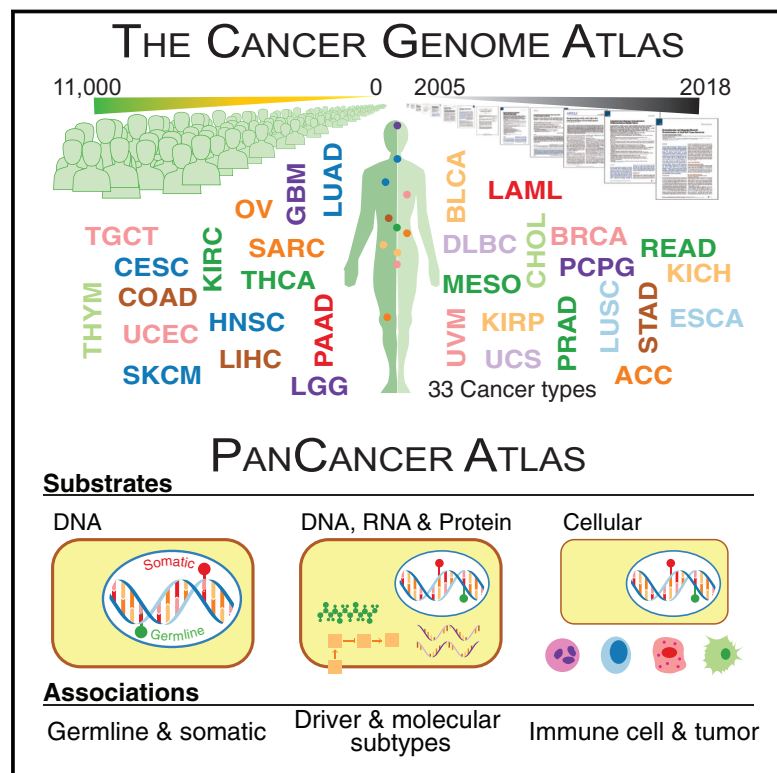


# Perspective on Oncogenic Processes at the End of the Beginning of Cancer Genomics

## Graphical Abstract



## Authors

Li Ding, Matthew H. Bailey, Eduard Porta-Pardo, ..., David A. Wheeler, Gad Getz, The Cancer Genome Atlas Research Network

## Correspondence

lding@wustl.edu (L.D.), wheeler@bcm.edu (D.A.W.), gadgetz@broadinstitute.org (G.G.)

## In Brief

A synthesized view on oncogenic processes based on PanCancer Atlas analyses highlights the complex impact of genome alterations on the signaling and multi-omic profiles of human cancers as well as their influence on tumor microenvironment.

## Highlights

- An overview of PanCancer Atlas analyses on oncogenic molecular processes
- Germline genome affects somatic genomic landscape in a pathway-dependent fashion
- Genome mutations impact expression, signaling, and multi-omic profiles
- Mutation burdens and drivers influence immune-cell composition in microenvironment



# Perspective on Oncogenic Processes at the End of the Beginning of Cancer Genomics

Li Ding,<sup>1,2,3,4,33,34,\*</sup> Matthew H. Bailey,<sup>1,2,33</sup> Eduard Porta-Pardo,<sup>5,6,33</sup> Vesteinn Thorsson,<sup>7</sup> Antonio Colaprico,<sup>8,9</sup> Denis Bertrand,<sup>10</sup> David L. Gibbs,<sup>7</sup> Amila Weerasinghe,<sup>1,2</sup> Kuan-lin Huang,<sup>1,2</sup> Collin Tokheim,<sup>11,12</sup> Isidro Cortés-Ciriano,<sup>13,14,15</sup> Reyka Jayasinghe,<sup>1,2</sup> Feng Chen,<sup>1,4</sup> Lihua Yu,<sup>16</sup> Sam Sun,<sup>17</sup> Catharina Olsen,<sup>8</sup> Jaegil Kim,<sup>18</sup> Alison M. Taylor,<sup>18,19</sup> Andrew D. Cherniack,<sup>18,19</sup> Rehan Akbani,<sup>20</sup> Chayaporn Suphavilai,<sup>10</sup> Niranjan Nagarajan,<sup>10</sup>

(Author list continued on next page)

<sup>1</sup>Department of Medicine, Washington University in St. Louis, St. Louis, MO 63110, USA

<sup>2</sup>McDonnell Genome Institute, Washington University in St. Louis, St. Louis, MO 63108, USA

<sup>3</sup>Department of Genetics, Washington University in St. Louis, St. Louis, MO 63110, USA

<sup>4</sup>Siteman Cancer Center, Washington University in St. Louis, St. Louis, MO 63110, USA

<sup>5</sup>Barcelona Supercomputing Centre, 08034 Barcelona, Spain

<sup>6</sup>Sanford Burnham Prebys Medical Discovery Institute, La Jolla, CA 92037, USA

<sup>7</sup>Institute for Systems Biology, Seattle, WA 98109, USA

<sup>8</sup>Machine Learning Group (MLG), Département d'Informatique, Université Libre de Bruxelles, 1050 Brussels, Belgium

<sup>9</sup>Department of Human Genetics, University of Miami, Miami, FL 33136, USA

<sup>10</sup>Computational and Systems Biology, Genome Institute of Singapore, Singapore, 13862

<sup>11</sup>Institute for Computational Medicine, Johns Hopkins University, Baltimore, MD 21218, USA

<sup>12</sup>Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21218, USA

<sup>13</sup>Harvard Medical School, Boston, MA 02115, USA

<sup>14</sup>Ludwig Center at Harvard, Boston, MA 02115, USA

<sup>15</sup>Department of Chemistry, University of Cambridge, Cambridge CB2 1EW, UK

<sup>16</sup>H3 Biomedicine Inc., Cambridge, MA 02139, USA

<sup>17</sup>Department of Radiation Oncology, Baylor College of Medicine, Houston, TX 77030, USA

<sup>18</sup>Broad Institute, Cambridge, MA 02142, USA

<sup>19</sup>Department of Medical Oncology, Dana-Farber Cancer Institute, 450 Brookline Avenue, Boston, MA 02215, USA

<sup>20</sup>Department of Bioinformatics and Computational Biology, University of Texas MD Anderson Cancer Center, Houston, TX 77498, USA

<sup>21</sup>Baskin School of Engineering, University of California, Santa Cruz, Santa Cruz, CA 95064, USA

(Affiliations continued on next page)

## SUMMARY

The Cancer Genome Atlas (TCGA) has catalyzed systematic characterization of diverse genomic alterations underlying human cancers. At this historic junction marking the completion of genomic characterization of over 11,000 tumors from 33 cancer types, we present our current understanding of the molecular processes governing oncogenesis. We illustrate our insights into cancer through synthesis of the findings of the TCGA PanCancer Atlas project on three facets of oncogenesis: (1) somatic driver mutations, germline pathogenic variants, and their interactions in the tumor; (2) the influence of the tumor genome and epigenome on transcriptome and proteome; and (3) the relationship between tumor and the microenvironment, including implications for drugs targeting driver events and immunotherapies. These results will anchor future characterization of rare and common tumor types, primary and relapsed tumors, and cancers across ancestry groups and will guide the deployment of clinical genomic sequencing.

## INTRODUCTION

In the nearly half century of the “War on Cancer,” prevention and treatment have progressed significantly, but many forms of the disease remain incurable. The advent of large-scale DNA sequencing ushered in new possibilities. Beginning with coding regions (Sjöblom et al., 2006), sequencing has sparked a revolution in cancer research. Genomic studies have identified numerous cancer driver genes (Kandoth et al., 2013; Lawrence et al., 2014) and germline variants that increase disease susceptibility (Lu et al., 2015). We increasingly understand the molecular determinants of oncogenesis, including tumor suppressor inactivation and pathway alteration. Significant progress has been made in identifying driver mutations (Porta-Pardo et al., 2017), assessing their druggability (Niu et al., 2016), disease subtyping (Waddell et al., 2015), prognosis (Cancer Genome Atlas Research Network et al., 2015), and residual disease detection (Martinez-Lopez et al., 2014).

Gene and protein expression are also key aspects. Studies have reported new fusions (Klijn et al., 2015), alternatively spliced transcripts (Oltean and Bates, 2014), expression-based stratification (Stricker et al., 2017), and implications of viral infection (Cao et al., 2016). Proteomic studies have made progress on subtyping (Lawrence et al., 2015), biomarker identification (Sogawa et al., 2016), and drug sensitivity and resistance (Ji



Joshua M. Stuart,<sup>21</sup> Gordon B. Mills,<sup>22</sup> Matthew A. Wyczalkowski,<sup>1,2</sup> Benjamin G. Vincent,<sup>23,24</sup> Carolyn M. Hutter,<sup>25</sup> Jean Claude Zenklusen,<sup>26</sup> Katherine A. Hoadley,<sup>23,27</sup> Michael C. Wendl,<sup>1,2,3</sup> Ilya Shmulevich,<sup>7</sup> Alexander J. Lazar,<sup>28</sup> David A. Wheeler,<sup>29,30,31,\*</sup> Gad Getz<sup>13,18,32,\*</sup> and The Cancer Genome Atlas Research Network

<sup>22</sup>Department of Systems Biology, University of Texas MD Anderson Cancer Center, Houston, TX 77498, USA

<sup>23</sup>Lineberger Comprehensive Cancer Center, University of North Carolina, Chapel Hill, NC 27599, USA

<sup>24</sup>Curriculum in Bioinformatics and Computational Biology, University of North Carolina, Chapel Hill, NC 27599, USA

<sup>25</sup>National Human Genome Research Institute, Bethesda, MD 20892, USA

<sup>26</sup>National Cancer Institute, Bethesda, MD 20892, USA

<sup>27</sup>Department of Genetics, University of North Carolina, Chapel Hill, NC 27599, USA

<sup>28</sup>Departments of Pathology, Genomic Medicine, and Translational Molecular Pathology, The University of Texas MD Anderson Cancer Center, Houston, TX 77498, USA

<sup>29</sup>Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX 77030, USA

<sup>30</sup>Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX 77030, USA

<sup>31</sup>Dan L Duncan Cancer Center, Baylor College of Medicine, Houston, TX 77030, USA

<sup>32</sup>Massachusetts General Hospital, Boston, MA 02114, USA

<sup>33</sup>These authors contributed equally

<sup>34</sup>Lead Contact

\*Correspondence: [lding@wustl.edu](mailto:lding@wustl.edu) (L.D.), [wheeler@bcm.edu](mailto:wheeler@bcm.edu) (D.A.W.), [gadgetz@broadinstitute.org](mailto:gadgetz@broadinstitute.org) (G.G.)

<https://doi.org/10.1016/j.cell.2018.03.033>

et al., 2017). Advancements have also been made in immune response (Biegging et al., 2014), infiltrate-based subtyping (Akbani et al., 2015), associations of PD-1/PD-L1 with prognosis (Danilova et al., 2016), interactions between immune reprogramming and angiogenesis (Tian et al., 2017), and immune cytolytic activity (Rooney et al., 2015). Each area shows enormous promise.

The era of the first large genome sequences was called the “end of the beginning” of genomics. It seems fitting to call the conclusion of The Cancer Genome Atlas (TCGA) the end of the beginning of cancer genomics. TCGA has systematized large-scale genomics-based cancer research, with its projects and data on 11,000 tumors from 33 cancer types having led to enormous advancements. The TCGA PanCancer Atlas project has a special focus on the oncogenic processes governing cancer development and progression, with its ten analysis working groups (AWGs) presenting their findings. Together we synthesized findings from consensus somatic mutation calling, fusion detection, splicing events, aneuploidy, image analysis, and the immune system in oncogenesis (Figure 1). Here, we concentrate on three themes: (1) interactions between somatic drivers and germline pathogenic variants; (2) links across genomic substrates, i.e., methylome, transcriptome, and proteome; and (3) tumor microenvironment and implications for targeted and immune therapies. We begin each section with an overview from AWG results and follow with additional analyses addressing questions not explored in individual AWG papers. The results of the PanCancer Atlas project will provide a foundation for subsequent phases of deeper, broader, and more sophisticated work that holds great promise for personalized cancer care.

## RESULTS

### Insights into Germline and Somatic Alterations

Previous TCGA studies often concentrated on focal copy-number alterations rather than chromosomal-level aneuploidy. The PanCancer Atlas Aneuploidy AWG systematically quantified aneuploidy (Taylor et al., 2018), correlated its degree with

genomic features, such as *TP53* status, mutational load, and level of lymphocytic infiltrate, and provided experimental evidence confirming some predictions.

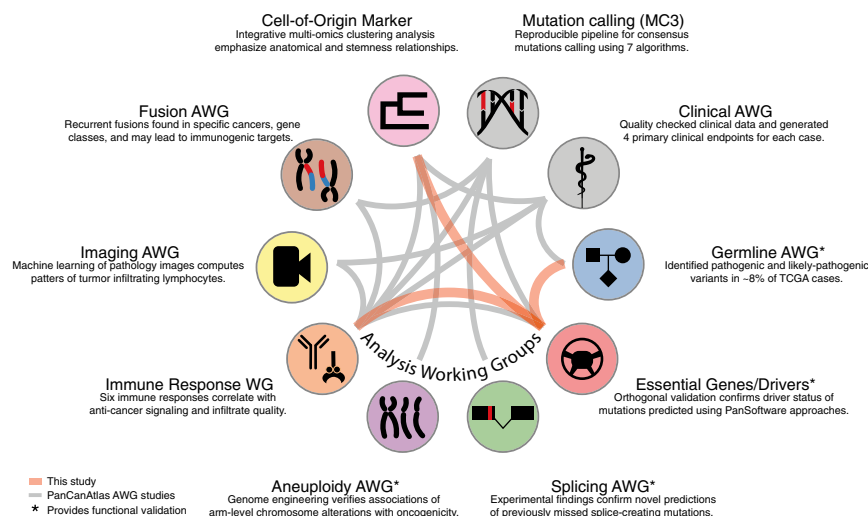
Gene fusions, which can drive overexpression or create fusion proteins, are another important class of drivers. The Fusion AWG systematically characterized fusions (Gao et al., 2018), finding that they are recurrent and disease defining in some neoplasms (e.g., *SS18/SSX1* or *SSX2* fusion in synovial sarcoma). In others, fusion drivers are present in small subsets of tumors (*ALK* or *ROS1* fusions in lung adenocarcinoma). The accompanying mutational events and how they differ among cancers provide functional insights (Gao et al., 2018).

Two other AWGs systematically characterized germline and somatic variants across 33 cancer types (Table S1) (Huang et al., 2018; Ellrott et al., 2018). They generated and analyzed 1.5 billion germline (Huang et al., 2018) and ~3.6 million somatic calls (Ellrott et al., 2018), making TCGA PanCancer Atlas the largest resource for investigating joint variant contributions to cancer. The germline group highlighted the two-hit hypothesis through loss of heterozygosity (LOH) and compound heterozygosity, rare copy-number events, and additional evidence supporting variant pathogenicity. The somatic dataset anchored a comprehensive analysis using 26 bioinformatic tools, identifying 299 driver genes and over 3,400 oncogenic mutations (Bailey et al., 2018). Similarly, the PanCancer Atlas Germline group identified >800 pathogenic or likely pathogenic germline variants in 99 predisposition genes affecting ~8% of all cases (Huang et al., 2018).

### Properties of Oncogenic Germline and Somatic Variants

Here, we used the 299 driver and 99 predisposition genes to study interactions of germline and somatic events in 9,389 samples (STAR Methods; Table S1). Many predisposition genes play roles in genome integrity (Figure 2A, green bars; Table S2). Alterations in these genes represent a higher fraction of germline variants (63%, 490/769) versus somatic drivers (14%, 8850/75825,  $p$  value =  $7e-151$  Fisher's Exact Test), highlighting the role of genome integrity in cancer predisposition. The remaining somatic alterations are largely from genes involved in cell cycle,

## PanCancer Atlas Oncogenic process



**Figure 1. Overview of the PanCancer Atlas Oncogenic Process Group**

PanCan Atlas studies use data from multiple working groups, with relationships shown by gray edges between associated studies. New connections described in this study are shown as orange edges.

### Germline/Somatic-Associated Microsatellite Instability Phenotypes

Many samples (250 out of 1,464) with non-synonymous somatic mutations in DNA mismatch repair (MMR) genes have high microsatellite instability (MSI) status (MSI-sensor score  $\geq 4$ ; [Figure 3C](#); [Table S5](#)) ([Niu et al., 2014](#)). Samples with germline pathogenic variants in MMR genes (18 out of 60) also have high MSI status. Notably, 16 of these 18 samples have both predisposition germline variants

and somatic mutations in MMR genes ([Table S2](#)), representing a population with potentially higher neoantigen load and response to checkpoint-blockade therapy. Indeed, samples with MSI-sensor scores  $\geq 4$  had higher expression of immune-response marker genes (*GZMA*, *PRF1*, *GZMK*, and *GZMH*) in the three cancer types with enough MSI high samples: colon adenocarcinoma and rectum adenocarcinoma (COADREAD), stomach adenocarcinoma (STAD), and uterine corpus endometrial carcinoma (UCEC) (two-sample Kolmogorov-Smirnov  $p < 0.01$ ; [Figure 3D](#)). This highlights the influence of mutations and MMR genes and the MSI phenotype in the immune response against tumors. Finally, using Moonlight we found several pathways that are differentially expressed depending on whether the mutations affecting *BRCA1* and/or *BRCA2* are somatic or germline ([Figures 3E, 3F, and S1](#)). For example, BRCA samples with somatic mutations in *BRCA1/2* downregulate genes involved in antigen processing and leukocyte cytotoxicity, whereas BRCA samples with germline *BRCA1/2* mutations downregulate genes involved in mitochondrial respiratory chain complex and metabolic pathways. The impact of *BRCA1/2* mutations may depend on both their somatic or germline status and the tissue of origin.

