parts of STPs. Such STPs, denoted $(g_s \Rightarrow g_t)$, are built using apriori information representing known gene-gene interactions from signaling pathways and biological processes. In our implementation, we used the Reactome database [19] as retrieved from Pathway Commons (version 10) [20], since the network information contained therein is comprehensive and well curated (see [S1 Text](#) for further details on network curation and summary statistics). Our approach may be implemented using other databases providing gene-gene interaction information to build STPs as well.

Let \mathcal{N} denote the family of directed pairs of genes from this database, annotated as regulator and target, including two kinds of links $g \rightarrow g'$ for which “ g controls state change of g' ” (notation: $(g \xrightarrow{state} g')$) or “ g controls the expression of g' ” (notation $g \xrightarrow{expr} g'$). We say that two genes g_s, g_t form a “source-target pair (STP)”, with notation $g_s \Rightarrow g_t$ if there exists a sequence of l intermediate genes g_1, \dots, g_l such that

$$g_s \longrightarrow g_1 \longrightarrow \dots \longrightarrow g_l \xrightarrow{expr} g_t$$

where the intermediate links are either \xrightarrow{state} or \xrightarrow{expr} and the last link is \xrightarrow{expr} . Such a sequence has $k = l + 1$ links, and the minimal number of links required to achieve the STP is called the length of $g_s \Rightarrow g_t$.

Let Λ_k^* denote the set of STP s of length k or less deduced from the pathway database. This set, which is tissue independent, includes a large number of pairs (more than 200,000 for $k = 3$). For computational efficiency, we take $k = 3$ in our experiments. For a more detailed analysis on the selection of $k = 3$, see [S1 Text](#) (Table A in [S1 Text](#)). Let $\Lambda^* = \Lambda_3^*$ from here on. To select pairs that are most relevant for a tissue, this set is reduced by applying a χ^2 test for independence, only keeping STPs ($g_s \Rightarrow g_t$) for which the independence between the events “ g_s DNA aberrant” and “ g_t RNA aberrant” is rejected at a 5% level by the test (without correction for multiple hypotheses, because we want to be conservative with this selection) using a dataset of tumor samples. Let Λ denote the set of remaining pairs (typically 5,000–10,000), which is therefore tissue dependent (see [Fig 1](#)).

We then let \mathcal{S} denote the set of sources in Λ , *i.e.*, the set of genes g such that there exists g' with $(g \Rightarrow g') \in \Lambda$ and, similarly, let \mathcal{T} be the set of all possible targets. We let $T(g)$ denote the set of all the targets of $g \in \mathcal{S}$, that is, $T(g) = \{g' \in \mathcal{T} : (g \Rightarrow g') \in \Lambda\}$ and $S(g')$ the set of all sources pointing to $g' \in \mathcal{T}$.

1.2.3 Paired aberrations. We can now define a family of binary random variables ($Z_\lambda, \lambda \in \Lambda$) of “Paired Aberrations” with $Z_\lambda = 1$ for STP $\lambda = (g_s \Rightarrow g_t)$ if and only if g_s is aberrant at the DNA level (either due to mutation or copy-number variation) and g_t is RNA-aberrant. That is $Z_\lambda = X_{g_s}^{dna} X_{g_t}^{rna}$, a product of binary variables. For $\lambda = (g_s \Rightarrow g_t)$, we will also use the notation $s(\lambda) = g_s$ and $t(\lambda) = g_t$ for the source and target in λ .

From this, we also define binary variables $Z_g^{(s)}$ indicating aberrations at the source level letting $Z_g^{(s)} = 1$ if and only if g participates in an aberrant STP as a source gene. Therefore,

$$Z_g^{(s)} = \max\{Z_\lambda : s(\lambda) = g\}, \text{ for } g \in \mathcal{S}. \tag{1}$$

Similarly, we consider aberrations at the target level letting

$$Z_g^{(t)} = \max\{Z_\lambda : t(\lambda) = g\}, \text{ for } g \in \mathcal{T}. \tag{2}$$

We will refer to the event $Z_g^{(s)} = 1$ as a “source aberration with target” for gene g and the event $Z_g^{(t)} = 1$ as an “target aberration with source” for gene g .

1.3 Coverings

1.3.1 Definition. We have defined three types of aberrations involving multiple genes (STP, source, target), indexed by three different sets of pairs, sources or targets: i) STP aberrations which involve one source and one target gene; ii) “source aberration with target” which involve one source gene and all its targets; iii) “target aberrations with source” which involve one target gene and all its sources. We will use the generic notation $(Z_i, i \in \mathcal{I})$ to refer to the variables associated to any one of them, so that \mathcal{I} is one of the index sets Λ, \mathcal{S} or \mathcal{T} and $Z_i = Z_\lambda, Z_g^{(s)}$ or $Z_g^{(t)}$, respectively. We identify small subsets of \mathcal{I} of “essential variables” for describing the stochastic behavior of Z , such that at least one aberration occurs with high probability.

Denote by (Ω, P) the probability space on which $\bar{Z} = (Z_i, i \in \mathcal{I})$ is defined. If $\alpha \in [0, 1]$ and J is a subset of \mathcal{I} , we will say that $\bar{Z}_J = (Z_j, j \in J)$ is a covering (or 1-covering) of Ω at level α if,

with probability larger than $1 - \alpha$, at least one of its variables is aberrant, *i.e.*,

$$P(\exists j \in J : Z_j = 1) \geq 1 - \alpha. \tag{3}$$

In other terms, Ω is covered (up to a subset of probability less than α) by the union of events $U_j = \{Z_j = 1\}, j \in J$. For simplicity, we will also refer to the index set, J , as a covering rather than the set of variables it indexes. More generally, we can define an r -covering at level α as a set J for which at least r of the variables in J are aberrant with probability $1 - \alpha$, *i.e.*,

$$P(|\{j \in J : Z_j = 1\}| \geq r) \geq 1 - \alpha. \tag{4}$$

1.3.2 Optimal coverings. We assume that a family of weights $(w_i, i \in \mathcal{I})$ is given and consider the function

$$\sigma(J) = \sum_{j \in J} w_j \tag{5}$$

representing the weighted size of J . (Although, in our experiments, we use $w_j = 1$ for all j , in which case $\sigma(J)$ is just the number of elements in J , we present a weighted version of the problem, which can be useful in some situations.) We define a minimal covering as any covering minimizing σ among all other coverings.

To rephrase this as an integer programming problem, we note that the subsets J of \mathcal{I} are in one-to-one correspondence with the set of all configurations $\psi = (\psi_j, j \in \mathcal{I})$, where $\psi_j = 1$ if $j \in J$ and 0 otherwise. The minimal covering problem can then be reformulated as minimizing $\sum_{j \in \mathcal{I}} w_j \psi_j$ subject to the existence of a random variable $Y: \Omega \rightarrow \{0, 1\}$ such that $P(Y = 1) \geq 1 - \alpha$ and

$$\sum_{j \in \mathcal{I}} \psi_j Z_j(\omega) \geq rY(\omega).$$

We have a finite sample of the distribution P , represented by a finite subset $\hat{\Omega}$ of Ω . We can approximate the covering problem by enforcing the constraints only for $\omega \in \hat{\Omega}$ and replacing $P(Y = 1) \geq 1 - \alpha$ by a sample fraction over $\hat{\Omega}$. We then determine a minimal r -covering at level α by minimizing

$$F(\psi, Y) = \sum_{j \in \mathcal{I}} w_j \psi_j \tag{6}$$

subject to the constraints

$$\left\{ \begin{array}{l} \forall \omega \in \hat{\Omega} : \sum_{j \in \mathcal{I}} \psi_j Z_j(\omega) \geq rY(\omega) \\ \sum_{\omega \in \hat{\Omega}} Y(\omega) \geq |\hat{\Omega}|(1 - \alpha) \end{array} \right. \tag{7}$$

Many optimal solutions are usual for equal w_j . Assume that one obtains several such sets $J^{(1)}, \dots, J^{(N)}$, all with same cardinality, and all providing coverings at level α of the considered population. (While it may be computationally prohibitive to compute all solutions, it is often possible to collect a large number of them.) These sets can be combined in at least two obvious

ways, namely via their union

$$J_{\text{all}} = \bigcup_{i=1}^N J^{(i)} \quad (8)$$

or their intersection

$$J_{\text{core}} = \bigcap_{i=1}^N J^{(i)}. \quad (9)$$

The latter (while not being a covering by itself) we particularly focus on as it captures essential abnormalities observed. Other possibilities include maximizing the sum of probabilities that each element is aberrant.

Once a subset J of variables is chosen (a covering, or a core), we obtain a representation of each sample ω as a binary vector $\bar{z}_j = (Z_j(\omega), j \in J)$, which should retain essential information from the whole *omics* profile associated to the sample. It has, in addition, a mechanistic interpretation, since each variable Z_j is associated to one or a group of STPs ($g \Rightarrow g'$). Because of the relatively small number of variables involved, all these events can be rendered together, using, for example, the visualization provided in Figs 2 or 4.

1.4 Measuring heterogeneity

We want to quantify the heterogeneity of a family of binary random variables $\bar{Z} = (Z_j, j \in J)$, defined on the probability space Ω , where J is a subset of \mathcal{I} (e.g., a covering). Similarly to the previous section, we assume that only a finite number of observations are available, represented by a finite subset $\hat{\Omega}$ of Ω . A natural measure for heterogeneity is the Shannon entropy [21, 22], that we need to estimate based on the finite random sample $(\bar{Z}(\omega), \omega \in \hat{\Omega} \subset \Omega)$. In our results, we focus on the entropy of \bar{Z} conditional to a specific cancer condition, phenotype or subtype.

Let $\mathfrak{S} = \mathfrak{S}_J$ denote the set of all binary configurations $\bar{z} = (z_j, j \in J)$, which has $2^{|J|}$ elements. Let $\pi(\bar{z}) = P(\bar{Z} = \bar{z})$, so that the Shannon entropy of \bar{Z} is

$$H(\pi) = - \sum_{\bar{z} \in \mathfrak{S}} \pi(\bar{z}) \log_2 \pi(\bar{z}). \quad (10)$$

The sample probability mass function of \bar{Z} is then given by

$$\hat{\pi}(\bar{z}) = \frac{N(\bar{z})}{N} \quad (11)$$

where $N(\bar{z}) = |\{\omega \in \hat{\Omega} : \bar{Z}(\omega) = \bar{z}\}|$ and $N = |\hat{\Omega}|$. One can plug these relative frequencies in the definition of the entropy to obtain the estimator

$$H(\hat{\pi}) = - \sum_{\bar{z} \in \mathfrak{S}} \hat{\pi}(\bar{z}) \log_2 \hat{\pi}(\bar{z}). \quad (12)$$

This estimator, however, significantly under-estimates the entropy for small and even moderate sample sizes, and several bias-correction methods have been introduced in the literature (see [23], from which (13) (see below) is obtained, and [24, 25]). This estimator computes the

entropy using the expression (in which ψ denotes the digamma function)

$$\hat{H}_c(\bar{Z}) = \log_2(e) \sum_{z \in \mathfrak{S}} \frac{N(\bar{z})}{N} \left(\log N - \psi(N(\bar{z})) - \frac{(-1)^{N(\bar{z})}}{N(\bar{z})(N(\bar{z}) + 1)} \right). \tag{13}$$

Still, this estimator is only accurate when the number of variables, $|J|$, is small, because the ratio $(2^{|J|} - 1)/N$ is the first order term in the expansion of the entropy bias [23, 26] in powers of $1/N$. In our experiment, we use $L = 4$, that is we estimate the entropy for at most 4 variables together. For a more detailed analysis on the selection of $L = 4$, see [S1 Text](#) (Table B in [S1 Text](#)). Since sets J of interest are typically larger, we estimate an upper-bound to the entropy in the following way.

