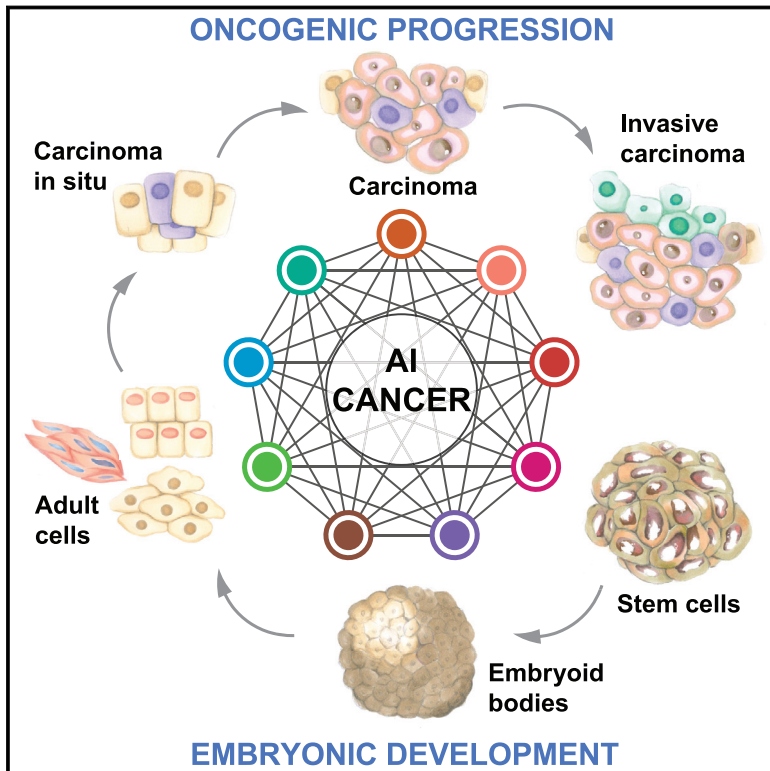


# Machine Learning Identifies Stemness Features Associated with Oncogenic Dedifferentiation

## Graphical Abstract



## Authors

Tathiane M. Malta, Artem Sokolov, Andrew J. Gentles, ..., Peter W. Laird, Houtan Noushmehr, Maciej Wiznerowicz

## Correspondence

hnoushm1@hfhs.org (H.N.),  
maciej.wiznerowicz@iimo.pl (M.W.)

## In Brief

Stemness features extracted from transcriptomic and epigenetic data from TCGA tumors reveal novel biological and clinical insight, as well as potential drug targets for anti-cancer therapies.

## Highlights

- Epigenetic and expression-based stemness indices measure oncogenic dedifferentiation
- Immune microenvironment content and PD-L1 levels associate with stemness indices
- Stemness index is increased in metastatic tumors and reveals intratumor heterogeneity
- Applying stemness indices reveals potential drug targets for anti-cancer therapies



# Machine Learning Identifies Stemness Features Associated with Oncogenic Dedifferentiation

Tathiane M. Malta,<sup>1,2,23</sup> Artem Sokolov,<sup>3,23</sup> Andrew J. Gentles,<sup>4</sup> Tomasz Burzykowski,<sup>5</sup> Laila Poisson,<sup>1</sup> John N. Weinstein,<sup>6</sup> Bożena Kamińska,<sup>7</sup> Joerg Huelsken,<sup>8</sup> Larsson Omberg,<sup>9</sup> Olivier Gevaert,<sup>4</sup> Antonio Colaprico,<sup>10,11</sup> Patrycja Czerwińska,<sup>12</sup> Sylwia Mazurek,<sup>12,13</sup> Lopa Mishra,<sup>14</sup> Holger Heyn,<sup>15</sup> Alex Krasnitz,<sup>16</sup> Andrew K. Godwin,<sup>17</sup> Alexander J. Lazar,<sup>6</sup> The Cancer Genome Atlas Research Network, Joshua M. Stuart,<sup>18</sup> Katherine A. Hoadley,<sup>19</sup> Peter W. Laird,<sup>20</sup> Houtan Noushmehr,<sup>1,2,23,\*</sup> and Maciej Wiznerowicz<sup>12,21,22,23,24,\*</sup>

<sup>1</sup>Henry Ford Health System, Detroit, MI 48202, USA

<sup>2</sup>University of São Paulo, Ribeirão Preto-SP 14049, Brazil

<sup>3</sup>Harvard Medical School, Boston, MA 02115, USA

<sup>4</sup>Stanford University, Palo Alto, CA 94305, USA

<sup>5</sup>Hasselt University, 3590 Diepenbeek, Belgium

<sup>6</sup>The University of Texas MD Anderson Cancer Center, Houston, TX 77030, USA

<sup>7</sup>Nencki Institute of Experimental Biology of PAS, 02093 Warsaw, Poland

<sup>8</sup>Swiss Federal Institute of Technology Lausanne (EPFL), CH-1015 Lausanne; Switzerland

<sup>9</sup>Sage Bionetworks, Seattle, WA 98109 USA

<sup>10</sup>Université Libre de Bruxelles, 1050 Bruxelles, Belgium

<sup>11</sup>Interuniversity Institute of Bioinformatics in Brussels (IB)<sup>2</sup>, 1050 Bruxelles; Belgium

<sup>12</sup>Poznań University of Medical Sciences, 61701 Poznań, Poland

<sup>13</sup>Postgraduate School of Molecular Medicine, Medical University of Warsaw, 02109 Warsaw, Poland

<sup>14</sup>George Washington University, Washington, D.C. 20052, USA

<sup>15</sup>Centre for Genomic Regulation (CNAG-CRG), 08003 Barcelona, Spain

<sup>16</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724, USA

<sup>17</sup>University of Kansas Medical Center, Kansas City, KS 66160, USA

<sup>18</sup>University of California, Santa Cruz, Santa Cruz, CA 95064, USA

<sup>19</sup>University of North Carolina, Chapel Hill, NC 27599, USA

<sup>20</sup>Van Andel Research Institute, Grand Rapids, MI 49503, USA

<sup>21</sup>Greater Poland Cancer Center, 61866 Poznań, Poland

<sup>22</sup>International Institute for Molecular Oncology, 60203 Poznań, Poland

<sup>23</sup>These authors contributed equally

<sup>24</sup>Lead Contact

\*Correspondence: [hnoushm1@hfhs.org](mailto:hnoushm1@hfhs.org) (H.N.), [maciej.wiznerowicz@iimo.pl](mailto:maciej.wiznerowicz@iimo.pl) (M.W.)

<https://doi.org/10.1016/j.cell.2018.03.034>

## SUMMARY

Cancer progression involves the gradual loss of a differentiated phenotype and acquisition of progenitor and stem-cell-like features. Here, we provide novel stemness indices for assessing the degree of oncogenic dedifferentiation. We used an innovative one-class logistic regression (OCLR) machine-learning algorithm to extract transcriptomic and epigenetic feature sets derived from non-transformed pluripotent stem cells and their differentiated progeny. Using OCLR, we were able to identify previously undiscovered biological mechanisms associated with the dedifferentiated oncogenic state. Analyses of the tumor microenvironment revealed unanticipated correlation of cancer stemness with immune checkpoint expression and infiltrating immune cells. We found that the dedifferentiated oncogenic phenotype was generally most prominent in metastatic tumors. Application of our stemness

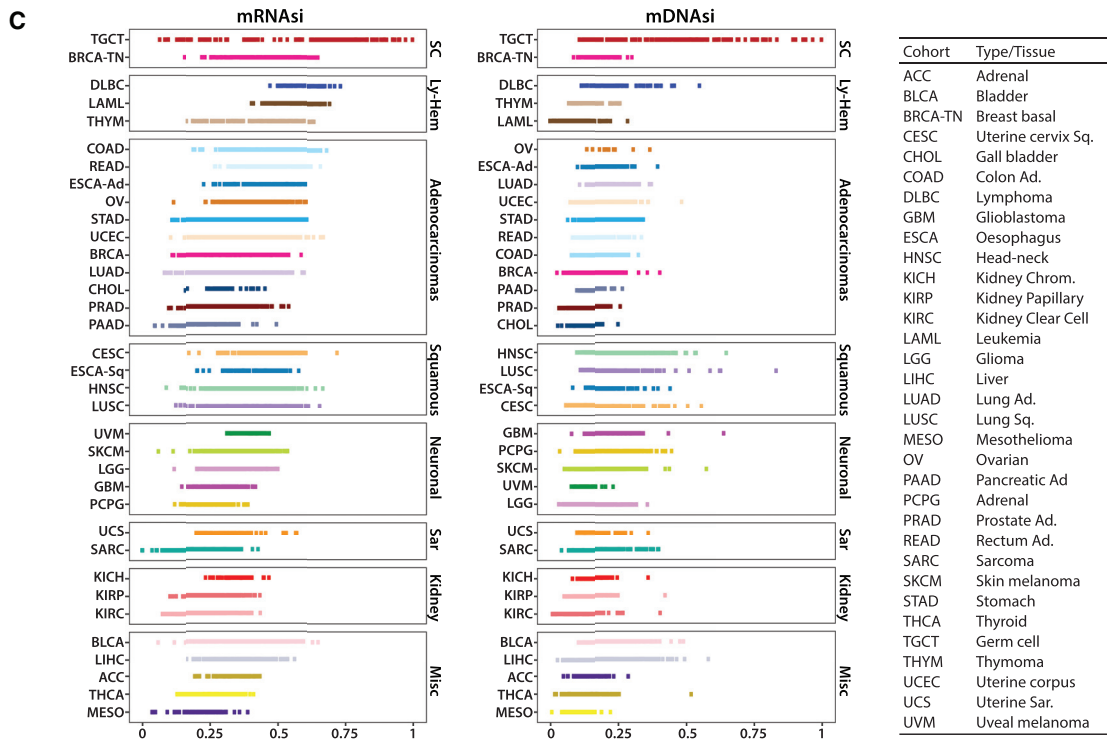
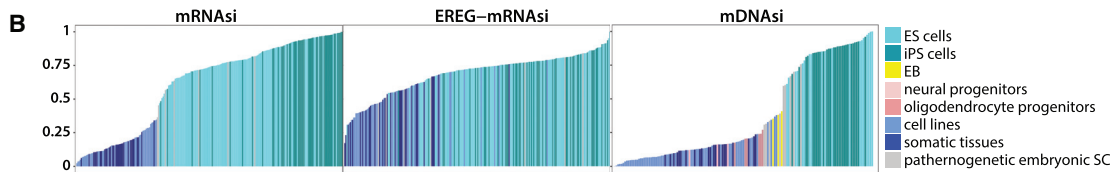
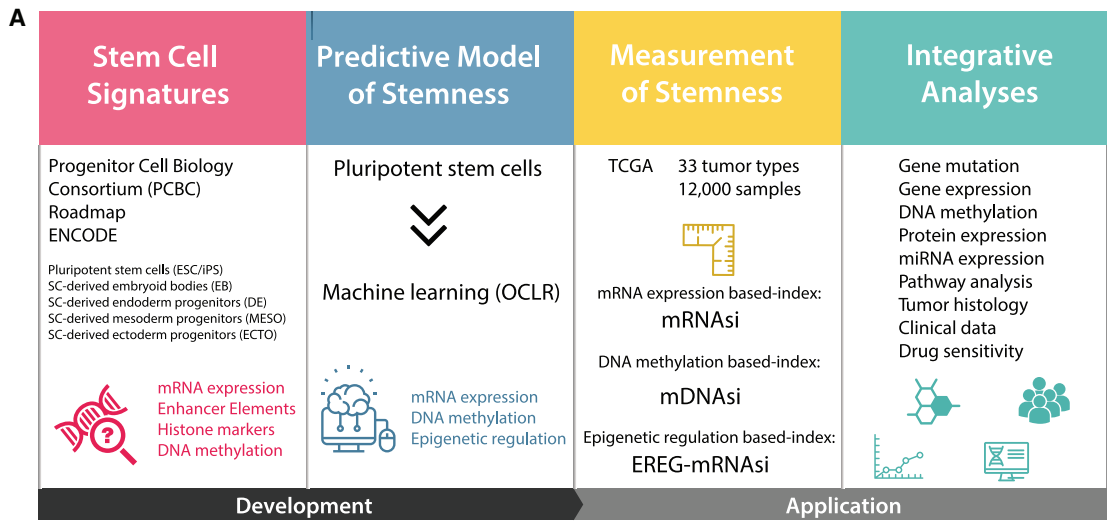
indices to single-cell data revealed patterns of intra-tumor molecular heterogeneity. Finally, the indices allowed for the identification of novel targets and possible targeted therapies aimed at tumor differentiation.

## INTRODUCTION

Stemness, defined as the potential for self-renewal and differentiation from the cell of origin, was originally attributed to normal stem cells that possess the ability to give rise to all cell types in the adult organism. Cancer progression involves gradual loss of a differentiated phenotype and acquisition of progenitor-like, stem-cell-like features. Undifferentiated primary tumors are more likely to result in cancer cell spread to distant organs, causing disease progression and poor prognosis, particularly because metastases are usually resistant to available therapies (Friedmann-Morvinski and Verma, 2014; Ge et al., 2017; Shibue and Weinberg, 2017; Visvader and Lindeman, 2012).

An increasing number of genomic, epigenomic, transcriptomic, and proteomic signatures have been associated with





**Figure 1. Development and Validation of the Stemness Indices**

(A) Overall methodology. Highlighted are data sources Progenitor Cell Biology Consortium (PCBC), Roadmap, and ENCODE databases; the OCLR machine-learning algorithm; and the resulting stemness indices mRNAsi, mDNAsi, and EREg-mRNAsi. The indices for each TCGA tumor sample were correlated with known cancer biology, tumor pathology, clinical information, and drug sensitivity.

(B) Stemness indices of the validation set derived using our stemness signature.