### DNA Damage Response Pathway

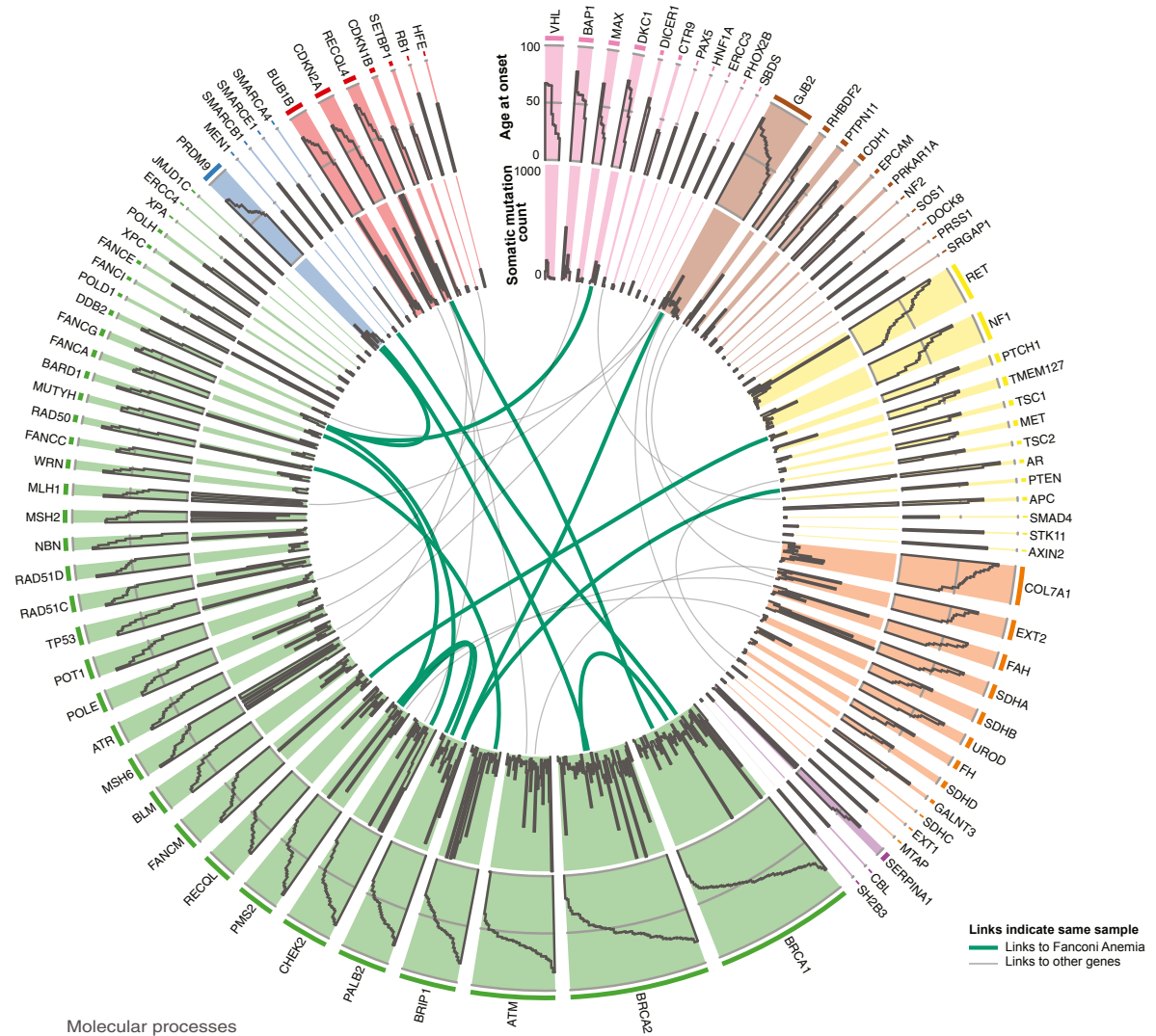
Most predisposition genes affecting genome integrity (64%, 23/36) belong to the Core DDR (DNA damage response) genes ([Knijnenburg et al., 2018](#)) ([Table S2](#)). Several show high germline variant counts, including *BRCA1*, *BRCA2*, *CHEK2*, *ATM*, *BRIP1*, *PALB2*, and *PMS2*. When considering germline and somatic mutations jointly, the most frequently mutated genes are *BRCA1* and *BRCA2*, together having 854 (571 samples) somatic and 153 (152 samples) germline mutations. We grouped samples with germline mutations, somatic, or no/low-impact mutations in these two genes by cancer type to establish associations between age of onset and somatic mutation load. Patients with germline *BRCA1/2* mutations develop cancer at younger ages compared to wild-type samples in OV, LUSC, and BRCA (false discovery rate [FDR]  $9.12e-6$ ,  $9.23e-3$ , and  $1.15e-2$ , respectively, t test). Mean age of diagnosis in patients with germline mutations is  $54.4 \pm 13.0$  years (standard deviation), compared to  $62.3 \pm 13.4$  years when the mutation is somatic across the pan-cancer cohort ( $p$  value =  $2.07e-10$ , 95% confidence interval [CI] =  $-10.27$ ,  $-5.57$ ); [Figure 3A](#); [Table S4](#)). As expected, germline or somatic variants associate with higher mutation load across cancer types ([Figure 3B](#)), being observed in OV samples with germline *BRCA1/2* mutations (FDR  $3e-3$ , t test) and BLCA, STAD somatic (FDR  $5.6e-3$ ,  $9.2e-6$ , t test).

and somatic mutations in MMR genes ([Table S2](#)), representing a population with potentially higher neoantigen load and response to checkpoint-blockade therapy. Indeed, samples with MSI-sensor scores  $\geq 4$  had higher expression of immune-response marker genes (*GZMA*, *PRF1*, *GZMK*, and *GZMH*) in the three cancer types with enough MSI high samples: colon adenocarcinoma and rectum adenocarcinoma (COADREAD), stomach adenocarcinoma (STAD), and uterine corpus endometrial carcinoma (UCEC) (two-sample Kolmogorov-Smirnov  $p < 0.01$ ; [Figure 3D](#)). This highlights the influence of mutations and MMR genes and the MSI phenotype in the immune response against tumors. Finally, using Moonlight we found several pathways that are differentially expressed depending on whether the mutations affecting *BRCA1* and/or *BRCA2* are somatic or germline ([Figures 3E, 3F, and S1](#)). For example, BRCA samples with somatic mutations in *BRCA1/2* downregulate genes involved in antigen processing and leukocyte cytotoxicity, whereas BRCA samples with germline *BRCA1/2* mutations downregulate genes involved in mitochondrial respiratory chain complex and metabolic pathways. The impact of *BRCA1/2* mutations may depend on both their somatic or germline status and the tissue of origin.

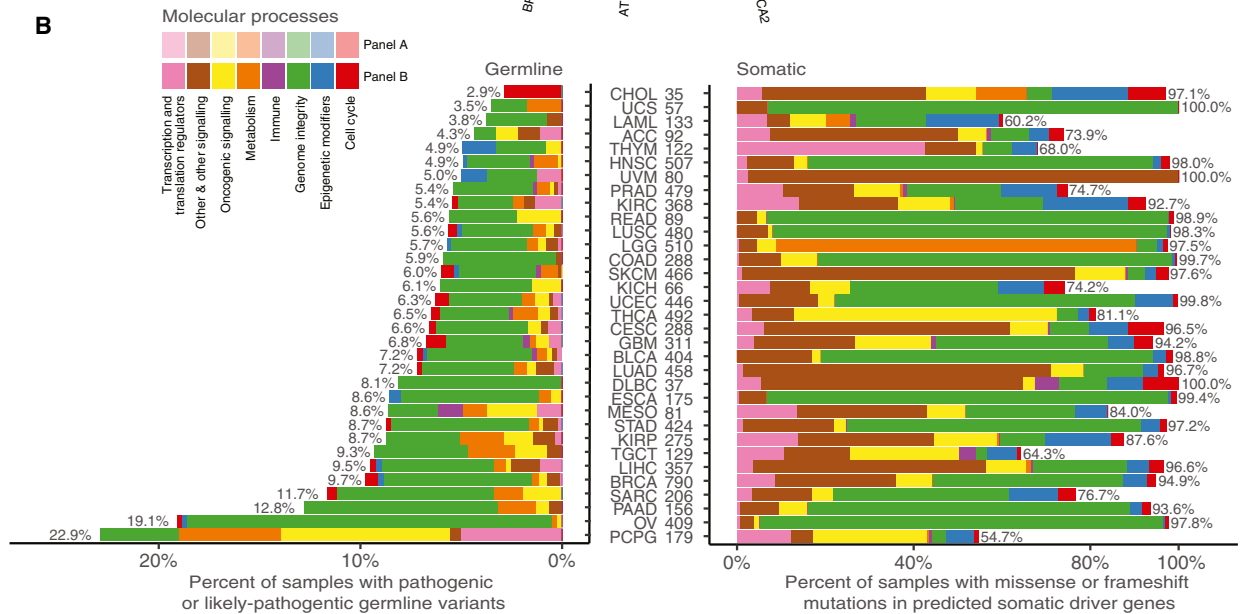
### Somatic-Somatic Interactions

Interactions among somatic driver genes, ranging from sequential dynamics to interactions of pathway and synthetic lethality, hold potential for therapeutic exploitation. We used the MC3 somatic mutation ([Ellrott et al., 2018](#)) dataset and the driver gene list ([Bailey et al., 2018](#)) to identify pairs of drivers that are mutually exclusive or tend to co-occur ([STAR Methods](#)). We found an extensive network of interactions (Cochran-Mantel-Haenszel test FDR  $< 0.1$ ; [Figure 4A](#); [Table S6](#)). *TP53* is the prime hub, co-occurring with *IDH1*, *ATRX*, *PPP2R1A*, *RB1*, and *CDKN2A* and mutually exclusive of *PIK3CA*, *HRAS*, *CTNNB1*, *ARID1A*, and *FGFR3*. As expected, driver genes and mutations that act via certain pathways/mechanisms show strong exclusivity, a primary example being *BRAF* and *HRAS/NRAS/KRAS*, all

A



B



(legend on next page)

of which affect the Ras signaling pathway. Other examples are pairs of homologous genes, such as *IDH1/IDH2* and *GNAQ/GNA11*, and interacting genes, such as *PIK3CA* and *PIK3R1*. These patterns held across virtually all 33 tumor types, indicating discovery of a key oncogenic relationship. We also observed exclusivity in specific tissues (Figure 4B), for example *BRAF*, *NRAS*, and *HRAS* in thyroid carcinoma (THCA) and *GNAQ* and *GNA11* in uveal melanoma.

At a larger scale, some cancer types require cooperation between gene networks. For example, in UCEC, there are two mutually exclusive networks, the first consisting of *TP53* and *PPP2R1A* (and occasionally *PTEN*) and the second *CTNNB1*, *PTEN*, and *CTCF*. This is consistent with previous descriptions of UCEC subtypes, with *TP53*-driven endometrial tumors having a copy-number high phenotype and *PTEN*-driven endometrial tumors being copy-number low or hypermutated (either via MSI and/or *POLE*). Finally, we observed cancer-specific somatic-somatic interactions. For instance, *TP53* and *KRAS* are mutually exclusive in COAD, READ, and LUAD (Table S6) but significantly co-occur in PAAD (Table S6). These observations highlight the importance of investigating both at the pan-cancer level and by tissue of origin (Park and Lehner, 2015).

### Insights into Interactions at -omics Levels

The tumor genome and transcriptome interact at multiple levels. For example, 1%–2% of genome mutations have detectable effects on splicing, with potential to alter the transcriptome and biochemical pathways (Wang and Cooper, 2007). Locally, *cis*-mutations can disrupt or activate splicing factor binding sites or splice sites. The Splicing AWG analyzed 8,656 TCGA tumors, finding that 1,964 mostly missense and synonymous mutations create novel splice junctions (Table S1) (Jayasinghe et al., 2018). They also produce neoantigens, often accompanied by an elevated immune response. Mutations in splice-governing genes result in large-scale abnormal splicing, providing potential biomarkers and therapeutic targets (Dvinge et al., 2016) and acting as proto-oncogenes or tumor suppressors (Yoshida et al., 2011). The Spliceosome Pathway AWG surveyed 33 tumor types for somatic mutations of over 400 splicing factor genes, identifying 119 genes with likely driver mutations (Seiler et al., 2018). They confirmed aberrant splicing of frequently mutated genes, suggesting that splicing de-regulation in cancer is broader than previously reported.

Integrating profiles from individual molecular platforms can provide insights into the molecular state of tumors and identify samples with shared regulation (sample clusters) across multiple assays. A recent analysis (Hoadley et al., 2018) performed clustering of individual platforms and subsequent clustering of cluster assignments (COCA) (Hoadley et al., 2014) on clusters derived from aneuploidy levels (10 clusters; 10,522 samples),

mRNA (25 clusters with at least 40 samples; 10,165 samples), miRNA (microRNA) (15 clusters; 10,170 samples), DNA methylation (25; 10,814), and reverse phase protein array (RPPA) (10; 7,858). They also performed integrative molecular subtyping with the iCluster method (Shen et al., 2009) in a joint analysis of aneuploidy, DNA methylation, mRNA, and miRNA levels across 9,759 tumor samples, identifying 28 iClusters. Consistent with previous multiplatform analyses (Hoadley et al., 2014), samples cluster primarily by tissue of origin.

### Cis- and Trans- Effects of Driver Mutations and Mutation Types

We analyzed the impact of somatic mutations in the *cis*-expression of driver genes. We grouped samples for each gene according to whether they contained frameshift or nonsense mutations (group I), missense (group II), or no mutations (group III). This analysis shows clear upregulation of cancer driver genes affected by missense mutations and downregulation of those affected by nonsense or frameshift mutations (Figures 4C and 4D; Table S7), consistent with previous findings (Hu et al., 2017; Alvarez et al., 2016). We observed reduced expression for tumor suppressors, such as *ATRX*, *BRCA1*, *NF1*, and *RB1*, and elevated expression of oncogenes, like *EGFR* and *KIT* (FDR < 0.1; Figure 4E). We highlight the top 15 genes showing significant expression differences between at least two of the three groups in at least one cancer type (Figures 4F, 4G, and S2). In most cases, the frameshift/nonsense group had significantly lower mRNA than the others, consistent with the hypothesis that they induce nonsense-mediated decay (NMD) (Lindeboom et al., 2016). The exception is *GATA3* in breast cancer, where samples with frameshift or nonsense mutations have higher mRNA levels (FDR = 4.54e–18 Welch's test; Figure 4G), likely because *GATA3* frameshift mutations can have gain-of-function, oncogenic effect (Mair et al., 2016). In cases such as *CASP8*, samples with missense mutations also overexpress the driver gene (FDR < 0.1; Figure 4G).

We used Moonlight to identify gene programs that are differentially expressed in each of the two mutated conditions when compared against non-mutated samples (Figure 4H; Method Details). Remarkably, several genes seem to affect different transcriptional programs, depending on the type of mutation affecting them. Following on the *GATA3* mutations in BRCA, samples with frameshift/nonsense mutations associate with downregulated genes related to microtubule dynamics or organization of cytoskeleton, an effect not seen in those with missense mutations. Similar effects also happen with *CDH1* in BRCA: samples with nonsense and frameshift mutations associate with upregulated genes involved in leukocyte migration but not in samples with missense *CDH1* mutations. The tissue of origin seems to also influence the transcriptional effects. For example, lower grade glioma (LGG) samples with any kind of

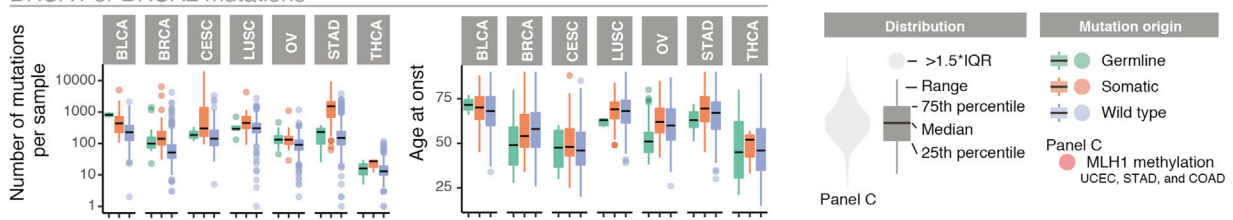
### Figure 2. Sequence-Level Evaluation of Samples with Pathogenic Germline Mutations

(A) Circos plot for each predisposition cancer gene. Width of each slice is proportional to germline-variant frequency. The outermost tier shows age at onset, while middle indicates total number of somatic mutations for each sample. Links designate one sample that has multiple pathogenic or likely pathogenic germline mutations and are green if one of the genes is from the Fanconi anemia pathway.

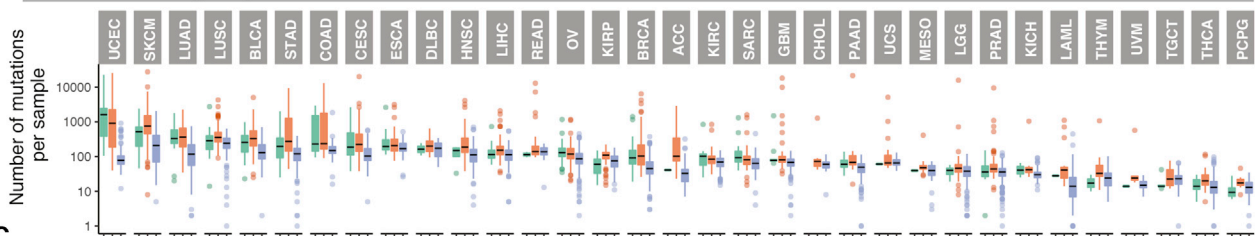
(B) Somatic and germline driver genes grouped into eight molecular process categories. On the x axis, germline and somatic proportions are plotted using number of samples as the denominator. Cancers are sorted by increasing germline contribution.

For a complete list of the TCGA cancer type abbreviations, please see <https://gdc.cancer.gov/resources-tcga-users/tcga-code-tables/tcga-study-abbreviations>.