Given two random variables X and Y , one always has $H(X, Y) \leq H(X) + H(Y)$. This implies that, if the set J is partitioned into subsets J_1, \dots, J_ℓ (i.e., $J = \bigcup_{h=1}^\ell J_h$ and $J_h \cap J_{h'} = \emptyset$ if $h \neq h'$), then

$$H(\bar{Z}) \leq H(\bar{Z}_{J_1}) + \dots + H(\bar{Z}_{J_\ell}) \tag{14}$$

where $\bar{Z}_{J_h} = (Z_j, j \in J_h)$, $h = 1, \dots, \ell$. We use the right-hand side as an upper-bound, determining the partition J_1, \dots, J_ℓ using the following greedy aggregating procedure:

1. Initialize the partition with singletons, i.e., $J_j = \{j\}$, $j \in J$, computing the estimated entropy $\hat{H}_c(\bar{Z}_j)$ of the binary variable \bar{Z}_j . Fix a maximal subset size, L .
2. Given a current decomposition J_1, \dots, J_ℓ , compute, for all pairs h, h' such that $|J_h \cup J_{h'}| \leq L$, the difference $\hat{H}_c(\bar{Z}_{J_h}) + \hat{H}_c(\bar{Z}_{J_{h'}}) - \hat{H}_c(\bar{Z}_{J_h \cup J_{h'}})$, remove the two sets $J_h, J_{h'}$ for which this difference is largest and replace them by their union (setting $\ell \rightarrow \ell - 1$).
3. If no pair h, h' satisfies $|J_h \cup J_{h'}| \leq L$, stop the procedure.

The obtained decomposition also provides a statistical model (denoted $\hat{\pi}_*$) approximating the distribution of \bar{Z} , namely the one for which $\bar{Z}_{J_1}, \dots, \bar{Z}_{J_\ell}$ are independent and the distribution of \bar{Z}_{J_i} is estimated using relative frequencies. To allow for comparisons between entropies evaluated for different sub-populations, we used this model within a Monte-Carlo simulation to estimate confidence intervals for $H(\pi)$. We generated $M = 1,000$ new N -samples of \bar{Z} (recall that N is the size of the original sample of \bar{Z} used to estimate the entropy), using the distribution $\hat{\pi}_*$, resulting in M new empirical distributions $\hat{\pi}_*^{(1)}, \dots, \hat{\pi}_*^{(M)}$ with associated corrected entropies $\hat{H}_c^{(1)}, \dots, \hat{H}_c^{(M)}$. Fixing a probability $\beta > 0$, we let $\hat{H}_c^{(\beta)}$ and $\hat{H}_c^{(1-\beta)}$ denote the β and $1 - \beta$ quantiles of the sample $\hat{H}_c^{(1)}, \dots, \hat{H}_c^{(M)}$ so that $\hat{H}_c^{(j)} - H(\hat{\pi}_*)$ belongs to $[\hat{H}_c^{(\beta)} - H(\hat{\pi}_*), \hat{H}_c^{(1-\beta)} - H(\hat{\pi}_*)]$ with probability $1 - 2\beta$. We use the same interval for the difference $\hat{H}_c - H(\pi)$, yielding the confidence interval for $H(\pi)$:

$$[\hat{H}_c + H(\hat{\pi}_*) - \hat{H}_c^{(1-\beta)}, \hat{H}_c + H(\hat{\pi}_*) - \hat{H}_c^{(\beta)}]. \tag{15}$$

1.5 Subtype analysis through partitioning

We assume here again a family of binary random variables $\bar{Z}_j = (Z_j, j \in J)$, where J is a tissue-dependent covering, observed through a finite sample $(\bar{Z}_j(\omega), \omega \in \hat{\Omega})$. We partition the sample space into disjoint subsets (S_1, \dots, S_ℓ) where each S_j is specified by a small number of events involving conjunctions or disjunctions of aberrations. This partition will be associated with the terminal nodes of a binary “coding tree” of limited depth d (e.g., $d = 5$), so that $\ell = 2^d$. To

each node in the tree we associate a subset S of $\hat{\Omega}$, and unless S is a terminal node, its children form a partition $S = S' \cup S''$, where $\omega \in S'$ is based on a certain splitting criterion.

While there are many ways to build such a tree-structured code, we opt for a decomposition for the tissue population that is as balanced as possible and unsupervised to compare the distributions between cancer subtypes (for a given tissue) with respect to a fixed partition. A sample is weighted inversely proportional to the size of its subtype, and at each node the event which balances the weight of the two daughter nodes is selected.

The standard choice for a binary tree are individual binary features, so events of the form $\{Z_{j(S)} = 1\}$, for a suitably chosen $j(S) \in J$. One could also use more complex splitting criteria, such as $\{Z_{j_1(S)} = 1 \text{ or } Z_{j_2(S)} = 1\}$, $\{Z_{j_1(S)} = 1 \text{ and } Z_{j_2(S)} = 1\}$ with $j_1(S), j_2(S) \in J$. (We have used both types of events in our experiments: two-gene disjunctions for trees based on source aberrations with targets and two-gene conjunctions for trees based on target aberrations with sources.) The stopping criterion is that either all samples at the node have identical configurations or a maximum depth has been reached.

2 Results

As described previously, we have delineated many gene pairs (STP, or “source-target pair”), with a binary random variable corresponding to each pair indicating whether a sample is source and target aberrant.

Given M STPs, there are 2^M possible “states” or “configurations” for each sample. We defined *cross-sample heterogeneity* as the entropy of the probability distribution P over configurations. Estimating the entropy of P is not feasible for modest sample sizes, since it requires estimating the probabilities of many rare events.

To overcome this computational barrier, the pool of STPs was substantially reduced using the notion of a “minimal set covering” in combinatorial optimization. In our case, the set to be covered is a population of cancer samples for a particular phenotype or subtype, a “covering” is a set of STPs for which, with high probability, cancer samples have at least one aberrant STP from the covering, and “minimal” means the smallest covering. All minimal coverings are necessarily of the same size, on the order of 10–100 for each tissue we study (breast, colon, liver, kidney, lung and prostate). In summary, our STPs are derived from Reactome, and then subsets of STPs of interest are identified for each cancer type based on the TCGA omics data.

Minimal set coverings are typically not unique. However, despite the differences between the coverings we can define a “core”, namely, the STPs that appear in *every* (minimal) solution. From a biological perspective, the core is a novel signature of the most salient events associated with tumors of a given type. We apply these concepts (STPs, cores and estimated entropies) to measuring cross-sample heterogeneity in tumor populations for a selection of tissues represented in TCGA data.

2.1 Source-target pairs

Based on the genes and interactions found in the Reactome [19], source-target pairs (STPs)($g_s \Rightarrow g_t$) are built, with the only parameter used being the maximum length of the directed chain from the source to the target (k ; see see Table C in [S1 Text](#)). There are then 272,237 valid STPs with 3,124 distinct source genes and 598 distinct target genes.

Our samples are those in TCGA with available matched mutation, extreme copy number variation (deleting or amplifying both copies), and mRNA expression data; the conversion of expression counts to aberration states was described in Section 1.2.1.

Nearly all samples exhibit at least one paired aberration.

2.2 Filtering

Given there are too many STPs to meaningfully analyze, we first filter based on rejecting the hypothesis that the existence of source and target aberrations are independent (see Section 1.2.2). The statistics of the STPs remaining after this filtering procedure for different tissue types are shown in see (see Table D in [S1 Text](#)). For example, for the 953 TCGA breast cancer samples, there are 17,261 valid STPs after the test for independence, with 2,130 source genes and 421 target genes.

Next we omit very rare events. For each tissue, and at each of the three levels (source, target, pair), we require each binary variable to be aberrant in at least 2% of the samples for that tissue. For details about the choice of 2%, see [S1 Text](#). The number of qualifying variables after the 2% filter was applied are given in Table F in [S1 Text](#).

For instance, for breast cancer samples, there are 4,026 STPs, 690 distinct source genes, and 256 distinct target genes after separately applying the 2% filter at each level.

2.3 Paired aberrations

[Table 1](#) shows examples of STPs $\lambda = (g \Rightarrow g')$ and their associated probabilities of aberration in the indicated tissue. For example, in the STPs shown for colon cancer in [Table 1](#), *APC* is the source gene, *AXIN2* is the target gene, and there exists a directed signaling path from *APC* to *AXIN2* of length at most three links (two intermediate genes) in Reactome such that the second-to-last link, namely the direct parent of *AXIN2*, controls the mRNA expression of *AXIN2*. This STP is aberrant in a given sample if *APC* is either mutated, deleted, or amplified *and* the mRNA expression of *AXIN2* is aberrant (with respect to baseline mRNA expression for *AXIN2*). In the case of *APC*, the DNA aberration is nearly always a mutation and *AXIN2* is over-expressed. See [Table E](#) in [S1 Text](#) for more information.

In [Table 1](#), probability $P(\text{DNA}\&\text{RNA})$ is the sample estimate, namely the fraction of colon samples for which *APC* is mutated *and* *AXIN2* is RNA-aberrant. Similarly, $P(\text{DNA})$ and $P(\text{RNA})$ stand for the marginal probabilities that the source is DNA-aberrant and the target RNA-aberrant, respectively. The conditional probabilities are then self-explanatory. For example, *APC* is mutated in 73.9% of our samples and in 79.1% of those samples *AXIN2* is RNA-aberrant. Multiplying these two probabilities gives the frequency of the joint occurrence (58.5%). Other STPs commonly found in colon samples include the four core STPs described in [Table 2](#).

The probabilities for $APC \Rightarrow AXIN2$ are atypically large. In particular, most pair probabilities are smaller than .575, generally of order 0.01–0.10 with a few above 0.3, usually involving main tumor drivers such as *PIK3CA* in breast cancer, and *TP53* and *KRAS* in lung cancer. Moreover, DNA aberrations tend to be considerably rarer than RNA aberrations, *i.e.*, the marginal source probabilities are generally far smaller than the marginal target probabilities. It is

Table 1. Examples of STPs. For each of the six tissues, one example of a common STP $\lambda = (g \Rightarrow g')$ is shown. $P(\text{DNA}\&\text{RNA})$ is our sample-based estimate of the probability that λ is an aberrant pair, namely, the fraction of samples of the indicated tissue for which the source gene g is DNA-aberrant and the target gene g' is RNA-aberrant. Similarly, $P(\text{DNA})$ (respectively, $P(\text{RNA})$) is the fraction of samples for which g is DNA-aberrant (resp., g' is RNA-aberrant), and $P(\text{RNA}|\text{DNA})$ is the (estimated) conditional probability that g' is RNA-aberrant given g is DNA-aberrant.