*(legend continued on next page)*

cancer stemness. Those molecular features are causally connected to particular oncogenic signaling pathways that regulate transcriptional networks that sustain the growth and proliferation of cancer cells (Ben-Porath et al., 2008; Eppert et al., 2011; Kim et al., 2010). Transcriptional and epigenetic dysregulation of cancer cells frequently leads to oncogenic dedifferentiation and acquisition of stemness features by altering core signaling pathways that regulate the phenotypes of normal stem cells (Bradner et al., 2017; Young, 2011). Cell-extrinsic mechanisms can also affect maintenance of the undifferentiated state, largely through epigenetic mechanisms. Tumors comprise a complex, diverse, integrated ecosystem of relatively differentiated cancer cells, cancer stem cells, endothelial cells, tumor-associated fibroblasts, and infiltrating immune cells, among other cell types. The microenvironment of a tumor, considered as a pathologically formed “organ,” is frequently characterized by hypoxia, as well as by abnormal levels of various cytokines, growth factors, and metabolites (Lyssiotis and Kimmelman, 2017). It provides numerous opportunities for cell-cell signals to modulate the epigenome and expression of stem-cell-like programs in cancer cells, frequently independent of their genetic backgrounds (Gin-gold et al., 2016).

Over the last decade, The Cancer Genome Atlas (TCGA) has illuminated the landscapes of primary tumors by generating comprehensive molecular profiles composed of genomic, epigenomic, transcriptomic, and (post-translational) proteomic characteristics (Hoadley et al., 2014; Tomczak et al., 2015), along with histopathological and clinical annotations. The resulting resource enabled us to analyze cancer stemness quite extensively in almost 12,000 samples of 33 tumor types.

First, we defined signatures to quantify stemness using publicly available molecular profiles from normal cell types that exhibit various degrees of stemness. By multi-platform analyses of their transcriptome, methylome, and transcription-factor binding sites using an innovative one-class logistic regression (OCLR) machine-learning algorithm (Sokolov et al., 2016), we obtained two independent stemness indices. One (mDNAsi) was reflective of epigenetic features; the other (mRNAsi) was reflective of gene expression. We then identified associations between the two stemness indices and novel oncogenic pathways, somatic alterations, and microRNA (miRNA) and transcriptional regulatory networks. Those features correlated with, and perhaps govern, cancer stemness in particular molecular subtypes of TCGA tumors. Importantly, higher values for stemness indices were associated with biological processes active in cancer stem cells and with greater tumor dedifferentiation, as reflected in histopathological grade. Metastatic tumor cells appeared more dedifferentiated phenotypically, probably contributing to their aggressiveness. We also found tumor heterogeneity at the single-cell level by measuring stemness in transcriptome profiles obtained from individual cancer cells. Using CIBERSORT to profile immune cell types in TCGA tumors, we

obtained insight into the interface of the immune system with stemness. Finally, we identified compounds specific to selected molecular targets and mechanisms that may eventually lead to novel treatments that trigger differentiation and exhaust the stemness potential of highly aggressive neoplasms.

## RESULTS

### DNA-Methylation- and mRNA-Expression-Based Stemness Classifiers

We analyzed publicly available non-tumor and tumor datasets for which transcriptomic and epigenomic molecular profiles were available (Figure 1A). We derived stemness indices using an OCLR algorithm trained on stem cell (ESC, embryonic stem cell; iPSC, induced pluripotent stem cell) classes and their differentiated ecto-, meso-, and endoderm progenitors. We chose OCLR because it does not penalize misclassification of stem-cell-derived progenitors at different stages of differentiation that still carry some of the undifferentiated features in their molecular profiles (its output was also validated against random forests in Figure S1A). OCLR-based transcriptomic and epigenetic signatures were applied to TCGA datasets to calculate the mRNAsi and mDNAsi. Each stemness index (si) ranges from low (zero) to high (one) stemness (Table S1). The tumor samples stratified by the indices were used for the integrative analyses.

### mRNA Expression-Based Stemness Index

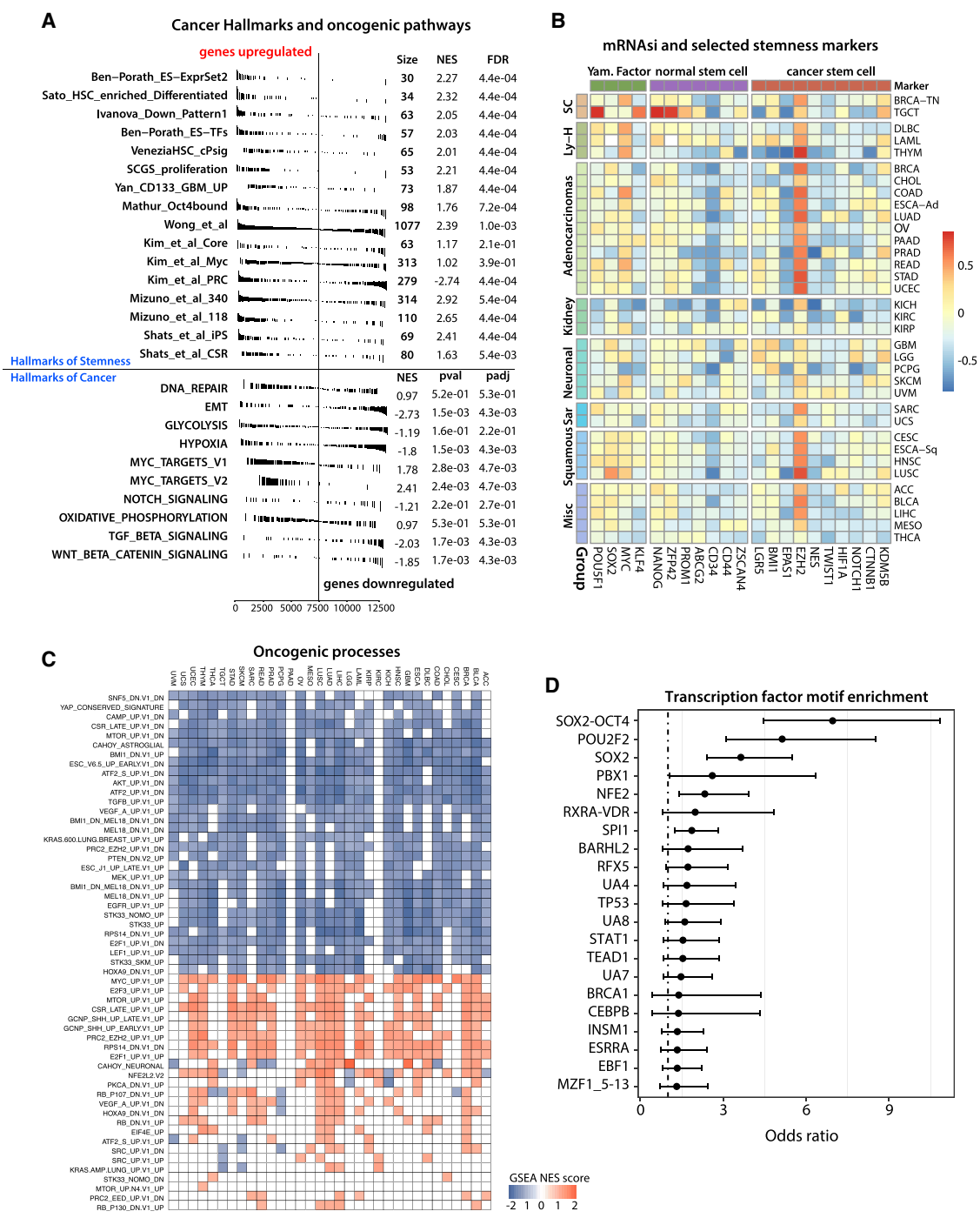
We validated the mRNAsi by applying it to an external dataset composed of both stem cells and somatic differentiated cells (Nazor et al., 2012) (Figure 1B) and by scoring molecular subtypes of breast cancers and gliomas that are characterized by different degrees of oncogenic dedifferentiation associated with pathology and clinical outcome (Figures S1B and S1C). All stem cell samples attained higher stemness index values than samples from differentiated cells. TCGA tumors display various degrees of cancer stemness as revealed by mRNAsi (Figure 1C [left]) and mDNAsi (Figure 1C [right]). Germ-cell tumors, basal breast cancer, and Ly-Hem cancers displayed highly dedifferentiated phenotypes in comparison to other tumor types.

Using gene set enrichment analysis (GSEA), we compared our signature to 16 gene sets that were associated with stemness in cancer and healthy cells in previous studies (Ben-Porath et al., 2008; Ivanova et al., 2006; Kim and Orkin, 2011; Mathur et al., 2008; Palmer et al., 2012; Sato et al., 2003; Venezia et al., 2004; Yan et al., 2011). These sets spanned 2,564 unique genes, with no 2 sets overlapping by more than 134 genes. In all cases, the published stemness gene sets were significantly enriched in mRNAsi (Figure 2A). We found that “cancer hallmark” gene sets were significantly enriched, as were MYC targets, which significantly contributed to the positive side of the signature (Hanahan and Weinberg, 2011). This is

(C) TCGA tumor types sorted by the stemness indices obtained from transcriptomic (mRNAsi) and epigenetic (mDNAsi) features; indices were scaled from 0 (low) to 1 (high). The TCGA tumor types were grouped based on their histology and cell of origin into stem cell-like (SC), lympho-hematopoietic (Ly-Hem), adenocarcinomas, squamous cell carcinomas (Squamous), neuronal lineage (Neuronal), sarcomas (Sar), kidney tumors (Kidney), and not belonging to any of the above (Misc) (Table S2).

See also Figures S1 and S2 and Tables S1 and S2.





**Figure 2. Biological Processes Associated with Cancer Stemness**

(A) Gene Set Enrichment Analysis showing RNA sequencing (RNA-seq)-based stemness signature evaluated in the context of gene sets representative for hallmarks of stemness and cancer.

(B) Correlation between mRNAi and mRNA expression for published hallmarks of stemness.

(C) Correlation between mRNAi and selected oncogenic processes.

(D) Association between the epigenomic-based stemness signature (EREG-mDNAi and EREG-mRNAi) and enrichment in the transcription factor binding sites.

See also Figure S2 and Table S2.

consistent with MYC being one of the transcription factors that drive pluripotency in ESCs (Young, 2011).

Wingless-related integration site (Wnt)/ $\beta$ -catenin and TGF- $\beta$  signaling pathways were significantly enriched on the negative side of the stemness signature. This negative enrichment does not imply absence of specific signals in cancer stem cells, but rather that this signaling is lower relative to stem-cell-derived progenitors, as captured by the signature weights. This is again consistent with other GSEA results, as both signaling pathways are known mediators of the epithelial-mesenchymal transition (EMT) mechanism (Gonzalez and Medici, 2014). We also computed the correlation of mRNAsi against mRNA expression of published pan-cancer EMT markers (Mak et al., 2016), which revealed significant correlations with for most tumors (Figure S2C). This is consistent with the biology of ESCs, which grow as epithelioid, polygonal cells *in vitro* and epithelial cancer precursors having stem-like properties. Importantly, most TCGA samples are primary tumors of an epithelial phenotype. Most skin melanoma cases come from lymph nodes, and this tumor type shows higher expression of vimentin, a key marker of a mesenchymal phenotype. mRNAsi is positively correlated with other core stem cell factors: *EZH2*, *OCT4*, and *SOX2* (Figure 2B and Table S2). Finally, Moonlight analysis of the oncogenic signatures from the Molecular Signatures Database (MSigDB) further validated our gene-expression-based index and confirmed engagement of MYC and *EZH2*, along with *E2F3*, *MTOR*, and *SHH* in driving oncogenic dedifferentiation (Figure 2C) (Colaprico et al., 2018).

### DNA Methylation-Based Stemness Index

We defined the mDNAsi using OCLR by combining (1) supervised classification between ESCs/iPSCs and their progenies, (2) stem cell signatures associated with pluripotency-specific genomic enhancer elements based on ChromHMM from Roadmap, and (3) ELMER, which uses DNA methylation to identify enhancer elements and correlates their state with the expression of nearby genes. 219 CpG probes (Figure S2A) were selected in training OCLR using the Progenitor Cell Biology Consortium (PCBC) datasets. By selecting probes previously defined to be active stemness-specific enhancers, we confirmed the ability of our approach to derive an mDNAsi. Since we focused exclusively on hypomethylated, functionally important CpG probes associated with stem cells, we further explored *cis*-activated genes.

We scored each TCGA sample using the mDNAsi and used an external dataset to confirm that stem cells had higher mDNAsi than differentiated samples (Figure 1B [left plot]). TCGA tumor types show different degrees of an inferred dedifferentiated phenotype (Figure 1C [right]). Within these, individual tumor samples show variation for cancer stemness. As anticipated, TCGA samples derived from the primary tumors show higher cancer stemness indices compared to non-tumor samples obtained from adjacent normal tissue of origin (Figure S1E [bottom]).

Most of our selected probes fell within non-promoter elements, yet the *SOX2*-*OCT4* transcription factor binding motif is one of the most highly enriched signatures within these regions. The *SOX2*-*OCT4* complex is a critical master regulator of pluripotency and stemness and is highly enriched in tumor samples with high mDNAsi (Figure 2D).