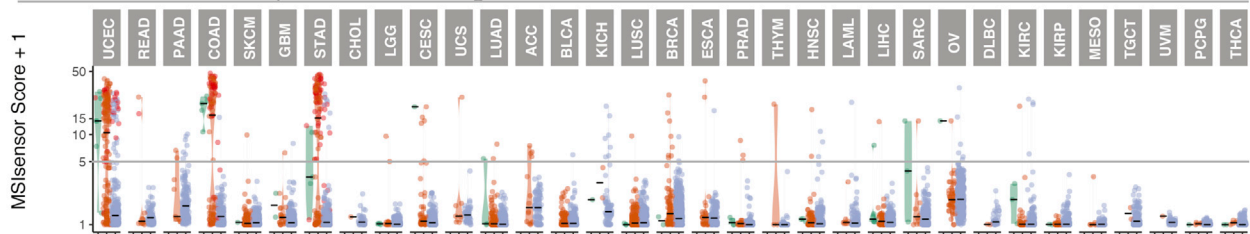
**A** *BRCA1* or *BRCA2* mutations



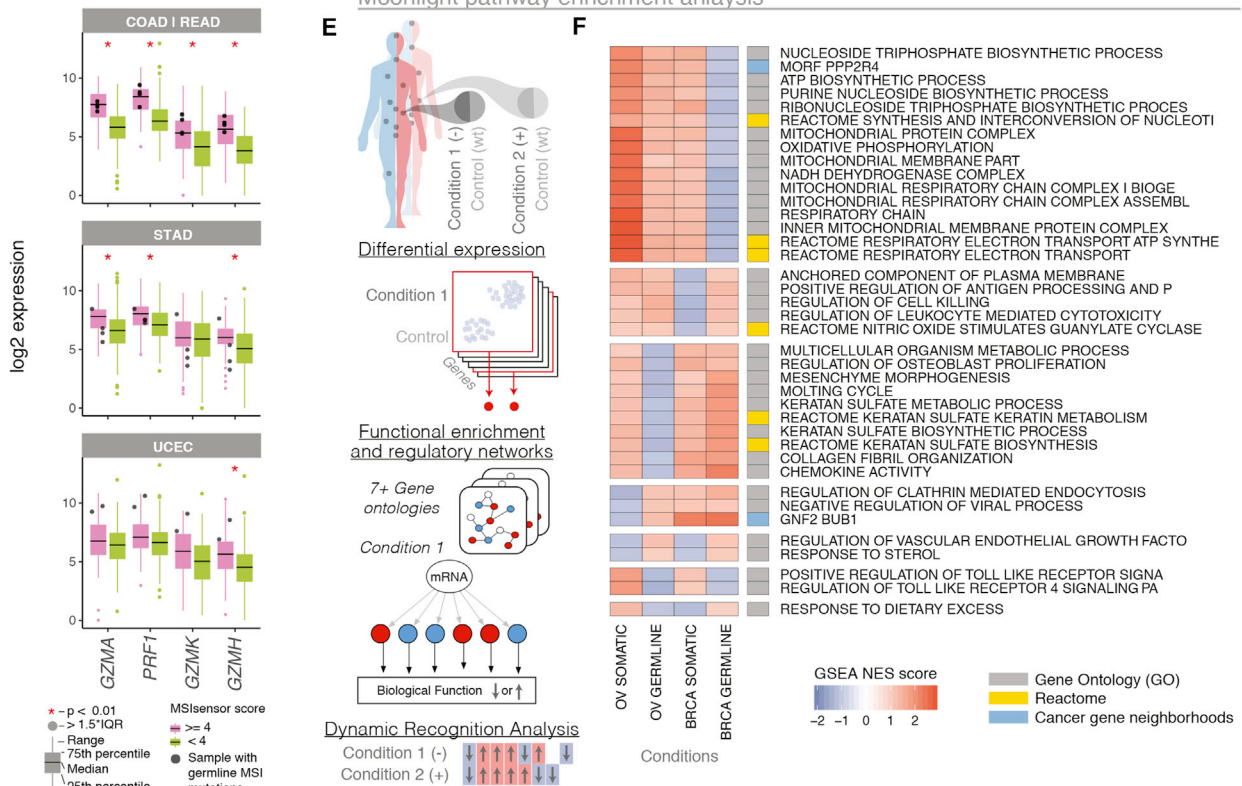
**B** Core DNA damage repair mutations



**C** MSIsensor score for samples with core MSI gene mutations



**D** Moonlight pathway enrichment analysis



(legend on next page)

*TP53* mutations associate with downregulated expression of leukocyte migration genes, but the expression of these genes remains unaltered in LIHC or BRCA samples with *TP53* mutations (Figure 4H). Overall, associations of driver mutations and the transcriptome of the cancer cell seem to be affected by both the original cell type and the type of driver gene mutation.

### Impacts of Genome Mutations on Transcriptomic Activities

Driver mutations often affect the expression of interacting genes and genes in the same pathway. We investigated this phenomenon by integrating protein interaction, transcriptomic, and mutation information using OncoIMPACT (Figure 5A). To reveal key deregulated oncogenic processes occurring in each cancer type, we calculated the fraction of patients for which an oncogenic process was associated with a driver mutation (Figure 5B). With few exceptions (e.g., KIRC), general tumorigenic processes, such as cell proliferation, death, signaling, and motility, are frequently deregulated across cancer types. These processes are mostly deregulated by *TP53*, *PTEN*, *KRAS*, and *PIK3CA*. Processes were more frequently deregulated in some cancers (e.g., head and neck squamous cell carcinoma [HNSC], skin cutaneous melanoma [SKCM], and breast invasive carcinoma [BRCA]). We also observed associations between oncogenic process and cancer types, e.g., Calcium signaling pathway deregulation and uveal melanoma (UVM), with frequent activating mutations in *GNA11* and *GNAQ* that are upstream members of the Calcium signaling pathway (Moore et al., 2016) and frequent deregulation of the Notch signaling pathway in bladder urothelial carcinoma (BLCA) due to inactivating driver mutations in this pathway (Rampias et al., 2014).

We also observed known pairs of significantly mutually exclusive mutated genes such as *TP53* and *PIK3CA* (Kandoth et al., 2013) and *KRAS* and *BRAF* (Loes et al., 2016) in cell death and MAPK signaling processes (Figure 5C; permutation test,  $p$  value  $< 10^{-5}$ ), suggesting that a single driver suffices to perturb these processes and that mutations in multiple drivers are functionally interchangeable in certain contexts. In heterogeneous tumors, this functional redundancy might serve as an important source of drug resistance and metastatic clones.

### Interactions between Different Molecular Layers

Having established the connections between driver events and the transcriptome, we investigated the relationship between

driver genes and the methylomic, transcriptomic, and proteomic profiles of tumors (Figure 6A). We used the clustering data from the Cell of origin AWG (Hoadley et al., 2018) to search for cluster combinations enriched in driver events (Figure 6B), identifying 40 genes associated with multiplatform clusters: *TP53*, *KRAS*, and *PIK3CA* mutations were enriched in ten or more multiplatform clusters, and *ARID1A*, *BRAF*, *CTNNB1*, *KMT2D*, *PTEN*, and *APC* mutations were significantly enriched in four or more clusters (Tables S8 and S9).

Interestingly, we found similar multiplatform clusters that differ in their associated genes. One notable case is comprised of LGG and glioblastoma multiforme (GBM) samples, which are predominantly covered by mRNA cluster 1 and RPPA cluster C1 but which differ markedly in their methylome profiles. *IDH1*-driven LGGs are in methylation cluster 1, where 330 of the 351 samples carried *IDH1* mutations, while *EGFR*-driven LGG and GBM are in methylation cluster 16 (Figure 6C). Another example is that *APC*- and *KRAS*-driven COAD/READ tumors are strongly enriched in mRNA cluster 15 and RPPA cluster C8 but separate in methylation clusters 10 and 11. Similar circumstances are observed for *PIK3CA*-driven BRCA tumors, which are enriched in mRNA and proteome clusters 23 and C6, respectively, but which can belong to methylation clusters 24 or 6 (Table S9).

Notably, we also found instances where specific driver genes differentiate among cluster combinations. For example, UCEC samples belong mostly to multiplatform clusters 4/18/C3 and 23/18/C3, which again differ only in methylation profile (Table S9). The first multi-cluster is enriched in *ARID1A*, *PTEN*, *CTNNB1*, and *PIK3CA* mutations and has fewer *TP53* mutations. The second cluster is conversely dominated by *TP53* and *PPP2R1A* mutations, indicating that differences in driver prevalences can be reflected in the methylation profile (Table S9). While multiplatform clusters are largely driven by tissue of origin (Figure 6D), they may also be affected by the mutations that drive tumor growth.

### Insights into Interactions in the Tumor Microenvironment

A third frontier involves interactions between cancer cells and the tumor microenvironment (TME), comprising stromal cells and the immune infiltrate. Results from the Immune Response Working Group (IRWG) (Thorsson et al., 2018) indicate that the TME can be characterized as belonging to one of six immune subtypes, namely wound healing (C1), IFN- $\gamma$  dominant (C2),

#### Figure 3. Evaluation of *BRCA1/BRCA2*, *DDR*, and *MSI* Genes Using Somatic and Germline Variation

- (A) Samples with *BRCA1* or *BRCA2* mutations are grouped by cancer type and stratified by somatic, germline, or wild-type status. Box-plots highlight mutations per sample (left) and age at onset (right). Outlier samples are plotted as points.
- (B) Box-plots for samples having mutations in DNA damage response genes grouped by cancer.
- (C) Violin plots of *MSI* sensor scores with samples grouped based on mutation status of *MSI* genes. Samples with *MLH1* promoter methylations status are shown in red.
- (D) Gene-expression differences for cytokine activators for three cancer types. Black dots are samples with predisposition germline mutation in *MSI* genes. Red stars highlight significant differences between groups.
- (E) Moonlight workflow shows how samples were stratified based on germline versus wild-type (condition 1) and somatic versus wild-type (condition 2) and integrated across pathways with genes that are labeled as differentially expressed. These were then compared using dynamic recognition analysis to identify patterns.
- (F) Normalized scores from gene set enrichment analysis for germline and somatic mutations in *BRCA1* and/or *BRCA2* only, as conditions of OV and BRCA cancer types. Only the first 50 characters of each pathway are shown (additional information in Figure S1).
- (A, B, and D) Boxplots indicate median *MSI* score with 25<sup>th</sup> and 75<sup>th</sup> percentile hinges and whiskers that extend to  $1.5 \times$  IQR.





inflammatory (C3), lymphocyte depleted (C4), immunologically quiet (C5), and TGF- $\beta$  dominant (C6) (Tables S8 and S10).

While immune signatures can infer levels of lymphocytic infiltrates in tumors, they provide no information on spatial distribution of the lymphocytes. The Imagine Analysis Working Group exploited high-resolution imaging of hematoxylin and eosin (H&E) to estimate tumor-associated lymphocyte densities and infiltration patterns across all samples from 13 of the 33 TCGA tumor types (Saltz et al., 2018). These data revealed relationships between degree of lymphocytic infiltrates measured by gene expression and feature extraction from imaging data using machine learning. Further correlations were made with cancer molecular subtypes, oncogenic events, and outcome, highlighting the power of the underutilized image resources of the TCGA.

### Impact of Driver Mutations on the Immune Communication Network

Here, we further study the relationship between specific driver events, composition of the immune infiltrate, and the signaling network among different cell types within distinct immune subtypes. The networks identified for each immune subtype (STAR Methods) might be relevant to identifying synergistic interventions between targeted drugs and immunotherapies.

*BRAF*-driven tumors have a higher proportion of CD8 T cells than *NRAS*-driven tumors (ANOVA  $p < 2e-5$  in both cases) (Figure 7A; Table S11) in the C3 immune subtype. Elevated CD8 T cell proportion, considered an important effector of checkpoint inhibition (Ji et al., 2012), correlates with better outcomes. We also identified a signaling loop involving CD8 T cells, *CD274* (PD-L1), and *PCDC1* (PD-1) (Method Details) in C3, where targeting *BRAF* and PD-L1 might have synergistic effects. The analysis also reveals an interesting network within the C5 subtype. Samples having mutations in *ATRX* or *TP53* have higher presence of macrophages and lower of CD8 (ANOVA  $p < 2e-8$  in both cases). Interestingly, these macrophages secrete HMGB1, which promotes proliferation and metastasis in glioma (Bassi et al., 2008), a prominent cancer type in C5.

Driver mutations in *KRAS/NRAS/HRAS* and *BRAF* V600 are among the most frequently predicted neoantigens in cancer (Thorsson et al., 2018) and could thus, as presented peptides, be directly steering immune response. Additionally, driver-gene mutations may impact the transcriptional regulation that guides

immune response. For example, *IDH1*-driven gliomas associate with lower levels of STAT1, which can decrease levels of immune infiltrate by ultimately decreasing the secretion of CXCL10, a critical chemokine for T cell trafficking in brain (Kohanbash et al., 2017). Also, models of transcriptional networks (Thorsson et al., 2018) implicate Ras family members and other driver genes in transcriptional control of genes affecting TME composition.

### Mutation Burden and Immune Fraction

Another way in which somatic mutations interact with the immune system is through neoantigens presented on class I or II major histocompatibility complex (MHC) proteins, which can activate immune cells. This has been studied by various PanCancer Atlas groups, describing splice-creating mutations and fusion events creating immunogenic neoantigens (Jaya-singhe et al., 2018; Gao et al., 2018) and neoantigens based on the derived HLA type and their predicted binding affinity (Thorsson et al., 2018).

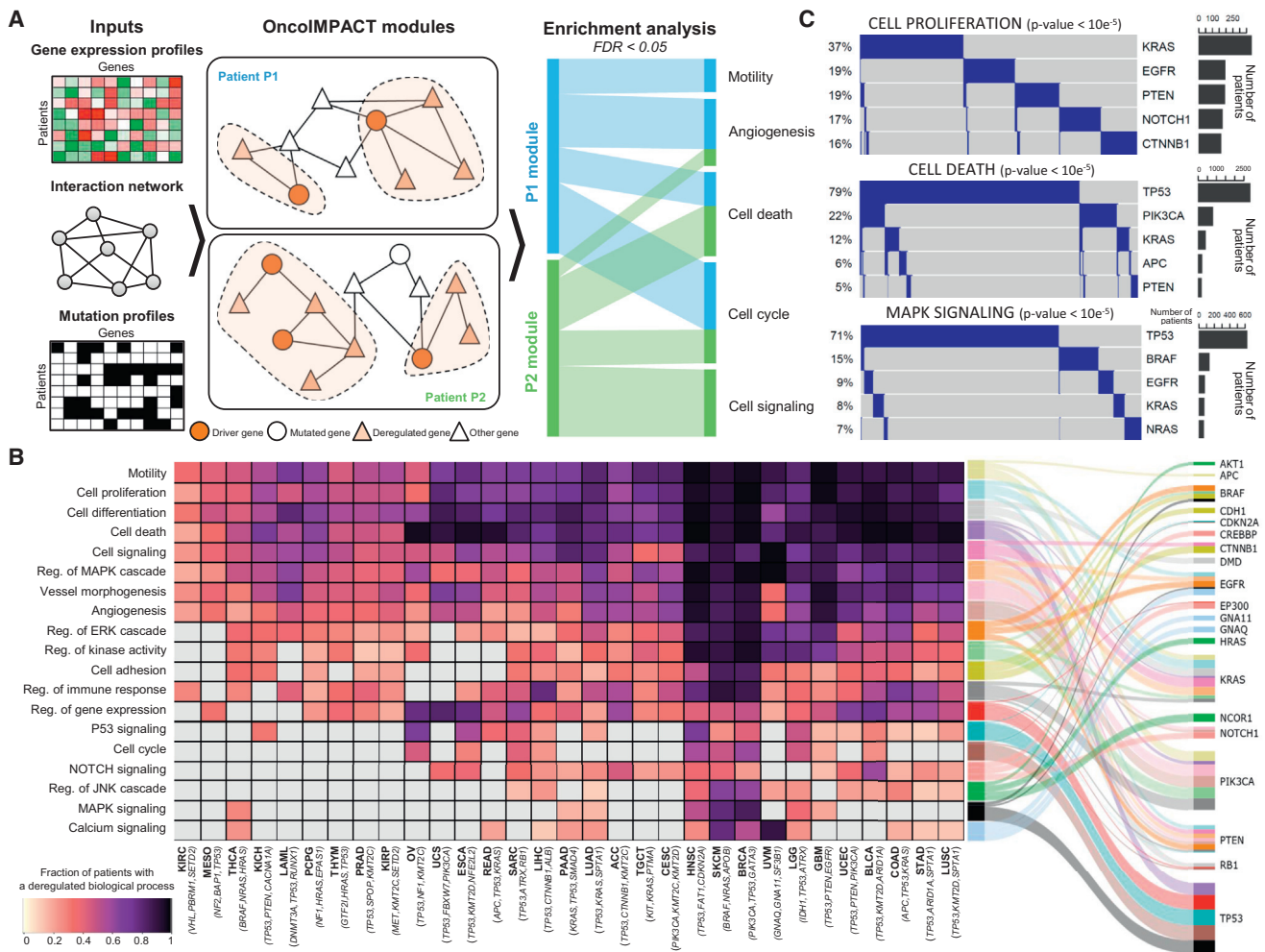
Using neoantigen predictions and immune infiltrate composition, we investigated associations between numbers of presented neoantigens and relative proportion of immune cells comprising immune subtypes (Table S12). These associations differ by immune subtype (Figure 7B). C2 has the greatest overall immune activity. Here, the CD8 T cell fraction increases with neoantigen load (FDR  $< 1e-15$ ; Figure 7C), suggesting that CD8 T cells may respond to neoantigen burden. CD4 T cell fraction and neutrophil fraction increase in relation to neoantigen burden in C3, perhaps reflective of the overall balanced immune response and good prognosis of C3 tumors (FDR  $< 1e-25$ ; Figure 7C). Macrophages have greater infiltration with neoantigen burden in C5, which contains many gliomas and for which TAMs (tumor-associated macrophages) support tumor growth (FDR  $< 5e-3$ ; Figure 7C).

### DISCUSSION

This study summarizes and expands the findings of the TCGA PanCancer Atlas project investigating oncogenic processes. The germline genome has far-ranging, pathway-dependent influences on the somatic landscape, often promoting somatic mutations. Interactions between driver genes and the transcriptome are context dependent, as is the impact of driver mutations in both *cis*- and *trans*-expression. Some oncogenic processes

#### Figure 4. Interactions between Somatic Driver Events

- (A) Mutual exclusivity and co-occurrence of driver events. Nodes sized according to degree and edges colored according to odds ratio of pairs of drivers: red for mutually exclusive (OR  $< 1$ ) and blue for co-occurrence (OR  $> 1$ ).
- (B) Tissue-specific interactions of driver events. Waterfall plots show whether each patient has clonal (dark purple), sub-clonal (light purple), or no driver mutation (gray). Each plot is flanked with a color corresponding to genes in (A).
- (C) Landscape of *cis*-expression changes shown for three mutation types, with FDR  $< 0.1$  considered significant.
- (D) Distribution of t values for gene-expression analyses, with FDR  $< 0.1$  considered significant.
- (E) *Cis*-effects of mutations in expression of driver genes. Gray violin plot depicts expression in all samples of driver gene in the tissue marked below each plot. Red boxes show expression of samples with any mutations in that gene; blue boxes show expression for samples with no mutation in that gene. Each dot represents a sample and is red if there is a copy-number alteration of the gene.
- (F) Same information as in (E), but separating samples according to frameshift and nonsense (green) versus missense mutations (orange). Selected genes show the top-15 t values when comparing between the missense and no-mutation groups (FDR  $< 0.1$ ).
- (G) Same as in (F), but genes selected by top-15 t values between nonsense/frameshift and no-mutations groups.
- (H) Moonlight scores for groups of mutations in driver genes in specific cancer types (y axis) and genes annotated with several gene ontology terms (x axis). Boxes colored red or blue if Moonlight Z-score is positive (overexpression of the biological function) or negative (downregulation), respectively. See also Figure S2.
- (E–G) Boxplots indicate median MSI score with 25<sup>th</sup> and 75<sup>th</sup> percentile hinges and whiskers that extend to 1.5  $\times$  IQR.



