



**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

## **The cancer genome atlas pan-cancer analysis project**

Downloaded from: <https://research.chalmers.se>, 2024-04-28 04:14 UTC

Citation for the original published paper (version of record):

Weinstein, J., Collisson, E., Mills, G. et al (2013). The cancer genome atlas pan-cancer analysis project. Nature Genetics, 45(10): 1113-1120. <http://dx.doi.org/10.1038/ng.2764>

N.B. When citing this work, cite the original published paper.

# The Cancer Genome Atlas Pan-Cancer analysis project

The Cancer Genome Atlas Research Network<sup>1</sup>, John N Weinstein<sup>2,3</sup>, Eric A Collisson<sup>4</sup>, Gordon B Mills<sup>3</sup>, Kenna R Mills Shaw<sup>5,6</sup>, Brad A Ozenberger<sup>7</sup>, Kyle Ellrott<sup>8,9</sup>, Ilya Shmulevich<sup>10</sup>, Chris Sander<sup>11</sup> & Joshua M Stuart<sup>8,9</sup>

The Cancer Genome Atlas (TCGA) Research Network has profiled and analyzed large numbers of human tumors to discover molecular aberrations at the DNA, RNA, protein and epigenetic levels. The resulting rich data provide a major opportunity to develop an integrated picture of commonalities, differences and emergent themes across tumor lineages. The Pan-Cancer initiative compares the first 12 tumor types profiled by TCGA. Analysis of the molecular aberrations and their functional roles across tumor types will teach us how to extend therapies effective in one cancer type to others with a similar genomic profile.

Cancer can take hundreds of different forms depending on the location, cell of origin and spectrum of genomic alterations that promote oncogenesis and affect therapeutic response.

Although many genomic events with direct phenotypic impact have been identified, much of the complex molecular landscape remains incompletely charted for most cancer lineages.

treated with targeted agents is increasing with the discovery of likely driver aberrations in most lung tumors<sup>16,17</sup>. Large-scale processes that shape cancer genomes have similarly been identified. Analyses of chromothripsis<sup>18</sup> and chromoplexy<sup>19</sup>, which involve the breakage and rearrangement of chromosomes at multiple loci, and kataegis<sup>20</sup>, which involves hypermutational processes associated with genomic rearrangements, are providing insights into tumor evolution (see Garraway and Lander<sup>21</sup> for a review).

## Analysis across tumor types

Increased numbers of tumor sample data sets enhance the ability to detect and analyze molecular defects in cancers. For example, driver genes can be pinpointed more precisely by narrowing regions affected by amplification and deletion to smaller segments of the chromosome using data on recurrent events across tumor types. The use of large cohorts has enabled DNA sequencing to uncover a list of recurrent genomic aberrations (mutations, amplifications, deletions, translocations, fusions and other structural variants), both known and novel, as common events across tumor types<sup>22</sup>. However, 'long tails' in the distributions of aberrations among samples have also been uncovered<sup>23</sup>. Indeed, a majority of the TCGA samples have distinct alterations not shared with other samples in their

## Molecular profiling of single tumor types

That cancer is fundamentally a genomic disease is now well established. Early on, large numbers of oncogenes were identified using functional assays on genetic material from tumors in positive-selection systems<sup>1-3</sup>, and a subset of tumor suppressor genes was identified by analyzing loss of heterozygosity<sup>4</sup>. More recently, systematic cancer genomics projects, including TCGA (**Box 1**), have applied emerging technologies to the analysis of specific tumor types. This disease-specific focus has identified novel oncogenic drivers and the genes contributing to functional change<sup>5-7</sup>, has established definitions of molecular subtypes<sup>8-12</sup> and has identified new biomarkers on the basis of genomic, transcriptomic, proteomic and epigenomic alterations. Some of these biomarkers have clinical implications<sup>13,14</sup>. For example, we now view ductal breast cancer as a collection of distinct diseases whose major subtypes (for example, luminal A, luminal B, HER2 and basal-like) are managed differently in the clinic; the outcomes for metastatic melanoma have improved as a result of therapeutic targeting of *BRAF*<sup>V600</sup> mutations<sup>15</sup>; and the fraction of lung cancers

<sup>1</sup>Full lists of members and affiliations appear at the end of the paper. <sup>2</sup>Department of Bioinformatics and Computational Biology, University of Texas MD Anderson Cancer Center, Houston, Texas, USA. <sup>3</sup>Department of Systems Biology, University of Texas MD Anderson Cancer Center, Houston, Texas, USA. <sup>4</sup>Department of Medicine, University of California, San Francisco, San Francisco, California, USA. <sup>5</sup>The Cancer Genome Atlas Program Office, Center for Cancer Genomics, National Cancer Institute, Bethesda, Maryland, USA. <sup>6</sup>University of Texas MD Anderson Cancer Center, Institute for Personalized Cancer Therapy, Houston, Texas, USA. <sup>7</sup>National Human Genome Research Institute, US National Institutes of Health, Bethesda, Maryland, USA. <sup>8</sup>Department of Biomolecular Engineering, University of California, Santa Cruz, Santa Cruz, California, USA. <sup>9</sup>Center for Biomolecular Science and Engineering, University of California, Santa Cruz, Santa Cruz, California, USA. <sup>10</sup>Institute for Systems Biology, Seattle, Washington, USA. <sup>11</sup>Computational Biology Center, Memorial Sloan-Kettering Cancer Center, New York, New York, USA. Correspondence should be addressed to J.M.S. (jstuart@ucsc.edu).

## BOX 1 TCGA: MISSION AND STRATEGY

Important information about the biological relevance of the molecular changes in cancer can be obtained through combined analysis of multiple different types of data.

For that reason, TCGA's principal aims are to generate, quality control, merge, analyze and interpret molecular profiles at the DNA, RNA, protein and epigenetic levels for hundreds of clinical tumors representing various tumor types and their subtypes. Cases that meet quality assurance specifications are characterized using technologies that assess the sequence of the exome, copy number variation (measured by SNP arrays), DNA methylation, mRNA expression and sequence, microRNA expression and transcript splice variation. Additional platforms applied to a subset of the tumors, including whole-genome sequencing and RPPAs, provide additional layers of data to complement the core genomic data sets and clinical data. By the end of 2015, the TCGA Research Network plans to have achieved the ambitious goal of analyzing the genomic, epigenomic and gene expression profiles of more than 10,000 specimens from more than 25 different tumor types.

TCGA has other, complementary aims as well: to promote the development and application of new technologies, to detect cancer-specific molecular alterations, to make data and results freely available to the scientific community, to develop tools and standard operating procedures that can serve other large-scale profiling projects and to build cadres of individuals (including experimentalists, computational biologists, statistical analysts, computer scientists and administrative staff) with the expertise to carry out such large-scale, team science projects. As of 24 July 2013, TCGA had mapped molecular patterns across 7,992 total cases representing 27 tumor types. The data, along with tools for exploring them, are publicly available at <http://www.cancergenome.nih.gov/>. Eight 'marker papers' (comprehensive initial publications on each of the tumor types) have been published so far<sup>8–12,14,27</sup>.

cohort. Despite the apparent uniqueness of each individual tumor in this regard, the set of molecular aberrations often integrates into known biological pathways that are shared by sets of tumor samples. In other cases, rare somatic mutations can be implicated as drivers by aggregating events across tumor types to improve the detection of patterns, for example, hotspot mutations in DNA segments that encode particular protein domains, leading to the identification of potential new drug targets.