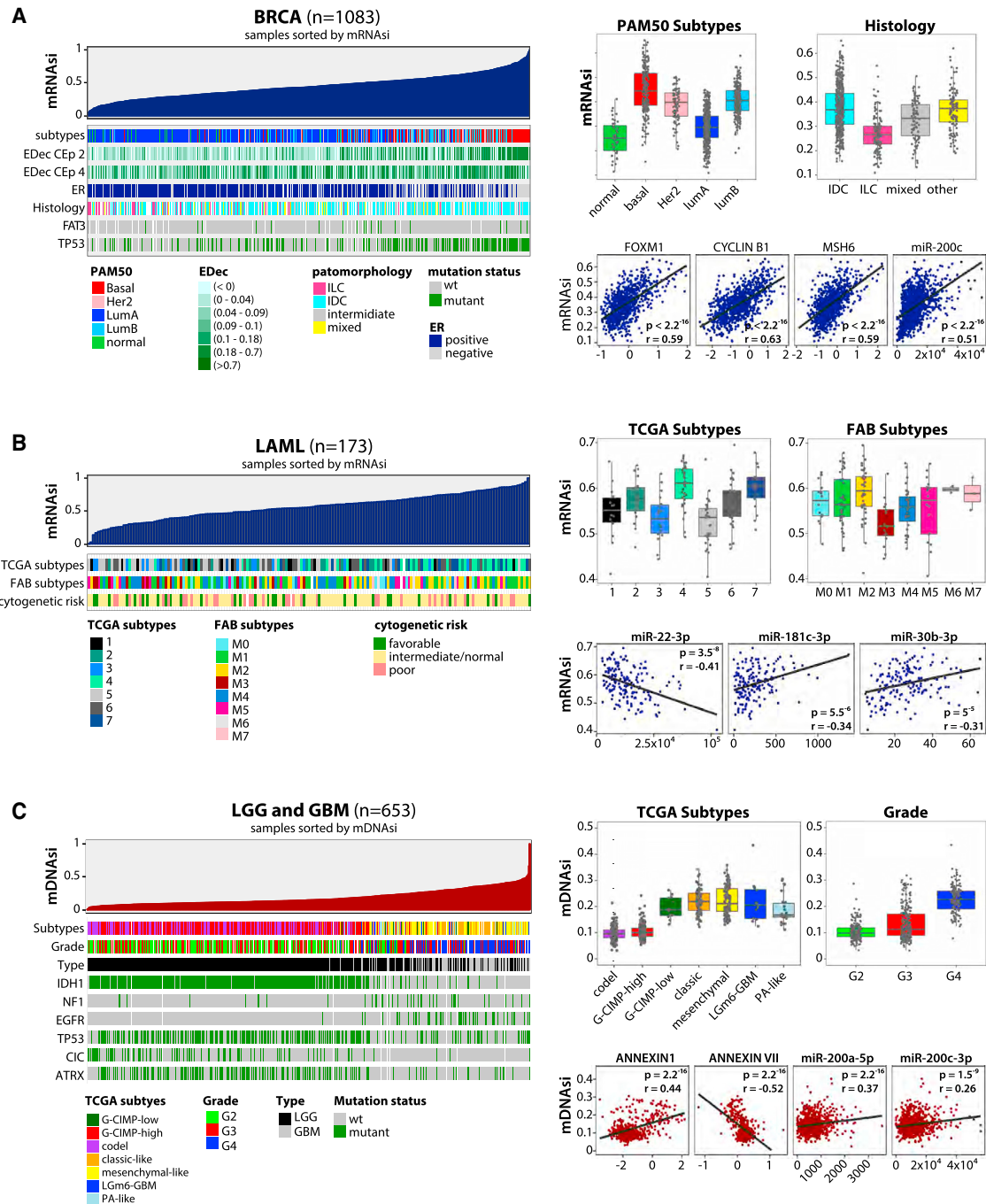
### Correlations of mRNAsi and mDNAsi

Since the inputs for mDNAsi and mRNAsi are not necessarily complementary, we explored stratification of glioma samples by the epigenetically regulated mRNAsi (EREG-mRNAsi), as stemness index generated using a set of stemness-related epigenetically regulated genes. The EREG-mRNAsi, based on both RNA expression and epigenetics, elucidates the discrepancy between mDNAsi and mRNAsi and shows a positive correlation with both indices (Figure S1F). Both mRNAsi and mDNAsi show good correspondence for a majority of tumors (Figures S1F and S2B). We observed major discrepancies in the case of brain lower grade glioma (LGG), thyroid carcinoma (THCA), and thymoma (THYM). For gliomas, mDNAsi is correlated positively with tumor pathology and clinical features, while mRNAsi shows a negative correlation. This result could arise from a high frequency of isocitrate dehydrogenase mutations (*IDH1/2*) mutations and resulting DNA hypermethylation.

### Stemness Index Can Stratify Recognized Undifferentiated Cancers

We examined breast invasive carcinoma (BRCA), acute myeloid leukemia (AML), and gliomas to study if the mRNAsi/mDNAsi predict stemness in poorly differentiated tumors. In BRCA, we found a strong association between the stemness index and known clinical and molecular features (Figure 3A [left]). The mRNAsi was highest in the basal subtype, known to exhibit an aggressive phenotype associated with an undifferentiated state. BRCA samples with high mRNAsi were more likely to be estrogen receptor (ER)-negative and enriched for *FAT3* and *TP53* mutations. We noted that high mRNAsi was associated with higher protein expression of *FOXM1*, *CYCLINB1*, and *MSH6*, as well as higher miRNA-200 family expression (Figure 3A [right]). Invasive lobular type of BRCA (ILC), characterized by better prognosis in comparison to invasive ductal carcinoma (IDC), has a lower mRNAsi (Figure 3A [right]). We also applied our indices to non-TCGA BRCA samples (Reynold et al., 2014) and found a similar correlation between mRNAsi and mDNAsi in those samples. Moreover, mRNAsi also stratified BRCA samples with distinct histology in this dataset (Figure S1B). Using datasets with estimated tumor cell type composition provided by the epigenetic deconvolution method (Onuchic et al., 2016), we found that both mRNAsi and mDNAsi were more highly correlated with malignant epithelial cells than with normal epithelial cells, suggesting that our indices identify distinct cancerous epithelial cell populations characterized by different features or degrees of stemness (Figure S1D).

We found an association between the mRNAsi, RNA expression subtypes previously defined by TCGA, and the French-American-British (FAB) classification of AML (Figure 3B). The mRNAsi showed the strongest correlation with the stage of myeloid differentiation of the AML samples. FAB subtypes M0 (undifferentiated), M1 (with minimal maturation), and M2 (with maturation) were characterized by high mRNAsi. In contrast, the M3 well-matured promyelocytic subtype, which is associated with benign chromosomal abnormalities and favorable clinical outcome, had low mRNAsi (Figure 3B [right upper]). High mRNAsi was associated with higher expression of



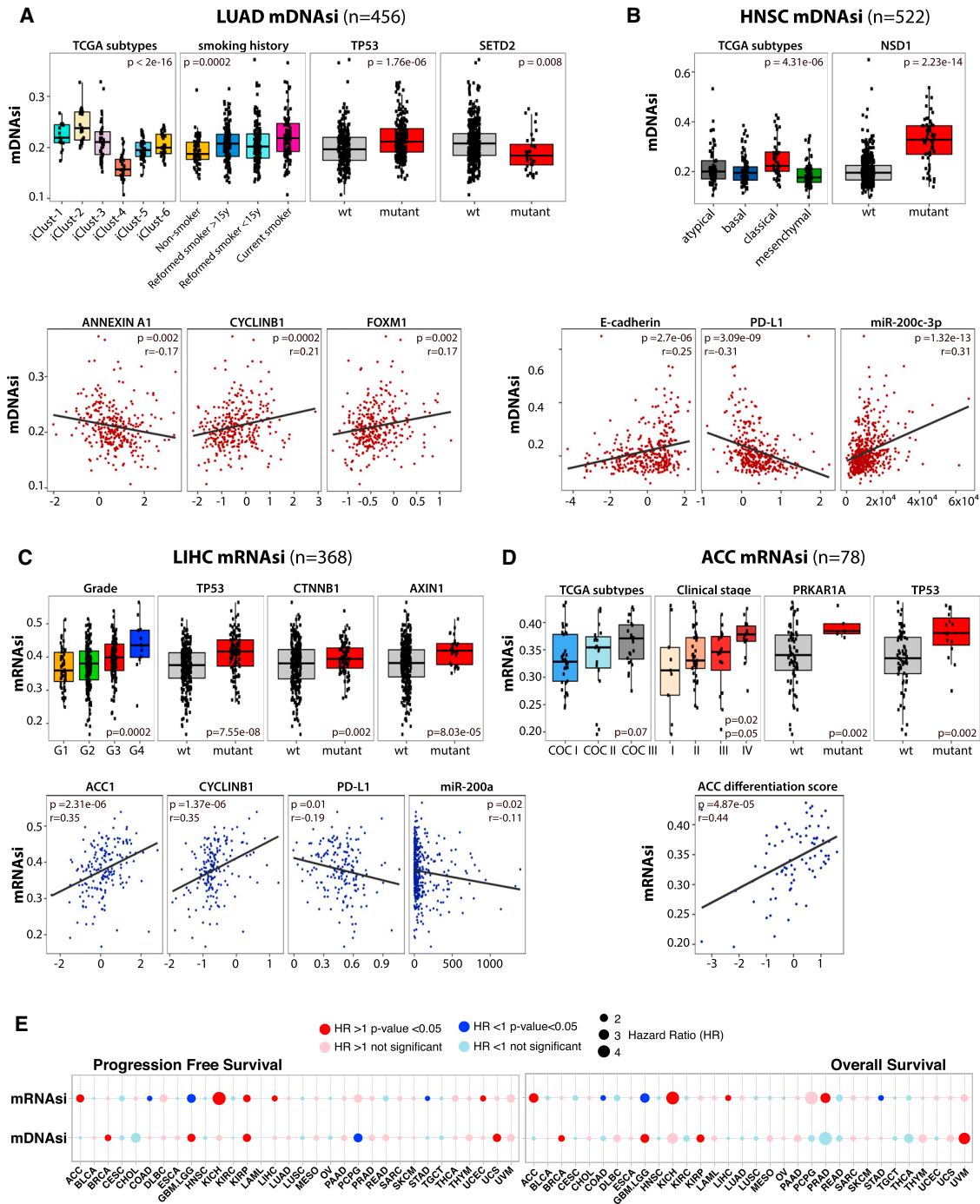
**Figure 3. Molecular and Clinical Features Associated with Stemness in Breast Cancer, Acute Myeloid Leukemia, and Gliomas**

(A) (Left) An overview of the association between known molecular and biological processes and stemness in BRCA. Columns represent samples sorted by mRNAsi from low to high (top row). Rows represent molecular and biological processes associated with mRNAsi. Rows named “EDec CEp 2 and 4” represent estimated cell type proportions. (Top right) Boxplots of mRNAsi in individual samples, stratified by molecular subtype and histology. (Bottom right) Correlation of mRNAsi and representative protein expression and microRNA.

(B) Similar to (A), association of mRNAsi in AML. (Top right) mRNAsi by mRNA-based molecular subtype and by FAB classification. (Bottom right) Correlation scores of mRNAsi and representative microRNA.

(C) As in (A) and (B), GBM and LGG sorted by mDNAsi. (Top right) mDNAsi by molecular subtype and grade. (Bottom right) Correlation scores of mDNAsi and representative protein expression and microRNA.

All molecular and clinical features shown are statistically significant. See also [Figures S1, S3, S4, and S5](#).



**Figure 4. Selected Molecular and Clinical Features Associated with the Stemness Indices in TCGA Tumors**

(A) Association of molecular and clinical features with stemness in LUAD. (Top) mDNAsi by integrative molecular subtypes, smoking history, and mutations of TP53 and SETD2. (Bottom) Correlation scores of mDNAsi and representative protein expression.

(B) Stemness in HNSC. (Top) mDNAsi stratified by molecular subtypes and mutation of NSD1. (Bottom) Correlation scores of mDNAsi and representative protein and microRNA expression.

(C) Stemness in LIHC. (Top) mRNAsi stratified by grade and mutations of TP53, CTNNB1, and AXIN1. (Bottom) Correlation scores of mRNAsi and representative protein and microRNA expression.

(legend continued on next page)

miR-181c-3p, miR-22-3p, and miR-30b-3p (Figure 3B [right bottom]).

We found a strong association between high mDNAsi, high pathologic grade, and recently published molecular subtypes of glioma (Figure 3C). mDNAsi was low in less-aggressive gliomas that are characterized by codel and glioma CpG methylator phenotype (G-CIMP)-high features and was highest in highly aggressive glioblastoma multiformes (GBMs) characterized by IDH mutations (G-CIMP low) and poor clinical outcome. Also, high mDNAsi is strongly associated with more aggressive classical and mesenchymal subtypes of GBM, suggesting that it can stratify tumors with distinct clinical outcomes. We also found that high mDNAsi was associated with mutations in *NF1* and *EGFR* and infrequent mutations in *IDH1*, *TP53*, *CIC*, and *ATRX* (Figure 3C [left]), with higher expression of ANNEXIN-A1 protein and lower expression of ANNEXIN-A7 and with expression of the miR-200 family (Figure 3C [right bottom]).

We obtained similar results on non-TCGA glioma samples for which both mRNA expression and DNA methylation data were available (Turcan et al., 2012) (Figure S1C). The negative correlation between mDNAsi and mRNAsi was restricted to LGG samples—specifically, the IDH mutant subtypes (G-CIMP high and codel). *IDH1* mutations are known to reduce cell differentiation, and high values of the mRNAsi in a subset of IDH mutant gliomas might capture this phenomenon (Lu et al., 2012).

### Pan-cancer Stemness Landscape

Next, we tested the ability of our indices to identify previously unexplored features of cancer stemness across all TCGA tumors. First, we performed an enrichment analysis by sorting all TCGA samples by stemness index for each tumor type and looking for associations with mutations and molecular and clinical features. The most salient associations of mRNAsi and mDNAsi are presented in Figure 4, while the following results of the comprehensive analyses are shown in the supplementary material: associations with mutations (Figure S3), associations with miRNA expression and protein abundance (Figure S4), associations with the tumor grading, and clinical outcome (Figure S5).

### Correlations of mRNAsi and mDNAsi with Mutations in Genes, miRNA, and Expression of Proteins

We found a strong association of mDNAsi with known molecular subtypes, with somatic mutations in *SETD2* and *TP53* genes, and with tobacco smoking status in lung adenocarcinoma (LUAD) (Figures 4A and S3). Current smokers and recently reformed smokers have higher mDNAsi than non-smokers or long-term reformed smokers. This suggests that the stemness of LUAD tumors might be activated in response to environmental stimuli such as smoking and might influence the aggressiveness of the tumor. We also found an association between mDNAsi and higher protein expression of CYCLINB1 and FOXM1, which is a pro-stemness transcription factor upstream of CYCLINB1 (Fig-

ure 4A [lower plots]). FOXM1 has been associated with dedifferentiation in pancreatic cancer cells (Bao et al., 2011), as well as tumor proliferation in the kidney (Xue et al., 2012) and ovarian (Wen et al., 2014) cancers. Our result suggests that it could be a driver of dedifferentiation and proliferation in breast and lung cancers. Stemness of LUAD tumors is also associated with lower expression of ANNEXIN-A1 (Figure 4A). ANNEXIN-A1 has been indicated as a differentiation marker in pancreatic (Bai et al., 2004) and urothelial (Kang et al., 2012) cancers; therefore, we suspect that the relationship between ANNEXIN and FOXM1 expression and tumor differentiation may extend to other tumor types (Figure S4C).