**Figure 5. Relationships between Oncogenic Processes and Driver Genes**

(A) Identifying processes deregulated by driver-gene modules using OncoIMPACT. Pathways associated with each module were identified using enrichment analysis (Method Details).

(B) Relationships among oncogenic processes, cancer types, and driver genes. Left: Heatmap shows fraction of samples with deregulated processes associated with sample-specific driver mutations. The three most frequently mutated driver genes are shown with each cancer type. Right: Graph of associations between processes and top three genes predicted to be responsible for their deregulation. Gray cells represent non-significant fraction of patients (binomial test, p value Bonferroni corrected > 0.05). Edge widths represent relative fraction of samples with deregulated processes associated to each driver gene.

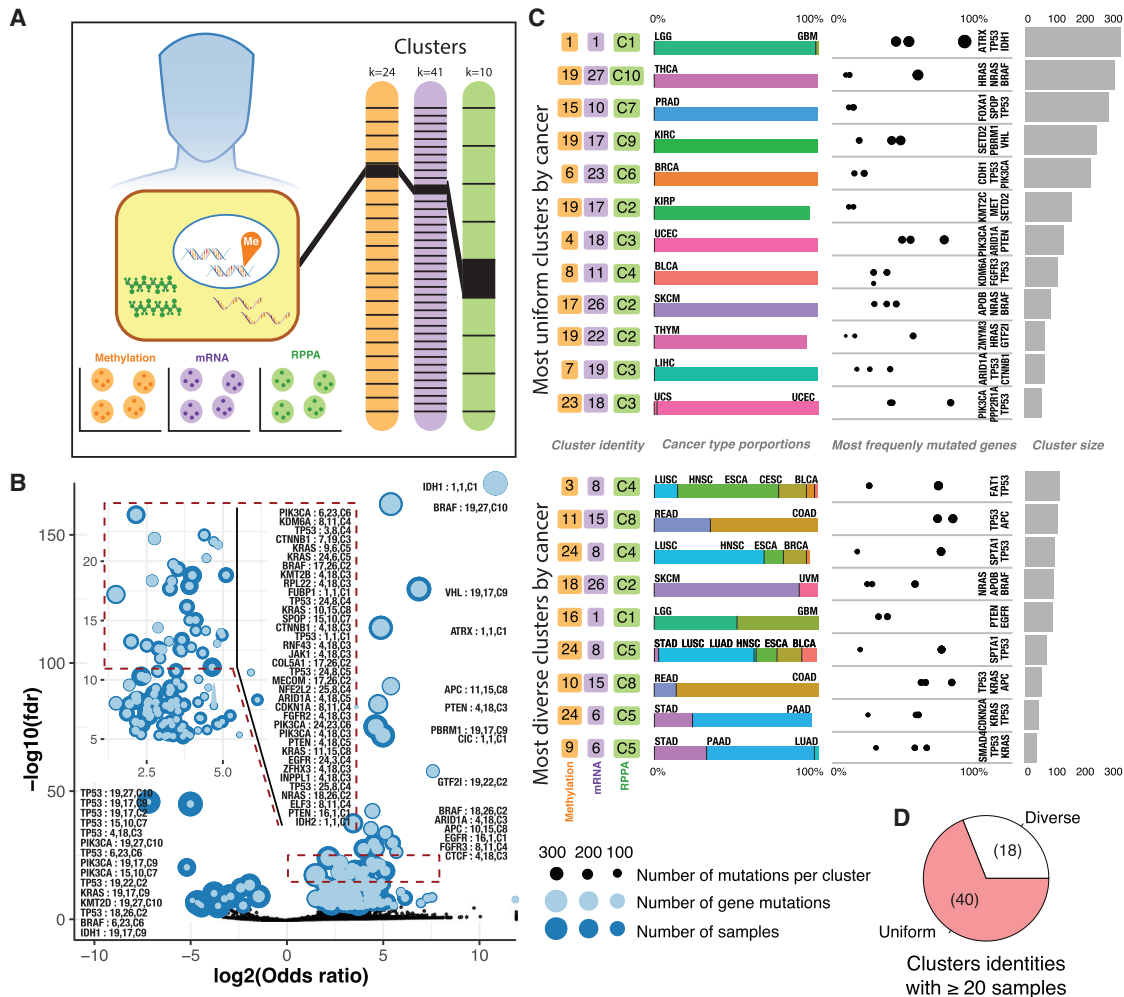
(C) Oncoprint of mutational profile of the five most mutated genes associated with deregulation of three biological processes. Left: Different samples harbor driver genes in a mutually exclusive manner, suggesting many samples have only one process driver gene. Right: Number of samples having driver gene mutated. p values are computed using R-exclusivity test (Method Details).

that tend to be deregulated in few cancer types, such as cell adhesion, are more related to specific genes rather than to prominent drivers. Findings also suggest that networks involving driver mutations, cell types, and cytokines might be used as blueprints for combining two or more immunomodulatory therapies (Tian et al., 2017) in selected tumors.

In summary, this work illuminates the complex milieu of oncogenic processes by integrating an enormous corpus of data obtained over the course of TCGA into organized themes. In effect, biomedical science is now graduating from studying the tumor in isolation to assessing it within its larger environmental context. The findings described here suggest drastic

changes in clinical practice and drug development. For example, molecular treatments will increasingly be developed with “multi-omics.” This strategy is being used to create small molecule inhibitors for druggable mutations (Drilon et al., 2017), mutation signatures (Davies et al., 2017), gene expression (Li et al., 2017), immunotherapeutic agents (Le et al., 2017), and vaccines (Ott et al., 2017). Bioinformatic systems will help efficiently design optimized treatment plans lurking within large combinatorial spaces with respect to dosage, efficacy, side effects, etc.

As we look to the future, there are many questions. For example, we are only beginning to realize that oncogenic



**Figure 6. Complexities of Multidimensional Molecular Evaluation**

(A) Clustering analysis was performed using three substrates: methylation, mRNA, and RPPA. Samples divided into 24 methylation clusters, 41 mRNA, and 10 RPPA clusters. Links show each tumor was given a unique cluster combination identifier.

(B) Gene-enrichment analysis for each cluster assignment is displayed as a volcano plot. Dashed square is enlarged in an inset. Overlapping dots show number of samples in the cluster assignment (dark blue) and the number of samples with a given mutation superimposed (light blue), jointly indicating the mutated proportion in that cluster.

(C) The 21 most gene-enriched cluster identities, with breakdown by tissue-type proportion and most frequently mutated gene from that cluster identity. Sample size for each identity appears in bar plot.

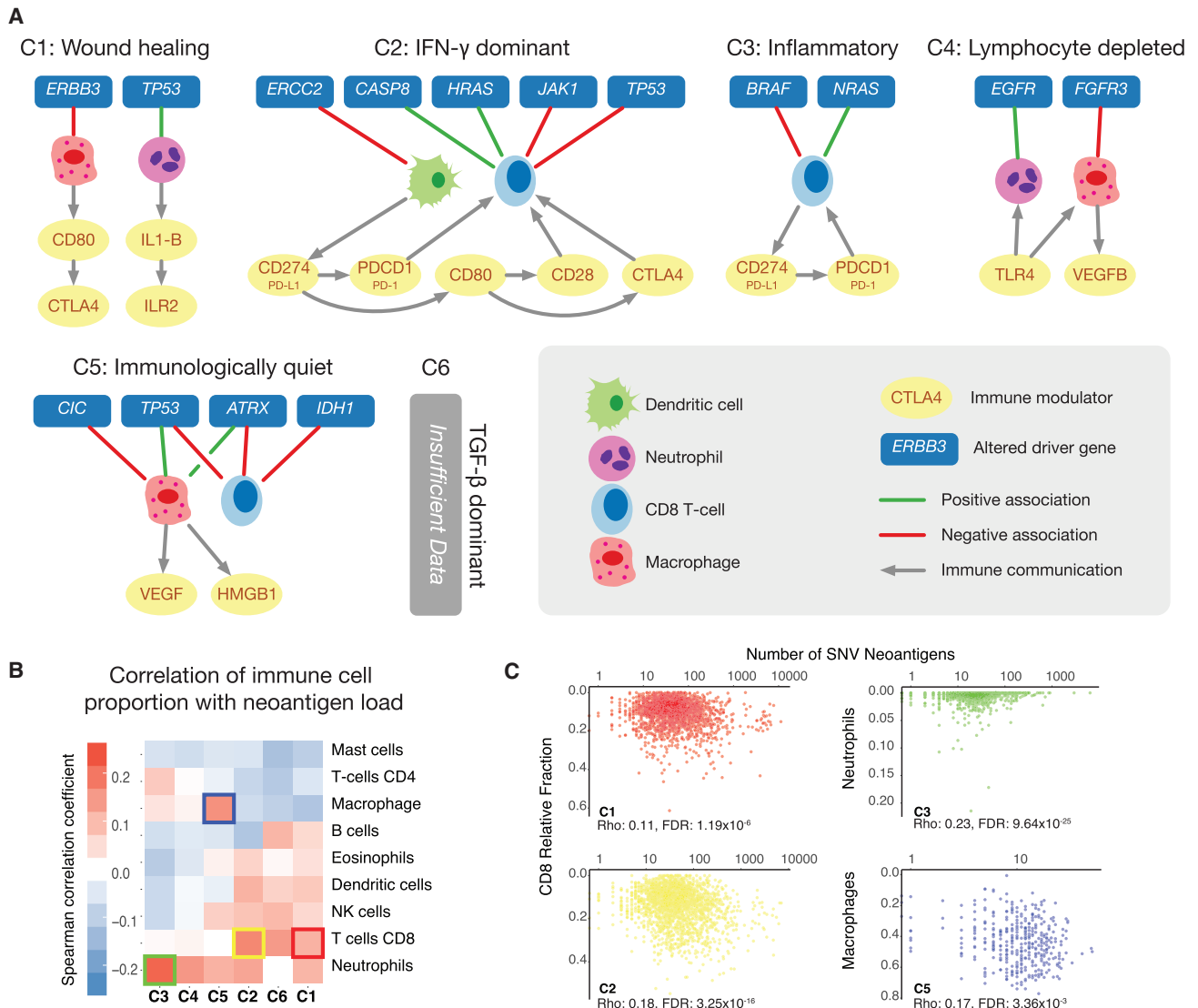
(D) The 58 cluster identities having  $\geq 20$  samples. Pie chart illustrates fraction of uniform clusters, where 90% of samples within a cluster are from a single cancer type.

mutations, such as *BRAF* V600E, frequently occur in healthy people (Martincorena et al., 2015). Could some somatic mutations be tolerated in normal development? If so, how does this impact our understanding of oncogenic mutations? TCGA data come mostly from primary tumors, yet patients usually succumb to metastases; can we find the alterations that drive this process? The next leaps to be taken by the Cancer Moonshot Initiative and Human Tumor Atlas Network (HTAN) will involve pre-cancer, primary, and metastatic tumors associated with treatment sensitivity or resistance and will advance the multidimensional mapping of human cancers over time for informing future cancer research and clinical decision-making.

**STAR★METHODS**

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE
- CONTACT FOR REAGENT AND RESOURCE SHARING
- EXPERIMENTAL MODEL AND SUBJECT DETAILS
- METHOD DETAILS
  - Germline variant calling
  - Somatic variant calling
  - Association testing between biological processes and germline or somatic *BRCA1/2* mutations



**Figure 7. Statistical Associations and Predicted Interactions within the Tumor Microenvironment**

(A) Networks of driver-gene events in distinct cancer-immune subtypes C1–C6 shown in each subpanel. Lines between events and immune cells are green if correlation between immune cell in samples with the driver event is positive and red if negative. Lines between cell types, ligands, and receptors denote interaction pairs known to occur in other contexts and for which there are concordant values across multiple tumor samples in the subtype.

(B) Heatmap shows Spearman correlation between number of predicted neoantigens in each sample of each immune subtype and proportion of different types of immune cells. Colored outline boxes are detailed in the next panel.

(C) In subtypes C1 and C2, proportion of CD8 T cells increases with burden of predicted neoantigens (left two plots). Correlation between number of neoantigens and Neutrophils in samples of C3 subtype (top right) and between number of neoantigens and fraction of macrophages in the TME in samples with C5 immune response (bottom right).

- Germline and somatic gene assignment to pathway analysis
- Detection of gene programs differentially expressed in samples with indels or nonsense mutations and missense mutations
- Identification of biological processes associated with cancer driver genes
- Integration for cell of origin clusters with mutations
- The cell-to-cell communication network
- **QUANTIFICATION AND STATISTICAL ANALYSIS**
  - Comparison of clinical and mutational impact of somatic and germline *BRCA1* and *BRCA2* variants
  - Comparison of clinical and mutational impact of somatic and germline DDR pathway alterations
  - Comparison of clinical and mutational impact of somatic and germline MSI pathway alterations
  - Correlation between MSI scores and expression of immune-related genes

- Mutation mutual exclusivity and co-occurrence analysis
- Association testing between different types of mutations and biological processes
- Correlation between driver events and immune cell types
- **DATA AND SOFTWARE AVAILABILITY**
  - Germline predisposition variant list
  - Driver gene list
  - Cell of origin transcript data
  - Expression and copy number data
  - Cancer Immune Subtypes
  - FANTOM5 network
  - Immune cellular fraction estimates
  - HLA typing and Predicting mutant peptide-MHC binding (neoantigens [pMHCs]) from SNVs
  - CIBERSORT
  - Moonlight
  - domainXplorer
  - OncoIMPACT
  - ABSOLUTE

#### SUPPLEMENTAL INFORMATION

Supplemental Information includes two figures and twelve tables and can be found with this article online at <https://doi.org/10.1016/j.cell.2018.03.033>.

#### ACKNOWLEDGMENTS

We thank patients who contributed to this study and the NCI Office of Cancer Genomics and acknowledge NIH grants U54 HG003273, U54 HG003067, U54 HG003079, U24 CA143799, U24 CA143835, U24 CA143840, U24 CA143843, U24 CA143845, U24 CA143848, U24 CA143858, U24 CA143866, U24 CA143867, U24 CA143882, U24 CA143883, U24 CA144025, U24 CA211006, and P30 CA016672.

#### AUTHOR CONTRIBUTIONS

L.D., G.G., and D.A.W. conceived the project. L.D. supervised the project. M.C.W., A.J.L., E.P.-P., M.H.B., S.S., A.W., K.H., V.T., A.C., D.B., R.J., F.C., L.Y., and L.D. drafted the manuscript. J.M.S., G.B.M., C.M.H., J.C.Z., D.A.W., G.G., and L.D. provided scientific input. M.H.B., M.A.W., and E.P.-P. produced figures. Analysis was performed by M.H.B., E.P.-P., K.H., A.C., C.O., I.C.-C., J.K., C.T., A.W., D.B., C.S., N.N., R.J., F.C., L.Y., K.A.H., R.A., V.T., D.L.G., I.S., B.G.V., and A.J.L. All authors approved submission.

#### DECLARATION OF INTERESTS

Michael Seiler, Peter G. Smith, Ping Zhu, Silvia Buonamici, and Lihua Yu are employees of H3 Biomedicine, Inc. Parts of this work are the subject of a patent application: WO2017040526 titled "Splice variants associated with neomorphic sf3b1 mutants." Shouyoung Peng, Anant A. Agrawal, James Palacino, and Teng Teng are employees of H3 Biomedicine, Inc. Andrew D. Cherniack, Ashton C. Berger, and Galen F. Gao receive research support from Bayer Pharmaceuticals. Gordon B. Mills serves on the External Scientific Review Board of AstraZeneca. Anil Sood is on the Scientific Advisory Board for Kiyatec and is a shareholder in BioPath. Jonathan S. Serody receives funding from Merck, Inc. Kyle R. Covington is an employee of Castle Biosciences, Inc. Preethi H. Gunaratne is founder, CSO, and shareholder of NextmiRNA Therapeutics. Christina Yau is a part-time employee/consultant at NantOmics. Franz X. Schaub is an employee and shareholder of SEngine Precision Medicine, Inc. Carla Grandori is an employee, founder, and shareholder of SEngine Precision Medicine, Inc. Robert N. Eisenman is a member of the Scientific

Advisory Boards and shareholder of Shenogen Pharma and Kronos Bio. Daniel J. Weisenberger is a consultant for Zymo Research Corporation. Joshua M. Stuart is the founder of Five3 Genomics and shareholder of NantOmics. Marc T. Goodman receives research support from Merck, Inc. Andrew J. Gentles is a consultant for Cibermed. Charles M. Perou is an equity stock holder, consultant, and Board of Directors member of BioClassifier and GeneCentric Diagnostics and is also listed as an inventor on patent applications on the Breast PAM50 and Lung Cancer Subtyping assays. Matthew Meyerson receives research support from Bayer Pharmaceuticals; is an equity holder in, consultant for, and Scientific Advisory Board chair for Origimed; and is an inventor of a patent for EGFR mutation diagnosis in lung cancer, licensed to LabCorp. Eduard Porta-Pardo is an inventor of a patent for domainXplorer. Han Liang is a shareholder and scientific advisor of Precision Scientific and Eagle Nebula. Da Yang is an inventor on a pending patent application describing the use of antisense oligonucleotides against specific lncRNA sequence as diagnostic and therapeutic tools. Yonghong Xiao was an employee and shareholder of TESARO, Inc. Bin Feng is an employee and shareholder of TESARO, Inc. Carter Van Waes received research funding for the study of IAP inhibitor ASTX660 through a Cooperative Agreement between NIDCD, NIH, and Astex Pharmaceuticals. Raunaq Malhotra is an employee and shareholder of Seven Bridges, Inc. Peter W. Laird serves on the Scientific Advisory Board for AnchorDx. Joel Tepper is a consultant at EMD Serono. Kenneth Wang serves on the Advisory Board for Boston Scientific, Microtech, and Olympus. Andrea Califano is a founder, shareholder, and advisory board member of DarwinHealth, Inc. and a shareholder and advisory board member of Tempus, Inc. Toni K. Choueiri serves as needed on advisory boards for Bristol-Myers Squibb, Merck, and Roche. Lawrence Kwong receives research support from Array BioPharma. Sharon E. Plon is a member of the Scientific Advisory Board for Baylor Genetics Laboratory. Beth Y. Karlan serves on the Advisory Board of Invitae.