Determining whether rare aberrations are drivers (oncogenic contributors) or just passengers (clonally propagated with neutral effect) and whether they are clinically actionable will require further functional evaluation as well as the analysis of additional tumors to increase power. The identification of more driver aberrations and acquired vulnerabilities for each individual tumor will undoubtedly boost personalized care. Developing treatments that target the ~140 drivers<sup>22</sup> validated so far, however daunting, appears possible; devising one-off therapies for the thousands of aberrations in the long tails will be much more challenging.

Although important general principles have emerged from decades of study<sup>24,25</sup>, until recently, most research on the molecular, pathological and clinical natures of cancers has been 'siloe'd' by tumor type<sup>26</sup>. One has only to glance at the directory of oncology departments in any major cancer center to realize that medical and surgical cancer care are, for the most

part, also divided by disease as defined by organ of origin. This framework has made sense for generations, but the results of molecular analysis are now calling this view into question; cancers of disparate organs have many shared features, whereas, conversely, cancers from the same organ are often quite distinct.

Important similarities among tumor subtypes from different organs have already been identified. For example, *TP53* mutations drive high-grade serous ovarian, serous endometrial and basal-like breast carcinomas, all of which share a global transcriptional signature involving the activation of similar oncogenic pathways<sup>10,27</sup>. Similarly, *ERBB2-HER2* is mutated and/or amplified in subsets of glioblastoma, gastric, serous endometrial, bladder and lung cancers. The result, at least in some cases, is responsiveness to HER2-targeted therapy, analogous to that previously observed for *HER2*-amplified breast cancer. Other commonalities across tumor types include inherited and somatic inactivation of the *BRCA1-BRCA2* pathway in both serous ovarian and basal-like breast cancers, microsatellite instability in colorectal and endometrial tumors, and the recently identified *POLE*-mediated ultramutator phenotype characterized by extremely high mutation rates, common to both colon and endometrial cancers<sup>12,27,28</sup>. Conversely, there are important cases in which the same genetic aberrations have very different effects depending on the organ within which they arise. A prime example is provided by the

NOTCH gene family, which is inactivated in some squamous cell cancers of the lung, head and neck<sup>29</sup>, skin<sup>30</sup> and cervix<sup>31</sup> but activated by mutation in leukemias<sup>32</sup>.

Such examples illustrate the importance of developing a comprehensive perspective across tumors, independent of histopathologic diagnosis; shared molecular patterns will enable etiologic and therapeutic discoveries in one disease that can be applied to another. Importantly, integrative interpretation of the data will help identify how the consequences of mutations vary across tissues, with important therapeutic implications. Relatively rare cancers, such as childhood malignancies, in particular stand to benefit from such an approach.

We know much more about the molecular details of major cancers than we did just a few years ago, but once a cancer is metastatic it remains incurable with few exceptions. Only time will tell whether the integration of molecular characteristics with data on histology, organ site and metastatic location will contribute to an improvement in patient outcomes. But the balance is shifting in this direction. Hence, the goal of the Pan-Cancer project is to identify and analyze aberrations in the tumor genome and phenotype that define cancer lineages as well as to identify aberrations that transcend particular lineages. This report outlines the scope of the project and introduces the first coordinated set of manuscripts to be published from the enterprise.

## The Pan-Cancer project

To gain analytical breadth—defining commonalities, differences and emergent themes across cancer types and organs of origin—TCGA launched the Pan-Cancer analysis project at a meeting held on 26–27 October 2012 in Santa Cruz, California. The Pan-Cancer project is a coordinated initiative whose goals are to assemble coherent, consistent TCGA data sets across tumor types, as well as across platforms, and then to analyze and interpret these data (Box 2). Within 2 months of the project's launch, a 'data freeze' was declared on the first 12 TCGA tumor types, each profiled using 6 different genomic, epigenomic, transcriptional and proteomic platforms (Fig. 1). Since that time, the aggregated data sets have been quality controlled, analyzed statistically and interpreted by a consortium of researchers, principally members of the TCGA Research Network.

The Pan-Cancer project lays the framework for an analytic process that, in the future, will include the integration of new tumor types and data from TCGA and other such enterprises. There are currently major consortium efforts in pediatric cancers (TARGET; Therapeutically

Applicable Research to Generate Effective Treatments) and adult cancers (ICGC; International Cancer Genomics Consortium), as well as smaller projects by research teams around the world. A critical component of such efforts will be the functional validation of aberrations in individual genes in team science efforts such as CTD<sup>2</sup> (Cancer Target Discovery and Development) and the elucidation of pathway and network relationships in programs such as the US National Cancer Institute's Integrative Cancer Biology Program.

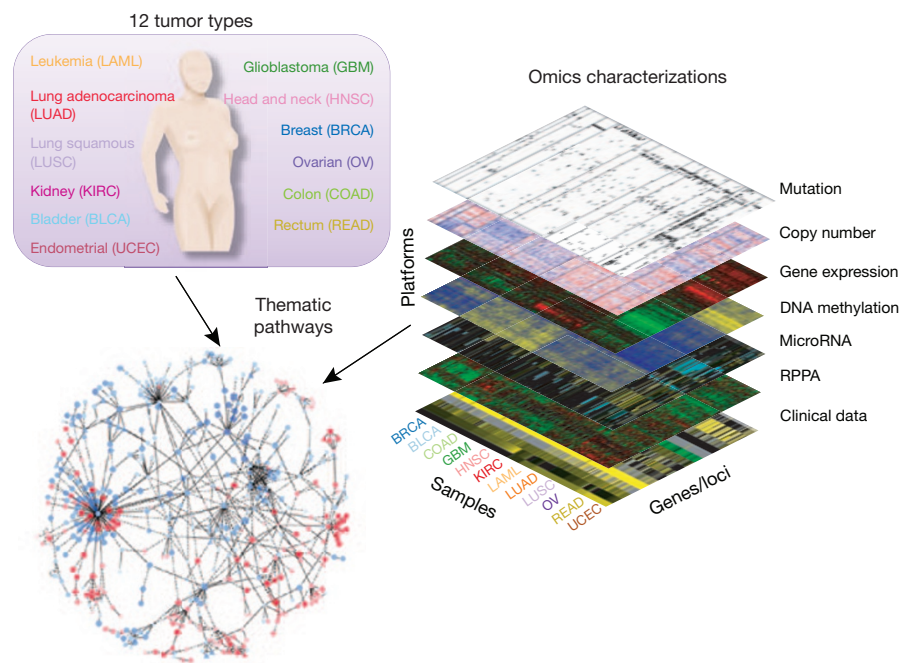
A number of investigations that go beyond the single-tumor perspective are being addressed in the collection of Pan-Cancer manuscripts. Examples of the kinds of questions addressed by these investigations are given below.

Can increases in statistical power help to distinguish new driver mutations from the background of passenger mutations as the sample size is increased by aggregating the 12 tumor types? Assembled Pan-Cancer data have, in fact, enabled the identification of new patterns of genomic drivers. New computational approaches that leverage cross-tumor principles of replication timing and gene expression correlated with background mutation rates now enable the identification of frequently mutated genes while eliminating many false-positive and false-negative calls made in several single-tumor-type projects<sup>33</sup>. Further, the power to identify multiple signals of positive selection has increased the ability to distinguish 'driver' from 'passenger' aberrations<sup>34</sup>.

What tissue associations underlie the major genomic structural changes in cancer? Improved methods for the analysis of structural variation of large chromosome segments have refined the ability to identify genomic and epigenetic regulators in multiple peak regions seen only by collating data across different cancer types. Tissue-associated patterns have now been established for the rate and timing of whole-genome duplication events<sup>35</sup>.