Tissue	Pair	$P(\text{DNA}\&\text{RNA})$	$P(\text{DNA})$	$P(\text{RNA})$	$P(\text{RNA} \text{DNA})$	$P(\text{DNA} \text{RNA})$
Breast	<i>PIK3CA</i> \Rightarrow <i>S100B</i>	0.316	0.356	0.838	0.888	0.377
Colon	<i>APC</i> \Rightarrow <i>AXIN2</i>	0.585	0.739	0.676	0.791	0.864
Kidney	<i>VHL</i> \Rightarrow <i>CA9</i>	0.482	0.485	0.967	0.994	0.498
Liver	<i>TP53</i> \Rightarrow <i>MYBL2</i>	0.308	0.319	0.814	0.965	0.379
Lung	<i>TP53</i> \Rightarrow <i>TOP2A</i>	0.529	0.535	0.923	0.988	0.573
Prostate	<i>PTEN</i> \Rightarrow <i>TWIST1</i>	0.161	0.216	0.654	0.745	0.246

<https://doi.org/10.1371/journal.pcbi.1008944.t001>

Table 2. Colon core STPs. There are four “core” STPs which appear in every minimal covering of the colon samples. $P(\text{DNA}\&\text{RNA})$ is the fraction of samples for which the source gene g is DNA-aberrant and target gene g' is RNA-aberrant; $P(\text{DNA})$ is the fraction of samples satisfying the source gene g is DNA-aberrant; $P(\text{RNA})$ is the fraction of samples with g' RNA-aberrant; $P(\text{RNA}|\text{DNA})$ is the fraction of DNA-aberrant samples for which g' is RNA-aberrant.

	$P(\text{DNA}\&\text{RNA})$	$P(\text{DNA})$	$P(\text{RNA})$	$P(\text{RNA} \text{DNA})$	$P(\text{DNA} \text{RNA})$
$APC \Rightarrow AXIN2$	0.585	0.739	0.676	0.791	0.864
$TP53 \Rightarrow PTPN12$	0.396	0.560	0.604	0.707	0.656
$PIK3CA \Rightarrow TNFRSF10B$	0.198	0.271	0.589	0.732	0.336
$MAML1 \Rightarrow PBX1$	0.034	0.039	0.401	0.875	0.084

<https://doi.org/10.1371/journal.pcbi.1008944.t002>

noteworthy that the conditional probability of a particular RNA aberration given a particular DNA aberration (as those in Table 1) is usually in the range 0.5–1, whereas the reverse is not the case: given a target gene is RNA-aberrant the probability of any particular gene serving as a source rarely exceeds 0.2 (see Tables G–K in S1 Text).

We have also defined separate source-level and target-level events in the sense of partially aberrant STPs; see Section 1.2.3 of Methods. Recall that $\{Z_g^{(s)} = 1\}$ represents the event that a given source gene g is DNA-aberrant and that there exists some target of g which is RNA-aberrant, denoted as “aberration with target” for “source aberration with downstream target aberration”. The probability of this event is denoted by $P(\text{DNA}\&\text{downstreamRNA})$; see Table 3 for some examples in Colon. Similarly, for the other direction, $\{Z_{g'}^{(t)} = 1\}$ is the event that some source gene renders $g \Rightarrow g'$ an aberrant STP. Tables 3 and 4 provide the probabilities and conditional probabilities for selected core genes at the source and target levels in colon; many other examples appear in Tables G–U in S1 Text.

Given a source gene g is aberrant, typically there is a strong likelihood that at least one of its targets g' is RNA-aberrant. These targets represent plausible downstream consequences of g

Table 3. Colon core source genes. There are five “core” source genes which appear in every minimal source covering of the colon samples. $P(\text{DNA})$ is the fraction of samples for which the indicated source gene is DNA-aberrant; $P(\text{DNA}\&\text{downstreamRNA})$ is the fraction of samples for which the indicated source gene is DNA-aberrant and there exists an RNA-aberrant gene among its targets. $P(\text{downstreamRNA}|\text{DNA})$ is the fraction of the samples with the indicated source gene DNA-aberrant for which there exists some RNA-aberrant gene among its targets.

	$P(\text{DNA}\&\text{downstreamRNA})$	$P(\text{DNA})$	$P(\text{downstreamRNA} \text{DNA})$
<i>APC</i>	0.585	0.739	0.791
<i>TP53</i>	0.560	0.560	1.000
<i>KRAS</i>	0.425	0.425	1.000
<i>LAMA5</i>	0.217	0.217	1.000
<i>MAML1</i>	0.034	0.039	0.875

<https://doi.org/10.1371/journal.pcbi.1008944.t003>

Table 4. Colon core target genes. There are six “core” target genes which appear in every minimal target covering of the colon samples. $P(\text{RNA})$ is the fraction of samples for which the indicated target gene is RNA-aberrant; $P(\text{RNA}\&\text{upstreamDNA})$ is the fraction of samples for which the indicated target gene is RNA-aberrant and there exists an DNA-aberrant gene among its sources. $P(\text{upstreamDNA}|\text{RNA})$ is the fraction of the samples with the indicated gene RNA-aberrant for which at least one of its sources is DNA-aberrant.

Target	$P(\text{RNA}\&\text{upstreamDNA})$	$P(\text{RNA})$	$P(\text{upstreamDNA} \text{RNA})$
<i>PERP</i>	0.710	0.807	0.880
<i>PDX1</i>	0.671	0.957	0.702
<i>AXIN2</i>	0.662	0.676	0.979
<i>SALL4</i>	0.638	0.918	0.695
<i>TNFRSF10B</i>	0.565	0.589	0.959
<i>MYBL2</i>	0.261	0.300	0.871

<https://doi.org/10.1371/journal.pcbi.1008944.t004>

having a DNA-aberration. The converse, however, is not valid; in particular, there are many targets g' for which there is no upstream DNA-aberrant source linked to g' . This makes sense since a gene can be RNA-aberrant for many reasons other than an upstream genetic aberration. In particular, the event driving the aberration of g' might be some perturbation not considered here, for example be epigenetic or fusion-related or as yet unrecognized.

2.4 Coverings

Recall that indexing a covering by source genes refers to leaving the particular aberrant target gene unspecified (indexing by targets is the opposite). The corresponding events were denoted in Methods by $\{Z_g^{(s)} = 1\}$ for a source gene and $\{Z_g^{(t)} = 1\}$ for a target gene.

As described in Section 1.3, minimal coverings composed of pairs, sources, or targets are all found with the same optimization program. For the pair and source levels, we calculate the optimal covering with the smallest possible $\alpha \geq 0$ and $r = 1$. (Here, the smallest α is such $1 - \alpha$ is the fraction of samples that have at least r aberrant STPs.)

At the target level, however, we select but $r = 3$, still using the smallest possible α ; that is, we attempt to cover tumor samples with at least three target aberrations (with source). This choice is justified by the higher frequency of RNA-aberrations in tumor samples.

Table 5 shows the optimal covering statistics at all levels for the six tissues of origin. For example, the minimum number of STPs (resp., sources, targets) necessary to cover the 953 breast cancer samples is 67 (resp., 60, 53), with the realized rate being 95% (resp., 96%, 96%) with $\alpha = 0.05$ (resp., 0.04, 0.04). In contrast, all colon samples can be covered with many fewer STPs, namely 11. In addition, the minimal covering size (size of solution) is usually largely determined by the incidence of aberrations in any given population, e.g., mutation rates. In particular, given two phenotypes A and B, if the samples of B are consistently more aberrant

Table 5. Statistics of optimal coverings. For each of the six tissues, this table provides basic information about the optimal coverings at all levels: STP, source with target, target with source. For instance, for breast cancer, there are 4,026 candidate STPs after both filters (rejecting source-target independence and 2% tissue sample frequency); the minimal covering size is 67 STPs; at least one of these 67 STPs is aberrant in 95.4% of the breast cancer samples; and there are 21 STPs which appear in every minimal covering.

Tissue	Samples	Covering Type	Quantity	Size of solution	Fraction of samples covered	Size of core set
Breast	953	STP	4,026	67	0.954	21
		Source	690	60	0.964	34
		Target	256	53	0.955	35
Colon	207	STP	1,195	11	1.000	4
		Source	525	10	1.000	5
		Target	226	15	0.995	6
Kidney	336	STP	347	26	0.827	12
		Source	133	28	0.854	21
		Target	176	60	0.890	45
Liver	360	STP	1,198	32	0.931	11
		Source	460	34	0.958	20
		Target	287	41	0.942	26
Lung	465	STP	3,154	27	0.985	10
		Source	908	25	0.989	19
		Target	350	29	0.985	26
Prostate	491	STP	430	53	0.686	32
		Source	211	53	0.743	42
		Target	160	72	0.699	66

<https://doi.org/10.1371/journal.pcbi.1008944.t005>

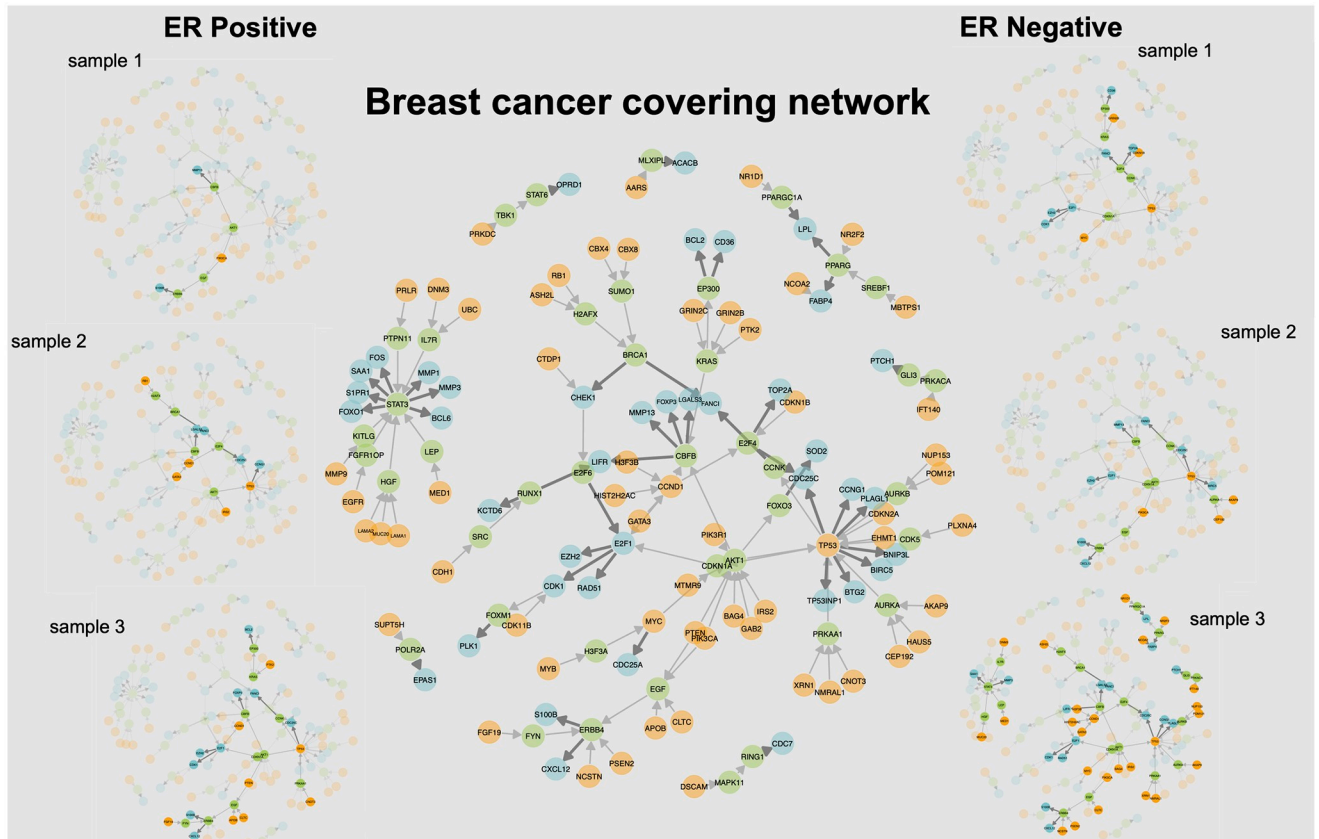


Fig 2. Networks of pair coverings in breast cancer. The network shown in the center depicts one covering of breast cancer samples by STPs, with source genes in orange, target genes in blue, and intermediary link genes in green. The thin and thick edges represent, respectively, the two types of relationships: “controls state change of” and “controls expression of” as designated in Reactome [19]. On the left are presented a selection of covering realizations for three ER-positive samples, where aberrant STPs are highlighted, while and on the right, three ER-negative samples are shown. The samples have different realizations over the covering network, and are ranked (top to bottom) by the number of events they exhibit. The sample networks demonstrate the inter-sample heterogeneity among the source and target realizations.