Received: November 17, 2017

Revised: February 20, 2018

Accepted: March 13, 2018

Published: April 5, 2018

#### REFERENCES

- Akbani, R., Akdemir, K.C., Aksoy, B.A., Albert, M., Ally, A., Amin, S.B., Arachchi, H., Arora, A., Auman, J.T., and Ayala, B. (2015). Genomic classification of cutaneous melanoma. *Cell* **161**, 1681–1696.
- Alvarez, M.J., Shen, Y., Giorgi, F.M., Lachmann, A., Ding, B.B., Ye, B.H., and Califano, A. (2016). Functional characterization of somatic mutations in cancer using network-based inference of protein activity. *Nat. Genet.* **48**, 838–847.
- Bailey, M.H., Tokheim, C., Porta-Pardo, E., Sengupta, S., Bertrand, D., Weerasinghe, A., Colaprico, A., Wendl, M.C., Kim, J., Reardon, B., et al. (2018). Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell* **173**. <https://doi.org/10.1016/j.cell.2018.02.060>.
- Bashashati, A., Haffari, G., Ding, J., Ha, G., Lui, K., Rosner, J., Huntsman, D.G., Caldas, C., Aparicio, S.A., and Shah, S.P. (2012). DriverNet: uncovering the impact of somatic driver mutations on transcriptional networks in cancer. *Genome Biol.* **13**, R124.
- Bassi, R., Giussani, P., Anelli, V., Colleoni, T., Pedrazzi, M., Patrone, M., Viani, P., Sparatore, B., Melloni, E., and Riboni, L. (2008). HMGB1 as an autocrine stimulus in human T98G glioblastoma cells: role in cell growth and migration. *J. Neurooncol.* **87**, 23–33.
- Beck, A.H., Espinosa, I., Edris, B., Li, R., Montgomery, K., Zhu, S., Varma, S., Marinelli, R.J., van de Rijn, M., and West, R.B. (2009). The macrophage colony-stimulating factor 1 response signature in breast carcinoma. *Clin Cancer Res.* **15**, 778–787.
- Bertrand, D., Chng, K.R., Sherbaf, F.G., Kiesel, A., Chia, B.K., Sia, Y.Y., Huang, S.K., Hoon, D.S., Liu, E.T., and Hillmer, A. (2015). Patient-specific driver gene prediction and risk assessment through integrated network analysis of cancer omics profiles. *Nucleic Acids Res.* **43**, e44.

- Biegging, K.T., Mello, S.S., and Attardi, L.D. (2014). Unravelling mechanisms of p53-mediated tumour suppression. *Nat. Rev. Cancer* *14*, 359.
- Calabrò, A., Beissbarth, T., Kuner, R., Stojanov, M., Benner, A., Asslaber, M., Ploner, F., Zatloukal, K., Samonigg, H., Poustka, A., et al. (2009). Effects of infiltrating lymphocytes and estrogen receptor on gene expression and prognosis in breast cancer. *Breast Cancer Res Treat.* *116*, 69–77.
- Cancer Genome Atlas Research Network, Brat, D.J., Verhaak, R.G., Aldape, K.D., Yung, W.K., Salama, S.R., Cooper, L.A., Rheinbay, E., Miller, C.R., Vitucci, M., et al. (2015). Comprehensive, integrative genomic analysis of diffuse lower-grade gliomas. *N. Engl. J. Med.* *372*, 2481–2498.
- Cao, S., Wendl, M.C., Wyczalkowski, M.A., Wylie, K., Ye, K., Jayasinghe, R., Xie, M., Wu, S., Niu, B., and Grubb, R., III. (2016). Divergent viral presentation among human tumors and adjacent normal tissues. *Sci. Rep.* *6*, 28294.
- Carter, S.L., Cibulskis, K., Helman, E., McKenna, A., Shen, H., Zack, T., Laird, P.W., Onofrio, R.C., Winckler, W., Weir, B.A., et al. (2012). Absolute quantification of somatic DNA alterations in human cancer. *Nat. Biotechnol.* *30*, 413–421.
- Chang, H.Y., Sneddon, J.B., Alizadeh, A.A., Sood, R., West, R.B., Montgomery, K., Chi, J.T., van de Rijn, M., Bolstein, D., and Brown, P.O. (2004). Gene expression signature of fibroblast serum response predicts human cancer progression: similarities between tumors and wounds. *PLoS Biol.* *2*, E7.
- Chapman, M.A., Lawrence, M.S., Keats, J.J., Cibulskis, K., Sougnez, C., Schinzel, A.C., Harview, C.L., Brunet, J.P., Ahmann, G.J., Adli, M., et al. (2011). Initial genome sequencing and analysis of multiple myeloma. *Nature* *471*, 467–472.
- Cibulskis, K., Lawrence, M.S., Carter, S.L., Sivachenko, A., Jaffe, D., Sougnez, C., Gabriel, S., Meyerson, M., Lander, E.S., and Getz, G. (2013). Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nat. Biotechnol.* *31*, 213–219.
- Colaprico, A., Olsen, C., Cava, C., Terkelsen, T., Silva, T.C., Olsen, A., Cantini, L., Bertoli, G., Zinoviyev, A., Barillot, E., et al. (2018). Moonlight: a tool for biological interpretation and driver genes discovery. *bioRxiv*. <https://doi.org/10.1101/265322>.
- Colaprico, A., Silva, T.C., Olsen, C., Garofano, L., Cava, C., Garolini, D., Sabetdot, T.S., Malta, T.M., Pagnotta, S.M., and Castiglioni, I. (2015). TCGAAbiolinks: an R/Bioconductor package for integrative analysis of TCGA data. *Nucleic Acids Res.* *44*, e71.
- Danilova, L., Wang, H., Sunshine, J., Kaunitz, G.J., Cottrell, T.R., Xu, H., Esandrio, J., Anders, R.A., Cope, L., and Pardoll, D.M. (2016). Association of PD-1/PD-L axis expression with cytolytic activity, mutational load, and prognosis in melanoma and other solid tumors. *Proc. Natl. Acad. Sci. USA* *113*, E7769–E7777.
- Davies, H., Glodzik, D., Morganello, S., Yates, L.R., Staaf, J., Zou, X., Ramakrishna, M., Martin, S., Boyault, S., Sieuwerts, A.M., et al. (2017). HRDetect is a predictor of BRCA1 and BRCA2 deficiency based on mutational signatures. *Nat. Med.* *23*, 517–525.
- Dees, N.D., Zhang, Q., Kandoth, C., Wendl, M.C., Schierding, W., Koboldt, D.C., Mooney, T.B., Callaway, M.B., Dooling, D., and Mardis, E.R. (2012). MuSiC: identifying mutational significance in cancer genomes. *Genome Res.* *22*, 1589–1598.
- Drilon, A., Siena, S., Ou, S.I., Patel, M., Ahn, M.J., Lee, J., Bauer, T.M., Farago, A.F., Wheler, J.J., Liu, S.V., et al. (2017). Safety and antitumor activity of the multitargeted pan-TRK, ROS1, and ALK inhibitor entrectinib: combined results from two phase I trials (ALKA-372-001 and STARTRK-1). *Cancer Discov.* *7*, 400–409.
- Dvinge, H., Kim, E., Abdel-Wahab, O., and Bradley, R.K. (2016). RNA splicing factors as oncoproteins and tumour suppressors. *Nat. Rev. Cancer* *16*, 413–430.
- Elliott, K., Bailey, M.H., Saksena, G., Covington, K.R., Kandoth, C., Stewart, C., Hess, J., Ma, S., McLellan, M., Sofia, H.J., et al. (2018). Scalable open science approach for mutation calling of tumor exomes using multiple genomic pipelines. *Cell Syst.* *6* <https://doi.org/10.1016/j.cels.2018.03.002>.
- Fan, Y., Xi, L., Hughes, D.S., Zhang, J., Zhang, J., Futreal, P.A., Wheeler, D.A., and Wang, W. (2016). MuSE: accounting for tumor heterogeneity using a sample-specific error model improves sensitivity and specificity in mutation calling from sequencing data. *Genome Biol.* *17*, 178.
- Foltz, S.M., Liang, W.-W., Xie, M., and Ding, L. (2017). MIRMMR: binary classification of microsatellite instability using methylation and mutations. *Bioinformatics* *33*, 3799–3801.
- Gao, Q., Liang, W.-W., Foltz, S.M., Mutharasu, G., Jayasinghe, R.G., Cao, S., Liao, W.-W., Reynolds, S.M., Wyczalkowski, M.A., Yao, L., et al. (2018). Driver fusions and their implications in the development and treatment of human cancers. *Cell Rep.* *23* <https://doi.org/10.1016/j.celrep.2018.03.050>.
- Hänzelmann, S., Castelo, R., and Guinney, J. (2013). GSVA: gene set variation analysis for microarray and RNA-seq data. *BMC Bioinformatics* *14*, 7.
- Hoadley, K.A., Yau, C., Wolf, D.M., Cherniack, A.D., Tamborero, D., Ng, S., Leiserson, M.D., Niu, B., McLellan, M.D., and Uzunangelov, V. (2014). Multiplatform analysis of 12 cancer types reveals molecular classification within and across tissues of origin. *Cell* *158*, 929–944.
- Hoadley, K.A., Yau, C., Hinoue, T., Wolf, D.M., Lazar, A.J., Drill, E., Shen, R., Taylor, A.M., Cherniack, A.D., Thorsson, V., et al. (2018). Cell-of-origin patterns dominate the molecular classification of 10,000 tumors from 33 types of cancer. *Cell* *173*. <https://doi.org/10.1016/j.cell.2018.03.022>.
- Hu, Z., Yau, C., and Ahmed, A.A. (2017). A pan-cancer genome-wide analysis reveals tumour dependencies by induction of nonsense-mediated decay. *Nat. Commun.* *8*, 15943.
- Huang, K., Mashl, R.J., Wu, Y., Ritter, D.I., Wang, J., Oh, C., Paczkowska, M., Reynolds, S., Wyczalkowski, M.A., Oak, N., et al. (2018). Pathogenic germline variants in 10,389 adult cancers. *Cell* *173*. <https://doi.org/10.1016/j.cell.2018.03.039>.
- Jayasinghe, R.G., Cao, S., Gao, Q., Wendl, M.C., Vo, N.S., Reynolds, S.M., Zhao, Y., Climente-González, H., Chai, S., Wang, F., et al. (2018). Systematic analysis of splice site-creating mutations in cancer. *Cell Rep.* *23* <https://doi.org/10.1016/j.celrep.2018.03.052>.
- Ji, R.R., Chasalow, S.D., Wang, L., Hamid, O., Schmidt, H., Cogswell, J., Alaparthi, S., Berman, D., Jure-Kunkel, M., Siemers, N.O., et al. (2012). An immune-active tumor microenvironment favors clinical response to ipilimumab. *Cancer Immunol. Immunother.* *61*, 1019–1031.
- Ji, Y., Wei, S., Hou, J., Zhang, C., Xue, P., Wang, J., Chen, X., Guo, X., and Yang, F. (2017). Integrated proteomic and N-glycoproteomic analyses of doxorubicin sensitive and resistant ovarian cancer cells reveal glycoprotein alteration in protein abundance and glycosylation. *Oncotarget* *8*, 13413–13427.
- Kanchi, K.L., Johnson, K.J., Lu, C., McLellan, M.D., Leiserson, M.D., Wendl, M.C., Zhang, Q., Koboldt, D.C., Xie, M., Kandoth, C., et al. (2014). Integrated analysis of germline and somatic variants in ovarian cancer. *Nat Commun* *5*, 3156.
- Kandoth, C., McLellan, M.D., Vandin, F., Ye, K., Niu, B., Lu, C., Xie, M., Zhang, Q., McMichael, J.F., Wyczalkowski, M.A., et al. (2013). Mutational landscape and significance across 12 major cancer types. *Nature* *502*, 333–339.
- Klijn, C., Durinck, S., Stawiski, E.W., Haverty, P.M., Jiang, Z., Liu, H., Degenhardt, J., Mayba, O., Gnadt, F., Liu, J., et al. (2015). A comprehensive transcriptional portrait of human cancer cell lines. *Nat. Biotechnol.* *33*, 306–312.
- Knijnenburg, T., Wang, L., Zimmermann, M., Chambwe, N., Gao, G., Cherniack, A., Fan, H., Shen, H., Way, G., Greene, C., et al. (2018). Genomic and molecular landscape of DNA damage repair deficiency across The Cancer Genome Atlas. *Cell Rep.* *23* <https://doi.org/10.1016/j.celrep.2018.03.076>.
- Koboldt, D.C., Zhang, Q., Larson, D.E., Shen, D., McLellan, M.D., Lin, L., Miller, C.A., Mardis, E.R., Ding, L., and Wilson, R.K. (2012). VarScan 2: somatic mutation and copy number alteration discovery in cancer by exome sequencing. *Genome Res.* *22*, 568–576.
- Kohanbash, G., Carrera, D.A., Shrivastav, S., Ahn, B.J., Jahan, N., Mazor, T., Chheda, Z.S., Downey, K.M., Watchmaker, P.B., Beppler, C., et al. (2017). Isocitrate dehydrogenase mutations suppress STAT1 and CD8+ T cell accumulation in gliomas. *J. Clin. Invest.* *127*, 1425–1437.

- Krzywinski, M., Schein, J., Birol, I., Connors, J., Gascoyne, R., Horsman, D., Jones, S.J., and Marra, M.A. (2009). Circos: an information aesthetic for comparative genomics. *Genome Res.* *19*, 1639–1645.
- Langfelder, P., and Horvath, S. (2008). WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics* *9*, 559.
- Larson, D.E., Harris, C.C., Chen, K., Koboldt, D.C., Abbott, T.E., Dooling, D.J., Ley, T.J., Mardis, E.R., Wilson, R.K., and Ding, L. (2012). SomaticSniper: identification of somatic point mutations in whole genome sequencing data. *Bioinformatics* *28*, 311–317.
- Lawrence, M.S., Stojanov, P., Mermel, C.H., Robinson, J.T., Garraway, L.A., Golub, T.R., Meyerson, M., Gabriel, S.B., Lander, E.S., and Getz, G. (2014). Discovery and saturation analysis of cancer genes across 21 tumour types. *Nature* *505*, 495–501.
- Lawrence, R.T., Perez, E.M., Hernández, D., Miller, C.P., Haas, K.M., Irie, H.Y., Lee, S.-I., Blau, C.A., and Villén, J. (2015). The proteomic landscape of triple-negative breast cancer. *Cell Rep.* *11*, 630–644.
- Le, D.T., Durham, J.N., Smith, K.N., Wang, H., Bartlett, B.R., Aulakh, L.K., Lu, S., Kemberling, H., Wilt, C., Luber, B.S., et al. (2017). Mismatch repair deficiency predicts response of solid tumors to PD-1 blockade. *Science* *357*, 409–413.
- Leiserson, M.D., Reyna, M.A., and Raphael, B.J. (2016). A weighted exact test for mutually exclusive mutations in cancer. *Bioinformatics* *32*, i736–i745.
- Li, B., and Dewey, C.N. (2011). RSEM: accurate transcript quantification from RNA-seq data with or without a reference genome. *BMC Bioinformatics* *12*, 323.
- Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* *25*, 1754–1760.
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., and Durbin, R.; 1000 Genome Project Data Processing Subgroup (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics* *25*, 2078–2079.
- Li, L., Karanika, S., Yang, G., Wang, J., Park, S., Broom, B.M., Manyam, G.C., Wu, W., Luo, Y., Basourakos, S., et al. (2017). Androgen receptor inhibitor-induced “BRCAness” and PARP inhibition are synthetically lethal for castration-resistant prostate cancer. *Sci. Signal.* *10*, eaam7479.
- Lindeboom, R.G., Supek, F., and Lehner, B. (2016). The rules and impact of nonsense-mediated mRNA decay in human cancers. *Nat. Genet.* *48*, 1112–1118.
- Liu, J., Lichtenberg, T., Hoadley, K.A., Poisson, L.M., Lazar, A.J., Cherniack, A.D., Kovach, A.J., Benz, C.C., Levine, D.A., Lee, A.V., et al. (2018). An Integrated TCGA Pan-Cancer Clinical Data Resource to drive high quality survival outcome analytics. *Cell* *173*. <https://doi.org/10.1016/j.cell.2018.02.052>.
- Lizio, M., Harshbarger, J., Shimoji, H., Severin, J., Kasukawa, T., Sahin, S., Abugessaisa, I., Fukuda, S., Hori, F., Ishikawa-Kato, S., et al. (2015). Gateways to the FANTOM5 promoter level mammalian expression atlas. *Genome Biol.* *16*, 22.
- Loes, I.M., Immervoll, H., Sorbye, H., Angelsen, J.H., Horn, A., Knappskog, S., and Lonning, P.E. (2016). Impact of KRAS, BRAF, PIK3CA, TP53 status and intraindividual mutation heterogeneity on outcome after liver resection for colorectal cancer metastases. *Int. J. Cancer* *139*, 647–656.
- Lu, C., Xie, M., Wendl, M.C., Wang, J., McLellan, M.D., Leiserson, M.D., Huang, K.-I., Wyczalkowski, M.A., Jayasinghe, R., and Banerjee, T. (2015). Patterns and functional implications of rare germline variants across 12 cancer types. *Nat Commun* *6*, 10086.
- Mair, B., Konopka, T., Kerzendorfer, C., Sleiman, K., Salic, S., Serra, V., Mueller, M.K., Theodorou, V., and Nijman, S.M. (2016). Gain-and loss-of-function mutations in the breast cancer gene GATA3 result in differential drug sensitivity. *PLoS Genet.* *12*, e1006279.
- Martincorena, I., Roshan, A., Gerstung, M., Ellis, P., Van Loo, P., McLaren, S., Wedge, D.C., Fullam, A., Alexandrov, L.B., and Tubio, J.M. (2015). High burden and pervasive positive selection of somatic mutations in normal human skin. *Science* *348*, 880–886.
- Martinez-Lopez, J., Lahuerta, J.J., Pepin, F., Gonzalez, M., Barrio, S., Ayala, R., Puig, N., Montalban, M.A., Paiva, B., Weng, L., et al. (2014). Prognostic value of deep sequencing method for minimal residual disease detection in multiple myeloma. *Blood* *123*, 3073–3079.
- Mashi, R.J., Scott, A.D., Huang, K.L., Wyczalkowski, M.A., Yoon, C.J., Niu, B., DeNardo, E., Yellapantula, V.D., Handsaker, R.E., Chen, K., et al. (2017). GenomeVIP: a cloud platform for genomic variant discovery and interpretation. *Genome Res.* *27*, 1450–1459.
- McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernysky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M., et al. (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* *20*, 1297–1303.
- Moore, A.R., Ceraudo, E., Sher, J.J., Guan, Y., Shoushtari, A.N., Chang, M.T., Zhang, J.Q., Walczak, E.G., Kazmi, M.A., Taylor, B.S., et al. (2016). Recurrent activating mutations of G-protein-coupled receptor CYSLTR2 in uveal melanoma. *Nat. Genet.* *48*, 675–680.
- Mularoni, L., Sabarinathan, R., Deu-Pons, J., Gonzalez-Perez, A., and López-Bigas, N. (2016). OncodriveFML: a general framework to identify coding and non-coding regions with cancer driver mutations. *Genome Biol.* *17*, 128.
- Newman, A.M., Liu, C.L., Green, M.R., Gentles, A.J., Feng, W., Xu, Y., Hoang, C.D., Diehn, M., and Alizadeh, A.A. (2015). Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods* *12*, 453–457.
- Nielsen, M., and Andreatta, M. (2016). NetMHCpan-3.0; improved prediction of binding to MHC class I molecules integrating information from multiple receptor and peptide length datasets. *Genome Med* *8*, 33.
- Niu, B., Scott, A.D., Sengupta, S., Bailey, M.H., Batra, P., Ning, J., Wyczalkowski, M.A., Liang, W.W., Zhang, Q., McLellan, M.D., et al. (2016). Protein-structure-guided discovery of functional mutations across 19 cancer types. *Nat. Genet.* *48*, 827–837.
- Niu, B., Ye, K., Zhang, Q., Lu, C., Xie, M., McLellan, M.D., Wendl, M.C., and Ding, L. (2014). MSIsensor: microsatellite instability detection using paired tumor-normal sequence data. *Bioinformatics* *30*, 1015–1016.
- Oltean, S., and Bates, D. (2014). Hallmarks of alternative splicing in cancer. *Oncogene* *33*, 5311.
- Ott, P.A., Hu, Z., Keskin, D.B., Shukla, S.A., Sun, J., Bozym, D.J., Zhang, W., Luoma, A., Giobbie-Hurder, A., Peter, L., et al. (2017). An immunogenic personal neoantigen vaccine for patients with melanoma. *Nature* *547*, 217–221.
- Park, S., and Lehner, B. (2015). Cancer type-dependent genetic interactions between cancer driver alterations indicate plasticity of epistasis across cell types. *Mol. Syst. Biol.* *11*, 824.
- Porta-Pardo, E., and Godzik, A. (2014). e-Driver: a novel method to identify protein regions driving cancer. *Bioinformatics* *30*, 3109–3114.
- Porta-Pardo, E., and Godzik, A. (2016). Mutation drivers of immunological responses to cancer. *Cancer Immunol. Res.* *4*, 789–798.
- Porta-Pardo, E., Kamburov, A., Tamborero, D., Pons, T., Grases, D., Valencia, A., Lopez-Bigas, N., Getz, G., and Godzik, A. (2017). Comparison of algorithms for the detection of cancer drivers at subgene resolution. *Nat. Methods* *14*, 782–788.
- Radenbaugh, A.J., Ma, S., Ewing, A., Stuart, J.M., Collisson, E.A., Zhu, J., and Haussler, D. (2014). RADIA: RNA and DNA integrated analysis for somatic mutation detection. *PLoS ONE* *9*, e111516.
- Rampias, T., Vgenopoulou, P., Avgeris, M., Polyzos, A., Stravodimos, K., Valavanis, C., Scorilas, A., and Klinakis, A. (2014). A new tumor suppressor role for the Notch pathway in bladder cancer. *Nat. Med.* *20*, 1199–1205.
- Reimand, J., and Bader, G.D. (2013). Systematic analysis of somatic mutations in phosphorylation signaling predicts novel cancer drivers. *Mol. Syst. Biol.* *9*, 637.
- Rooney, M.S., Shukla, S.A., Wu, C.J., Getz, G., and Hacohen, N. (2015). Molecular and genetic properties of tumors associated with local immune cytolytic activity. *Cell* *160*, 48–61.
- Saltz, J.H., Gupta, R., Hou, L., Kurc, T., Singh, P., Nguyen, V., Samaras, D., Shroyer, K.R., Zhao, T., Batiste, R., et al. (2018). Spatial organization and molecular correlation of tumor-infiltrating lymphocytes using deep learning