What pathways emerge as critical and potentially actionable when all mutational events across many tissues are considered together? New classes of mutations, such as those in chromatin-remodeling genes, are emerging as cancer drivers identified only by (i) collecting less frequent events across tumor types, (ii) integrating event types such as mutations, copy number changes and epigenetic silencing, (iii) combining multiple algorithms to identify predicted drivers<sup>34</sup> and (iv) aggregating genes using gene networks and pathways<sup>36</sup>.

Can an increase in the number of samples enhance analysis of the co-occurrence and mutual exclusivity of gene aberrations and



**Figure 1** Integrated data set for comparing and contrasting multiple tumor types. The TCGA Pan-Cancer project assembled data from thousands of patients with primary tumors occurring in different sites of the body, covering 12 tumor types (top left) including glioblastoma multiforme (GBM), lymphoblastic acute myeloid leukemia (LAML), head and neck squamous carcinoma (HNSC), lung adenocarcinoma (LUAD), lung squamous carcinoma (LUSC), breast carcinoma (BRCA), kidney renal clear-cell carcinoma (KIRC), ovarian carcinoma (OV), bladder carcinoma (BLCA), colon adenocarcinoma (COAD), uterine cervical and endometrial carcinoma (UCEC) and rectal adenocarcinoma (READ). Six types of omics characterization were performed creating a 'data stack' (right) in which data elements across the platforms are linked by the fact that the same samples were used for each, thus maximizing the potential of integrative analysis. Use of the data enables the identification of general trends, including common pathways (bottom left), revealing master regulatory hubs activated (red) or deactivated (blue) across different tissue types.

improve the ability to distinguish driver aberrations from passengers? A bird's-eye view of genomic and epigenomic events yields a 'fate map' of the alternative routes to carcinogenesis in a decision tree that spans tissue boundaries<sup>37</sup>.

Can molecular subtypes be delineated to disentangle tissue-specific from tissue-independent components of disease? Analyses of the epigenome, transcriptome and proteome show a strong influence of tissue on the state of altered pathways in tumor cells. For instance, analysis of the gene expression landscape reinforces the dominant tissue dependence of altered pathways and complements simultaneous profiling of over a hundred proteins important in cancer<sup>38</sup>. Using all of the tumor types together allows for any tumor-specific signals to be subtracted from the data sets. Intriguingly, subtracting tissue-specific signal from DNA microarray gene expression data sets identifies signatures of immune stromal influence that transcend tumor type boundaries (R. Verhaak, personal communication). Further, events that are common across lineages become apparent in a cross-tumor analysis<sup>38</sup>. Examples are the

hormonal dependencies of breast, ovarian and endometrial cancers and a common 'squamous cell' signature across head and neck, lung, cervical and bladder cancers.

Which events actionable in one tumor lineage are also actionable in another tumor lineage, potentially increasing the range of indications for specific targeted therapeutics? A systematic evaluation of machine-learning approaches is needed to highlight methodological principles for predicting patient outcomes using integrated information across tissues (H. Liang, personal communication).

### Limitations of analysis across tumor types

Several data integration challenges place unavoidable limitations on cross-tumor analysis at the current time. A key challenge is the integration of data that have been generated on different platforms or updates of the same platform, as technologies improve. In the Pan-Cancer studies, for example, there have been transitions to much higher density DNA methylation arrays, use of different exome capture technologies, addition of RNA sequencing to microarray-based RNA



## BOX 2 COORDINATION OF DATA AND RESULTS

The first goal of the Pan-Cancer Analysis Working Group was to assemble data from the separate disease projects to build a well-coordinated joint data set spanning multiple tumor types. A data freeze (21 December 2012) based on six different genomic and epigenomic characterization platforms was made available as the panCan12 data set to all analysis groups. Twelve tumor types (GBM, OV, BRCA, LUSC, LUAD, COAD, READ, KIRC, UCEC, BLCA, HNSC and LAML; see **Fig. 1** for definitions) were selected on the basis of data maturity, adequate sample size and publication or submission for publication of primary analyses. The Pan-Cancer 12 data set includes a total of 5,074 tumor samples, of which 93% had been assessed for genomic, epigenomic, and gene and protein expression data on at least one platform each (**Table 1**). The essential purpose of such a joint data set is twofold: to increase the statistical power to detect functional genomic determinants of disease and to identify both tissue-specific aspects of cancer and intrinsic molecular commonalities across tumor types.

The Pan-Cancer analysis project started as an informal collaboration among members of the TCGA Research Network but then quickly expanded to include many other interested researchers. Ensuring standardization and consistency of the data and annotations across multiple platforms and clinical data elements was a necessity for the project. To coordinate analyses across this large group of researchers, formal pipelines were created to establish a coherent working base of data and results.

The process of TCGA data generation and Pan-Cancer analysis was as follows (**Fig. 2**). First, tumor and germline samples were obtained from a large number of tissue source sites and processed by the Biospecimen Core Resource (with sample selection according to criteria established for each tumor type and with

extensive quality controls) to generate purified DNA, RNA and protein preparations. The preparations were sent to Genome Characterization Centers (GCCs) and Genome Sequencing Centers (GSCs) for molecular profiling, and the resulting data were deposited at the TCGA Data Coordinating Center (DCC) to provide a primary source of data at four levels of data processing. Seven Genome Data Analysis Centers (GDACs), along with analysts at the GCCs and GSCs and the external research community, shared analysis and interpretation of the data, coordinating activities through face-to-face meetings and regular (usually weekly) teleconferences.

A data freeze was created by pulling higher levels of interpreted data (Level 3) from the DCC into a coordinating repository called Synapse created by Sage Bionetworks. To create a coherent data set, a sample 'white list' was created by synchronizing flagged samples with the DCC on the basis of annotations and criteria from the individual disease working groups. The Pan-Cancer project leveraged the TCGA infrastructure for sample acquisition, sample processing and data generation on individual tumor types, as well as for the production of derived data sets and a variety of analysis results assembled in the Broad Institute's Firehose system. Assembled robust, self-consistent data sets across all 12 Pan-Cancer tumor types were deposited into Synapse. The Synapse system implements mechanisms for tracking provenance and metadata, stable digital object identifiers (DOIs) for data referencing and flexible methods for data access, either through a wiki-like web-based environment or programmatically through application programming interfaces (APIs). The panCan12 data sets and selected results are available under the Synapse accession doi:10.7303/syn300013.

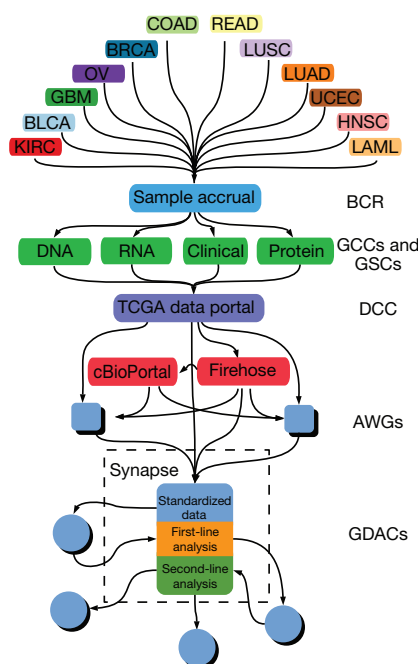
characterization and increases in the quality and number of antibodies available for reverse-phase proteomic arrays (RPPAs). A series of analyses of batch effects has been carried out to assess systematic and platform-

specific biases (R. Akbani, personal communication). However, more work is needed to establish best practices for minimizing unwanted batch effects while preserving biological signals.