<https://doi.org/10.1371/journal.pcbi.1008944.g002>

than those of A, then the minimal B covering will be smaller. More comprehensive statistics for all tissues are given in Table F in [S1 Text](#); as seen in the column labeled “No. of solutions”, there are in general a great many instances of minimal coverings.

[Fig 2](#) shows one such tissue-level covering obtained for breast cancer as a graphical network with nodes representing genes forming the STPs. The source and target genes are shown in orange and blue respectively, while genes representing the intermediary links are shown in green. Note that while the union of the coverings may also be visualized in a similar manner, it contains many more STPs that make readability of the resulting graph difficult; therefore we have opted to show only individual coverings here. [Figs C—G in S1 Text](#) depict the networks associated with the coverings obtained for the other types of cancer.

These visual representations allow us to go beyond lists of names and numbers and begin to interpret coverings in biological terms and incorporate mechanism (see [Discussion](#)). For instance, in the breast network shown in [Fig 2](#), several important breast cancer genes (e.g., *STAT3* [27], *TP53* [28], *BRCA1* [29], and *ERBB4* [30]) all form important hubs through which multiple sources and targets in the covering link according to Reactome. Similarly, the network figures for the remaining networks show similar positioning for many important cancer genes: *NOTCH1* [31] in liver, *NOTCH3* [32, 33] and *EGFR* [34] in lung are some other

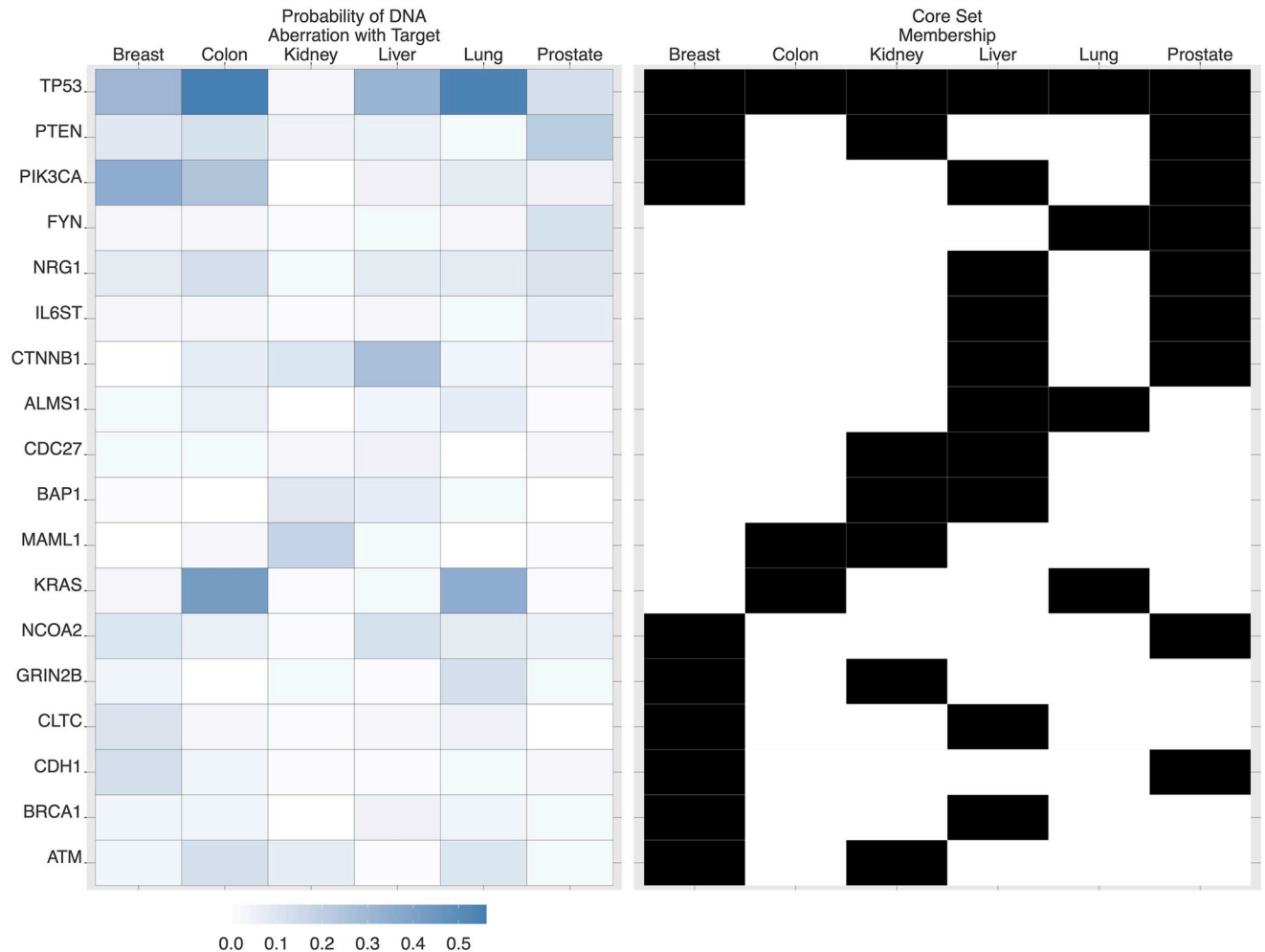


Fig 3. Core set across tissues at source level. There are 18 source genes which appear in the core set of at least two tissues. For instance, gene *TP53* is a core gene for all six tissues, and genes *PTEN* and *PIK3CA* are core genes for three tissues. The color in the heatmap on the left represents the probability that the corresponding source gene is DNA-aberrant and there exists an RNA-aberrant target gene (thereby forming an aberrant *source-target pair*). On the right, black marks indicate the membership of each gene to the corresponding core set for each tumor type.

<https://doi.org/10.1371/journal.pcbi.1008944.g003>

examples. Finally, *KRAS*, *TP53* [28] and *STAT3* [27] make an appearance in multiple cancers. See Figs C–G in [S1 Text](#).

For a given tissue and fixed covering level, about 30%–60% of the genes appearing in any covering in fact appear in all coverings, referred to as the *core set* (see [Tables 2–4](#) for colon, and in [Tables G–U](#) in [S1 Text](#) for the other tissues). The STP $TP53 \Rightarrow PTPN12$ is aberrant in 39.6% of the colon samples (see [Table 2](#)), the source gene *KRAS* in 42.5% of samples (see [Table 3](#)), and the target gene *PERP* in 80.7% of samples (see [Table 4](#)). From [Table 3](#) we see that there is some aberrant target for *every sample* for which *KRAS* is DNA-aberrant in colon; hence the probability that *KRAS* is DNA-aberrant *and* there is a matching target gene is again 42.5%. Finally, from the target covering we see that targets gene *PDX1* is RNA-aberrant in 95.7% of colon samples (see [Table 4](#)) but only 70.2% of samples for which *PDX1* is RNA-aberrant have some corresponding corresponding upstream DNA-aberrant source gene.

[Fig 3](#) shows 18 core source genes across multiple tissues. *TP53* is a core source gene shared by all 6 tissues, and is DNA-aberrant in more than 60% of colon cancers, and also a large

Table 6. Probabilities of source aberration with downstream target for breast cancer subtypes. For PAM50 subtypes of breast cancer, the heatmap represents the probabilities that the indicated gene is a DNA-aberrant source gene with some downstream RNA-aberrant target. The sources are selected from the set of core genes for coverings of the given tissue; the selection criterion is that the probability of a DNA-aberration is high for at least one of the subtypes for that tissue. Core sources with varying probabilities present interesting candidates for discrimination between subtypes. For example, the DNA-aberration frequency of *TP53* is much higher in the HER2-enriched and Basal-like subtypes than in Luminal A and Luminal B, whereas an aberration in *PIK3CA* is less frequent among basal-like samples than among the other subtypes.

	Luminal A	Luminal B	HER2-enriched	Basal-like
<i>TP53</i>	0.123	0.273	0.655	0.770
<i>MED1</i>	0.055	0.116	0.545	0.046
<i>PIK3CA</i>	0.438	0.314	0.418	0.172
<i>CLTC</i>	0.068	0.190	0.255	0.057
<i>PTEN</i>	0.050	0.074	0.073	0.253
<i>PTK2</i>	0.105	0.190	0.182	0.310
<i>GATA3</i>	0.151	0.190	0.055	0.253
<i>NCSTN</i>	0.068	0.132	0.127	0.264
<i>NCOA2</i>	0.064	0.149	0.218	0.126

<https://doi.org/10.1371/journal.pcbi.1008944.t006>

percentage in other cancers. All core source genes across multiple tissues are shown in Fig L in [S1 Text](#). In Fig M in [S1 Text](#), we show all core target genes across multiple tissues. For instance, *CDC25C* appears in core set of breast, lung and prostate cancer, and the probability that *CDC25* is RNA aberrant and there exist an upstream aberrant source is nearly 0.8.

2.5 Subtype coverings

Having computed the tissue-level coverings, we examine them with respect to certain phenotypes of interest, including the PAM50 subtypes in breast cancer [17], smoking history in lung cancer [35], Gleason grade in prostate cancer [16], and and the CRIS-classes in colon cancer [18]. We observe a large range of aberration frequencies among subtypes. [Table 6](#) shows the probabilities of DNA-aberration (with targets) for PAM50 subtypes, with genes selected from the core set of source breast cancer coverings; [Tables 7–9](#) show similar selections of sources.

We observe potential discriminating *sources* between subtypes. For example, *TP53* has a much lower likelihood of aberration for the luminal subtypes in comparison to the basal-like and HER2-enriched subtypes. Similar observations can be made among the subtypes of other cancers (see [Tables 7–9](#)). Finally, such patterns persist for target-level analyses and are presented in [Tables V–I](#) in [S1 Text](#).

A comparison between subtypes can also be captured as a graphical network, as shown in [Fig 4](#). Similarly, [Fig H](#) in [S1 Text](#) presents the breast covering with the size of the nodes representing the source (with target) and target (with source) aberration probabilities for the molecular subtypes considered. Finally, similar networks are also presented in [Fig I](#) in [S1 Text](#) for lung cancer with respect to smoking history, and in [Fig J](#) in [S1 Text](#) for Gleason grade in for prostate cancer.