Analyses of head and neck squamous cell carcinoma (HNSC) samples revealed that high indices are correlated with *NSD1* mutation, E-cadherin protein expression, miR-200-3p, and previously identified classical molecular subtypes (Figure 4B). *NSD1* mutation was recently linked in HNSC tumors to blockade of cellular differentiation and promotion of oncogenesis (Papillon-Cavanagh et al., 2017). Interestingly, miR-200 family members have been implicated in cancer initiation and metastasis, as well as self-renewal of healthy stem cells (Gregory et al., 2008; Tellez et al., 2011). HNSC tumors with high mDNAsi have reduced programmed death ligand 1 (PD-L1) protein level (Figure 4B).

In liver hepatocellular carcinoma (LIHC) samples, we found an association between mRNAsi and high pathological grade (Figure 4C). Negative associations between mRNAsi and the probability of overall survival (OS) or progression-free survival (PFS) were detected (Figures 4E and S5C). In contrast to the majority of tumor types, LIHC samples with high mRNAsi have low expression of miR-200 family members (Figure 4C). The miR-200 family is known to be associated with progression of hepatocellular carcinoma (Tsai et al., 2017; Wong et al., 2015), and the miR-200b-ZEB1 circuit has been suggested as a master regulator of stemness in these cancers (Tsai et al., 2017). We found associations of mRNAsi with higher CYCLINB1 and ACC1 and with lower PD-L1 and ANNEXIN-A1 protein expression in LIHC (Figure 4C). ACC1 was associated with pathomorphological markers of LIHC aggressiveness (vascular invasion and poor differentiation), and its upregulation was correlated with poor OS and disease recurrence in hepatocellular carcinoma patients (Wang et al., 2016). LIHC samples with high mRNAsi were associated with specific genomic alterations (e.g., *TP53*, *CTNNB1*, *AXIN1*).

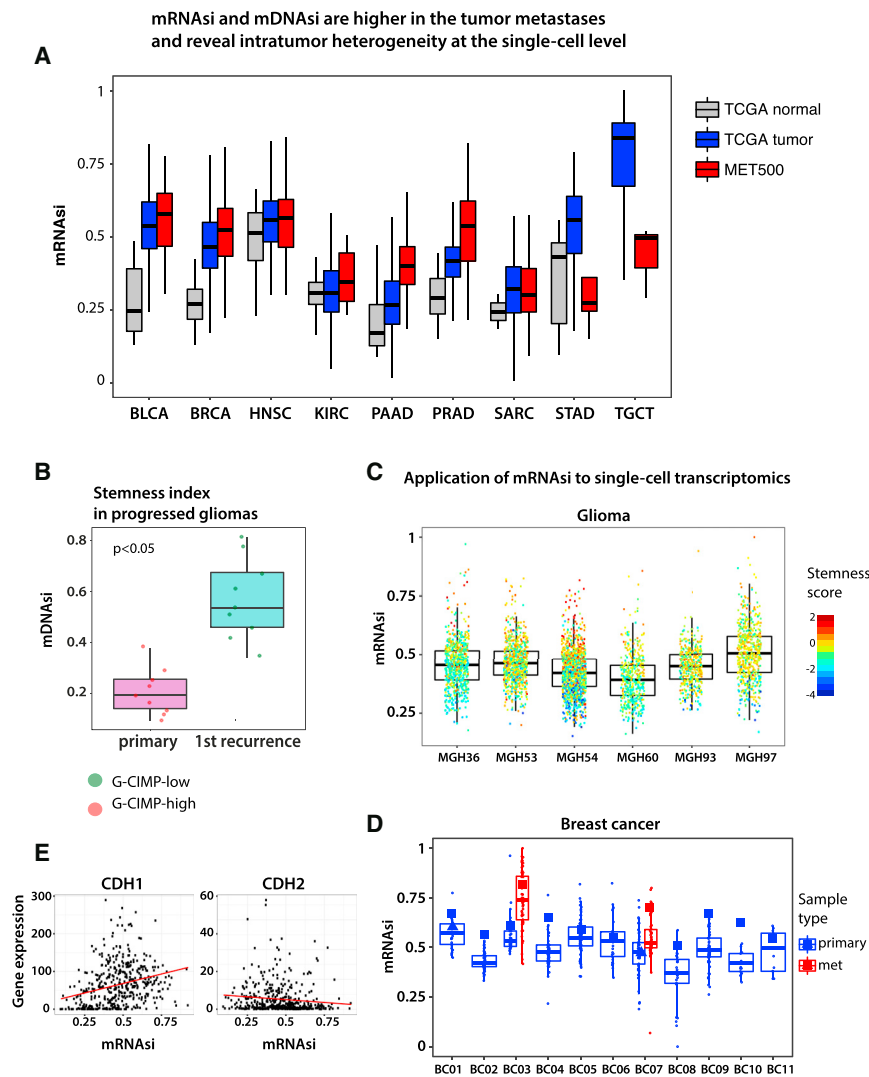
Detailed analyses of adrenocortical carcinoma (ACC) samples revealed an association between high mRNAsi and defined molecular subtypes (Zheng et al., 2016), clinical stage, and mutations in *PRKAR1A* and *TP53* genes (Figure 4D). We found a positive correlation between mRNAsi and adrenal differentiation score that is based on expression of 25 genes that are important for adrenal function (Zheng et al., 2016) (Figure 4D).

(D) Stemness in ACC. (Top) mRNAsi stratified by mRNA molecular subtypes, clinical stage, and mutations of *PRKAR1A* and *TP53*. (Bottom) Correlation scores of mRNAsi and adrenal differentiation score.

(E) Cox proportional hazards model analysis. (Left) Progression-free survival. (Right) Overall survival. Hazard ratio greater than one denotes a trend toward higher stemness index with worse outcome.

See also Figures S3, S4, and S5.





**Figure 5. Analysis of Cancer Stemness in the Context of Metastatic State and Intratumor Heterogeneity**

(A) mRNAsi is higher in cancer metastases in comparison to the TCGA primary tumors.

(B) mDNAsi is higher in recurrent glioma samples compared to the primary glioma occurrence from the same patient. G-CIMP, glioma CpG methylator phenotype.

(C and D) Application of mRNAsi to a single-cell transcriptome of gliomas and breast cancer reveal intratumor heterogeneity and various degrees of the oncogenic dedifferentiation.

(E) Correlation of mRNAsi and mRNA expression of CDH1 (epithelial marker) and CDH2 (mesenchymal marker) in the cancer metastases.

clinical factors. We found a positive correlation between previously published glioma subtypes and mDNAsi, suggesting that mDNAsi might recapitulate prognostic molecular subtypes (Figure 3C). The discordance between the mRNAsi and the mDNAsi for gliomas may be explained in part by the dominant genomic alteration associated with the LGG tumor type. Roughly 80% of LGG tumors carry an *IDH1/2* mutation and, as demonstrated by our group and others, confer a genome-wide hypermethylator phenotype (G-CIMP) (Noushmehr et al., 2010; Turcan et al., 2012). Given that the mDNAsi is driven primarily by low methylation levels associated with the stemness phenotype, the LGG tumors might resemble non-stem like phenotypes, which are predominantly hypermethylated. The subgroup of G-CIMP with the lowest overall DNA methylation levels

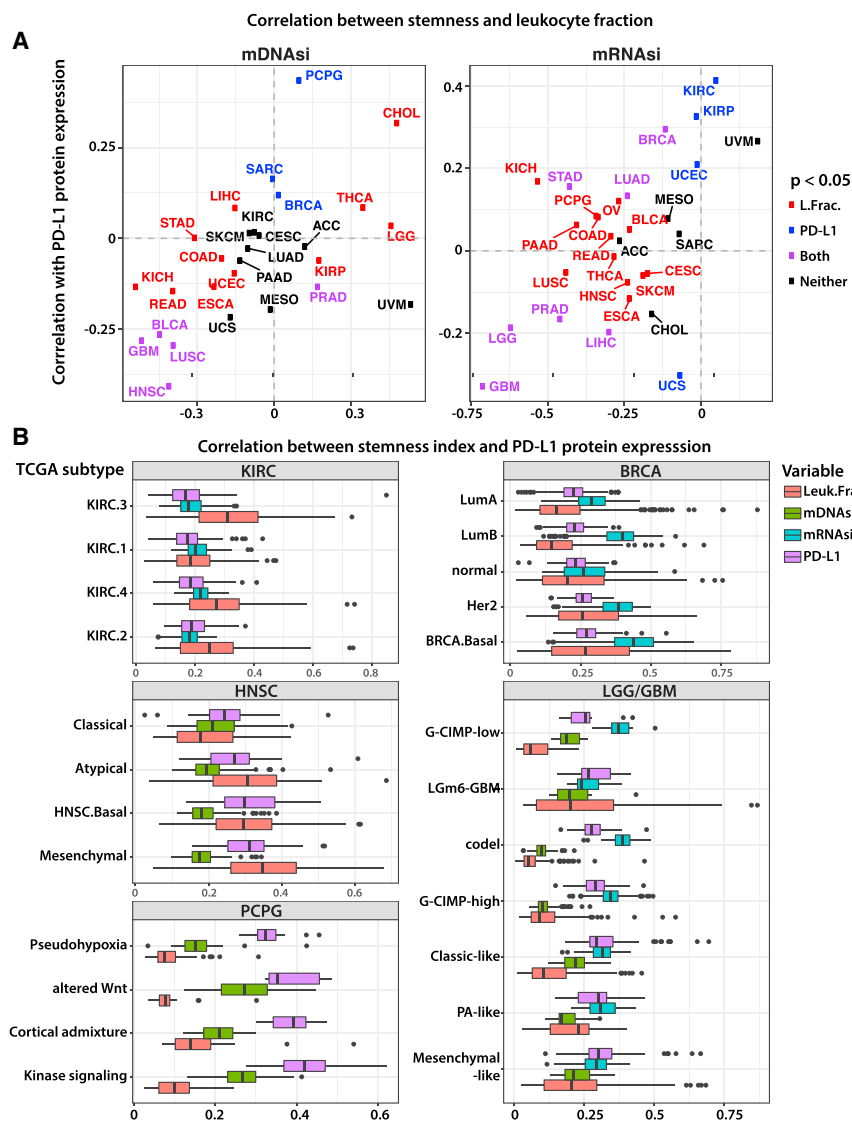
### Stemness Indices Are Correlated with Tumor Pathology and Predictive of Clinical Outcome

We observed a positive correlation between tumor histology and pathology grading and both stemness indices for the majority of the TCGA cases (Figures 3A, 3C, 4C, 4D, S1B, S5A, and S5B). For mRNAsi, the most significant correlations were found for BRCA (IDC and ILC), cervical squamous cell carcinoma and endocervical adenocarcinoma (CESC), LIHC, pancreatic adenocarcinoma (PAAD), and uterine corpus endometrial carcinoma (UCEC) (Figure S5A). Interestingly, mRNAsi shows low values in GBM and stomach adenocarcinoma (STAD). On the other hand, mDNAsi strongly stratifies glioma by the pathology grade, culminating with the highest value for GBM (Figure S5B). The reversed values of mDNAsi and mRNAsi in case of gliomas were also evident in the clinical data analyses. An adverse association between the mRNAsi and survival was detected (Figure 4E), which was significant for OS and PFS after adjusting for clinical factors (Figures S5C). In contrast, the mDNAsi had no significant association with OS and PFS after correcting for

(G-CIMP low) is associated with the worst outcomes. Compared to G-CIMP-high tumors, G-CIMP-low tumors are known to be more proliferative, express cell-cycle-related genes, and have various stem-cell-like genomic features (Ceccarelli et al., 2016).

### Cancer Stemness Indices Are Higher in Tumor Metastases and Reveal Intratumor Heterogeneity

The TCGA samples are derived mostly from primary tumors, except for skin melanoma, for which tissues are mostly metastatic lymph nodes. We used the mRNAsi to interrogate the MET500 dataset comprising expression profiles from 500 metastatic samples obtained from 22 different organs (Robinson et al., 2017). In most cases, mRNAsi was significantly higher in metastatic samples compared to primary TCGA tumors (Figure 5A). Prostate and pancreatic adenocarcinoma metastases had the most dedifferentiated phenotypes and are also more aggressive and resistant to therapies in contrast to primary tumors. Weaker association with the mRNAsi was due to a small number of available samples ( $n < 20$ ). Interestingly, testicular germ cell tumors



**Figure 6. Association of Stemness Index with Immune Microenvironment**

(A) mDNAsi and mRNAsi in the context of immune microenvironment. Each panel shows the Spearman correlation between the stemness index and PD-L1 protein expression plotted against Spearman correlation between the same stemness index and total leukocyte fraction, as estimated from DNA methylation data.

(B) Highlight of tumor types that exhibit strong correlation between stemness and PD-L1 expression or total leukocyte fraction. See also Figure S6.

had higher stemness index in breast cancer (Figure 5D). Interestingly, the negative correlation of EMT signature and stemness that we observed in TCGA primary tumors was also found in metastatic samples (Figure 5E).

**Stemness Index Evaluated in the Context of Immune Response**

We found that for many tumors, higher stemness indices are associated with a reduced leukocyte fraction and lower PD-L1 expression (Figure 6A). For mDNAsi, the most distinctive negative correlations were found in the PanCan-12 squamous cluster (LUSC, lung squamous cell carcinoma; HNSC; BLCA, bladder urothelial carcinoma) (Hoadley et al., 2014) and in GBM (Figures 6A [left panel] and S6B). For the mRNAsi, the highest negative correlation values were seen in GBM/LGG, prostate adenocarcinoma (PRAD), LIHC, and uterine carcinosarcoma (UCS) tumors (Figures 6A [right panel] and S6A). We expect that such tumors will be less susceptible to immune

(TGCTs) present the less differentiated phenotypes in primary tumors when compared to distant metastases. Primary TGCT tumor cells have high mRNAsi and may differentiate when metastasizing to distant organs. A similar trend was observed for STAD.