The nature and quality of available clinical data vary widely by cancer type. Differences in these data limit the ability to establish one-size-fits-all norms for the comparison of demographic information, histopathologic characterization, behavioral context and clinical outcomes. For example, the Pan-Cancer survival data are relatively robust for serous ovarian cancer because of its poor prognosis but are still immature for breast and endometrial cancers, as (thankfully) most patients with these cancers do better for longer periods of time. Certain data elements are routinely collected only when they are anticipated to be relevant (for example,

the smoking history of patients with lung, bladder and head and neck cancers). Clear viral etiologies have been identified in several solid tumor types, including head and neck cancer, cervical cancer, Kaposi's sarcoma and hepatocellular carcinoma. However, a pan-cancer analysis of the infectious etiologies of other cancers could not be conducted at present because infection status was recorded for only some tumors and tumor types (as an optional data element). Finally, tumor stage and grade are not easily comparable across different tumor types because, for good reason, each tumor type has its own system. This set of challenges to cross-tumor analysis highlights the fact that current clinical practice is largely conducted according to classification by tissue or organ.

Statistically speaking, care must be taken to ensure that the increased sample size



**Figure 2** Data coordination for the Pan-Cancer TCGA project. Data were collected by the Biospecimen Collection Resource (BCR) from 12 different tumor types and characterized on 6 major platforms by the Genome Characterization Centers and Genomic Sequencing Centers (GCCs and GSCs). Data sets were deposited in the TCGA Data Coordinating Center (DCC) from which they were then distributed to the Broad Institute's Firehose and the Memorial Sloan-Kettering Cancer Center's cBioPortal for various automated processing pipelines. Analysis Working Groups (AWGs) conducted focused analyses on individual tumor types. Results from the DCC, Firehose and AWGs were collected and stored in Sage Bionetworks' Synapse database system to create a data freeze. Genome data analysis centers (GDACs) accessed and deposited both data and results through Synapse to coordinate distributed analyses.

**Table 1** Data freeze used by the Pan-Cancer project as defined on 21 December 2012

Cancer	RPPA <sup>a</sup>	DNA methylation <sup>b</sup>	Copy number <sup>c</sup>	Mutation <sup>d</sup>	microRNA <sup>e</sup>	Expression <sup>f</sup>
LUSC	195	358	345	178	332	227
READ	130	162	164	69	143	71
GBM	214	405	578	290	501	495
LAML	NA	194	198	197	187	179
HNSC	212	310	310	277	309	303
BLCA	54	126	126	99	121	96
KIRC	423	457	457	417	442	431
UCEC	200	512	511	248	497	333
LUAD	237	431	357	229	365	355
OV	332	592	577	316	454	581
BRCA	408	888	887	772	870	817
COAD	269	420	422	155	407	192
Total	2,674	4,855	4,932	3,247	4,628	4,080

Tabulated are the numbers of unique tumor samples available for each tumor type (rows) and each measurement platform (columns). NA, not available.

<sup>a</sup>Reverse-phase protein arrays measuring protein and phosphoprotein abundance. <sup>b</sup>DNA methylation at CpG islands. <sup>c</sup>Microarray-based measurement of copy number. <sup>d</sup>Samples subjected to whole-exome sequencing to determine single-nucleotide and structural variants. <sup>e</sup>Sequencing of microRNAs. <sup>f</sup>RNA sequencing and microarray gene expression analysis.

achieved by cross-cancer comparison does not lead to increased false-negative rates for discovery (for example, by 'diluting out' an important mutation specific to one disease) or false-positive rates (for example, by compounding false positives known to result from current single-tumor investigations<sup>33</sup>).

Tumor lineage has an important role in the observed patterns of co-aberrations and gene expression profiles that indicate different consequences of seemingly similar events, for example, involving the same gene(s) or amplicon(s). Likewise, new methods for accurately probing cross-tumor trends will need to account explicitly for differences across tissues in mutation rates, copy number changes on the focal and arm-level scales, and the prevalence of other co-occurring events in the genetic and epigenetic backgrounds.

Despite these challenges, the collection of Pan-Cancer publications presented here represents a landmark in the continuing effort to understand the common and contrasting biologies of cancers from a molecular perspective. Still, major questions amenable to further cross-tumor investigations remain (**Box 3**), and the techniques used to compare different tumors will undoubtedly improve with use, time and further collaborative efforts.

### Future directions

The Pan-Cancer project represents one of the first of what will surely be many efforts to coordinate analysis across the molecular landscape of cancer, especially as additional tumor types are investigated in large numbers. Further increasing the number of samples per tumor type and the variety of tumor types will improve our ability to detect rare driver events in heterogeneous tumor samples. But the true

power will come from a detailed analysis across tumor types—with links to high-quality clinical outcomes and eventual experimental validation and clinical trials to test the hypotheses that emerge. Technologies such as laser capture microdissection and cell sorting will improve the ability to distinguish whether omic signals

arise from malignant or stromal cells. Histone profiling, protein analysis based on mass spectrometry and deconvolution of tumor heterogeneity through single-cell sequencing are examples of technologies expected to add important new dimensions of information. Continued efforts to identify the progenitor

### BOX 3 EXAMPLES OF ADDITIONAL MAJOR QUESTIONS AMENABLE TO FURTHER PAN-CANCER ANALYSES

- What is the spectrum of nucleotide- and dinucleotide-level changes associated with different carcinogenic etiologies (for example, tobacco, pathogens or inflammation) operating in different parts of the body?
- Will integration of additional data sources, including additional tumor types from TCGA and other projects, increase the power of analysis in useful ways?
- How can characterization based on molecular changes complement pathological analysis for classification of cancers into tumor lineages with potentially different clinical management?
- Can molecular profiles effectively categorize cancers for therapeutic decision-making?
- Are there predictive expression-based signatures for genomic events that transcend tissues, reflecting pathways disrupted by the alterations?
- Will comprehensive protein analysis through emerging mass spectrometry approaches in the Clinical Proteomic Tumor Analysis Consortium (CPTAC) and other efforts extend the power of the genomic, transcriptomic and proteomic analyses in TCGA?
- Will whole genome analysis demonstrate the influence of mobile elements, mutations in non-coding regions and connections to constitutional risk-associated loci?
- How are changes in protein families distributed across different tumor types?
- Are aberrations in specific protein domains or pathways distributed differentially across tumor lineages?
- Beyond the known examples, including in cervical, head and neck, esophageal and hepatocellular cancers, can we identify other cancer types that show virally mediated initiation?
- Are bacteria associated with different cancer lineages (as *Fusobacteria* are in colorectal cancer<sup>43</sup>)?
- Can the answers to any of these questions help in the design of novel therapies and clinical trials, with the ultimate goal of improving patient outcomes?

cells of tumors will enable universal properties to be distinguished from parochial ones. Clone-level and other types of studies may identify even more connections among tumor types. Longitudinal genomic studies on primary resected tumors paired with their local recurrences and/or metastases will be undertaken by large consortium efforts, which have heretofore been restricted to primary disease and have lacked information about response to treatment. The characteristics of primary tumors may change markedly when they metastasize to distant sites, particularly bone and brain. Analysis of metastasis across tumor types will therefore be highly informative.

The power of cross-tumor analysis will increase as technologies for monitoring individual tumor cells at high resolution come into play. Now that the price of genome sequencing has fallen, the next Pan-Cancer enterprise will be able to analyze large numbers of whole-genome sequences across tumor types. Whole-genome analysis will complement the current studies by shedding light on mutational processes in the noncoding parts of the genome, which have not been as well explored so far. This expanded analysis will bring focus to disruptions in promoter and enhancer sites and aberrations in noncoding RNAs, as well as the genomic integration processes at work in tumor evolution that result from mobile endogenous and exogenous DNA elements such as retrotransposons and viruses. Whole-genome sequencing will create a backdrop against which genome-wide association studies can relate inherited predispositions to particular forms of cancer. Systems-oriented approaches, based on relevant pathways and networks, will add to the therapeutic opportunities that arise from the wealth of data. Experimental follow-up will be critical to assess the functional consequences and therapeutic liabilities of these new findings.