We also compared coverings of subtypes controlling for population sizes. For each of the phenotypes under a given sub-typing, an equal number of samples were selected and coverings for all these samples simultaneously were obtained. Then we examined the proportion from each subtype that was covered, repeating over multiple sampling iterations ([Fig 5](#)). A general pattern of more pathological phenotypes having higher coverage proportions can be observed throughout these results (see [Fig K](#) in [S1 Text](#) for further results). The more malignant phenotypes tend to have larger aberration probabilities. This corresponds with the observation that the size of the covering obtained for a subtype while sampling equal numbers from each group

Table 7. Probabilities of source aberration with downstream target for lung cancer subtypes. For smoking history based categories of lung cancer, the heatmap represents the probabilities that the indicated gene is a DNA-aberrant source gene with some downstream RNA-aberrant target. The sources are selected from the set of core genes for coverings of the given tissue; the selection criterion is that the probability of a DNA-aberration is high for at least one of the subtypes for that tissue. *TP53* and *KRAS* are both more frequently DNA-aberrant (with some downstream RNA-aberrant target) among smokers than non-smokers whereas *EGFR* is a more aberrant source among non-smokers.

	Smoker	Recently reformed	Reformed	Non smoker
<i>TP53</i>	0.581	0.452	0.500	0.259
<i>EGFR</i>	0.093	0.129	0.233	0.370
<i>KRAS</i>	0.395	0.371	0.350	0.148
<i>ANK2</i>	0.140	0.274	0.117	0.037
<i>STK11</i>	0.140	0.290	0.183	0.074
<i>SPTA1</i>	0.372	0.355	0.317	0.185

<https://doi.org/10.1371/journal.pcbi.1008944.t007>

Table 8. Probabilities of source aberration with downstream target for colon cancer subtypes. For CRIS-class subtypes of colon cancer, the heatmap represents the probabilities that the indicated gene is a DNA-aberrant source gene with some downstream RNA-aberrant target. The sources are selected from the set of core genes for coverings of the given tissue; the selection criterion is that the probability of a DNA-aberration is high for at least one of the subtypes for that tissue.

	CRIS-A	CRIS-B	CRIS-C	CRIS-D	CRIS-E
<i>APC</i>	0.348	0.217	0.686	0.750	0.710
<i>TP53</i>	0.261	0.478	0.771	0.571	0.774
<i>KRAS</i>	0.500	0.478	0.057	0.393	0.581
<i>LAMA5</i>	0.370	0.261	0.229	0.107	0.194

<https://doi.org/10.1371/journal.pcbi.1008944.t008>

indicated larger covering solutions obtained for more benign subtypes in comparison to more malignant ones.

2.6 Measuring heterogeneity

We applied the approach described in Section 1.4 to assess the relative heterogeneity of different cancer phenotypes and subtypes in the analyzed tissues. Note that we are analyzing heterogeneity within phenotypes at a population level, so that our measurements are primarily about the variability across tumors in this population. Without single-cell data, one cannot evaluate variability within a tumor, although it is likely that a higher variability at this level should also trigger larger heterogeneity at the population level. It is worth noting, however, that the analytical framework we propose here can be easily extended to single cell data once paired molecular measurements will become available in the future.

Table 9. Probabilities of source aberration with downstream target for prostate cancer subtypes. For primary gleason grade subtypes of prostate cancer, the heatmap represents the probabilities that the indicated gene is a DNA-aberrant source gene with some downstream RNA-aberrant target. The sources are selected from the set of core genes for coverings of the given tissue; the selection criterion is that the probability of a DNA-aberration is high for at least one of the subtypes for that tissue.

	3	4	5
<i>TP53</i>	0.071	0.143	0.286
<i>MYC</i>	0.036	0.094	0.245
<i>PTEN</i>	0.133	0.245	0.347
<i>ZFH3</i>	0.056	0.102	0.184
<i>FGF17</i>	0.102	0.139	0.163

<https://doi.org/10.1371/journal.pcbi.1008944.t009>

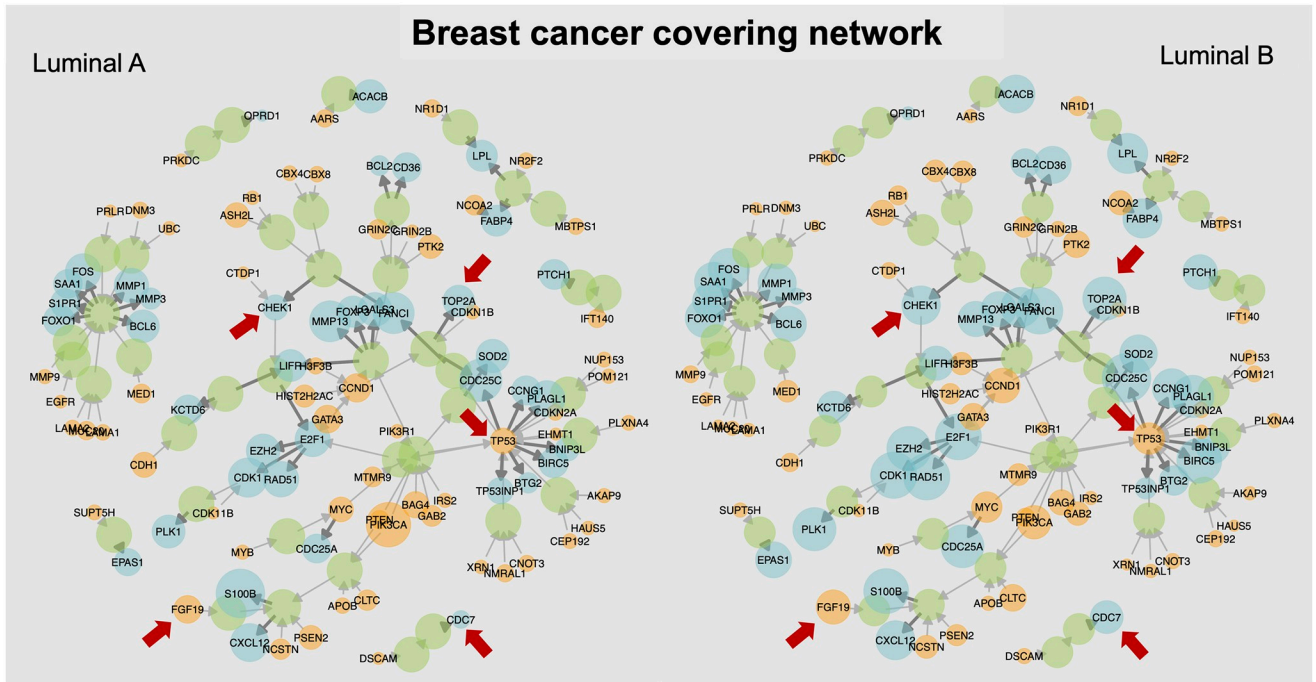


Fig 4. Comparison of one covering network for luminal breast cancer subtypes. The probabilities of DNA aberration (with targets) and RNA aberration (with sources) over the Luminal A and Luminal B populations of breast cancer are depicted by the size of each node in the network, which corresponds to one possible covering. The red arrows indicate some sources and target genes that have noticeable differences in the respective probabilities between the two luminal subtypes (e.g., *TP53*, *CHEK1*, *PIK3CA*, and *TOP2A*, also see Table 6).

<https://doi.org/10.1371/journal.pcbi.1008944.g004>

We base our analysis on coverings estimated on source aberration with targets and on target aberration with sources (see Tables II and III in S1 Text). In all cases, coverings are obtained for each tissue of origin separately, and entropy estimates are computed after restricting the data to samples exhibiting each considered cancer phenotype (e.g., the breast cancer molecular subtypes, smoking status in lung cancer, and so on. .).

At the source level, the general trend is that heterogeneity estimates increase with increasing disease severity. In prostate cancer, for instance, entropy grows with Gleason sum, primary Gleason grade, and with tumor status, while no clear ordering is observable for lymph node status (see Table 10). Similar observations can be also made for tumors originating in other tissues (see Table II in S1 Text). In breast cancer, the entropy for ER positive tumors is less than that for ER negative ones, and it also increases with tumor size, and with more aggressive molecular subtypes (i.e., Luminal A < Luminal B < HER 2 < Basal, with a small overlap between confidence intervals for Luminal B and HER 2). For lung, samples from patients with recent smoking history (reformed for less than 15 years or current smokers) have a higher entropy than those with either ancient or no history.

At the target level, a similar trend of increasing heterogeneity with increasing disease severity is observed: in prostate cancer for all variables considered (Gleason sum, primary Gleason grade, tumor stage, size, and lymph node status), in kidney cancer for tumor stage and tumor size, and in breast cancer for the molecular subtypes (with Luminal A samples exhibiting the lowest heterogeneity while Luminal B the highest). Finally, in the tumor types originating in the other tissues, we observe large overlaps between confidence intervals, and no obvious and clear trends emerged across cancer subtypes. Complete summaries for this analysis can be found in Table III in S1 Text.