Using another dataset, we found that mDNAsi was significantly higher in glioma samples obtained at first recurrence in contrast to primary gliomas (Figure 5B). Our results reveal significant dedifferentiation of glioma cancer cells that contribute to glioma recurrence, which is frequently associated with poor prognosis and resistance to treatment (de Souza et al., 2018).

By taking advantage of single-cell transcriptome datasets, we used mRNAsi to probe tumor heterogeneity for oncogenic dedifferentiation of individual cancer cells (Chung et al., 2017; Tirosch et al., 2016). We revealed high variation of stemness in the glioma and breast primary tumors. Individual glioma cells showed higher variation of oncogenic dedifferentiation in comparison to breast cancer cells (Figure 5C). Single cells from metastases

checkpoint blockade treatments due to insufficient immune cell infiltration or preexisting downregulation of the PD-L1 pathway, which makes further inhibition ineffective. Our findings are consistent with previous reports showing a strong correlation between PD-L1 protein expression and infiltration of CD8+ cytotoxic lymphocytes (Zaretsky et al., 2016).

We further explored correlations between stemness and immune microenvironment variables in the context of molecular subtypes of tumors. Figure 6B highlights several tumor types with the strongest (positive or negative) correlations. Except for kidney renal clear cell carcinoma (KIRC), the association between stemness and PD-L1 expression and leukocyte fraction is readily apparent from the increasing and decreasing trends of individual variables across the molecular subtypes. For example, we found mesenchymal tumors to have the highest PD-L1 expression levels, the most significant leukocyte fractions, and the lowest mDNAsi compared to other HNSC subtypes, suggesting potential susceptibility to checkpoint

blockade inhibitors. The use of immunotherapy for HNSC tumors is under active investigation (Economopoulou et al., 2016; Fuereder, 2016) with the recent FDA approval of pembrolizumab; however, whether the effectiveness of therapy is limited to specific HNSC molecular subtypes is not clear from those reports.

To assess other relationships between stemness and tumor microenvironment, we computed correlations between stemness indices and individual types of immune cells. By applying CIBERSORT, we scored 22 immune cell types for their relative abundance in TCGA tumor samples. These cell types included natural killer (NK) cells, monocytes, macrophages, dendritic and mast cells, eosinophils, and neutrophils. We also obtained absolute estimates by scaling their relative abundance by overall leukocyte infiltration in each tumor as determined by ESTIMATE applied to DNA methylation data. For any given TCGA sample, we calculated the correlation between mDNAs/mRNAsi and the estimated fraction of individual immune cell types. In addition to individual immune subpopulation fractions, we considered the functional activation of distinct cells by measuring the difference between activated and resting fractions of NK cells, CD4+ T cells, and macrophages. This approach was motivated by recent observations that activation of peripheral CD4+ T cells triggered by immunotherapy is responsible for the specific killing of tumor cells (Spitzer et al., 2017).

Although the squamous cluster tumors had a negative correlation between stemness and the fraction of CD4+ T cell populations, the activation state of the CD4+ T cells was higher in dedifferentiated tumors. This finding is consistent with our observation that PD-L1 protein expression is lower in these tumors, suggesting again that immune checkpoint blockade might be ineffective, and an additional mechanism of immune evasion may be operative. The opposite trend is present in thymomas, where PD-L1 protein expression and the fraction of the CD4+ T cell population are positively correlated with tumor dedifferentiation. Likewise, the activation state of CD4+ T cells is lower in dedifferentiated tumors, suggesting that they might be more susceptible to immunotherapy treatments (Figures S6A and S6B).

### Connectivity Map Analysis Identifies Potential Compounds/Inhibitors Capable of Targeting the Stemness Signature

We employed the Connectivity Map (CMap), a data-driven, systematic approach for discovering associations among genes, chemicals, and biological conditions, to search for candidate compounds that might target pathways associated with stemness. We found enrichment for compounds associated with stemness in at least three cancer types (Figure 7A). 5 compounds are significantly enriched in more than 10 cancer types and have been reported to inhibit stemness-related tumorigenicity: the dopamine receptor antagonists thioridazine and prochlorperazine (Cheng et al., 2015; Lu et al., 2015; Dolma et al., 2016), the Wnt signaling inhibitor pyrvinium (Xu et al., 2016), the HSP90 inhibitor tanespimycin, and the protein synthesis inhibitor puromycin. Further, the telomerase inhibitor gossypol induced apoptosis and growth inhibition of cancer stem cells (CSCs) (Volate et al., 2010), and histone deacetylase inhibitors such as trichostatin A (SAHA) reduced glioblastoma stem cell

growth (Chiao et al., 2013). According to our analysis, pyrvinium and puromycin could inhibit stemness in LUAD. We found several candidates with recognized anti-CSC activity for HNSC, including the aforementioned compounds. For LIHC, thioridazine, a prospective inhibitor of lung cancer stem cells (Yue et al., 2016), pyrvinium, puromycin, prochlorperazine, and others are potential compounds targeting undifferentiated tumors (Figure 7).

CMap mode-of-action (MoA) analysis of the 74 compounds revealed 56 mechanisms of action shared by the above compounds (Figure 7B and Table S4B). Five compounds (fluspirilene, pimozide, prochlorperazine, thioridazine, and trifluoperazine) shared the MoA of dopamine receptor antagonist. We observed that entinostat, trichostatin-a, and vorinostat shared MoA as HDAC inhibitors, and LY-294002, zaprinast, zardaverine share MoA as phosphodiesterase inhibitors.

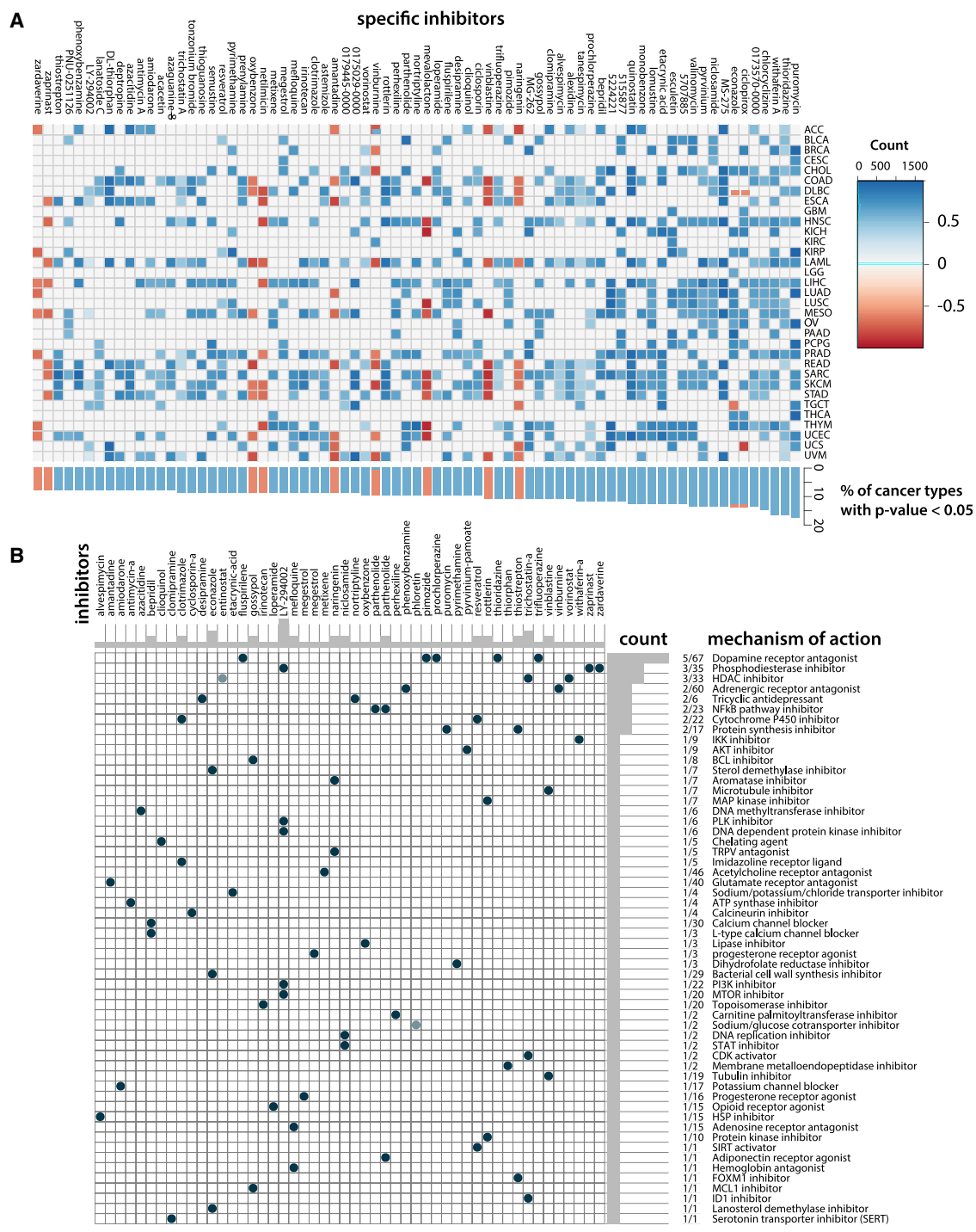
CMap target analysis revealed 212 distinct drug-target genes shared by the mentioned compounds (Figure S7 and Table S4C). Eight genes are targets of five different compounds—namely, DRD2 (8 drugs), HTR2A (7 drugs), HRH1 (6 drugs), ADRA1A (5 drugs), CALM1 (5 drugs), CHRM3 (5 drugs), HTR1A (5 drugs), and HTR2C (5 drugs).

Recent polypharmacology studies suggest the need to design compounds that act on multiple genes or molecular pathways. In this study, we observed similar mechanisms of action among different compounds, suggesting that selective therapies can target the undifferentiated phenotypes for selected cancer types.

## DISCUSSION

This study is based on integrated analysis of cancer stemness in almost 12,000 primary human tumors of 33 different cancer types. We interrogated TCGA data for mutations, DNA methylation, expression of mRNA and miRNA, expression and post-translational modification of proteins, histopathological grade, and clinical outcome. Applying CIBERSORT, we gained insight into the tumor microenvironment and composition of immune cell infiltrates. By applying a machine-learning algorithm to molecular datasets from normal stem cells and their progeny, we developed two different molecular metrics of stemness and then used them to assess epigenomic and transcriptomic features of TCGA cancers according to their grade of oncogenic dedifferentiation. Ultimately, the analyses led us to potentially actionable targets (and their MoAs) as candidates for possible differentiation therapy of solid tumors and metastases. Our approach could be applied to longitudinal study of samples from primary, recurrent, and metastatic cancers, and gene expression signatures identified in the tumor samples can be used to interrogate CMap to suggest actionable targets and inhibitors for further analysis.

To the best of our knowledge, this is the first study in which molecular PCBC datasets comprised of stem cells and defined populations of their differentiated progeny have been leveraged to develop a classification tool and machine-learning algorithm for analysis of a spectrum of human malignancies. A number of cancer stemness scores, based on genes that are differentially expressed between CSCs and non-CSCs, have been published



**Figure 7. Correlation of Cancer Stemness With Drug Resistance: Connectivity Map Analysis**

(A) Heatmap showing enrichment score (positive in blue, negative in red) of each compound from the CMap for each cancer type. Compounds are sorted from right to left by descending number of cancer type significantly enriched.

(B) Heatmap showing each compound (perturbagen) from the CMap that shares mechanisms of action (rows) and sorted by descending number of compound with shared mechanisms of action.

See also [Figure S7](#) and [Tables S3](#) and [S4](#).

and are relevant to clinical outcomes in AML (Eppert et al., 2011; Gentles et al., 2010; Ng et al., 2016). In those studies, gene sets enriched in ESCs (e.g., targets of NANOG, OCT4, SOX2, and c-MYC) were frequently overexpressed in poorly differentiated tumors compared with well-differentiated ones. In breast cancers, those gene sets were associated with high-grade estrogen receptor-negative, basal-like tumors and poor clinical outcome (Ben-Porath et al., 2008). Another web-based tool, StemChecker, uses a curated set of 49 published stemness signatures defined by gene expression, RNAi screens, transcription factor binding sites, text mining of the literature, and other computational approaches. But it has been tested only for pancreatic ductal adenocarcinoma. In that case, high expression of stemness genes correlated with poor prognosis (Pinto et al., 2015). All previous studies were transcriptome-based and limited to a narrow set of genes and a small number of tumor types.

In the present study, we found oncogenic dedifferentiation to be associated with several characteristics: mutations in genes that encode oncogenes and epigenetic modifiers, perturbations in specific mRNA/miRNA transcriptional networks, and deregulation of signaling pathways. Cancer stemness also appeared to involve core expression of *MYC*, *OCT4*, *SOX2*, and other genes involved in the regulatory circuitry that underlies normal and malignant self-renewal potential. Our indices derived from mRNA expression and DNA methylation signatures reliably stratified tumors of known stemness phenotype. High mRNAsi was associated with basal breast carcinomas but also Her2 and lumB subtypes that are more aggressive than the hormone-dependent lumA group. In contrast, high mDNAsi was strongly associated with high-grade glioblastomas, poor OS, and PFS. The association between stemness signatures and adverse outcome for some tumor types, including gliomas, may reflect malignant cell origins or the impact of their microenvironment.

Dedifferentiated cells can arise from different sources: from long-lived stem or progenitor cells that accumulate mutations in oncogenic pathways or via dedifferentiation from non-stem cancer cells that convert to CSCs through deregulation of developmental and/or non-developmental pathways. It is important to distinguish between the inherent stemness of CSCs and dedifferentiation induced by the tumor microenvironment. However, addressing that issue would require further validation beyond the scope of this study using other genomic datasets and/or laboratory experiments.