### From many tumors to the individual

The hope is that investigations across tumor type such as the Pan-Cancer project will ultimately inform clinical decision-making. We hope such studies will enable the discovery of novel therapeutic agents that can be tested clinically—perhaps in novel adaptive, biomarker-based clinical trials that cross boundaries between tumor types. Toward this end, TCGA Pan-Cancer data sets have been made available publicly in one location. Although coordination remains a challenge, the data sets comprise an unequalled resource for the integrative analysis of cancer in its many forms.

A key challenge is the development of clinical trial strategies for connecting subsets of tumors from different tissues in terms of molecular signatures. Recent analyses of pharmacological profiling experiments across a diverse panel of cancer cell lines has suggested that common genetic alterations can sometimes predict response to therapy across multiple cell lineages<sup>39–42</sup>. Biomarker-based design of clinical trials can increase statistical power, greatly decreasing the size, expense and duration of clinical trials.

The number and size of omic data sets on cancer available to the research community for mining and exploring continue to expand rapidly, and computational tools to derive insights into the fundamental causes of cancer are becoming more powerful. It is important to note that the full potential of the enterprise will be realized only over time and with broader efforts. Still, the collection of TCGA Pan-Cancer publications represents a significant contribution to a new period of discovery in cancer research.

### ACKNOWLEDGMENTS

We thank J. Zhang for the administrative coordination of TCGA Pan-Cancer Analysis Working Group activities, C. Perou and K. Hoadley for contributions to Figure 1, and A. Margolin, D. Wheeler, M. Meyerson and L. Ding for comments on early drafts of the manuscript. The study was funded by the National Cancer Institute and the National Human Genome Research Institute.

### COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.



This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0/>.

1. Soda, M. *et al. Nature* **448**, 561–566 (2007).
2. Parada, L.F., Tabin, C.J., Shih, C. & Weinberg, R.A. *Nature* **297**, 474–478 (1982).

3. Payne, G.S., Bishop, J.M. & Varmus, H.E. *Nature* **295**, 209–214 (1982).
4. Baker, S.J. *et al. Science* **244**, 217–221 (1989).
5. Tomlins, S.A. *et al. Science* **310**, 644–648 (2005).
6. Davies, H. *et al. Nature* **417**, 949–954 (2002).
7. Mardis, E.R. *et al. N. Engl. J. Med.* **361**, 1058–1066 (2009).
8. Cancer Genome Atlas Research Network. *Nature* **455**, 1061–1068 (2008).
9. Cancer Genome Atlas Research Network. *Nature* **474**, 609–615 (2011).
10. Cancer Genome Atlas Network. *Nature* **490**, 61–70 (2012).
11. Cancer Genome Atlas Research Network. *Nature* **489**, 519–525 (2012).
12. Cancer Genome Atlas Network. *Nature* **487**, 330–337 (2012).
13. Perou, C.M. *et al. Nature* **406**, 747–752 (2000).
14. Cancer Genome Atlas Research Network. *Nature* **499**, 43–49 (2013).
15. Chapman, P.B. *et al. N. Engl. J. Med.* **364**, 2507–2516 (2011).
16. Paez, J.G. *et al. Science* **304**, 1497–1500 (2004).
17. Takeuchi, K. *et al. Nat. Med.* **18**, 378–381 (2012).
18. Stephens, P.J. *et al. Cell* **144**, 27–40 (2011).
19. Baca, S.C. *et al. Cell* **153**, 666–677 (2013).
20. Alexandrov, L.B. *et al. Nature* **500**, 415–421 (2013).
21. Garraway, L.A. & Lander, E.S. *Cell* **153**, 17–37 (2013).
22. Vogelstein, B. *et al. Science* **339**, 1546–1558 (2013).
23. Wheeler, D.A. & Wang, L. *Genome Res.* **23**, 1054–1062 (2013).
24. Hanahan, D. & Weinberg, R.A. *Cell* **144**, 646–674 (2011).
25. Hanahan, D. & Weinberg, R.A. *Cell* **100**, 57–70 (2000).
26. McDermott, U. & Settleman, J. *J. Clin. Oncol.* **27**, 5650–5659 (2009).
27. Kandoth, C. *et al. Nature* **497**, 67–73 (2013).
28. Palles, C. *et al. Nat. Genet.* **45**, 136–144 (2013).
29. Stransky, N. *et al. Science* **333**, 1157–1160 (2011).
30. Wang, N.J. *et al. Proc. Natl. Acad. Sci. USA* **108**, 17761–17766 (2011).
31. Zagouras, P., Stifani, S., Blaumueller, C.M., Carcangiu, M.L. & Artavanis-Tsakonas, S. *Proc. Natl. Acad. Sci. USA* **92**, 6414–6418 (1995).
32. Weng, A.P. *et al. Science* **306**, 269–271 (2004).
33. Lawrence, M.S. *et al. Nature* **499**, 214–218 (2013).
34. Tamborero, D. *et al. Sci. Rep.* doi:10.1038/srep02650 (26 September 2013).
35. Zack, T.I. *et al. Nat. Genet.* doi:10.1038/ng.2760 (26 September 2013).
36. Hofree, M., Shen, J.P., Carter, H., Gross, A. & Ideker, T. *Nat. Methods* doi:10.1038/nmeth.2651 (15 September 2013).
37. Ciriello, G. *et al. Nat. Genet.* doi:10.1038/ng.2762 (26 September 2013).
38. Li, J. *et al. Nat. Methods* doi:10.1038/nmeth.2650 (15 September 2013).
39. Barretina, J. *et al. Nature* **483**, 603–607 (2012).
40. Garnett, M.J. *et al. Nature* **483**, 570–575 (2012).
41. Weinstein, J.N. *Nature* **483**, 544–545 (2012).
42. Heiser, L.M. *et al. Proc. Natl. Acad. Sci. USA* **109**, 2724–2729 (2012).
43. Kostic, A.D. *et al. Genome Res.* **22**, 292–298 (2012).