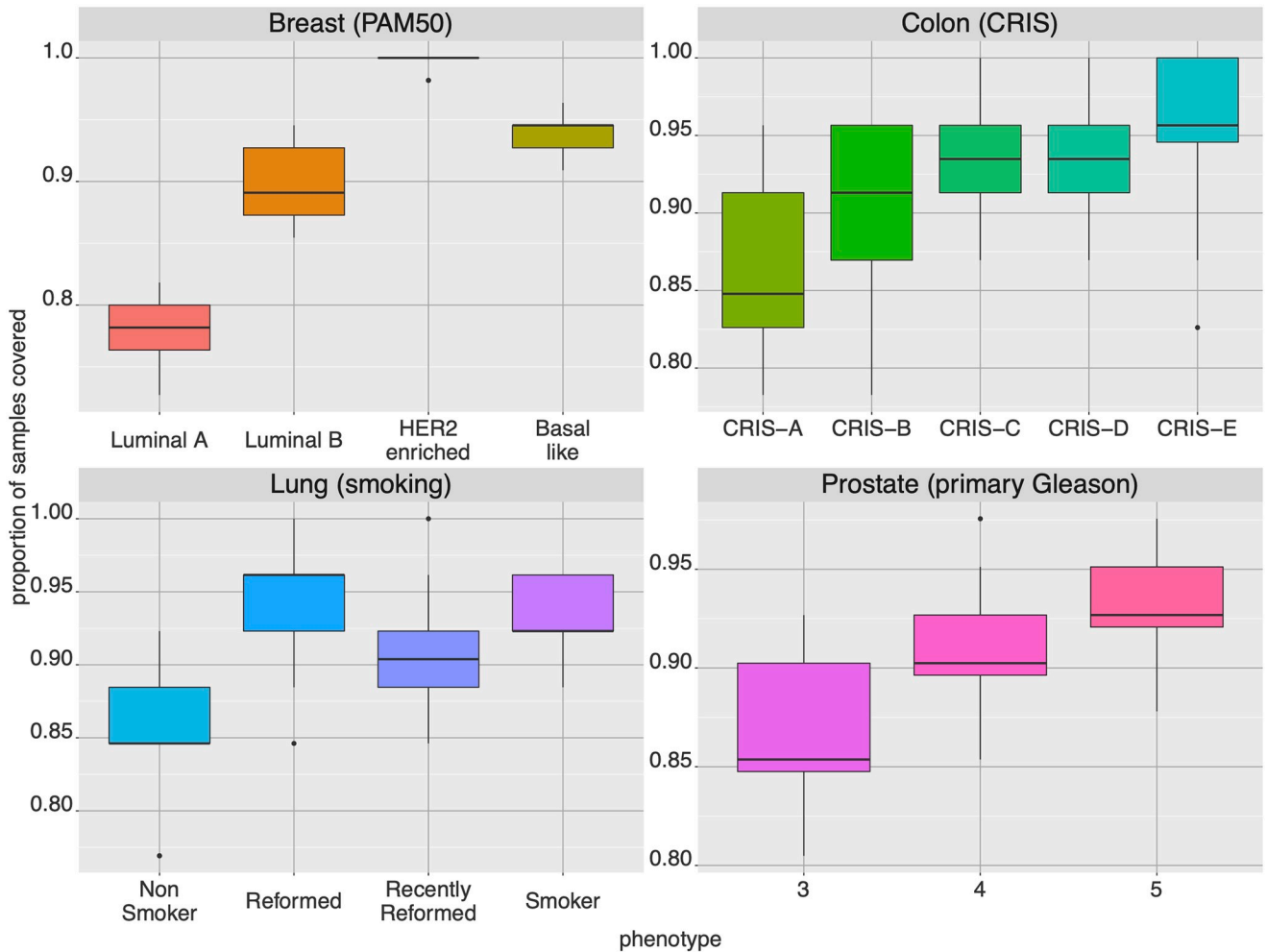


Fig 5. Rates of covering assembly. For each of four tissues (breast, colon, lung and prostate), several phenotypes are compared based on the proportion of samples actually covered when requesting 90% coverage or more for the given tissue by the optimization procedure. The boxplots represent the results of 20 iterations of normalizing for sample size among the phenotypes by random sampling. In general, coverings for more aggressive phenotypes assemble faster.

<https://doi.org/10.1371/journal.pcbi.1008944.g005>

2.7 Partitioning

We applied the approach described in Section 1.5 to all six cancer types. Let T denote a coding tree with terminal nodes $\{t_1, t_2, \dots, t_j\}$. Recall that for each subtype, the resulting histogram (number of samples per bin) is, by construction, as balanced as possible for the whole population. It is then easy to visualize the histograms conditional on tumor sub-populations (samples for a given subtype), and assess differences across subtypes or phenotypes in this representation (See Fig 6, for ER status in breast cancer based on target aberration with sources).

For instance, the event chosen at the root node is whether both *MMP1* and *E2F1* are RNA aberrant, each with some upstream DNA aberrant source. Such samples take the left branch whereas all others go to the right. The number in the terminal node t is the overall number of tumor samples reaching t (regardless of subtypes). The two histograms show the numbers of samples collected at terminal nodes for the ER positive and the ER negative sub-populations. In particular, the two distributions are significantly different (permutation test $p \leq 0.0001$). Other coding trees for other cancer subtypes are showed in Figs N–P in S1 Text.

Table 10. Entropy estimation at source level. Entropy estimates for source aberrations with target for prostate Gleason sum, primary Gleason grade, tumor status, and lymph-node status. *N* is the total number of samples available in the given subtype.

Subtype	Value	N	Entropy	Conf. Interval
All		491	12.16	[11.67, 12.68]
Gleason sum	6	45	8.05	[6.82, 8.98]
Gleason sum	7	244	10.80	[10.10, 11.45]
Gleason sum	8	63	11.33	[10.07, 12.41]
Gleason sum	9	135	14.21	[13.24, 15.06]
Primary Gleason grade	3	196	9.45	[8.68, 10.12]
Primary Gleason grade	4	245	12.33	[11.65, 13.01]
Primary Gleason grade	5	49	16.60	[15.07, 17.88]
Tumor Status	T2	186	10.73	[9.92, 11.45]
Tumor Status	T3-T4	298	12.88	[12.26, 13.52]
Lymph Node Status	Negative	342	12.14	[11.55, 12.71]
Lymph Node Status	Positive	77	12.82	[11.66, 13.76]

<https://doi.org/10.1371/journal.pcbi.1008944.t010>

3 Discussion

Our approach in introducing this framework is rooted in the biological tenant that cancer is driven by genetic and genomic alterations that alter normal cell behavior through the modification of downstream gene expression programs governing cell proliferation, cell identity, and cell differentiation. The representation of a cancer profile is binary and integrative, centered on STPs which encode paired DNA and RNA aberrations. The set of possible pairings is fixed, based on signaling pathways and regulatory network topology; the STP is aberrant in an individual profile if the first, or “source”, gene carries the specified DNA aberration and the gene expression of the second, or “target”, gene diverges from the normal baseline. Note that we are not assuming that the DNA alterations are decidedly always “drivers” of cancer or that the source-target links are necessarily causative; rather, given an aberrant STP, we regard the aberrant target as a putative effect of the upstream source aberration, which itself provides a putative explanation for the aberration of the target.

Whereas we do not deal directly with the driver-passenger distinction, *i.e.*, we treat all mutations in the same way, most of the genes that emerge from our analysis are known drivers, particularly the core covering. This is further expected given our requirement that candidate DNA aberrations have at least 2% incidence in the cancer population. That said, given the heterogeneous nature of cancer as a disease, it is likely that the list of currently accepted, known cancer drivers is not exhaustive. In this perspective, our approach—requiring additional constraints for defining aberrant STPs—might also be helpful in prioritizing previously unknown driver events from the large set of potential candidates.

We apply STPs, together with integer programming, to extract parsimonious sets of important aberrations in the tumor populations. In addition to minimal coverings with paired aberrations, the same algorithm can be applied to find the minimal coverings by DNA aberrations alone for which there is some downstream effect (the source gene appears in some aberrant STPs), and vice-versa for coverings by RNA-aberrant target genes (*i.e.*, gene expression alterations plausibly associated with some upstream DNA aberrations).

UNCOVER by [36] utilizes a similar computational method, attempting to distinguish driver and passenger somatic event sets with a high degree of mutual exclusivity between alterations. As posed in [36], the problem is NP-hard [37] and only small sets of genes, usually of

3.1 Cores

In general the minimal coverings are not unique; of particular interest are the core STPs which are those belonging to *every* minimal solution. The restriction to coverings, especially to cores, massively reduces the number of considered aberrations in a given tumor population, making it mathematically feasible to quantify and measure tumor heterogeneity at the population level in a natural, information-theoretic way.

Overall, we were able to identify well-known cancer aberrations, as well as to uncover novel potential molecular circuits involved in this disease. The presence of tumor suppressor *TP53* in all cores confirm the well-known notion that this DNA alteration is the most frequent and (possibly) important aberration across multiple tumor types of different lineage [43]. Other aberrations involving key cancer drivers, are discussed in our supplementary results section (see [S1 Text](#)).

No single gene appears in the core target set of all tissue types. In fact, only one gene, *FABP4*, lies in at least four out of six tissue-specific core signatures. Irrespective of source, target, or pair level, all core genes belong to signaling pathways commonly disrupted in the analyzed cancer types. For instance, paradigmatic examples of affected cancer pathways emerging from our analyses are *Ras* and *Wnt* in colon cancer and the *PI3K* and *mTor* pathways in breast and prostate cancers.

More unexpected and presumably novel genes of interest include for instance, *GRIN2B*—a gene encoding for a subunit of a N-methyl-D-aspartate (NMDA) receptor family member—is a source-level core gene for both breast and kidney cancer. Despite a relatively low incidence of DNA aberration (2.1% and 3.6% respectively in kidney and breast cancer), this gene was always associated with the divergent expression of a downstream RNA in both tumor types. Our findings, along with the previously reported promoter hyper-methylation observed in gastric [44], esophageal [45], and lung cancer [46], collectively suggest this gene might play a role as tumor suppressor.

Another interesting example is *PTK6* which encodes a cytoplasmic protein kinase also known as breast cancer kinase. This is a core target gene in breast, lung, and prostate cancer, with high probabilities of RNA aberration and upstream DNA aberration in breast and lung (17.2% and 50.1% respectively).

Possibly, effects of different upstream DNA alterations can propagate and converge on downstream targets to explain their aberrant expression. Thus *PTK6* could represent a suitably “unifying” target for treatment, despite the heterogeneous set of mutations observed in the patient population. Interestingly, inhibition of *PTK6* has been proposed for treatment in triple negative breast cancer [47] and *PTEN*-null prostate cancer [48].

In addition, our analyses also point to specific interaction pairs, further underscoring the importance of adopting a network view that goes beyond “hubs”, individual genes, and known cancer driver, when interpreting the core sets. [Table 2](#) include a number of pairs for colon cancer that can be directly mapped to specific signaling pathways. The *APC* \Rightarrow *AXIN2* pair participates into Wnt signaling, while the *MAML1* \Rightarrow *PBX1* pair is part of the Notch3 signaling network. Both these pathways are known to regulate the homeostasis of the colonic epithelium, and their alterations are well documented in colon cancer [49].

Finally, interesting differences between cancer subtypes and phenotypes emerged when we analyzed STP coverings and tumor heterogeneity at population level. For instance, in lung cancer, source-level paired aberrations involving *KRAS* were most strongly associated with smoking, whereas those involving *EGRF* showed an opposite trend, consistent with well-established patterns [50]. Similarly, *KRAS* aberrations were virtually absent in the CRIS-C colon cancer subtype, which was in turn enriched for aberrations involving *TP53*, as previously described

[18]. *CHEK1* in breast cancer is also a notable regulator of the response to DNA damage, which is over-expressed in triple-negative breast cancer (TNBC) and has therefore been proposed as a potential target for treatment [51]. Notably, we were not only able to confirm *CHEK1* aberration (with over-expression) in basal-like tumors (which are enriched for TNBC), but also reveal this aberration in the luminal B subtype, suggesting a possible vulnerability of this more aggressive type of breast cancer.

3.2 Heterogeneity

Regarding significant differences in the computed entropy (of the joint distribution of aberrations) between tumor groups, larger entropy estimates were typically associated with more severe disease phenotypes. This suggests that there is more diversity or variation in the DNA and RNA profiles of sub-populations of patients with more aggressive disease phenotypes. Such heterogeneity observed at the population level probably reflects the variability present at the individual level—*i.e.*, the intra-tumoral heterogeneity, stemming from genetic, epigenetic, and cellular variation—which is a well-known factor impacting clinical outcome and therapy response [52, 53]. Also notably, aggressive phenotypes are covered more efficiently than are less aggressive ones. This may be reflective of more aggressive phenotypes accumulating many more aberrations over time.

Large inter-patient variations in the genetic aberration profiles have been reported for individuals with the same diagnosis. Due to this heterogeneity combined with small sample sizes, it is difficult to associate specific changes in gene expression with specific aberrations in cancer genomes. In a recent pan-cancer study [11], the authors use matched whole-genome DNA and RNA sequencing data for 1,188 patients in order to identify co-occurrent DNA and RNA aberrations focusing on fusions, copy number changes, and mutation-driven aberrant splicing, followed by putative causal or mechanistic explanations. In contrast we leverage mechanistic constraints via prior, independent, biological information.