Both stemness indices were lowest in normal cells, increased in primary tumors, and highest in metastases, consistent with the idea that tumor progression generally involves oncogenic dedifferentiation. Interestingly, we observed negative associations between stemness and EMT gene signatures. The relationship between EMT and stemness remains a hotly debated topic, with several studies showing that EMT is necessarily associated with stemness (Fabregat et al., 2016). However, most TCGA data are obtained from primary tumors, which exist in a pre-EMT state, since EMT is strongly associated with tumor progression and with metastasis for many tumor types. Cancer cells in many primary solid tumors are basically epithelial regardless of their degrees of dedifferentiation, but some cells in such con-

texts could acquire mesenchymal characteristics either by accumulating additional mutations or by undergoing epigenetic changes shaped by the tumor microenvironment. Those mesenchymal cells can traverse the underlying tissue, enter the bloodstream, and seed distant organs, where they reacquire an epithelial phenotype to form metastatic tumors.

We observed epithelial phenotypes and increased stemness index in molecular profiles of tumor-type-matched metastatic samples in the MET500 cohort. This portends an association between dedifferentiation and spread of tumor cells to distant organs. The observation is further supported by high mDNAsi in samples from recurrent gliomas. It appears that tumor growth *de novo*, or at recurrence/metastasis, is associated with an increased stemness phenotype. Decreased mRNAsi levels seen in TGCT suggest its possible differentiation as a germ cell tumor type induced by the microenvironment of liver or lung parenchyma, the organs it most often colonizes. Clinically, in general, tumor progression is associated with greater aggressiveness and resistance to therapy of almost all types.

The mRNAsi was high for individual primary glioma and breast cancer cells. Interestingly, when applied to transcriptomic profiles obtained from analysis of single cancer cells in bulk tumors, stemness indices revealed a high degree of intratumor heterogeneity with respect to dedifferentiation phenotype. The heterogeneity was greater in gliomas than in breast cancer, suggesting that intratumor environment, including stromal cells, hypoxia, and infiltration of immune cells, may play a role in shaping CSC niches and affect cancer cell developmental plasticity. Further molecular analyses of cancer cells stratified by the stemness phenotype would provide novel insights into the biology of primary tumors.

We found that for a number of tumor types (GBM, LUSC, HNSC, and BLCA), higher mDNAsi was associated with reduced leukocyte fraction and/or lower PD-L1 expression. Such tumors are expected to be less susceptible to immune checkpoint blockade, due either to insufficient immune cell infiltration of tumors or to inherent downregulation of the PD-L1 pathway. Both factors can render immune checkpoint immunotherapy ineffective. The interaction between PD-L1 on cancer cells and PD1 receptor on T cells helps cancer cells elude the immune system by preventing activation of cytotoxic T cells in lymph nodes and subsequent recruitment of other immune cell types to the tumor site (Chen and Mellman, 2013). The presence of tumor-infiltrating lymphocytes and/or PD-L1 expression correlates with aggressiveness in gastrointestinal stromal tumors (Bertucci et al., 2015) and breast carcinomas (Polónia et al., 2017). Common features shared between cancer cells and stem cells in the context of the immune response are being highlighted by a growing number of studies showing that vaccination with ESCs or iPSCs can raise specific immune response against cancer cells (Kooreman et al., 2018). That finding may indicate that both cell populations use protein networks that, in tumors, result in uncontrolled self-renewal and dedifferentiated phenotypes histopathologically defined by loss of architecture specific to the tissue of origin. We speculate that the indices described here may help predict the efficacy of stem-cell-based immunotherapies and contribute to the identification of patients who will respond to such therapies.



We interrogated CMap using the gene expression signatures from tumor samples with the highest and lowest mRNAsi levels. Surprisingly, perhaps, the CMap analysis, which is based on only a limited number of treated cell lines, very precisely selected drugs that have been shown to affect cancer stem cells with specificity. These translational analyses may ultimately pave the way for implementation of differentiation therapies for solid tumors.

Here, we have also shown that cancer hallmarks can be extracted from datasets on cells with defined phenotypes and used to train machine-learning methods applicable to index molecular profiles of cancer. Our mRNAsi and mDNAsi can be translated into stemness scores (e.g., STEM50) that stratify tumors based on their dedifferentiation features, thus providing biomarkers for prediction of patient outcomes and response to differentiation therapies.

By defining new metrics of cancer stemness and using them to interrogate TCGA datasets, our results provide a comprehensive characterization of dedifferentiation as new and significant hallmarks of cancer. The strengths of the approach are that it leverages features of dedifferentiated cells across a spectrum of tumor types that reflect tumor pathology and, in some cases, clinical outcome. This study also provides strategies for integrated analysis of cancer genomics based on machine-learning methods trained on molecular profiles obtained from cells with defined phenotypes. The findings based on those methods may advance the development of objective diagnostics tools for quantitating cancer stemness in clinical tumors, perhaps leading eventually to new biomarkers that predict tumor recurrence, guide treatment selection, or improve responses to therapy.

## 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](#)
  - DNA Methylation Data
  - RNA Expression Data
  - Stemness Index Derived Using OCLR
  - DNA Methylation Stemness Signatures
  - Stemness versus Molecular and Clinical Features
  - Stemness versus Clinical Predictors
  - Compounds Targeting with Cancer Stemness
- [QUANTIFICATION AND STATISTICAL ANALYSIS](#)
- [DATA AND SOFTWARE AVAILABILITY](#)

## SUPPLEMENTAL INFORMATION

Supplemental Information includes seven figures and four tables and can be found with this article online at <https://doi.org/10.1016/j.cell.2018.03.034>.

## ACKNOWLEDGMENTS

We thank Marcin Cieřlik from Michigan Center for Translational Pathology at University of Michigan for providing the MTE500 dataset. This work was supported by the following grants: 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, and P30 CA016672; NCI grants 5R01CA180778, 3U24CA143858, 1U24CA210990, 5U54HG006097, 1U24CA210949, and 1U24CA210950; NIGMS grant 5R01GM109031; the Henry Ford Cancer Institute's Early Career Investigator Award grant A20054 to T.M.M.; Sao Paulo Research Foundation (FAPESP) grants 2014/02245-3 and 2016/01975-3 to T.M.M. and H.N.; FAPESP grants 2014/08321-3, 2015/07925-5, 2016/01389-7, 2016/10436-9, 2016/06488-3, 2016/12329-5, and 2016/15485-8, and Henry Ford Hospital grant A30935 to H.N.; Spanish Institute of Health Carlos III grant CP14/00229; Mary K. Chapman Foundation gift "Chapman Foundation Fund for Bioinformatics," CPRIT grant RP13039, and the Michael & Susan Dell Foundation grant "The Lorraine Dell Program in Bioinformatics" to J.N.W.; and the Polish Science Foundation Welcome grant 2010/3-3 to M.W.

## AUTHOR CONTRIBUTIONS

The Cancer Genome Atlas Research Network contributed collectively to this study. The contributions of other authors are as follows: epigenetic-derived stemness index, T.M.M., P.W.L., and H.N.; mRNA-expression-derived stemness index, A.S.; methodology, T.M.M., A.S., and H.N.; integrative analysis, T.M.M., A.S., and H.N.; clinical analysis, T.B. and L.P.; drug analysis, A.C.; data interpretation, T.M.M., A.S., J.H., B.K., H.H., A.J.G., A.C., L.M., A.J.L., A.K., A.K.G., J.M.S., K.A.H., P.C., O.G., S.M., P.W.L., H.N., and M.W.; data curation, L.O.; writing, T.M.M., A.S., J.W., B.K., J.H., A.C., J.N.W., H.N., and M.W.; visualization, T.M.M., A.S., H.N., and M.W.; overall concept and coordination, T.M.M., A.S., P.W.L., H.N., and M.W.

## 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 Ci-bermed. 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: September 5, 2017

Revised: January 30, 2018

Accepted: March 14, 2018

Published: April 5, 2018

## REFERENCES

- Agarwal, S., Graepel, T., Herbrich, R., Har-Peled, S., and Roth, D. (2005). Generalization Bounds for the Area Under the ROC Curve. *J. Mach. Learn. Res.* 6, 393–425.
- Bai, X.-F., Ni, X.-G., Zhao, P., Liu, S.-M., Wang, H.-X., Guo, B., Zhou, L.-P., Liu, F., Zhang, J.-S., Wang, K., et al. (2004). Overexpression of annexin 1 in pancreatic cancer and its clinical significance. *World J. Gastroenterol.* 10, 1466–1470.
- Bao, B., Wang, Z., Ali, S., Kong, D., Banerjee, S., Ahmad, A., Li, Y., Azmi, A.S., Miele, L., and Sarkar, F.H. (2011). Over-expression of FoxM1 leads to epithelial-mesenchymal transition and cancer stem cell phenotype in pancreatic cancer cells. *J. Cell. Biochem.* 112, 2296–2306.
- Ben-Porath, I., Thomson, M.W., Carey, V.J., Ge, R., Bell, G.W., Regev, A., and Weinberg, R.A. (2008). An embryonic stem cell-like gene expression signature in poorly differentiated aggressive human tumors. *Nat. Genet.* 40, 499–507.
- Benjamini, Y., and Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B* 57, 289–300.
- Bertucci, F., Finetti, P., Mamessier, E., Pantaleo, M.A., Astolfi, A., Ostrowski, J., and Birnbaum, D. (2015). PDL1 expression is an independent prognostic factor in localized GISt. *Oncolimmunology* 4, e1002729.
- Bradner, J.E., Hnisz, D., and Young, R.A. (2017). Transcriptional addiction in cancer. *Cell* 168, 629–643.
- Ceccarelli, M., Barthel, F.P., Malta, T.M., Sabedot, T.S., Salama, S.R., Murray, B.A., Morozova, O., Newton, Y., Radenbaugh, A., Pagnotta, S.M., et al.; TCGA Research Network (2016). Molecular profiling reveals biologically discrete subsets and pathways of progression in diffuse glioma. *Cell* 164, 550–563.
- Chen, D.S., and Mellman, I. (2013). Oncology meets immunology: the cancer-immunity cycle. *Immunity* 39, 1–10.
- Cheng, H.W., Liang, Y.H., Kuo, Y.L., Chuu, C.P., Lin, C.Y., Lee, M.H., Wu, A.T., Yeh, C.T., Chen, E.I., Whang-Peng, J., et al. (2015). Identification of thioridazine, an antipsychotic drug, as an antiglioblastoma and anticancer stem cell agent using public gene expression data. *Cell Death Dis* 6, e1753.
- Chiao, M.T., Cheng, W.Y., Yang, Y.C., Shen, C.C., and Ko, J.L. (2013). Suberoylanilide hydroxamic acid (SAHA) causes tumor growth slowdown and triggers autophagy in glioblastoma stem cells. *Autophagy* 9, 1509–1526.
- Chung, W., Eum, H.H., Lee, H.-O., Lee, K.-M., Lee, H.-B., Kim, K.-T., Ryu, H.S., Kim, S., Lee, J.E., Park, Y.H., et al. (2017). Single-cell RNA-seq enables comprehensive tumour and immune cell profiling in primary breast cancer. *Nat. Commun.* 8, 15081.
- Colaprico, A., Silva, T.C., Olsen, C., Garofano, L., Cava, C., Garolini, D., Sabedot, T.S., Malta, T.M., Pagnotta, S.M., Castiglioni, I., et al. (2016). TCGAAbiolinks: an R/Bioconductor package for integrative analysis of TCGA data. *Nucleic Acids Res.* 44, e71.
- 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*, doi 10.1101/265322.
- Daily, K., Ho Sui, S.J., Schriml, L.M., Dexheimer, P.J., Salomonis, N., Schroll, R., Bush, S., Keddache, M., Mayhew, C., Lotia, S., et al. (2017). Molecular, phenotypic, and sample-associated data to describe pluripotent stem cell lines and derivatives. *Sci. Data* 4, 170030.
- Davis, S., Bilke, S., Triche Jr, T., and Bootwalla, M. (2015). methylumi: Handle Illumina methylation data. R Package Version 2.18.0.
- de Souza, C.F., Sabedot, T.S., Malta, T.M., Stetson, L., Morozova, O., Sokolov, A., Laird, P.W., Wiznerowicz, M., Iavarone, A., Snyder, J., et al. (2018). Distinct DNA Methylation Shift in a Subset of Glioma CpG Island Methylator Phenotype (G-CIMP) during Tumor Recurrence. *Cell Rep.* Published online April 10, 2018. <https://doi.org/10.1016/j.celrep.2018.03.107>.
- Dolma, S., Selvadurai, H.J., Lan, X., Lee, L., Kushida, M., Voisin, V., Whetstone, H., So, M., Aviv, T., Park, N., et al. (2016). Inhibition of dopamine receptor D4 impedes autophagic flux, proliferation, and survival of glioblastoma stem cells. *Cancer Cell* 29, 859–873.
- Economopoulou, P., Perisanidis, C., Giotakis, E.I., and Psyrri, A. (2016). The emerging role of immunotherapy in head and neck squamous cell carcinoma (HNSCC): anti-tumor immunity and clinical applications. *Ann. Transl. Med.* 4, 173.
- Eppert, K., Takenaka, K., Lechman, E.R., Waldron, L., Nilsson, B., van Galen, P., Metzeler, K.H., Poepl, A., Ling, V., Beyene, J., et al. (2011). Stem cell gene expression programs influence clinical outcome in human leukemia. *Nat. Med.* 17, 1086–1093.
- Fabregat, I., Malfettone, A., and Soukupova, J. (2016). New Insights into the Crossroads between EMT and Stemness in the Context of Cancer. *J. Clin. Med.* 5.
- Friedmann-Morvinski, D., and Verma, I.M. (2014). Dedifferentiation and reprogramming: origins of cancer stem cells. *EMBO Rep.* 15, 244–253.
- Fuereder, T. (2016). Immunotherapy for head and neck squamous cell carcinoma. *Memo* 9, 66–69.
- Ge, Y., Gomez, N.C., Adam, R.C., Nikolova, M., Yang, H., Verma, A., Lu, C.P.-J., Polak, L., Yuan, S., Elemento, O., and Fuchs, E. (2017). Stem cell lineage infidelity drives wound repair and cancer. *Cell* 169, 636–650.
- Gentles, A.J., Plevritis, S.K., Majeti, R., and Alizadeh, A.A. (2010). Association of a leukemic stem cell gene expression signature with clinical outcomes in acute myeloid leukemia. *JAMA* 304, 2706–2715.
- Gentles, A.J., Newman, A.M., Liu, C.L., Bratman, S.V., Feng, W., Kim, D., Nair, V.S., Xu, Y., Khuong, A., Hoang, C.D., et al. (2015). The prognostic landscape of genes and infiltrating immune cells across human cancers. *Nat. Med.* 21, 938–945.
- Gevaert, O., Villalobos, V., Sikic, B.I., and Plevritis, S.K. (2013). Identification of ovarian cancer driver genes by using module network integration of multi-omics data. *Interface Focus* 3, 20130013.
- Gingold, J., Zhou, R., Lemischka, I.R., and Lee, D.-F. (2016). Modeling Cancer with Pluripotent Stem Cells. *Trends Cancer* 2, 485–494.
- Gonzalez, D.M., and Medici, D. (2014). Signaling mechanisms of the epithelial-mesenchymal transition. *Sci. Signal.* 7, re8.
- Gregory, P.A., Bert, A.G., Paterson, E.L., Barry, S.C., Tsykin, A., Farshid, G., Vadas, M.A., Khew-Goodall, Y., and Goodall, G.J. (2008). The miR-200 family and miR-205 regulate epithelial to mesenchymal transition by targeting ZEB1 and SIP1. *Nat. Cell Biol.* 10, 593–601.
- Hanahan, D., and Weinberg, R.A. (2011). Hallmarks of cancer: the next generation. *Cell* 144, 646–674.
- Hoadley, K.A., Yau, C., Wolf, D.M., Cherniack, A.D., Tamborero, D., Ng, S., Leiserson, M.D.M., Niu, B., McLellan, M.D., Uzunangelov, V., et al.; Cancer Genome Atlas Research Network (2014). Multiplatform analysis of 12 cancer types reveals molecular classification within and across tissues of origin. *Cell* 158, 929–944.
- Ivanova, N., Dobrin, R., Lu, R., Kotenko, I., Levorse, J., DeCoste, C., Schafer, X., Lun, Y., and Lemischka, I.R. (2006). Dissecting self-renewal in stem cells with RNA interference. *Nature* 442, 533–538.
- Kang, W.-Y., Chen, W.-T., Huang, Y.-C., Su, Y.-C., and Chai, C.-Y. (2012). Overexpression of annexin 1 in the development and differentiation of urothelial carcinoma. *Kaohsiung J. Med. Sci.* 28, 145–150.