- on pathology images. *Cell Rep.* 23 <https://doi.org/10.1016/j.celrep.2018.03.086>.
- Scrucca, L., Fop, M., Murphy, T.B., and Raftery, A.E. (2016). mclust 5: clustering, classification and density estimation using Gaussian finite mixture models. *R J* 8, 289–317.
- Seiler, M., Peng, S., Agrawal, A.A., Palacino, J., Teng, T., Zhu, P., Smith, P.G., The Cancer Genome Atlas Research Network, Buonamici, S., Yu, L., et al. (2018). Somatic mutational landscape of splicing factor genes and their functional consequences across 33 cancer types. *Cell Rep.* 23 <https://doi.org/10.1016/j.celrep.2018.01.088>.
- Shen, R., Olshen, A.B., and Ladanyi, M. (2009). Integrative clustering of multiplatform genomic data types using a joint latent variable model with application to breast and lung cancer subtype analysis. *Bioinformatics* 25, 2906–2912.
- Silva, T.C., Colaprico, A., Olsen, C., D'Angelo, F., Bontempi, G., Ceccarelli, M., and Noushmehr, H. (2016). TCGA Workflow: Analyze cancer genomics and epigenomics data using Bioconductor packages. *F1000Res.* 5 <https://doi.org/10.12688/f1000research.8923.2>.
- Siragusa, E., Weese, D., and Reinert, K. (2013). Fast and accurate read mapping with approximate seeds and multiple backtracking. *Nucleic Acids Res.* 41, e78.
- Sjöblom, T., Jones, S., Wood, L.D., Parsons, D.W., Lin, J., Barber, T.D., Mandelker, D., Leary, R.J., Ptak, J., Silliman, N., et al. (2006). The consensus coding sequences of human breast and colorectal cancers. *Science* 314, 268–274.
- Sogawa, K., Takano, S., Iida, F., Satoh, M., Tsuchida, S., Kawashima, Y., Yoshitomi, H., Sanda, A., Kadera, Y., Takizawa, H., et al. (2016). Identification of a novel serum biomarker for pancreatic cancer, C4b-binding protein alpha-chain (C4BPA) by quantitative proteomic analysis using tandem mass tags. *Br. J. Cancer* 115, 949–956.
- Stricker, T.P., Brown, C.D., Bandlamudi, C., McNeerney, M., Kittler, R., Montoya, V., Peterson, A., Grossman, R., and White, K.P. (2017). Robust stratification of breast cancer subtypes using differential patterns of transcript isoform expression. *PLoS Genet.* 13, e1006589.
- Szolek, A., Schubert, B., Mohr, C., Sturm, M., Feldhahn, M., and Kohlbacher, O. (2014). OptiType: precision HLA typing from next-generation sequencing data. *Bioinformatics* 30, 3310–3316.
- Tamborero, D., Gonzalez-Perez, A., and Lopez-Bigas, N. (2013). Oncodrive-CLUST: exploiting the positional clustering of somatic mutations to identify cancer genes. *Bioinformatics* 29, 2238–2244.
- Tatlow, P.J., and Piccolo, S.R. (2016). A cloud-based workflow to quantify transcript-expression levels in public cancer compendia. *Sci Rep* 6, 39259.
- Taylor, A.M., Shih, J., Ha, G., Gao, G.F., Zhang, X., Berger, A.C., Schumacher, S.E., Wang, C., Hu, H., Liu, J., Lazar, A.J., The Cancer Genome Atlas Research Network, Cherniack, A.D., Beroukhi, R., and Meyerson, M. (2018). Genomic and functional approaches to understanding cancer aneuploidy. *Cancer Cell* 33. <https://doi.org/10.1016/j.ccell.2018.03.007>.
- Teschendorff, A.E., Gomez, S., Arenas, A., El-Ashry, D., Schmidt, M., Gehrman, M., and Caldas, C. (2010). Improved prognostic classification of breast cancer defined by antagonistic activation patterns of immune response pathway modules. *BMC Cancer* 10, 604.
- Thorsson, V., Gibbs, D.L., Brown, S.D., Wolf, D., Bortone, D.S., Yang, T.-H.O., Porta-Pardo, E., Gao, G., Plaisier, C.L., Eddy, J.A., et al. (2018). The immune landscape of cancer. *Immunity* 48. <https://doi.org/10.1016/j.immuni.2018.03.023>.
- Tian, L., Goldstein, A., Wang, H., Ching Lo, H., Sun Kim, I., Welte, T., Sheng, K., Dobrolecki, L.E., Zhang, X., Putluri, N., et al. (2017). Mutual regulation of tumour vessel normalization and immunostimulatory reprogramming. *Nature* 544, 250–254.
- Tokheim, C.J., Papadopoulos, N., Kinzler, K.W., Vogelstein, B., and Karchin, R. (2016). Evaluating the evaluation of cancer driver genes. *Proc. Natl. Acad. Sci. USA* 113, 14330–14335.
- Waddell, N., Pajic, M., Patch, A.M., Chang, D.K., Kassahn, K.S., Bailey, P., Johns, A.L., Miller, D., Nones, K., Quek, K., et al. (2015). Whole genomes redefined the mutational landscape of pancreatic cancer. *Nature* 518, 495–501.
- Wang, G.-S., and Cooper, T.a. (2007). Splicing in disease: disruption of the splicing code and the decoding machinery. *Nat. Rev. Genet.* 8, 749–761.
- Wang, K., Singh, D., Zeng, Z., Coleman, S.J., Huang, Y., Savich, G.L., He, X., Mieczkowski, P., Grimm, S.A., Perou, C.M., et al. (2010). MapSplice: Accurate mapping of RNA-seq reads for splice junction discovery. *Nucleic Acids Res.* 38, e178.
- Wilkerson, M.D., and Hayes, D.N. (2010). ConsensusClusterPlus: a class discovery tool with confidence assessments and item tracking. *Bioinformatics* 26, 1572–1573.
- Wolf, D.M., Lenburg, M.E., Yau, C., Boudreau, A., and van 't Veer, L.J. (2014). Gene co-expression modules as clinically relevant hallmarks of breast cancer diversity. *PLoS ONE* 9, e88309.
- Ye, K., Wang, J., Jayasinghe, R., Lameijer, E.W., McMichael, J.F., Ning, J., McLellan, M.D., Xie, M., Cao, S., Yellapantula, V., et al. (2016). Systematic discovery of complex insertions and deletions in human cancers. *Nat. Med.* 22, 97–104.
- Yoshida, K., Sanada, M., Shiraishi, Y., Nowak, D., Nagata, Y., Yamamoto, R., Sato, Y., Sato-Otsubo, A., Kon, A., Nagasaki, M., et al. (2011). Frequent pathway mutations of splicing machinery in myelodysplasia. *Nature* 478, 64–69.

## STAR★METHODS

### KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Deposited Data		
Public MC3 MAF	Ellrott et al., 2018	<a href="https://gdc.cancer.gov/about-data/publications">https://gdc.cancer.gov/about-data/publications</a>
TCGA Clinical data	Liu et al., 2018	<a href="https://gdc.cancer.gov/about-data/publications">https://gdc.cancer.gov/about-data/publications</a>
Germline genes used	Huang et al., 2018	Table S2; <a href="https://gdc.cancer.gov/about-data/publications">https://gdc.cancer.gov/about-data/publications</a>
Pan-Immune clusters and immune infiltrates.	Thorsson et al., 2018	Table S12; <a href="https://gdc.cancer.gov/about-data/publications">https://gdc.cancer.gov/about-data/publications</a>
Cell-of-Origin cluster	Hoadley et al., 2018	Table S8; <a href="https://gdc.cancer.gov/about-data/publications">https://gdc.cancer.gov/about-data/publications</a>
DNA Damage Response Genes	Knijnenburg et al., 2018	Table S2; <a href="https://gdc.cancer.gov/about-data/publications">https://gdc.cancer.gov/about-data/publications</a>
Essential Genes/Drivers genes used	Bailey et al., 2018	Table S2; <a href="https://gdc.cancer.gov/about-data/publications">https://gdc.cancer.gov/about-data/publications</a>
Software and Algorithms		
domainXplorer	Porta-Pardo and Godzik, 2016	<a href="https://github.com/eduardporta/domainXplorer">https://github.com/eduardporta/domainXplorer</a>
MSIsensor	Niu et al., 2014	<a href="https://github.com/ding-lab/msisensor">https://github.com/ding-lab/msisensor</a>
Moonlight	Colaprico, et al. 2018	<a href="https://www.bioconductor.org/packages/devel/bioc/vignettes/MoonlightR/inst/doc/Moonlight.html">https://www.bioconductor.org/packages/devel/bioc/vignettes/MoonlightR/inst/doc/Moonlight.html</a>
OncolMPACT	Bertrand et al. 2015	<a href="https://github.com/CSB5/OncolMPACT">https://github.com/CSB5/OncolMPACT</a>
ABSOLUTE	Carter et al. 2012	<a href="http://archive.broadinstitute.org/cancer/cga/ABSOLUTE">http://archive.broadinstitute.org/cancer/cga/ABSOLUTE</a>
GSVA	Hänzelmann et al., 2013	<a href="https://bioconductor.org/packages/release/bioc/html/GSVA.html">https://bioconductor.org/packages/release/bioc/html/GSVA.html</a>
FANTOM5	Lizio et al., 2015	<a href="http://fantom.gsc.riken.jp/5/">http://fantom.gsc.riken.jp/5/</a>
CIBERSORT	Newman et al., 2015	<a href="http://cibersort.stanford.edu/index.php">http://cibersort.stanford.edu/index.php</a>
Clue (iCluster)	Shen et al., 2009	<a href="https://www.mskcc.org/departments/epidemiology-biostatistics/biostatistics/icluster">https://www.mskcc.org/departments/epidemiology-biostatistics/biostatistics/icluster</a>

### CONTACT FOR REAGENT AND RESOURCE SHARING

Further information and requests for resources and reagents should be directed to and will be fulfilled by the Lead Contact, Li Ding ([lding@wustl.edu](mailto:lding@wustl.edu)).

### EXPERIMENTAL MODEL AND SUBJECT DETAILS

For this research we used data collected by The Cancer Genome Atlas. Under the direction of the National Cancer Institute (NCI) and the National Human Genome Research Institute (NHGRI), TCGA collected both tumor and non-tumor biospecimens from more than 10,000 human samples with informed consent under that authorization of local Institutional Review Boards (<https://cancergenome.nih.gov/abouttcga/policies/informedconsent>). These steps ensured that patients were exposed to no unnecessary risks and that the resulting research is legal, ethical, and well designed. Mutation and clinical data (including age and sex) used for this manuscript are deposited by the GDC (<https://gdc.cancer.gov/about-data/publications/pancanatlas>).

### METHOD DETAILS

#### Germline variant calling

TCGA sequence information was obtained from the database of Genotypes and Phenotypes (dbGaP). Data from paired tumor and germline samples were independently aligned to human reference GRCh37-lite using BWA (Li and Durbin, 2009) v0.5.9 and de-duplicated using Picard 1.29. GenomeVIP (Mashl et al., 2017) was used to orchestrate germline calling using the following tools. Germline single nucleotide variants (SNVs) were identified using Varscan (Koboldt et al., 2012) version 2.3.8 (default parameters, except where  $-\text{min-var-freq}$  0.10,  $-\text{p}$  value 0.10,  $-\text{min-coverage}$  3,  $-\text{strand-filter}$  1) operating on an mpileup stream produced by samtools (Li et al., 2009) version 1.2 (default parameters, except where  $-\text{q}$  1  $-\text{Q}$  13) and GATK (McKenna et al., 2010) version 3.5 using the haplotype caller in single-sample mode with duplicate or unmapped reads removed and calls with quality threshold of 10 retained.

Germline indels were identified using Varscan and GATK, both as configured as above, along with Pindel (Ye et al., 2016) version 0.2.5b8. We specified an insert size of 500 whenever this information was not present in the BAM header. Variants were limited to coding regions of full length transcripts obtained from Ensembl release 70 plus two additional base pairs flanking each exon that cover splice donor/acceptor sites. The union of GATK and VarScan SNVs was processed through our in-house false-positive filter (Kanchi et al., 2014). We included indels called by at least two out of the three callers (GATK, Varscan, Pindel) and high-confidence, Pindel-unique calls (at least 30x coverage and 20% VAF). The combined indels set was again processed through our false-positive filter (default parameters, except where-min-homopolymer 10-min-var-freq 0.2-min-var-count = 6). The entire process is described in more detail in (Huang et al., 2018). For germline and somatic variant comparison we restricted our data to the overlap of samples with at least one mutations in the MC3 MAF after restricting variants as outlined below. This overlap removed one gene from the germline predisposition list (CYLD).

### Somatic variant calling

A publicly available MAF file (syn7824274, <https://gdc.cancer.gov/about-data/publications/mc3-2017>) was compiled by the TCGA MC3 Working Group and annotated with filter flags to highlight potential artifacts and discrepancies (Ellrott et al., 2018). A host of possible artifacts were flagged, including strand-bias, contamination, Oxo-guanine artifacts, and low normal read depth. If a mutation escaped flagging and was called by 2 or more variant calling tools, it was labeled a 'PASS'. We restricted analysis to PASS calls, except for samples from OV and LAML, which were early entrants in TCGA that were whole genome amplified (WGA). Of the 412 OV and 141 LAML samples in our dataset, 347 (84%) and 141 (100%), respectively, had artificial variants induced by WGA. In order to maintain sample sizes and uniformity in mutation calling, we did not filter mutations containing only 'wga' filter tags from these two cancer types. Seven bioinformatic tools were applied, five for Single Nucleotide Variants (SNV) and three for short Insertion Deletion (INDEL) events, with Varscan 2 providing both types of analysis. This list is comprised of MuTect (Cibulskis et al., 2013), VarScan2 (Koboldt et al., 2012), Indelocator (Chapman et al., 2011), Pindel (Ye et al., 2016), SomaticSniper (Larson et al., 2012), RADIA (Radabaugh et al., 2014), and MuSE (Fan et al., 2016). The final call set was filtered to identify cohort level artifacts and was subject to extensive variant, subject, and cohort level QC. In total, 22,485,627 putative variants were identified and 2,907,335 high confidence mutations were retained after filtering.

### Association testing between biological processes and germline or somatic BRCA1/2 mutations

Additionally, Moonlight (Colaprico et al., 2018) analysis was considered to incorporate multiple molecular levels to identify differentially expressed genes in the context of biological pathways (Figures 3 and S1). For this analysis samples with germline predisposition variants in the *BRCA1* and/or *BRCA2* were considered for OV and BRCA. Similarly if a sample harbored somatic missense, frameshift, nonsense, splice site, or in-frame in *BRCA1* or *BRCA2*, that sample was aggregated into the somatic group. If a sample had both germline and somatic mutations, it was not considered for this comparison. A full table of GSEA results is publicly available at <https://github.com/ibsquare/MoonlightOP> "Moonlight\_GSEA\_NES\_results\_Rebut\_v3."