Genome Characterization Center: Kyle Chang<sup>12</sup>, Chad J Creighton<sup>12</sup>, Caleb Davis<sup>12</sup>, Lawrence Donehower<sup>12</sup>, Jennifer Drummond<sup>12</sup>, David Wheeler<sup>12</sup>, Adrian Ally<sup>13</sup>, Miruna Balasundaram<sup>13</sup>, Inanc Birol<sup>13–15</sup>, Yaron S N Butterfield<sup>13</sup>, Andy Chu<sup>13</sup>, Eric Chuah<sup>13</sup>, Hye-Jung E Chun<sup>13</sup>, Noreen Dhalla<sup>13</sup>, Ranabir Guin<sup>13</sup>, Martin Hirst<sup>13</sup>, Carrie Hirst<sup>13</sup>, Robert A Holt<sup>13</sup>, Steven J M Jones<sup>13</sup>, Darlene Lee<sup>13</sup>, Haiyan I Li<sup>13</sup>, Marco A Marra<sup>13</sup>, Michael Mayo<sup>13</sup>, Richard A Moore<sup>13</sup>, Andrew J Mungall<sup>13</sup>, A Gordon Robertson<sup>13</sup>, Jacqueline E Schein<sup>13</sup>, Payal Sipahimalani<sup>13</sup>, Angela Tam<sup>13</sup>, Nina Thiessen<sup>13</sup>, Richard J Varhol<sup>13</sup>, Rameen Beroukhim<sup>16</sup>, Ami S Bhatt<sup>16,17</sup>, Angela N Brooks<sup>16,18</sup>, Andrew D Cherniack<sup>16</sup>, Samuel S Freeman<sup>16</sup>, Stacey B Gabriel<sup>16</sup>, Elena Helman<sup>16,19</sup>, Joonil Jung<sup>16</sup>, Matthew Meyerson<sup>16,17</sup>, Akinyemi I Ojesina<sup>16,17</sup>, Chandra Sekhar Pedamallu<sup>16,17</sup>, Gordon Saksena<sup>16</sup>, Steven E Schumacher<sup>16,20</sup>, Barbara Tabak<sup>16,20</sup>, Travis Zack<sup>16,21</sup>, Eric S Lander<sup>16</sup>, Christopher A Bristow<sup>22</sup>, Angela Hadjipanayis<sup>23</sup>, Psalm Haseley<sup>24</sup>, Raju Kucherlapati<sup>25</sup>, Semin Lee<sup>24</sup>, Eunjung Lee<sup>24</sup>, Lovelace J Luquette<sup>24</sup>, Harshad S Mahadeshwar<sup>22</sup>, Angeliki Pantazi<sup>23</sup>, Michael Parfenov<sup>23</sup>, Peter J Park<sup>24,26,27</sup>, Alexei Protopopov<sup>22</sup>, Xiaojia Ren<sup>23</sup>, Netty Santoso<sup>23</sup>, Jonathan Seidman<sup>23</sup>, Sahil Seth<sup>22</sup>, Xingzhi Song<sup>22</sup>, Jiabin Tang<sup>22</sup>, Ruibin Xi<sup>24,28,29</sup>, Andrew W Xu<sup>24</sup>, Lixing Yang<sup>24</sup>, Dong Zeng<sup>22</sup>, J Todd Auman<sup>30</sup>, Saianand Balu<sup>31</sup>, Elizabeth Buda<sup>32</sup>, Cheng Fan<sup>31</sup>, Katherine A Hoadley<sup>31</sup>, Corbin D Jones<sup>32,33</sup>, Shaowu Meng<sup>31</sup>, Piotr A Mieczkowski<sup>34</sup>, Joel S Parker<sup>31,34</sup>, Charles M Perou<sup>31,34,35</sup>, Jeffrey Roach<sup>36</sup>, Yan Shi<sup>31</sup>, Grace O Silva<sup>31,34</sup>, Donghui Tan<sup>34</sup>, Umadevi Veluvolu<sup>34</sup>, Scot Waring<sup>31,32</sup>, Matthew D Wilkerson<sup>34</sup>, Junyuan Wu<sup>31</sup>, Wei Zhao<sup>31,34</sup>, Tom Bodenheimer<sup>31</sup>, D Neil Hayes<sup>31,37</sup>, Alan P Hoyle<sup>31</sup>, Stuart R Jeffreys<sup>31</sup>, Lisle E Mose<sup>31</sup>, Janae V Simons<sup>31</sup>, Mathew G Soloway<sup>31</sup>, Stephen B Baylin<sup>38</sup>, Benjamin P Berman<sup>39</sup>, Moiz S Bootwalla<sup>39</sup>, Ludmila Danilova<sup>40</sup>, James G Herman<sup>38</sup>, Toshinori Hinoue<sup>39</sup>, Peter W Laird<sup>39</sup>, Suhn K Rhie<sup>39</sup>, Hui Shen<sup>39</sup>, Timothy Triche Jr<sup>39</sup>, Daniel J Weisenberger<sup>39</sup>, Scott L Carter<sup>16</sup>, Kristian Cibulskis<sup>16</sup>, Lynda Chin<sup>16,22</sup>, Jianhua Zhang<sup>22</sup>, Gad Getz<sup>16,41,42</sup>, Carrie Sougnez<sup>16</sup> & Min Wang<sup>12</sup>

Genome Data Analysis Center: Gordon Saksena<sup>16</sup>, Scott L Carter<sup>16</sup>, Kristian Cibulskis<sup>16</sup>, Lynda Chin<sup>16,22</sup>, Jianhua Zhang<sup>22</sup>, Gad Getz<sup>16,41,42</sup>, Huyen Dinh<sup>12</sup>, Harsha Vardhan Doddapaneni<sup>12</sup>, Richard Gibbs<sup>12</sup>, Preethi Gunaratne<sup>12,43</sup>, Yi Han<sup>12</sup>, Divya Kalra<sup>12</sup>, Christie Kovar<sup>12</sup>, Lora Lewis<sup>12</sup>, Margaret Morgan<sup>12</sup>, Donna Morton<sup>12</sup>, Donna Muzny<sup>12</sup>, Jeffrey Reid<sup>12</sup>, Liu Xi<sup>12</sup>, Juok Cho<sup>16</sup>, Daniel DiCara<sup>16</sup>, Scott Frazer<sup>16</sup>, Nils Gehlenborg<sup>16,24</sup>, David I Heiman<sup>16</sup>, Jaegil Kim<sup>16</sup>, Michael S Lawrence<sup>16</sup>, Pei Lin<sup>16</sup>, Yingchun Liu<sup>16</sup>, Michael S Noble<sup>16</sup>, Petar Stojanov<sup>16,17</sup>, Doug Voet<sup>16</sup>, Hailei Zhang<sup>16</sup>, Lihua Zou<sup>16</sup>, Chip Stewart<sup>16</sup>, Brady Bernard<sup>10</sup>, Ryan Bressler<sup>10</sup>, Andrea Eakin<sup>10</sup>, Lisa Iype<sup>10</sup>, Theo Knijnenburg<sup>10</sup>, Roger Kramer<sup>10</sup>, Richard Kreisberg<sup>10</sup>, Kalle Leinonen<sup>10</sup>, Jake Lin<sup>10</sup>, Yuexin Liu<sup>48</sup>, Michael Miller<sup>10</sup>, Sheila M Reynolds<sup>10</sup>, Hector Rovira<sup>10</sup>, Ilya Shmulevich<sup>10</sup>, Vesteinn Thorsson<sup>10</sup>, Da Yang<sup>44</sup>, Wei Zhang<sup>44</sup>, Samirkumar Amin<sup>45</sup>, Chang-Jiun Wu<sup>22</sup>, Chia-Chin Wu<sup>22</sup>, Rehan Akbani<sup>2</sup>, Kenneth Aldape<sup>46</sup>, Keith A Baggerly<sup>2</sup>, Bradley Broom<sup>2</sup>, Tod D Casasent<sup>2</sup>, James Cleland<sup>2,46</sup>, Chad Creighton<sup>47</sup>, Deepti Dodda<sup>2</sup>, Mary Edgerton<sup>48</sup>, Leng Han<sup>2</sup>, Shelley M Herbrich<sup>2</sup>, Zhenlin Ju<sup>2</sup>, Hoon Kim<sup>2</sup>, Seth Lerner<sup>49</sup>, Jun Li<sup>2</sup>, Han Liang<sup>2</sup>, Wenbin Liu<sup>2</sup>, Philip L Lorenzi<sup>2</sup>, Yiling Lu<sup>3</sup>, James Melott<sup>2</sup>, Gordon B Mills<sup>3</sup>, Lam Nguyen<sup>2,46</sup>, Xiaoping Su<sup>2</sup>, Roeland Verhaak<sup>2</sup>, Wenyi Wang<sup>2</sup>, John N Weinstein<sup>2,3</sup>, Andrew Wong<sup>2,46</sup>, Yang Yang<sup>2,50</sup>, Jun Yao<sup>51</sup>, Rong Yao<sup>2</sup>, Kosuke Yoshihara<sup>2</sup>, Yuan Yuan<sup>2,52</sup>, Alfred K Yung<sup>51</sup>, Nianxiang Zhang<sup>2</sup>, Siyuan Zheng<sup>2</sup>, Michael Ryan<sup>3,46</sup>, David W Kane<sup>2,53</sup>, B Arman Aksoy<sup>11</sup>, Giovanni Ciriello<sup>11</sup>, Gideon Dresdner<sup>11</sup>, Jianjiong Gao<sup>11</sup>, Benjamin Gross<sup>11</sup>, Anders Jacobsen<sup>11</sup>, Andre Kahles<sup>11</sup>, Marc Ladanyi<sup>54,55</sup>, William Lee<sup>11</sup>, Kjong-Van Lehmann<sup>11</sup>, Martin L Miller<sup>11</sup>, Ricardo Ramirez<sup>11</sup>, Gunnar Rättsch<sup>11</sup>, Boris Reva<sup>11</sup>, Chris Sander<sup>11</sup>, Nikolaus Schultz<sup>11</sup>, Yasin Senbabaoglu<sup>11</sup>, Ronglai Shen<sup>56</sup>, Rileen Sinha<sup>11</sup>, S Onur Sumer<sup>11</sup>, Yichao Sun<sup>11</sup>, Barry S Taylor<sup>4,56,57</sup>, Nils Weinhold<sup>11</sup>, Suzanne Fei<sup>58</sup>, Paul Spellman<sup>58</sup>, Christopher Benz<sup>59</sup>, Daniel Carlin<sup>8,9</sup>, Melissa Cline<sup>8,9</sup>, Brian Craft<sup>8,9</sup>, Kyle Ellrott<sup>8,9</sup>, Mary Goldman<sup>8,9</sup>, David Haussler<sup>8,9,60</sup>, Singer Ma<sup>8,9</sup>, Sam Ng<sup>8,9</sup>, Evan Paull<sup>8,9</sup>, Amie Radenbaugh<sup>8,9</sup>, Sofie Salama<sup>8,9,60</sup>, Artem Sokolov<sup>8,9</sup>, Joshua M Stuart<sup>8,9</sup>, Teresa Swatloski<sup>8,9</sup>, Vladislav Uzunangelov<sup>8,9</sup>, Peter Waltman<sup>8,9</sup>, Christina Yau<sup>59</sup>, Jing Zhu<sup>8,9</sup> & Stanley R Hamilton<sup>44</sup>