- Kim, J., and Orkin, S.H. (2011). Embryonic stem cell-specific signatures in cancer: insights into genomic regulatory networks and implications for medicine. *Genome Med.* 3, 75.
- Kim, J., Woo, A.J., Chu, J., Snow, J.W., Fujiwara, Y., Kim, C.G., Cantor, A.B., and Orkin, S.H. (2010). A Myc network accounts for similarities between embryonic stem and cancer cell transcription programs. *Cell* 143, 313–324.
- Kooreman, N.G., Kim, Y., de Almeida, P.E., Termglinchan, V., Diecke, S., Shao, N.-Y., Wei, T.-T., Yi, H., Dey, D., Nelakanti, R., et al. (2018). Autologous iPSC-Based Vaccines Elicit Anti-tumor Responses In Vivo. *Cell Stem Cell*. Published online February 8, 2018. <https://doi.org/10.1016/j.stem.2018.01.016>.
- Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M., Yen, A., Heravi-Moussavi, A., Kheradpour, P., Zhang, Z., Wang, J., Ziller, M.J., et al.; Roadmap Epigenomics Consortium (2015). Integrative analysis of 111 reference human epigenomes. *Nature* 518, 317–330.
- Lamb, J., Crawford, E.D., Peck, D., Modell, J.W., Blat, I.C., Wrobel, M.J., Lerner, J., Brunet, J.-P., Subramanian, A., Ross, K.N., et al. (2006). The Connectivity Map: using gene-expression signatures to connect small molecules, genes, and disease. *Science* 313, 1929–1935.
- Lu, C., Ward, P.S., Kapoor, G.S., Rohle, D., Turcan, S., Abdel-Wahab, O., Edwards, C.R., Khanin, R., Figueroa, M.E., Melnick, A., et al. (2012). IDH mutation impairs histone demethylation and results in a block to cell differentiation. *Nature* 483, 474–478.
- Lu, M., Li, J., Luo, Z., Zhang, S., Xue, S., Wang, K., Shi, Y., Zhang, C., Chen, H., and Li, Z. (2015). Roles of dopamine receptors and their antagonist thioridazine in hepatoma metastasis. *Oncol. Targets Ther* 8, 1543–1552.
- Lyssiotis, C.A., and Kimmelman, A.C. (2017). Metabolic interactions in the tumor microenvironment. *Trends Cell Biol.* 27, 863–875.
- Mak, M.P., Tong, P., Diao, L., Cardnell, R.J., Gibbons, D.L., William, W.N., Skoulidis, F., Parra, E.R., Rodriguez-Canales, J., Wistuba, I.I., et al. (2016). A Patient-Derived, Pan-Cancer EMT Signature Identifies Global Molecular Alterations and Immune Target Enrichment Following Epithelial-to-Mesenchymal Transition. *Clin. Cancer Res.* 22, 609–620.
- Mathur, D., Danford, T.W., Boyer, L.A., Young, R.A., Gifford, D.K., and Jaenisch, R. (2008). Analysis of the mouse embryonic stem cell regulatory networks obtained by ChIP-chip and ChIP-PET. *Genome Biol.* 9, R126.
- Nazor, K.L., Altun, G., Lynch, C., Tran, H., Harness, J.V., Slavin, I., Garitaonandia, I., Müller, F.-J., Wang, Y.-C., Boscolo, F.S., et al. (2012). Recurrent variations in DNA methylation in human pluripotent stem cells and their differentiated derivatives. *Cell Stem Cell* 10, 620–634.
- Ng, S.W.K., Mitchell, A., Kennedy, J.A., Chen, W.C., McLeod, J., Ibrahimova, N., Arruda, A., Popescu, A., Gupta, V., Schimmer, A.D., et al. (2016). A 17-gene stemness score for rapid determination of risk in acute leukaemia. *Nature* 540, 433–437.
- Noushmehr, H., Weisenberger, D.J., Diefes, K., Phillips, H.S., Pujara, K., Berman, B.P., Pan, F., Pelloski, C.E., Sulman, E.P., Bhat, K.P., et al.; Cancer Genome Atlas Research Network (2010). Identification of a CpG island methylator phenotype that defines a distinct subgroup of glioma. *Cancer Cell* 17, 510–522.
- Onuchic, V., Hartmaier, R.J., Boone, D.N., Samuels, M.L., Patel, R.Y., White, W.M., Garovic, V.D., Oesterreich, S., Roth, M.E., Lee, A.V., and Milosavljevic, A. (2016). Epigenomic Deconvolution of Breast Tumors Reveals Metabolic Coupling between Constituent Cell Types. *Cell Rep.* 17, 2075–2086.
- Palmer, N.P., Schmid, P.R., Berger, B., and Kohane, I.S. (2012). A gene expression profile of stem cell pluripotentiality and differentiation is conserved across diverse solid and hematopoietic cancers. *Genome Biol.* 13, R71.
- Papillon-Cavanagh, S., Lu, C., Gayden, T., Mikael, L.G., Bechet, D., Karamboulas, C., Ailles, L., Karamchandani, J., Marchione, D.M., Garcia, B.A., et al. (2017). Impaired H3K36 methylation defines a subset of head and neck squamous cell carcinomas. *Nat. Genet.* 49, 180–185.
- Pinto, J.P., Kalathur, R.K., Oliveira, D.V., Barata, T., Machado, R.S.R., Machado, S., Pacheco-Leyva, I., Duarte, I., and Futschik, M.E. (2015). StemChecker: a web-based tool to discover and explore stemness signatures in gene sets. *Nucleic Acids Res.* 43 (W1), W72–W77.
- Polónia, A., Pinto, R., Cameselle-Teijeiro, J.F., Schmitt, F.C., and Paredes, J. (2017). Prognostic value of stromal tumour infiltrating lymphocytes and programmed cell death-ligand 1 expression in breast cancer. *J. Clin. Pathol.* 70, 860–867.
- R Development Core Team (2017). R: A Language and Environment for Statistical Computing (R Foundation for Statistical Computing).
- Reyngold, M., Turcan, S., Giri, D., Kannan, K., Walsh, L.A., Viale, A., Drobnjak, M., Vahdat, L.T., Lee, W., and Chan, T.A. (2014). Remodeling of the methylation landscape in breast cancer metastasis. *PLoS ONE* 9, e103896.
- Robinson, D.R., Wu, Y.-M., Lonigro, R.J., Vats, P., Cobain, E., Everett, J., Cao, X., Rabban, E., Kumar-Sinha, C., Raymond, V., et al. (2017). Integrative clinical genomics of metastatic cancer. *Nature* 548, 297–303.
- Royston, P., and Altman, D.G. (1994). Regression Using Fractional Polynomials of Continuous Covariates: Parsimonious Parametric Modelling. *J. R. Stat. Soc. Ser. C. Appl. Stat* 43, 429.
- Salomonis, N., Dexheimer, P.J., Omberg, L., Schroll, R., Bush, S., Huo, J., Schriml, L., Ho Sui, S., Keddache, M., Mayhew, C., et al. (2016). Integrated Genomic Analysis of Diverse Induced Pluripotent Stem Cells from the Progenitor Cell Biology Consortium. *Stem Cell Reports* 7, 110–125.
- Sato, N., Sanjuan, I.M., Heke, M., Uchida, M., Naef, F., and Brivanlou, A.H. (2003). Molecular signature of human embryonic stem cells and its comparison with the mouse. *Dev. Biol.* 260, 404–413.
- Sergushichev, A. (2016). An algorithm for fast preranked gene set enrichment analysis using cumulative statistic calculation. *bioRxiv*, doi 10.1101/060012.
- Shibue, T., and Weinberg, R.A. (2017). EMT, CSCs, and drug resistance: the mechanistic link and clinical implications. *Nat. Rev. Clin. Oncol.* 14, 611–629.
- Silva, T.C., Colaprico, A., Olsen, C., Bontempi, G., Ceccarelli, M., Berman, B.P., and Noushmehr, H. (2017). TCGAAbiLinksGUI: A graphical user interface to analyze cancer molecular and clinical data. *bioRxiv*, doi 10.1101/147496.
- Sokolov, A., Paull, E.O., and Stuart, J.M. (2016). One-class detection of cell states in tumor subtypes. *Pac. Symp. Biocomput.* 27, 405–416.
- Spitzer, M.H., Carmi, Y., Reticker-Flynn, N.E., Kwek, S.S., Madhiredy, D., Martins, M.M., Gherardini, P.F., Prestwood, T.R., Chabon, J., Bendall, S.C., et al. (2017). Systemic immunity is required for effective cancer immunotherapy. *Cell* 168, 487–502.
- StataCorp (2013). Stata Statistical Software: Release 13 (College Station, TX: StataCorp LP).
- Sturm, D., Witt, H., Hovestadt, V., Khuong-Quang, D.-A., Jones, D.T.W., Kovermann, C., Pfaff, E., Tönjes, M., Sill, M., Bender, S., et al. (2012). Hotspot mutations in H3F3A and IDH1 define distinct epigenetic and biological subgroups of glioblastoma. *Cancer Cell* 22, 425–437.
- Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S., and Mesirov, J.P. (2005). Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA* 102, 15545–15550.
- Subramanian, A., Narayan, R., Corsello, S.M., Peck, D.D., Natoli, T.E., Lu, X., Gould, J., Davis, J.F., Tubelli, A.A., Asiedu, J.K., et al. (2017). A next generation connectivity map: L1000 platform and the first 1,000,000 profiles. *Cell* 171, 1437–1452.
- Tellez, C.S., Juri, D.E., Do, K., Bernauer, A.M., Thomas, C.L., Damiani, L.A., Tessema, M., Leng, S., and Belinsky, S.A. (2011). EMT and stem cell-like properties associated with miR-205 and miR-200 epigenetic silencing are early manifestations during carcinogen-induced transformation of human lung epithelial cells. *Cancer Res.* 71, 3087–3097.
- Therneau, T.M., and Grambsch, P.M. (2000). Estimating the survival and hazard functions. In *Modeling Survival Data: Extending the Cox Model* (New York, NY: Springer New York), pp. 7–37.
- Tirosh, I., Venteicher, A.S., Hebert, C., Escalante, L.E., Patel, A.P., Yizhak, K., Fisher, J.M., Rodman, C., Mount, C., Filbin, M.G., et al. (2016). Single-cell