### Germline and somatic gene assignment to pathway analysis

Assignment of genes to specific pathways was performed to provide a landscape of frequently mutated biological processes across 33 cancer types. Primarily genes were classified into 24 unique categories which combined the drivers and essentiality working group classification supplemented by Kegg pathway designations provided by Moonlight. These pathways included: apoptosis, cell cycle, chromatin SWI/SNF complex, chromatin histone modifiers, chromatin other, epigenetics DNA modifiers, genome integrity, histone modification, immune signaling, MAPK signaling, metabolism, NFKB signaling, NOTCH signaling, other, other signaling, PI3K signaling, protein homeostasis/ubiquitination, RNA abundance, RTK signaling, splicing, TGFB signaling, TOR signaling, Transcription factor, and Wnt/B-catenin signaling. This was then further reduced to the 8 molecular processes shown on Figure 2. Of note, one germline predisposition gene was missing from the Circos figure (Krzywinski et al., 2009), *CYLD* due to missing somatic data for a single sample.

In order to calculate the prominent molecular process in each tumor type, a single process was assigned to each sample. This was calculated as follows. If a sample did not carry a predisposing germline variant or missense/frameshift mutation in a driver gene then it was merely added to the denominator of that cancer type. Otherwise, if a sample carried a mutations in a germline and/or somatic driver gene, each driver mutation was compared to the ranked order molecular processes based on the cancer type as a whole. For example, if the top molecular processes, by frequency, for LGG were ranked metabolism, genome integrity, and oncogenic signaling, and a sample only carried mutations in both a metabolic gene and a genome integrity gene, then that sample would be classified for the highest rank of that particular cancer.

### Detection of gene programs differentially expressed in samples with indels or nonsense mutations and missense mutations

Cancer Genome Atlas (TCGA) cohort were available in Genomic Data Commons (GDC) Data Portal and were used in this study in September 2017. We focused on these 16 cancer types because the top 15 cases of cancer-gene combinations for two groups (30 combinations in total) from the frameshift / missense from the significant *cis*-expression associations RNA-seq raw counts of 7668 cases as legacy archive, and using the reference of hg19 were downloaded, normalized and filtered using the R/Bioconductor

package TCGAbiolinks version 2.5.9 (Colaprico et al., 2015) using GDCprepare for tumor types (level 3, and platform “IlluminaHiSeq\_RNASeqV2”) using data.type as “Gene expression quantification” and file.type as “results.” This allowed us to extract the raw signal for expression of a gene for each case following the TCGA pipeline used to create Level 3 expression data from RNA Sequence data that uses MapSplice (Wang et al., 2010) to do the alignment and RSEM to perform the quantification (Li and Dewey, 2011). Integrative analysis using mutation, clinical and gene expression were performed following our recent TCGA’s workflow (Silva et al., 2016).

For this study we used TCGAbiolinks version 2.7.6 and MoonlightR Version 1.2.0 in October 2017 with the following parameters: (i) for Differential Phenotype Analysis (DPA) we filtered out differentially expressed genes with  $\text{fdr.cut} = 0.01$  and  $\text{logFC.cut} = 1$ , (ii) for Functional Enrichment Analysis (FEA) we considered significantly enriched biological processes (BP) by each signature of DEGs with a Fisher Test FDR less than 0.01, (iii) for Gene regulatory network (GRN) the pairwise mutual information was computed using entropy estimates from k-nearest ( $k = 3$ ) neighbor distances filtering out non-significant interactions using a permutation test ( $\text{nboot} = 100$ ,  $\text{nGenesPerm} = 1000$ ), (iv) Upstream Regulator Analysis (URA) was performed considering the output of previous steps with  $\text{nCores} = 64$ . Hierarchical cluster analysis using a complete linkage method to finds similar cluster of biols was applied to generate the heatmap (Figure 4H) sorted by each cancer type. A full list of Moonlight significance scores are publicly available at <https://github.com/ibsquare/MoonlightOP> “Moonlight\_FrameShift\_Missense\_SupplementalData.”

We used Moonlight (Colaprico et al., 2018) to find pathways and biological processes that show differences in the expression levels of their genes based on the presence and type of mutations in driver genes. We had three groups: WT, missense and frameshift/nonsense. Samples with both types of mutations, missense and frameshift/nonsense were excluded from this analysis.

### Identification of biological processes associated with cancer driver genes

OncolMPACT (Bertrand et al., 2015) integrates genomic and transcriptomic profiles using a gene interaction network model to discern patient-specific drivers based on their “phenotypic” effect. We used this tool to predict patient-specific modules of deregulated genes associated with mutational driver genes. Modules are constructed by: 1) identifying phenotype genes defined as significantly deregulated genes associated with a driver mutation (deregulated in  $\geq 5\%$  of patients, permutation test,  $\text{FDR} < 0.1$ ) for a particular cancer type, 2) aggregating patient specific modules by linking driver genes to the phenotypes genes using the protein interaction network. For each cancer type, deregulated genes of a patient were identified by calculating the  $\log_2$  fold-change between the patient gene expression value and the cancer type median gene expression value. After obtaining the gene modules predicted by OncolMPACT based on patients’ transcriptomic and mutational profiles (SNV, indels and CNA), we selected, for each patient, the largest module containing at least one driver gene from the PanCancer Atlas oncogenic process working group cancer driver genes list. Genes affected by a focal amplification/deletion were filtered out from the modules, as their change in expression may be associated with the copy number change. Biological processes associated with each module were identified by using enrichment analysis on MSigDB’s GO\_BP and KEGG\_PATHWAY gene lists (Fisher exact test,  $\text{FDR} < 0.05$ ). Patient-specific predictions were then combined at the cancer type level to obtain the fraction of patients for which an oncogenic process was associated with a driver mutation. To control for Type 1 errors introduced by the FDR threshold (0.05 of the predictions are expected to be false positive), we performed a binomial test for each fraction reported (expected frequency 0.05) and filtered out any fraction with a Bonferroni corrected p values  $> 0.05$ . The total number of samples used in this analysis was 6,224 (samples from DLBC and CHOL were excluded due to their small module sizes).

Additionally, we tested if the five most frequently mutated driver genes were significantly mutually exclusive in each oncogenic process using the R-exclusivity test (Leiserson et al., 2016). For each oncogenic process, we constructed a mutation matrix where rows are driver genes and columns are samples. We then counted the number of samples harboring mutually exclusive driver mutations and performed a permutation test by maintaining frequencies of all five driver genes. The reported p value is based on the number of permuted matrices (100,000) showing higher numbers of samples harboring mutually exclusive driver mutations. The full table of results from this analysis can be located at [https://github.com/CSB5/OncolMPACT/blob/development/TCGA\\_PAN\\_CAN\\_ANALYSIS/gene\\_list\\_driver.csv](https://github.com/CSB5/OncolMPACT/blob/development/TCGA_PAN_CAN_ANALYSIS/gene_list_driver.csv).

### Integration for cell of origin clusters with mutations

Sample and cluster information was provided in the private communication with the cell-of-origin group for 3 additional molecular levels, methylation, mRNA, and reverse phase protein array (RPPA). These sets had varying samples sizes based on data quality and availability (Table S8). These 3 level identifiers were concatenated to create a new cluster identification number that was utilized for down stream analysis and investigation. From the data provided we identified 166 samples with a single sample in the classifier. Samples with missense, indel, or splice site mutations (considered drivers for this analysis) in any of the 299 genes identified by the PanCancer Atlas drivers group were merged in by sample and a gene enrichment analysis was performed comparing clusters sizes (by sample) to the number of samples with a driver mutation.  $\text{FDR} \leq 0.05$  was considered significant. We also determined what fraction of the cluster ids originate from a single tissue of origin. To address this, we implemented a simple heuristic to estimate cluster homogeneity. We define cluster homogeneity as those clusters with  $\geq 20$  samples that have  $\geq 90\%$  of the samples from a single cancer type (Figure 6D). 58/414 cluster have 20 or more samples, of which, 69% are homogeneous (40/58), however there are a number of clusters that capture more universal molecular patterns and are shared across cancer types.

### The cell-to-cell communication network

A network of documented ligand-receptor, cell-receptor, and cell-ligand pairs was retrieved from the FANTOM5 resource at ([http://fantom.gsc.riken.jp/5/suppl/Ramilowski\\_et\\_al\\_2015/](http://fantom.gsc.riken.jp/5/suppl/Ramilowski_et_al_2015/)). Because CIBERSORT cell types are more granular than immune cells in FANTOM5, CIBERSORT abundance estimates were aggregated by summing to yield estimates for FANTOM5 immune cell abundances, as defined above. This network was augmented with additional known interactions of immunomodulators, and only ligand-receptor edges that contained at least one cell or one immune modulator were retained, yielding a 'scaffold' of possible interactions.

From the scaffold of possible interactions, interactions were identified that could be playing a role within the TME in each subtype as follows. Cellular fractions were binned into tertiles (low, medium, high), as were gene expression values for ligands and receptors, yielding ternary values for all 'nodes' in the network. The binning was performed over all TCGA samples. In subsequent processing, nodes and edges were treated uniformly in processing, without regard to type (cell, ligand, receptor). From the scaffold, interactions predicted to take place in the TME were identified *first* by a criterion for the nodes to be included ('present' in the network), *then* by a criterion for inclusion of edges. For nodes, if at least 66% of samples within a subtype map to mid or high value bins, the node is entered into the subtype-network. An edge present in the scaffold network between any two nodes is then evaluated for inclusion. A contingency table is populated for the ternary values of the two nodes, over all samples in the subtype, and a concordance versus discordance ratio ("concordance score") is calculated for the edge in terms of the values of  $((\text{high,high})+(\text{low,low})) / ((\text{low,high})+(\text{high,low}))$ . Edges were retained with concordance score  $> 2.9$ , set based on evaluation of quantile distributions (Table S11). Additional details in (Thorsson et al., 2018).

## QUANTIFICATION AND STATISTICAL ANALYSIS

### Comparison of clinical and mutational impact of somatic and germline *BRCA1* and *BRCA2* variants

We grouped samples according to whether they had *BRCA1* and/or *BRCA2* germline, somatic or no mutations. We then compared the number of somatic mutations (Ellrott et al., 2018) in each group using a t test. We also used the clinical data (<https://www.synapse.org/#!Synapse:syn4983466.1>) to compare the age at onset of each group using also a Welch's two sample t test compared to wild type. Samples with both, germline and somatic *BRCA1/2* mutations were included in both categories. These results are reported in Table S4 and distinguishable with the column header AnalysisGrouping (Figure 3A).

### Comparison of clinical and mutational impact of somatic and germline DDR pathway alterations

We grouped samples according to whether they had germline, somatic or no mutations in the core DDR pathway (Figure 3B). This pathway consists of 80 genes according to genes from the Pathways DDR AWG (Table S2). The number of mutations was compared using Welch's two sample t test compared to wild type. Samples with both, germline and somatic in DDR genes mutations were included in both categories. These results are reported in Table S4 and distinguishable with the column header AnalysisGrouping.

### Comparison of clinical and mutational impact of somatic and germline MSI pathway alterations

We grouped the samples as in Figure 3C, but using the MSI pathway definition instead, which consists of 33 genes (Table S2). We used MSIsensor (Niu et al., 2014) to determine the MSI score of each sample and compared the scores in each group using a Welch's two sample t test compared to wild type (Table S3). In addition to stratifying our analysis by mutation status in MSI and germline predisposition genes, promoter methylation status for MLH1 was appended to UCEC, COAD, and STAD and was obtained from MIRMR (Foltz et al., 2017).

### Correlation between MSI scores and expression of immune-related genes

We grouped samples according to whether they had high or low MSI scores (MSIsensor score  $\geq 4$  and MSIsensor score  $< 4$  respectively). Then we compared the  $\log_2$  expression of immune-related genes (*GZMA*, *PRF1*, *GZMK* and *GZMH*) in both groups using both Student's t test and a two sample Kolmogorov–Smirnov test (KS-test). We limited our analysis to those cancer types because there were sufficient number of MSIhigh samples: UCEC, STAD and COADREAD. We used the KS-test significance of p value  $< 0.01$  for (Figure 2D). All groups indicated as significant also showed significance using the t test except when comparing *GZMH* abundance in UCEC (t test p value = 0.49; KS-test p value = 0.003).

### Mutation mutual exclusivity and co-occurrence analysis

We performed a mutually exclusivity/co-occurring mutation analysis of samples between all official pairs (258/299) of consensus driver genes from (Bailey et al., 2018), which included splice site mutations, but excluded non-coding and silent mutations. The analysis was run at the gene level. We used a two-sided exact Cochran-Mantel-Haenszel test (mantelhaen.test R function) to identify significant patterns for each individual cancer type and for the PanCancer set as a whole, with multiple test correction of FDR  $< 0.1$ . The covariate stratum for this test used mutation burden and the identity of the cancer type for the PanCancer analysis. Mutation burden was dichotomized at a 500 mutations threshold based on an even split of the minimum hypermutated sample threshold (1,000 mutations per sample). This was intended to control for spurious co-occurrence inferences induced by samples with very high mutation burden. Odds ratios of greater or less than one indicate tendencies toward co-occurrence and mutual exclusivity, respectively. Note that in the tissue-specific analyses, this amounts to the tables being 2x2x2 (Gene1 / Gene2 / Mutation burden) whereas in the

Pancan analysis they are 2x2x66 (Gene1 / Gene2 /Tissue + mutation burden). We corrected for multiple hypotheses using the Benjamini-Hochberg FDR method, reporting all gene pairs having a FDR < 0.1.

### Association testing between different types of mutations and biological processes

We conducted this analysis on the extended consensus driver list of 299 genes, grouping the associated samples for each cancer type into three categories; (i) samples having only frameshift indels or nonsense mutations (FSN), (ii) those having only missense mutations (MIS), and (iii) those having no mutations (WT). Samples with both types of mutations, missense and frameshift/nonsense, were not included in this analysis. For each combination of cancer type and gene, we compiled subsets of samples for these three categories. Any cancer-gene combination not having at least five samples in each of the three categories was excluded for lack of power.

RNA-Seq gene expression data were obtained for each sample category for the above cancer-gene combinations. All RSEM value sets were transformed into normal distributions with Box-Cox transformations, after which Z-Scores were calculated. For a given cancer type, gene, and respective subsets of samples (distinguished by mutation category), Welch's t-Test was performed to assess the significance of the difference of expression distributions between the test subset and the subset of wild-type samples from the same cancer type and gene. Here, the t-statistic is

$$t = \frac{X_1 - X_2}{\sqrt{\frac{S_1^2}{N_1} + \frac{S_2^2}{N_2}}}$$

where,  $X_i$ ,  $S_i$ , and  $N_i$  are the respective sample mean, standard deviation, and tally of the  $i$ th distribution. Welch's test is especially appropriate, since we do not always find equal variances or sample numbers between the distributions. The t-scores and degrees of freedom generated by the t test were used to perform a two-tailed significance test against the t-distributions. The distribution of t-scores and their corresponding significance status is depicted in [Figure 4](#). The results from this analysis are reported in [Table S7](#), and separated by Mutated (any non-silent mutation) and "Frame\_Shift\_And\_Nonsense" or "Missense\_Only" under the column header "AnalysisGrouping." These two groups ("Mutated" and "Frame\_Shift\_And\_Nonsense"/"Missense\_Only") were tested independent of each other. Additionally, we have included results by expanding our analysis to all non-silent mutations and show the top results in [Figure S2](#).

### Correlation between driver events and immune cell types

We focused our analysis on the set of 299 driver genes and > 3400 driver mutations from ([Bailey et al., 2018](#)). We considered that a sample had a driver event if it carried a frameshift or truncating mutation, or a missense mutation detected by at least 2 different signals of oncogenicity ([Bailey et al., 2018](#)). In order to reduce the issues related to multiple-testing we analyzed only driver events present in 10 or more samples. We considered both individual driver mutations and entire driver genes that met these criteria.

Then, for each of the six immune subtypes ([Thorsson et al., 2018](#)) we checked for a correlation between the presence of the driver event and the quantity of different immune cells in the tumor microenvironment. The quantification of immune cells is described in "Immune Fraction Estimates" below. Then, we used domainXplorer to identify driver events that correlate with the presence of different immune cell types ([Porta-Pardo and Godzik, 2016](#)). Briefly, domainXplorer uses a linear correlation model that accounts for different variables that might bias the results, such as the tissue of origin or the number of mutations in the tumor sample. The model is:

$$CF = \beta_0 + \beta_1 T + \beta_2 N + \beta_3 D$$

where  $CF$  is the cell fraction of each sample,  $T$  is the tissue of origin for each sample,  $N$  the total number of mutations in the sample and  $D$  is a binary variable showing whether the sample has a certain driver event or not. To correct for multiple testing, the Benjamini-Hochberg method was applied to p values of the  $D$  factor from the ANOVA test of each driver event ([Table S11](#)).

## DATA AND SOFTWARE AVAILABILITY

### Germline predisposition variant list

The list of germline variants was obtained from [Huang et al. \(2018\)](#). While the details on how to obtain the final 1,461 germline variants are explained in detail in the manuscript, in brief the group first selected for cancer-relevant pathogenic variants, based on whether they were found in the curated cancer variant database or in the curated cancer predisposition gene list, and their associated ClinVar trait. This resulted in 1,678 variants for manual review using the Integrative Genomics Viewer (IGV). For candidate germline variants having the same genomic change as somatic mutations, we further filtered for the germline variants that may have originated from contaminated adjacent normal samples by eliminating variants called from adjacent normal, the VAF in normal < 30%, and co-localizing with any known somatic mutation.

### Driver gene list

The list of driver genes was obtained from [Bailey et al. \(2018\)](#). The details about how this list was created are further detailed in that manuscript, but in brief, the Driver AWG combined the predictions of 8 different tools comprising algorithms based on mutation frequency (MuSiC2 [[Dees et al., 2012](#)] and MutSig2CV [[Lawrence et al., 2014](#)]), features (20/20 [[Tokheim et al., 2016](#)], CompositeDriver [<https://github.com/khuranalab/CompositeDriver>] and OncodriveFML [[Mularoni et al., 2016](#)]), clustering (OncodriveCLUS [[Tamborero et al., 2013](#)]), and externally defined regions (e-Driver [[Porta-Pardo and Godzik, 2014](#)] and ActiveDriver [[Reimand and Bader, 2013](#)]).

The preliminary total of 2,101 potential driver genes was identified by taking the union of genes predicted by the eight driver-gene discovery tools. They refined this list by calculating, for each gene predicted in each cancer type, a consensus score that compensated for outlier results and correlation among tools. The consensus score was defined as a weighted sum of the number of tools that predicted the gene to be a driver in each cancer type (see Gene Discovery Weighting Strategy). They required a minimum of two tools to agree, where both could not be outliers (score  $\geq 1.5$ ).