One finding with some commonality is the significant correlation between DNA and RNA alterations, observed in [11]. Results on our end can be seen from the values of $P(\text{upstreamDNA}|\text{RNA})$ (see Table 4 and Tables Q—U in S1 Text). In our framework, STPs can and usually do involve different partners and therefore our probabilities are not strictly comparable to those of PCAWG. Thus for us, “RNA” in $P(\text{upstreamDNA}|\text{RNA})$ refers to RNA-aberration in a fixed target gene and “upstream DNA” means that the target gene is linked (forms an STP) with *some* DNA-aberrant gene. Therefore a larger set of explanations is available for a given RNA aberration.

A multivariate statistical approach is seen in [10], where the authors first predict gene expression from (phospho)protein expression and gene-specific transcription factor (TF) binding sites using affinity regression, then predict TF and protein activities from somatic changes. Biological analysis centers on specific genes and pathways, notably the dysregulating effect on TFs of activating mutations in the *PIK3CA* pathway. This pathway also emerges as pivotal in our results: indeed, the most common STPs in breast cancer are $PIK3CA \Rightarrow S100B$ and $PIK3CA \Rightarrow MMP13$ (see Table 1); *PIK3CA* is one of only three core source genes appearing in at least three tissues (see Fig 3); and in breast cancer, *PIK3CA* is virtually certain to have a downstream RNA aberrant target (see Table L in S1 Text). Whereas the methods here and in [10] are largely non-overlapping, the spotlight falls on many of the same DNA-RNA associations.

In [54], sample-specific Bayesian inference is applied to somatic genomic alteration and differential gene expression data to identify driver genes. Characterizations of cancer types are then obtained by summarizing the discovered relationships at the sample and the population

levels. Whereas this program bears similarities with ours the objectives and methodology are quite different: our analysis is top-down, based on applying a known network to directly characterize a tumor population with a relatively concise set of paired genomic-transcriptomic relationships, and designed for quantification of inter-tumor heterogeneity. In contrast, the approach in [54] is model-driven and the networks are learned.

3.3 Limitations and extensions

At the DNA level, we have only considered non-synonymous somatic mutations and extreme copy number variations. Specifically, we did not consider any annotation for mutations (*e.g.*, specific base changes) beyond population frequencies. We have also limited our annotation of downstream effects to the presence or absence of deviation of RNA expression from a baseline (normal) population, not accounting for the direction of the aberration from baseline (*i.e.*, up-regulation or down-regulation). The simplicity of such a binarized representation has enabled new findings and allows for some analysis of mechanism. Examination of source level paired aberrations facilitate the process of annotating mutations with an unknown effect (so-called variants of uncertain significance). Mutations that are recovered at the population level together with copy number losses can be presumed to be inactivating, and vice-versa.

Further consistency constraints may be imposed for a deeper analysis of the biology. For example, an STP may appear to be “inconsistent” if the source gene is duplicated or has an activating mutation and yet the target gene is down-aberrant, assuming the intermediate genes do not further modulate signaling propagation through the network. Such situations should clearly not be excluded. Needless to say there are many other explanations for such observations, *e.g.*, methylation to take but one example. Indeed, there are many cases of aberrant target genes which do not appear in any STP *i.e.*, for which there is no putative explanation in terms of upstream mutations and copy number variations. Uncovering a mechanistically coherent picture of the upstream-downstream synergy would evidently require incorporating additional types of data (*e.g.*, gene fusions, histone modifications and changes in methylation), other sources of transcriptional dysregulation (*e.g.*, expression of microRNAs) and other downstream effects, such as post-transcriptional changes in regulation and aberrant protein structure and concentration. Without such data, making assumptions about consistency among the catalogued and detected anomalies would result in damaging over-simplifications.

The covering signatures vary considerably from one tissue or subtype to another. For instance, among our six tissues and at any level (source, target, or pair) the core set of features (those shared by all minimal coverings) is the smallest in colon and accounts for all colon samples in TCGA, whereas substantially larger signatures were necessary in other tumors (*e.g.*, in breast), and some populations could not even be largely covered (*e.g.*, in prostate) regardless of the number of features. On one hand, a plausible bias here is that better and more refined network information is available for some cancer types than others. For instance, colon cancer has served for years as a model of tumorigenesis, and a wealth of data is available to derive “realistic” signaling pathways and regulatory networks compared with cancer types studied to a lesser extent. On the other hand other DNA alterations, beside mutations and copy number changes, can drive tumorigenesis and may be necessary to efficiently “cover” a cancer population. A prototypical example is that the exceptionally large coverings in prostate cancer may be due to the absence of data on gene fusions; in fact, over a half of the tumors could be accounted for by a small subset of such alterations (*e.g.*, the fusion between *TMPRSS2* and *ERG*, or other *ETS* family genes [55]).

Finally, the theoretical framework we have developed is based on the “Regulators \rightarrow Targets” paradigm and would support the incorporation of additional *omics* information. In fact,

gene fusion, epigenetic measurements, epigenomic states, enhancer expression, and so forth, could be simply integrated to generate an expanded repertoire of STPs, as could proteomics or metabolomics serve as additional downstream “targets”. To this end, epigenetics events (*e.g.*, methylation status, chromatin modification marks, and so on) could be easily integrated with DNA aberrations at the “source” level, while protein levels could be combined with RNA measurements. In both cases, a set of mechanistic rules would be required to integrate the different data types. For instance, a specific “source” gene in a pair could be defined as “aberrant” if it is mutated, OR it is deleted, OR it is hyper-methylated, and so on. Similarly, a “target” gene could be deemed “aberrant” based on biologically justified rules for combining protein and RNA data. Importantly, from the computational point of view, adding further modalities would not change the number of constraints in the optimization.

4 Conclusions

We have described an integrated analysis of DNA and RNA aberrations, which is grounded in cancer biology and enabled by a highly simplified summary representation of the complex and heterogeneous landscape of aberrations in cancer populations. The summary is a collection of STPs, each linking a particular DNA aberration with a downstream RNA expression change, and derived automatically from a “covering” algorithm in combinatorial optimization. Beside recapitulating many known alterations, our collection of STPs flags potentially important aberrations and interactions which might go unrecognized using simple frequency criteria, given the accumulation of low frequency events at the population level. This integrated representation could facilitate discriminating cancer drivers from passenger aberrations, and suggest potential novel therapeutic targets for further functional studies. Furthermore, this representation allows for a rigorous quantitative estimate of heterogeneity in a cancer population and across distinct tumor phenotypes, which would not be otherwise feasible. Indeed, in order to quantify heterogeneity beyond a simple listing of possibilities, it is necessary to assign likelihoods to these possibilities and their co-occurrences, in which case the entropy of the distributions over the possible combinations is the natural measure. The heterogeneity differences observed between distinct cancer phenotypes, along with the interactions among paired aberrations, suggest that our approach can represent an alternative to standard statistical filtering to identify important features for predictive model building and machine learning application in cancer. Finally, our analytical framework provides a highly efficient and innovative computational tool for harnessing the expanding data on tumor samples emerging from large consortia projects.

Supporting information

S1 Text. The file contains all supplementary figures and tables referenced from the main paper. Specifically, this supporting file contains two subsections. In the Supplementary Table subsection, Tables A—B show results pertaining to the selection of parameters k and L . Tables C—D show the basic statistics of interactions before and after filters. Tables F—U illustrate the statistics of “Optimal Covering” and core set with associated probabilities for 5 tissues at 3 levels. Tables V—I display the divergence probabilities at different cancer subtypes. Tables II and III show the entropy analysis across distinct tissues. In Supplementary Fig subsection, Fig A shows probabilities relating to target cores, and Fig B shows effects of filter thresholds on fraction of samples covered. Figs C—G show pair covering network across different cancers. Figs H—J show annotated network for cancer subtypes. Figs L—M show the complete core set across tissues. Figs N—P display coding trees for cancer subtypes.
(PDF)

S1 Data. The compressed archive contains 6 “DNABinary.txt.gz” files which correspond to binary DNA aberration matrices for all tumor types, 6 “RNABinary.txt.gz” files which corresponds to binary RNA aberration matrices for all tumor types, and “full_signature.xlsx” which contains one full signature for each tumor type at each level (pair, source, target). (ZIP)

Acknowledgments

We thank Drs. Giovanni Parmigiani, Nathan Price, Diego Fernando Sanchez, and Eddie Luidy-Imada for helpful discussions.

Author Contributions

Conceptualization: Qian Ke, Wikum Dinalankara, Laurent Younes, Donald Geman, Luigi Marchionni.

Data curation: Qian Ke, Wikum Dinalankara, Luigi Marchionni.

Formal analysis: Qian Ke, Wikum Dinalankara, Laurent Younes.

Funding acquisition: Donald Geman, Luigi Marchionni.

Investigation: Qian Ke, Wikum Dinalankara, Laurent Younes, Donald Geman, Luigi Marchionni.

Methodology: Qian Ke, Laurent Younes, Donald Geman, Luigi Marchionni.

Project administration: Luigi Marchionni.

Supervision: Laurent Younes, Donald Geman, Luigi Marchionni.

Visualization: Qian Ke, Wikum Dinalankara, Luigi Marchionni.

Writing – original draft: Qian Ke, Wikum Dinalankara, Laurent Younes, Donald Geman, Luigi Marchionni.

Writing – review & editing: Qian Ke, Wikum Dinalankara, Laurent Younes, Donald Geman, Luigi Marchionni.

References

1. Cancer Genome Atlas Research Network, Weinstein JN, Collisson EA, Mills GB, Shaw KRM, Ozenberger BA, et al. The Cancer Genome Atlas Pan-Cancer analysis project. *Nat Genet.* 2013; 45(10):1113–20. <https://doi.org/10.1038/ng.2764> PMID: 24071849
2. 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–D947. <https://doi.org/10.1093/nar/gky1015> PMID: 30371878
3. Zehir A, Benayed R, Shah RH, Syed A, Middha S, Kim HR, et al. Mutational landscape of metastatic cancer revealed from prospective clinical sequencing of 10,000 patients. *Nature medicine.* 2017; 23(6):703. <https://doi.org/10.1038/nm.4333> PMID: 28481359
4. ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analysis of whole genomes. *Nature.* 2020; 578(7793):82–93. <https://doi.org/10.1038/s41586-020-1969-6>
5. Vogelstein B, Kinzler KW. The path to cancer—three strikes and you’re out. *N Engl J Med.* 2015; 373(20):1895–1898. <https://doi.org/10.1056/NEJMp1508811>
6. Tomasetti C, Marchionni L, Nowak MA, Parmigiani G, Vogelstein B. Only three driver gene mutations are required for the development of lung and colorectal cancers. *Proc Natl Acad Sci U S A.* 2015; 112(1):118–23. <https://doi.org/10.1073/pnas.1421839112>
7. Sever R, Brugge JS. Signal transduction in cancer. *Cold Spring Harbor perspectives in medicine.* 2015; 5(4):a006098. <https://doi.org/10.1101/cshperspect.a006098>