Sequencing Center: Gad Getz<sup>16,41,42</sup>, Carrie Sougnez<sup>16</sup>, Scott Abbott<sup>61</sup>, Rachel Abbott<sup>61</sup>, Nathan D Dees<sup>61</sup>, Kim Delehaanty<sup>61</sup>, Li Ding<sup>61–63</sup>, David J Dooling<sup>61</sup>, Jim M Eldred<sup>61</sup>, Catrina C Fronick<sup>61</sup>, Robert Fulton<sup>61</sup>, Lucinda L Fulton<sup>61</sup>, Joelle Kalicki-Veizer<sup>61</sup>, Krishna-Latha Kanchi<sup>61</sup>, Cyriac Kandoth<sup>61</sup>, Daniel C Koboldt<sup>61</sup>, David E Larson<sup>61</sup>, Timothy J Ley<sup>61,64</sup>, Ling Lin<sup>61</sup>, Charles Lu<sup>61</sup>, Vincent J Magrini<sup>61</sup>, Elaine R Mardis<sup>61,63,65</sup>, Michael D McLellan<sup>61</sup>, Joshua F McMichael<sup>61</sup>, Christopher A Miller<sup>61</sup>, Michelle O'Laughlin<sup>61</sup>, Craig Pohl<sup>61</sup>, Heather Schmidt<sup>61</sup>, Scott M Smith<sup>61</sup>, Jason Walker<sup>61</sup>, John W Wallis<sup>61</sup>, Michael C Wendl<sup>61,65,66</sup>, Richard K Wilson<sup>61,63,65</sup>, Todd Wylie<sup>61</sup> & Qunyan Zhang<sup>61,65</sup>



**Data Coordinating Center:** Robert Burton<sup>67</sup>, Mark A Jensen<sup>53</sup>, Ari Kahn<sup>53</sup>, Todd Pihl<sup>53</sup>, David Pot<sup>53</sup> & Yunhu Wan<sup>53</sup>

**Tissue Source Site:** Douglas A Levine<sup>68</sup>

**Biospecimen Core Resource Center:** Aaron D Black<sup>69</sup>, Jay Bowen<sup>69</sup>, Jessica Frick<sup>69</sup>, Julie M Gastier-Foster<sup>69,70</sup>, Hollie A Harper<sup>69</sup>, Carmen Helsel<sup>69</sup>, Kristen M Leraas<sup>69</sup>, Tara M Lichtenberg<sup>69</sup>, Cynthia McAllister<sup>69</sup>, Nilsa C Ramirez<sup>69,70</sup>, Samantha Sharpe<sup>69</sup>, Lisa Wise<sup>69</sup> & Erik Zmuda<sup>69</sup>

**National Cancer Institute/National Human Genome Research Institute Project Team:** Stephen J Chanock<sup>71</sup>, Tanja Davidsen<sup>71</sup>, John A Demchok<sup>71</sup>, Greg Eley<sup>72</sup>, Ina Felau<sup>71</sup>, Brad A Ozenberger<sup>7</sup>, Margi Sheth<sup>71</sup>, Heidi Sofia<sup>7</sup>, Louis Staudt<sup>71</sup>, Roy Tarnuzzer<sup>71</sup>, Zhining Wang<sup>71</sup>, Liming Yang<sup>71</sup> & Jiashan Zhang<sup>71</sup>

**Collaborators:** Larsson Omberg<sup>73</sup>, Adam Margolin<sup>73</sup>, Benjamin J Raphael<sup>74</sup>, Fabio Vandin<sup>74</sup>, Hsin-Ta Wu<sup>74</sup>, Mark D M Leiserson<sup>74</sup>, Stephen C Benz<sup>75</sup>, Charles J Vaske<sup>75</sup>, Houtan Noushmehr<sup>76,77</sup>, Theo Knijnenburg<sup>10</sup>, Denise Wolf<sup>78</sup>, Laura Van 't Veer<sup>78</sup>, Eric A Collisson<sup>4</sup>, Dimitris Anastassiou<sup>79</sup>, Tai-Hsien Ou Yang<sup>79</sup>, Nuria Lopez-Bigas<sup>80,81</sup>, Abel Gonzalez-Perez<sup>80</sup>, David Tamborero<sup>80</sup>, Zheng Xia<sup>82</sup>, Wei Li<sup>82</sup>, Dong-Yeon Cho<sup>83</sup>, Teresa Przytycka<sup>83</sup>, Mark Hamilton<sup>84</sup>, Sean McGuire<sup>84</sup>, Sven Nelander<sup>85,86</sup>, Patrik Johansson<sup>85,86</sup>, Rebecka Jörnsten<sup>87</sup>, Teresia Kling<sup>85,86</sup> & Jose Sanchez<sup>87</sup>