To maximize the coverage of the analysis and ensure the accuracy of the final list, they reviewed previous findings in 31 individual cancer types and PanCancer-12 from TCGA. For cancer types not yet having a TCGA publication, they consulted with the relevant analysis working groups (LIHC, TGCT, UVM, SARC, PAAD, and THYM). They included in the final consensus list all those genes that were previously described as drivers by experts in the cancer-specific analysis of TCGA datasets and that were also identified by at least one of the eight algorithms, even if they did not meet the consensus score threshold ( $\geq 1.5$ ). Then, to limit false positives in the expanded list, they applied linear discriminant analysis, removing 45 genes from the consensus they detected as likely false positives.

Finally, given the limitations of a systematic approach, they additionally manually rescued 41 genes based on supportive evidence from the following sources: hypermutator phenotype related genes (since they excluded hypermutated samples in our systematic discovery), established cancer genes from LAML because of low quality variant calling originating from tumor contamination of the normal samples, genes supported by omic network tools: OncoIMPACT ([Bertrand et al., 2015](#)) and DriverNet ([Bashashati et al., 2012](#)). Addition of genes to the final list was subjected to expert manual curation.

### Cell of origin transcript data

The PanCancer Atlas Cell Origin manuscript provided us with cluster data for 3 additional substrates: methylation, mRNA, and RPPA ([Table S9](#)). This overview supports the notion that cancers should be classified by their molecular characteristics and can effectively identify molecular subgroup patterns. Methylation data used unsupervised clustering of 10,814 tumors using Ward's method to cluster the distance matrix computed with the Jaccard index. This resulted in 25 number of clusters. Unsupervised consensus clustering using Consensus Cluster Plus ([Wilkerson and Hayes, 2010](#)) was performed on RSEM (mRNA normalized expression) for 10,165 samples and 15,363 genes and resulted in 43 clusters (25 with at least 40 samples). And finally, reverse phase protein arrays (RPPA) was also clustered using Pearson's correlation coefficient as the distance metric and Ward's method as the linkage function, which resulted in 10 clusters.

### Expression and copy number data

Gene expression and copy number information for each sample were retrieved from the Genomic Data Commons unless indicated otherwise in specific sections of [STAR Methods](#).

### Cancer Immune Subtypes

To characterize the commonality and diversity of intratumoral immune states, we scored 160 published immune expression signatures on all available TCGA PanCancerAtlas tumor samples and performed cluster analysis to identify similarity modules of multiple immune signature sets. The 160 immune expression signatures were selected based on extensive literature search, utilizing diverse resources considered to be reliable and comprehensive based on expert opinions of immuno-oncologists. 83 signatures were derived in the context of immune response studies in cancer and the remaining 77 are of general validity for immunity. TCGA RNA-seq values from the PanCancer Atlas normalized gene expression matrix were scored for each of the 160 identified gene expression signatures using single-sample gene set enrichment (ssGSEA) analysis, using the R package GSVA. Clusters of similar signature scores were identified by weighted gene correlation network analysis (WGCNA) ([Langfelder and Horvath, 2008](#)). Based on the WGCNA analysis, five immuno-oncology-related immune expression signatures: activation of macrophages/monocytes ([Beck et al., 2009](#)), overall lymphocyte infiltration (dominated by T and B cells) ([Calabrò et al., 2009](#)), TGF- $\beta$  response ([Teschendorff et al., 2010](#)), IFN- $\gamma$  response ([Wolf et al., 2014](#)), and wound healing ([Chang et al., 2004](#)), robustly reproduced co-clustering of the immune signature sets, and were selected to perform cluster analysis of all cancer types, with the exception of hematologic neoplasias (acute myeloid leukemia, LAML; diffuse large B cell lymphoma, DLBC; and thymoma, THYM). Clustering of tumor samples scored on these five signatures was performed using model based clustering, using the mclust R package ([Scrucca et al., 2016](#)), with the number of clusters, K, determined by maximization of Bayesian Information Criterion (BIC). Maximal BIC was found with a six cluster solution, and the six resulting clusters C1-C6 (with 2416, 2591, 2397, 1157, 385 and 180 cases, respectively) were characterized by a distinct distribution of scores over the five representative signatures, and effectively categorized each TCGA sample as belonging to one of six cancer "immune subtypes," namely Wound Healing (C1), IFN- $\gamma$  Dominant (C2), Inflammatory (C3), Lymphocyte Depleted (C4), Immunologically Quiet (C5), or TGF- $\beta$  Dominant (C6). Additional details in ([Thorsson et al., 2018](#); [Tables S11 and S12](#)).

### FANTOM5 network

A network of documented ligand-receptor, cell-receptor, and cell-ligand pairs was retrieved from the FANTOM5 resource at ([http://fantom.gsc.riken.jp/5/suppl/Ramilowski\\_et\\_al\\_2015/](http://fantom.gsc.riken.jp/5/suppl/Ramilowski_et_al_2015/)).

### Immune cellular fraction estimates

The relative fraction of 22 immune cell types within the leukocyte compartment were estimated by applying CIBERSORT (Newman et al., 2015) to TCGA RNASeq data (Table S12). As several key immune genes used in the signatures are absent from TCGA GAF (Generic Annotation File) Version 3.0, we applied CIBERSORT to a re-quantification of the TCGA data using Kallisto and the Gencode GTF, which includes the missing genes. A version of the entire TCGA RNA-seq data normalized to Gencode with Kallisto was computed on the ISB Cancer Genomics Cloud by Steve Piccolo's group at BYU (<https://osf.io/gqrz9/wiki/home/>) (Tatlow and Piccolo, 2016). In this study, the 22 CIBERSORT values were aggregated into 9 overall cell types as follows

```
Mast.cells = Mast.cells.resting + Mast.cells.activated,
Dendritic.cells = Dendritic.cells.resting + Dendritic.cells.activated,
Macrophage = Macrophages.M0 + Macrophages.M1 + Macrophages.M2, NK.cells = NK.cells.resting+NK.cells.activated,
B.cells = B.cells.naive + B.cells.memory,
T.cells.CD4 = T.cells.CD4.naive+T.cells.CD4.memory.resting+T.cells.CD4.memory.activated
Neutrophils = Neutrophils,
Eosinophils = Eosinophils,
T.cells.CD8 = T.cells.CD8
```

Additional details in (Thorsson et al., 2018), where this particular combination is referred to as “Aggregate 2.”

### HLA typing and Predicting mutant peptide-MHC binding (neoantigens [pMHCs]) from SNVs

HLA class I typing of samples (raw RNA-Seq from 8872 samples and aligned reads from 715 samples) was performed on the Seven Bridges Cancer Genomics Cloud using a Common Workflow Language (CWL) description of the OptiType tool (version 1.2) (Szolek et al., 2014). The aligned RNA-Seq samples were first converted to raw sequences using a CWL description of the Picard SamtoFastq tool (version 1.140). The reads from each raw RNA-Seq sample were first aligned to the HLA class I database using a CWL description of the yara aligner (version 0.9.9) (Siragusa et al., 2013) with its error rate parameter set to 3%. Next, the CWL description of OptiType was used to compute the HLA class I types for the sample. Potential neoantigenic peptides were identified using NetMHCpan v3.0 (Nielsen and Andreatta, 2016), based on HLA types. For each sample, all pairs of MHC and minimal mutant peptide were input into NetMHCpan v3.0 using default settings. NetMHCpan will automatically extract all 8-11-mer peptides from a minimal peptide sequence and predict binding for each peptide-MHC pair. After computation, the results were parsed to only retain peptides which included the mutated position. Peptides containing amino acid mutations were identified as potential antigens on the basis of a predicted binding to autologous MHC ( $IC_{50} < 500$  nM) and detectable gene expression meeting an empirically determined threshold of 1.6 transcripts-per-million (TPM). This threshold was selected in order to divide the bimodal distribution in the expression data. Additional details in (Thorsson et al., 2018)

### CIBERSORT

CIBERSORT (cell-type identification by estimating relative subsets of RNA transcripts, Newman et al., 2015) uses a set of 22 immune cell reference profiles to derive a base (signature) matrix which can be applied to mixed samples to determine relative proportions of immune cells. It can be accessed at <https://cibersort.stanford.edu>.

### Moonlight

Moonlight (Colaprico et al., 2018) is a new methodology available as R bioconductor package, (<https://bioconductor.org/packages/release/bioc/html/MoonlightR.html>, DOI: 10.18129/B9.bioc.MoonlightR) that does not only identify driver genes playing a dual role (e.g., tumor suppressor genes (TSGs) in one cancer type and oncogenes (OCGs) in another), but also helps in elucidating the biological processes underlying their specific roles.

For this study we used MoonlightR Version 1.2.0 in July 2017 with the following parameters: (i) for DPA we filtered out differentially expressed genes with  $fdr.cut = 0.01$  and  $logFC.cut = 1$ , (ii) for FEA we considered significantly enriched biological processes by each signature of DEGs with a Fisher Test FDR less than 0.01, (iii) for GRN the pairwise mutual information was computed using entropy estimates from k-nearest ( $k = 3$ ) neighbor distances filtering out non-significant interactions using a permutation test ( $nboot = 100$ ,  $nGenesPerm = 1000$ ), (iv) URA was performed considering the output of previous steps with  $nCores = 64$ , (v) First we retrieved a list of validated OCGs and TSGs from the Catalogue of somatic mutations in cancer (COSMIC). The list consists of 84 OCGs, 55 TSGs, 17 dual role genes and 439 genes without validated role. Second PRA was performed considering the URA output as input for the random forest learning approach together with the list of known OCGs and TSGs (COSMIC) used to construct the training set and using a permutation test with  $nrand = 1000$  for obtaining p values filtered by  $FDR = 0.01$ .



**domainXplorer**

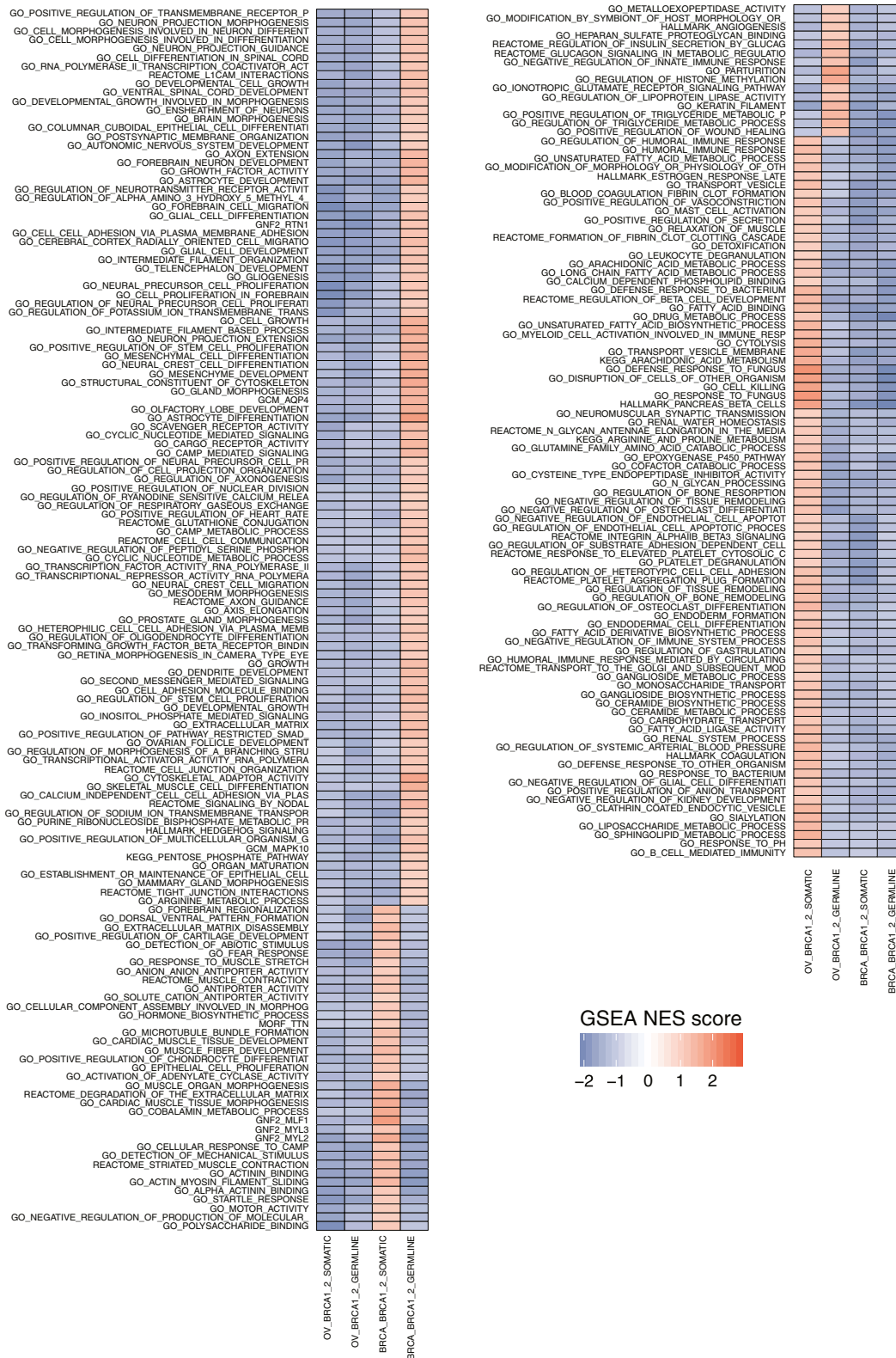
This pipeline identifies events that show statistically significant correlations with the presence of immune cells in the tumor microenvironment (Porta-Pardo and Godzik, 2016). It accounts for several potentially confounding factors, such as the presence of neo-antigens. It can be accessed at <https://github.com/eduardporta/domainXplorer>.

**OncolIMPACT**

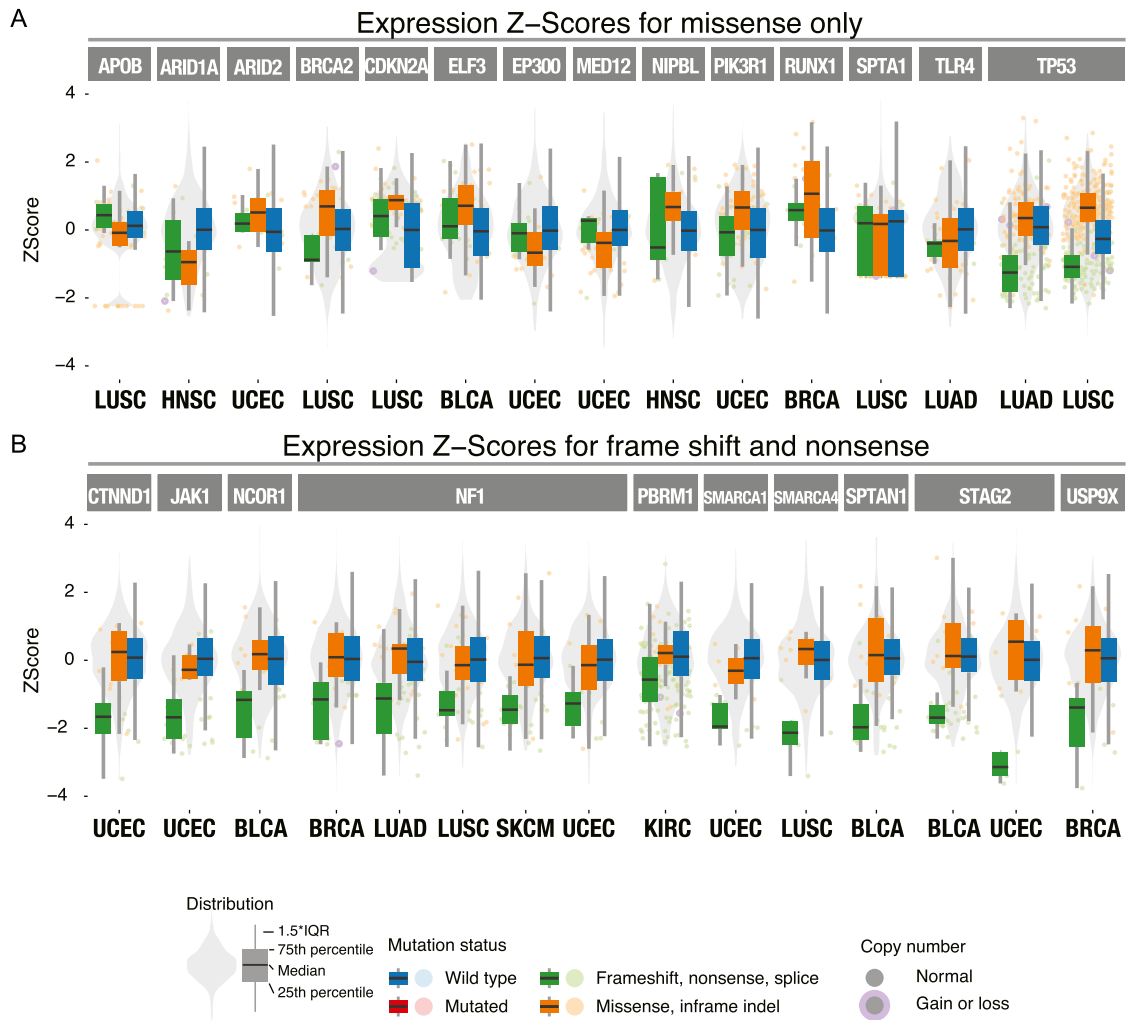
Integrates genomic and transcriptomic profiles using a gene interaction network model to discern patient-specific drivers based on their “phenotypic” effect. It can be accessed at <https://github.com/CSB5/OncolIMPACT>.

**ABSOLUTE**

We used ABSOLUTE (Carter et al., 2012) calls to infer whether each mutation was clonal or sub-clonal. ABSOLUTE optimizes/solves a mixture model for the observed allelic fraction for each mutation (i.e., the mutated reads could have arisen from 1 copy, 2 copies, 3 copies, etc. or from a subclonal population). We defined ‘clonal’ as all mutations that were predicted only as clonal by ABSOLUTE (n = 910,138 out of a total 1,451,623 mutations, 62%). It can be accessed at <http://software.broadinstitute.org/cancer/software/genepattern/modules/docs/ABSOLUTE>.



**Figure S1. Moonlight Analysis of Enriched Pathways for Samples with Germline or Somatic Mutations in *BRCA1* or *BRCA2*, Related to Figure 3**  
 Shown here are the extended set of pathways not shown in Figure 3.



**Figure S2. Alternative Grouping for *Cis*-expression Differences, Related to Figure 4**

(A and B) For Figure 4, only missense mutations and frameshift indels were considered. The top 15 *t* values using an extended definition of missense mutations to include in-frame indels (A). Frameshift/nonsense mutations, here, include splice-site mutations (B).