8. Haigis KM, Cichowski K, Elledge SJ. Tissue-specificity in cancer: The rule, not the exception. *Science*. 2019; 363(6432):1150–1151. <https://doi.org/10.1126/science.aaw3472>
9. 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. <https://doi.org/10.1038/nature12213> PMID: 23770567
10. Osmanbeyoglu HU, Toska E, Chan C, Baselga J, Leslie CS. Pancancer modelling predicts the context-specific impact of somatic mutations on transcriptional programs. *Nature communications*. 2017; 8:14249.
11. PCAWG Transcriptome Core Group, Calabrese C, Davidson NR, Demircioğlu D, Fonseca NA, He Y, et al. Genomic basis for RNA alterations in cancer. *Nature*. 2020; 578(7793):129–136. <https://doi.org/10.1038/s41586-020-1970-0> PMID: 32025019
12. Costello JC, Heiser LM, Georgii E, Gönen M, Menden MP, Wang NJ, et al. A community effort to assess and improve drug sensitivity prediction algorithms. *Nature biotechnology*. 2014; 32(12):1202. <https://doi.org/10.1038/nbt.2877> PMID: 24880487
13. 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):371–385. <https://doi.org/10.1016/j.cell.2018.02.060> PMID: 29625053
14. Dinalankara W, Ke Q, Xu Y, Ji L, Pagane N, Lien A, et al. Digitizing omics profiles by divergence from a baseline. *Proceedings of the National Academy of Sciences*. 2018; 115(18):4545–4552. <https://doi.org/10.1073/pnas.1721628115> PMID: 29666255
15. Hristov BH, Singh M. Network-based coverage of mutational profiles reveals cancer genes. *Cell systems*. 2017; 5(3):221–229. <https://doi.org/10.1016/j.cels.2017.09.003>
16. Humphrey PA. Gleason grading and prognostic factors in carcinoma of the prostate. *Modern pathology*. 2004; 17(3):292–306. <https://doi.org/10.1038/modpathol.3800054>
17. Parker JS, Mullins M, Cheang MCU, Leung S, Voduc D, Vickery T, et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *J Clin Oncol*. 2009; 27(8):1160–7. <https://doi.org/10.1200/JCO.2008.18.1370> PMID: 19204204
18. Isella C, Brundu F, Bellomo SE, Galimi F, Zanella E, Porporato R, et al. Selective analysis of cancer-cell intrinsic transcriptional traits defines novel clinically relevant subtypes of colorectal cancer. *Nature communications*. 2017; 8(1):1–16. <https://doi.org/10.1038/ncomms15107> PMID: 28561063
19. Jassal B, Matthews L, Viteri G, Gong C, Lorente P, Fabregat A, et al. The reactome pathway knowledge-base. *Nucleic Acids Res*. 2020; 48(D1):D498–D503. <https://doi.org/10.1093/nar/gkz1031> PMID: 31691815
20. Cerami EG, Gross BE, Demir E, Rodchenkov I, Babur O, Anwar N, et al. Pathway Commons, a web resource for biological pathway data. *Nucleic Acids Res*. 2011; 39(Database issue):D685–90. <https://doi.org/10.1093/nar/gkq1039> PMID: 21071392
21. Shannon CE. A mathematical theory of communication. *The Bell system technical journal*. 1948; 27(3):379–423. <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>
22. Cover TM, Thomas JA. *Elements of information theory*. John Wiley & Sons; 2012.
23. Schürmann T. Bias analysis in entropy estimation. *Journal of Physics A: Mathematical and General*. 2004; 37(27):L295.
24. Grassberger P. Finite sample corrections to entropy and dimension estimates. *Physics Letters A*. 1988; 128(6-7):369–373. [https://doi.org/10.1016/0375-9601\(88\)90193-4](https://doi.org/10.1016/0375-9601(88)90193-4)
25. Grassberger P. Entropy estimates from insufficient samplings. *arXiv preprint physics/0307138*. 2003;.
26. Miller G. Note on the bias of information estimates. *Information theory in psychology: Problems and methods*. 1955; p. 95–100.
27. Huynh J, Chand A, Gough D, Ernst M. Therapeutically exploiting STAT3 activity in cancer—using tissue repair as a road map. *Nature Reviews Cancer*. 2019; 19(2):82–96. <https://doi.org/10.1038/s41568-018-0090-8>
28. Olivier M, Hollstein M, Hainaut P. TP53 mutations in human cancers: origins, consequences, and clinical use. *Cold Spring Harbor perspectives in biology*. 2010; 2(1):a001008.
29. Semmler L, Reiter-Brennan C, Klein A. BRCA1 and breast cancer: a review of the underlying mechanisms resulting in the tissue-specific tumorigenesis in mutation carriers. *Journal of breast cancer*. 2019; 22(1):1–14. <https://doi.org/10.4048/jbc.2019.22.e6>
30. Segers VFM, Dugaucquier L, Feyen E, Shakeri H, De Keulenaer GW. The role of ErbB4 in cancer. *Cell Oncol (Dordr)*. 2020; 43(3):335–352.
31. Viatour P, Ehmer U, Saddic LA, Dorrell C, Andersen JB, Lin C, et al. Notch signaling inhibits hepatocellular carcinoma following inactivation of the RB pathway. *Journal of Experimental Medicine*. 2011; 208(10):1963–1976. <https://doi.org/10.1084/jem.20110198> PMID: 21875955

32. Aster JC, Pear WS, Blacklow SC. The varied roles of Notch in cancer. *Annual Review of Pathology: Mechanisms of Disease*. 2017; 12:245–275. <https://doi.org/10.1146/annurev-pathol-052016-100127>
33. Nowell CS, Radtke F. Notch as a tumour suppressor. *Nature Reviews Cancer*. 2017; 17(3):145. <https://doi.org/10.1038/nrc.2016.145>
34. Sigismund S, Avanzato D, Lanzetti L. Emerging functions of the EGFR in cancer. *Molecular oncology*. 2018; 12(1):3–20. <https://doi.org/10.1002/1878-0261.12155>
35. Pfeifer GP, Denissenko MF, Olivier M, Tretyakova N, Hecht SS, Hainaut P. Tobacco smoke carcinogens, DNA damage and p53 mutations in smoking-associated cancers. *Oncogene*. 2002; 21(48):7435–7451. <https://doi.org/10.1038/sj.onc.1205803>
36. Basso RS, Hochbaum DS, Vandin F. Efficient algorithms to discover alterations with complementary functional association in cancer. *PLoS computational biology*. 2019; 15(5):e1006802. <https://doi.org/10.1371/journal.pcbi.1006802>
37. Aho AV, Hopcroft JE. The design and analysis of computer algorithms. Pearson Education India; 1974.
38. Pasquier N, Bastide Y, Taouil R, Lakhal L. Discovering frequent closed itemsets for association rules. In: *International Conference on Database Theory*. Springer; 1999. p. 398–416.
39. Alves R, Rodriguez-Baena DS, Aguilar-Ruiz JS. Gene association analysis: a survey of frequent pattern mining from gene expression data. *Briefings in Bioinformatics*. 2010; 11(2):210–224. <https://doi.org/10.1093/bib/bbp042>
40. Uno T, Kiyomi M, Arimura H. LCM ver. 3: collaboration of array, bitmap and prefix tree for frequent itemset mining. In: *Proceedings of the 1st international workshop on open source data mining: frequent pattern mining implementations*; 2005. p. 77–86.
41. Geistlinger L, Csaba G, Küffner R, Mulder N, Zimmer R. From sets to graphs: towards a realistic enrichment analysis of transcriptomic systems. *Bioinformatics*. 2011; 27(13):i366–i373. <https://doi.org/10.1093/bioinformatics/btr228>
42. Vaske CJ, Benz SC, Sanborn JZ, Earl D, Szeto C, Zhu J, et al. Inference of patient-specific pathway activities from multi-dimensional cancer genomics data using PARADIGM. *Bioinformatics*. 2010; 26(12):i237–i245. <https://doi.org/10.1093/bioinformatics/btq182> PMID: 20529912
43. Levine AJ, Oren M. The first 30 years of p53: growing ever more complex. *Nat Rev Cancer*. 2009; 9(10):749–58. <https://doi.org/10.1038/nrc2723>
44. Liu JW, Kim MS, Nagpal J, Yamashita K, Poeta L, Chang X, et al. Quantitative hypermethylation of NMDAR2B in human gastric cancer. *Int J Cancer*. 2007; 121(9):1994–2000. <https://doi.org/10.1002/ijc.22934> PMID: 17620329
45. Kim MS, Yamashita K, Baek JH, Park HL, Carvalho AL, Osada M, et al. N-methyl-D-aspartate receptor type 2B is epigenetically inactivated and exhibits tumor-suppressive activity in human esophageal cancer. *Cancer Res*. 2006; 66(7):3409–18. <https://doi.org/10.1158/0008-5472.CAN-05-1608> PMID: 16585162
46. Tamura H, Suzuki M, Moriya Y, Hoshino H, Okamoto T, Yoshida S, et al. Aberrant methylation of N-methyl-D-aspartate receptor type 2B (NMDAR2B) in non-small cell carcinoma. *BMC Cancer*. 2011; 11:220. <https://doi.org/10.1186/1471-2407-11-220>
47. Ito K, Park SH, Nayak A, Byerly JH, Irie HY. PTK6 Inhibition Suppresses Metastases of Triple-Negative Breast Cancer via SNAIL-Dependent E-Cadherin Regulation. *Cancer Res*. 2016; 76(15):4406–17. <https://doi.org/10.1158/0008-5472.CAN-15-3445>
48. Wozniak DJ, Kajdacsy-Balla A, Macias V, Ball-Kell S, Zenner ML, Bie W, et al. PTEN is a protein phosphatase that targets active PTK6 and inhibits PTK6 oncogenic signaling in prostate cancer. *Nat Commun*. 2017; 8(1):1508. <https://doi.org/10.1038/s41467-017-01574-5> PMID: 29142193
49. Bertrand FE, Angus CW, Partis WJ, Sigounas G. Developmental pathways in colon cancer: crosstalk between WNT, BMP, Hedgehog and Notch. *Cell Cycle*. 2012; 11(23):4344–51. <https://doi.org/10.4161/cc.22134>
50. Herbst RS, Heymach JV, Lippman SM. Lung cancer. *N Engl J Med*. 2008; 359(13):1367–80. <https://doi.org/10.1056/NEJMra0802714>
51. Marzio A, Puccini J, Kwon Y, Maverakis NK, Arbini A, Sung P, et al. The F-Box Domain-Dependent Activity of EMI1 Regulates PARPi Sensitivity in Triple-Negative Breast Cancers. *Mol Cell*. 2019; 73(2):224–237.e6. <https://doi.org/10.1016/j.molcel.2018.11.003> PMID: 30554948
52. Jamal-Hanjani M, Quezada SA, Larkin J, Swanton C. Translational implications of tumor heterogeneity. *Clin Cancer Res*. 2015; 21(6):1258–66. <https://doi.org/10.1158/1078-0432.CCR-14-1429>
53. Marusyk A, Janiszewska M, Polyak K. Intratumor Heterogeneity: The Rosetta Stone of Therapy Resistance. *Cancer Cell*. 2020; 37(4):471–484. <https://doi.org/10.1016/j.ccell.2020.03.007>

54. 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 computational biology*. 2019; 15(7):e1007088. <https://doi.org/10.1371/journal.pcbi.1007088> PMID: 31276486
55. Tomlins SA, Bjartell A, Chinnaiyan AM, Jenster G, Nam RK, Rubin MA, et al. ETS gene fusions in prostate cancer: from discovery to daily clinical practice. *Eur Urol*. 2009; 56(2):275–86. <https://doi.org/10.1016/j.eururo.2009.04.036> PMID: 19409690