<sup>12</sup>Human Genome Sequencing Center, Baylor College of Medicine, Houston, Texas, USA. <sup>13</sup>Canada's Michael Smith Genome Sciences Centre, BC Cancer Agency, Vancouver, British Columbia, Canada. <sup>14</sup>Department of Medical Genetics, University of British Columbia, Vancouver, British Columbia, Canada. <sup>15</sup>School of Computing Science, Simon Fraser University, Burnaby, British Columbia, Canada. <sup>16</sup>The Eli and Edythe L. Broad Institute of the Massachusetts Institute of Technology and Harvard University, Cambridge, Massachusetts, USA. <sup>17</sup>Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, Massachusetts, USA. <sup>18</sup>Dana-Farber Cancer Institute, Boston, Massachusetts, USA. <sup>19</sup>Harvard-MIT Division of Health Sciences & Technology, Cambridge, Massachusetts, USA. <sup>20</sup>Department of Cancer Biology, Dana-Farber Cancer Institute, Boston, Massachusetts, USA. <sup>21</sup>Biophysics Program, Harvard University, Boston, Massachusetts, USA. <sup>22</sup>Institute for Applied Cancer Science, Department of Genomic Medicine, University of Texas MD Anderson Cancer Center, Houston, Texas, USA. <sup>23</sup>Department of Genetics, Harvard Medical School, Boston, Massachusetts, USA. <sup>24</sup>The Center for Biomedical Informatics, Harvard Medical School, Boston, Massachusetts, USA. <sup>25</sup>Harvard Medical School-Partners HealthCare Center for Genetics and Genomics, Boston, Massachusetts, USA. <sup>26</sup>Division of Genetics, Brigham and Women's Hospital, Boston, Massachusetts, USA. <sup>27</sup>Informatics Program, Children's Hospital, Boston, Massachusetts, USA. <sup>28</sup>School of Mathematical Sciences, Peking University, Beijing, China. <sup>29</sup>Center for Statistical Science, Peking University, Beijing, China. <sup>30</sup>Eshelman School of Pharmacy, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA. <sup>31</sup>Lineberger Comprehensive Cancer Center, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA. <sup>32</sup>Carolina Center for Genome Sciences, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA. <sup>33</sup>Department of Biology, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA. <sup>34</sup>Department of Genetics, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA. <sup>35</sup>Department of Pathology and Laboratory Medicine, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA. <sup>36</sup>Research Computing Center, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA. <sup>37</sup>Department of Internal Medicine, Division of Medical Oncology, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA. <sup>38</sup>Cancer Biology Division, The Sidney Kimmel Comprehensive Cancer Center at Johns Hopkins University, Baltimore, Maryland, USA. <sup>39</sup>USC Epigenome Center, University of Southern California Keck School of Medicine, Los Angeles, California, USA. <sup>40</sup>The Sidney Kimmel Comprehensive Cancer Center at Johns Hopkins University, Baltimore, Maryland, USA. <sup>41</sup>Cancer Center, Massachusetts General Hospital, Boston, Massachusetts, USA. <sup>42</sup>Department of Pathology, Massachusetts General Hospital, Boston, Massachusetts, USA. <sup>43</sup>Department of Biology and Biochemistry, University of Houston, Houston, Texas, USA. <sup>44</sup>University of Texas MD Anderson Cancer Center, Houston, Texas, USA. <sup>45</sup>Department of Genomic Medicine, University of Texas MD Anderson Cancer Center, Houston, Texas, USA. <sup>46</sup>In Silico Solutions, Fairfax, Virginia, USA. <sup>47</sup>Baylor College of Medicine, Houston, Texas, USA. <sup>48</sup>Department of Pathology, University of Texas MD Anderson Cancer Center, Houston, Texas, USA. <sup>49</sup>Department of Urology, Baylor College of Medicine, Houston, Texas, USA. <sup>50</sup>Division of Biostatistics, University of Texas Health Science Center at Houston, School of Public Health, Houston, Texas, USA. <sup>51</sup>Neuro-Oncology Department, University of Texas MD Anderson Cancer Center, Houston, Texas, USA. <sup>52</sup>Graduate Program in Structural and Computational Biology and Molecular Biophysics, Baylor College of Medicine, Houston, Texas, USA. <sup>53</sup>SRA International, Inc., Fairfax, Virginia, USA. <sup>54</sup>Department of Pathology and Human Oncology, Memorial Sloan-Kettering Cancer Center, New York, New York, USA. <sup>55</sup>Pathogenesis Program, Memorial Sloan-Kettering Cancer Center, New York, New York, USA. <sup>56</sup>Department of Epidemiology and Biostatistics, Memorial Sloan-Kettering Cancer Center, New York, New York, USA. <sup>57</sup>Helen Diller Family Comprehensive Cancer Center, University of California, San Francisco, San Francisco, California, USA. <sup>58</sup>Molecular & Medical Genetics, Oregon Health & Science University, Portland, Oregon, USA. <sup>59</sup>Buck Institute for Research on Aging, Novato, California, USA. <sup>60</sup>Howard Hughes Medical Institute, University of California, Santa Cruz, Santa Cruz, California, USA. <sup>61</sup>The Genome Institute, Washington University, St. Louis, Missouri, USA. <sup>62</sup>Department of Medicine, Washington University, St. Louis, Missouri, USA. <sup>63</sup>Siteman Cancer Center, Washington University, St. Louis, Missouri, USA. <sup>64</sup>Division of Oncology, Washington University, St. Louis, Missouri, USA. <sup>65</sup>Department of Genetics, Washington University, St. Louis, Missouri, USA. <sup>66</sup>Department of Mathematics, Washington University, St. Louis, Missouri, USA. <sup>67</sup>SAIC-Frederick, Inc., Frederick, Maryland, USA. <sup>68</sup>Gynecology Service, Department of Surgery, Memorial Sloan-Kettering Cancer Center, New York, New York, USA. <sup>69</sup>The Research Institute at Nationwide Children's Hospital, Columbus, Ohio, USA. <sup>70</sup>The Ohio State University, Columbus, Ohio, USA. <sup>71</sup>National Cancer Institute, US National Institutes of Health, Bethesda, Maryland, USA. <sup>72</sup>Scintentis, Statham, Georgia, USA. <sup>73</sup>Sage Bionetworks, Seattle, Washington, USA. <sup>74</sup>Department of Computer Science, Brown University, Providence, Rhode Island, USA. <sup>75</sup>Five3 Genomics, LLC, Santa Cruz, California, USA. <sup>76</sup>Department of Genetics, Ribeirão Preto Medical School, University of São Paulo, São Paulo, Brazil. <sup>77</sup>Center for Integrative Systems Biology (CISBi), NAP/USP, São Paulo, Brazil. <sup>78</sup>Department of Laboratory Medicine, University of California, San Francisco, San Francisco, California, USA. <sup>79</sup>Department of Electrical Engineering, Columbia University, New York, New York, USA. <sup>80</sup>Research Unit on Biomedical Informatics, Department of Experimental and Health Sciences, Universitat Pompeu Fabra, Barcelona, Spain. <sup>81</sup>Catalan Institution for Research and Advanced Studies (ICREA), Barcelona, Spain. <sup>82</sup>Dan L. Duncan Cancer Center, Baylor College of Medicine, Houston, Texas, USA. <sup>83</sup>National Center for Biotechnology Information, National Library of Medicine, US National Institutes of Health, Bethesda, Maryland, USA. <sup>84</sup>Department of Molecular and Cellular Biology, Baylor College of Medicine, Houston, Texas, USA. <sup>85</sup>Institution for Immunology, Genetics and Pathology, Uppsala University, Uppsala, Sweden. <sup>86</sup>Science for Life Laboratory, Uppsala University, Uppsala, Sweden. <sup>87</sup>Mathematical Sciences, University of Gothenburg/Chalmers, Gothenburg, Sweden.