- RNA-seq supports a developmental hierarchy in human oligodendrogloma. *Nature* 539, 309–313.
- Tomczak, K., Czerwińska, P., and Wiznerowicz, M. (2015). The Cancer Genome Atlas (TCGA): an immeasurable source of knowledge. *Contemp. Oncol. (Pozn.)* 19 (1A), A68–A77.
- Tsai, S.-C., Lin, C.-C., Shih, T.-C., Tseng, R.-J., Yu, M.-C., Lin, Y.-J., and Hsieh, S.-Y. (2017). The miR-200b-ZEB1 circuit regulates diverse stemness of human hepatocellular carcinoma. *Mol. Carcinog.* 56, 2035–2047.
- Turcan, S., Rohle, D., Goenka, A., Walsh, L.A., Fang, F., Yilmaz, E., Campos, C., Fabius, A.W.M., Lu, C., Ward, P.S., et al. (2012). IDH1 mutation is sufficient to establish the glioma hypermethylator phenotype. *Nature* 483, 479–483.
- Venezia, T.A., Merchant, A.A., Ramos, C.A., Whitehouse, N.L., Young, A.S., Shaw, C.A., and Goodell, M.A. (2004). Molecular signatures of proliferation and quiescence in hematopoietic stem cells. *PLoS Biol.* 2, e301.
- Visvader, J.E., and Lindeman, G.J. (2012). Cancer stem cells: current status and evolving complexities. *Cell Stem Cell* 10, 717–728.
- Volate, S.R., Kawasaki, B.T., Hurt, E.M., Milner, J.A., Kim, Y.S., White, J., and Farrar, W.L. (2010). Gossypol induces apoptosis by activating p53 in prostate cancer cells and prostate tumor-initiating cells. *Mol. Cancer Ther* 9, 461–470.
- Wang, M.-D., Wu, H., Fu, G.-B., Zhang, H.-L., Zhou, X., Tang, L., Dong, L.-W., Qin, C.-J., Huang, S., Zhao, L.-H., et al. (2016). Acetyl-coenzyme A carboxylase alpha promotion of glucose-mediated fatty acid synthesis enhances survival of hepatocellular carcinoma in mice and patients. *Hepatology* 63, 1272–1286.
- Wen, N., Wang, Y., Wen, L., Zhao, S.-H., Ai, Z.-H., Wang, Y., Wu, B., Lu, H.-X., Yang, H., Liu, W.-C., and Li, Y. (2014). Overexpression of FOXM1 predicts poor prognosis and promotes cancer cell proliferation, migration and invasion in epithelial ovarian cancer. *J. Transl. Med.* 12, 134.
- Wickham, H. (2009). *ggplot2* (Springer Nature).
- Wong, C.-M., Wei, L., Au, S.L.-K., Fan, D.N.-Y., Zhou, Y., Tsang, F.H.-C., Law, C.-T., Lee, J.M.-F., He, X., Shi, J., et al. (2015). MiR-200b/200c/429 subfamily negatively regulates Rho/ROCK signaling pathway to suppress hepatocellular carcinoma metastasis. *Oncotarget* 6, 13658–13670.
- Xu, L., Zhang, L., Hu, C., Liang, S., Fei, X., Yan, N., Zhang, Y., and Zhang, F. (2016). WNT pathway inhibitor pyrvinium pamoate inhibits the self-renewal and metastasis of breast cancer stem cells. *Int. J. Oncol.* 48, 1175–1186.
- Xue, Y.J., Xiao, R.H., Long, D.Z., Wang, X.N., Zhang, G.N., Yuan, Y.H., Wu, G.Q., Yang, J., Wu, Y.T., Xu, H., et al. (2012). Overexpression of FoxM1 is associated with tumor progression in patients with clear cell renal cell carcinoma. *J. Transl. Med* 10, 200.
- Yan, X., Ma, L., Yi, D., Yoon, J.G., Diercks, A., Foltz, G., Price, N.D., Hood, L.E., and Tian, Q. (2011). A CD133-related gene expression signature identifies an aggressive glioblastoma subtype with excessive mutations. *Proc. Natl. Acad. Sci. USA* 108, 1591–1596.
- Yao, L., Shen, H., Laird, P.W., Farnham, P.J., and Berman, B.P. (2015). Inferring regulatory element landscapes and transcription factor networks from cancer methylomes. *Genome Biol.* 16, 105.
- Young, R.A. (2011). Control of the embryonic stem cell state. *Cell* 144, 940–954.
- Yue, H., Huang, D., Qin, L., Zheng, Z., Hua, L., Wang, G., Huang, J., and Huang, H. (2016). Targeting Lung Cancer Stem Cells with Antipsychological Drug Thioridazine. *BioMed Res. Int.* 2016, 6709828.
- Zaretsky, J.M., Garcia-Diaz, A., Shin, D.S., Escuin-Ordinas, H., Hugo, W., Hu-Lieskovan, S., Torrejon, D.Y., Abril-Rodriguez, G., Sandoval, S., Barthly, L., et al. (2016). Mutations Associated with Acquired Resistance to PD-1 Blockade in Melanoma. *N. Engl. J. Med.* 375, 819–829.
- Zheng, S., Cherniack, A.D., Dewal, N., Moffitt, R.A., Danilova, L., Murray, B.A., Lerario, A.M., Else, T., Knijnenburg, T.A., Ciriello, G., et al.; Cancer Genome Atlas Research Network (2016). Comprehensive Pan-Genomic Characterization of Adrenocortical Carcinoma. *Cancer Cell* 29, 723–736.

## STAR★METHODS

### KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Software and Algorithms		
Workflow to reproduce the stemness index	This paper	<a href="https://bioinformaticsfmrp.github.io/PanCanStem_Web/">https://bioinformaticsfmrp.github.io/PanCanStem_Web/</a>
CIBERSORT	Gentles et al., 2015	<a href="https://precog.stanford.edu/">https://precog.stanford.edu/</a>
R 3.3.1	R Development Core Team, 2017	<a href="https://www.R-project.org">https://www.R-project.org</a>
ggplot2 (v2.2.1)	Wickham, 2009	<a href="https://cran.r-project.org/web/packages/ggplot2/index.html">https://cran.r-project.org/web/packages/ggplot2/index.html</a>
gelnet (v1.2.1)	Sokolov et al., 2016	<a href="https://cran.r-project.org/web/packages/gelnet/index.html">https://cran.r-project.org/web/packages/gelnet/index.html</a>
GSEA	Subramanian et al., 2005	<a href="https://software.broadinstitute.org/gsea/index.jsp">https://software.broadinstitute.org/gsea/index.jsp</a>
TCGAbiolinks (v2.4.3)	Colaprico et al., 2016	<a href="http://bioconductor.org/packages/release/bioc/html/TCGAbiolinks.html">http://bioconductor.org/packages/release/bioc/html/TCGAbiolinks.html</a>
ELMER (v1.4.1)	Yao et al., 2015	<a href="http://bioconductor.org/packages/ELMER/">http://bioconductor.org/packages/ELMER/</a>
fgsea (v1.2.1)	Sergushichev, 2016	<a href="http://bioconductor.org/packages/release/bioc/html/fgsea.html">http://bioconductor.org/packages/release/bioc/html/fgsea.html</a>
Methylumi (v2.20.0)	Davis et al., 2015	<a href="http://bioconductor.org/packages/release/bioc/html/methylumi.html">http://bioconductor.org/packages/release/bioc/html/methylumi.html</a>
MoonlightR (v1.2.0)	Colaprico et al., 2018	<a href="http://bioconductor.org/packages/release/bioc/html/MoonlightR.html">http://bioconductor.org/packages/release/bioc/html/MoonlightR.html</a>
Amaretto	Gevaert et al., 2013	<a href="http://med.stanford.edu/gevaertlab/software.html">http://med.stanford.edu/gevaertlab/software.html</a>
STATA (v13)	StataCorp, 2013	<a href="https://www.stata.com/">https://www.stata.com/</a>
Deposited Data		
TCGA data	NIH Genomic Data Commons (GDC)	<a href="https://gdc.cancer.gov/about-data/publications/PanCanStemness-2018">https://gdc.cancer.gov/about-data/publications/PanCanStemness-2018</a>
Progenitor Cell Biology Consortium (PCBC)	Daily et al., 2017; Salomonis et al., 2016	<a href="https://www.synapse.org/#!Synapse:syn1773109/wiki/54962">https://www.synapse.org/#!Synapse:syn1773109/wiki/54962</a>
Chromatin State (ChromHMM)	Kundaje et al., 2015	<a href="http://www.roadmapepigenomics.org">http://www.roadmapepigenomics.org</a>
Stem Cells Validation set	Nazor et al., 2012	mRNA expression (GEO: GSE30652) and DNA methylation (GEO: GSE30654)
Glioma validation set	Sturm et al., 2012	mRNA expression (GEO: GSE36245) and DNA methylation (GEO: GSE36278)
Glioma validation set	Turcan et al., 2012	GEO: GSE30339
BRCA validation set	Reyngold et al., 2014	GEO: GSE59000
Deconvolution of breast cancer (BRCA)	Onuchic et al., 2016	<a href="http://genboree.org/theCommons/projects/edec">http://genboree.org/theCommons/projects/edec</a>
MET500 - Metastatic solid tumors	Robinson et al., 2017	Database of Genotypes and Phenotypes (dbGaP) accession number phs000673.v2.p1
Gliomas Single Cell RNA expression	Tirosh et al., 2016	GEO: GSE70630
BRCA Single Cell RNA expression	Chung et al., 2017	GEO: GSE75688

### CONTACT FOR REAGENT AND RESOURCE SHARING

Further information and requests for resources should be directed to and will be fulfilled by the Lead Contact, Maciej Wiznerowicz ([maciej.wiznerowicz@iimo.pl](mailto:maciej.wiznerowicz@iimo.pl)).

### EXPERIMENTAL MODEL AND SUBJECT DETAILS

Clinical and molecular data were collected from the NIH Genomic Data Commons (GDC) of 11,392 participants from The Cancer Genome Atlas PanCancer Atlas cohort (<https://gdc.cancer.gov/about-data/publications/PanCanStemness-2018>).



## METHOD DETAILS

### DNA Methylation Data

A total of 9,627 PanCancer TCGA samples across 33 different tumor types were available for DNA methylation using the robust Illumina HumanMethylation 450 (HM450) platform. TCGA samples included primary (8,471), recurrent (41), and metastatic tumor (394) tissues and a set of 721 non-tumor tissues.

Level 3 data were downloaded from TCGA Data Portal using TCGAbiolinks functions GDCquery, GDCdownload and GDCprepare importing into R (<http://www.r-project.org>) for further analysis (Colaprico et al., 2016).

DNA methylation level 3 data are  $\beta$ -values that were calculated from pre-processed raw data using the methylumi Bioconductor package (Davis et al., 2015). Pre-processing steps included background correction, dye-bias normalization, and calculation of  $\beta$ -values and detection p values.  $\beta$ -values range from zero to one, with zero indicating no DNA methylation and one indicating complete DNA methylation. A detection p value compares the signal intensity difference between the analytical probes and a set of negative control probes on the array. Any data point with a corresponding p value greater than 0.01 is deemed not statistically significantly different from background and is thus masked as “NA” in TCGA level 3 data. The data levels and the files contained in each data level package are on the NIH Genomic Data Commons (GDC).

In addition to TCGA data, we used a dataset of 99 human stem/progenitor cells from the Progenitor Cell Biology Consortium (PCBC) (<https://www.synapse.org/#!Synapse:syn1773109>) to define stem cell signatures (Daily et al., 2017; Salomonis et al., 2016). PCBC samples were profiled using the Illumina HumanMethylation 450 (HM450) platform and consisted of 4 embryonic stem cells (ESC), 40 induced pluripotent stem cells (iPSC), 22 stem cell (SC)-derived embryoid bodies (EB), 11 SC-derived mesoderm day 5 (MESO, mesothelioma), 11 SC-derived ectoderm (ECTO), and 11 SC-derived definitive endoderm (DE). We downloaded raw IDAT files from PCBC Genomic Data Commons and processed the data according to the TCGA standard level 3 protocol described above.

### RNA Expression Data

PanCancer TCGA RNA sequence level 3 normalized data were downloaded from the GDC Data Portal using TCGAbiolinks functions GDCquery, GDCdownload and GDCprepare importing into R (<http://www.r-project.org>) for further analysis (Colaprico et al., 2016). A total of 10,852 samples across 33 tumor types were available, including primary (9,702), recurrent (45) and metastatic tumor (395) tissues and a set of 710 non-tumor tissues.

We also downloaded PCBC RNA sequence data from the PCBC Synapse Portal (<https://www.synapse.org/#!Synapse:syn1773109>), consisting 16 ESC, 77 iPSC, 66 SC-derived EB, 29 SC-derived MESO, 29 SC-derived ECTO, and 36 SC-derived DE (Daily et al., 2017; Salomonis et al., 2016).

### Stemness Index Derived Using OCLR

To calculate a stemness index (si) based on mRNA expression or DNA methylation, we built a predictive model using one-class logistic regression (OCLR) (Sokolov et al., 2016) on the pluripotent stem cell samples (ESC and iPSC) from the PCBC dataset (Daily et al., 2017; Salomonis et al., 2016).

For mRNA expression-based signatures, to ensure compatibility with the TCGA PanCancer cohort, we first mapped the gene names from Ensembl IDs to Human Genome Organization (HUGO), dropping any genes that had no such mapping. The resulting training matrix contained 12,945 mRNA expression values measured across all available PCBC samples. For DNA methylation-based signatures, we used each of the signatures (probe set) described below.

We mean-centered the data, then applied OCLR to just the samples labeled SC (which included both ESC and iPSC). We chose to use the one-class framework because of its robustness in the absence of the a “negative” class. The PCBC data does not have data for fully differentiated cells, and progenitor cell types might exhibit some of the stemness signals.

Once the signature is obtained, it can be applied to score new samples. For RNA expression data, we computed Spearman correlations between the model’s weight vector and the new sample’s expression profile. We advocate for the use of Spearman correlation over the more traditional dot product operation because it is more robust with respect to potential cross-dataset batch effects that may arise. For DNA methylation data, which follow the beta distribution, the samples were scored using the standard application of a linear model:  $f(x) = w^T x + b$ .

We validated our approach using leave-one-out cross-validation by withholding each SC sample in turn. A separate signature was then trained on all other SC samples and used to score the withheld sample as well as all the non-SC samples. The performance was measured using the area under the curve (AUC) metric, which can be interpreted as the probability that the model correctly ranks a positive sample above a negative (Agarwal et al., 2005). In our cross-validation experiment, every withheld SC sample was scored higher than all the non-SC samples, yielding an overall AUC of 1.0.

We performed additional validation of the stemness signature by applying it to an external dataset composed of pluripotent stem cells (ESC and iPSC), somatic cells (17 distinct tissue types and several primary cell lines of diverse origin), and hydatidiform mole samples (Nazor et al., 2012). The mRNA expression data for the study were downloaded from GEO (GEO: GSE30652) as were DNA methylation data (GEO: GSE30654). We observed that all of the SC samples were correctly scored above all of the somatic samples

by both platforms (Figure 1B). This is particularly striking for mRNA expression, because mRNA expression in study by the Nazor et al. was measured using microarrays, whereas the signature was trained using RNA-seq data.

Having validated the signature by using cross-validation and external SC data, we then applied it to score the TCGA PanCancer cohort using the same Spearman correlation (RNA expression) or linear model (DNA methylation) operators. The indices were subsequently mapped to the [0, 1] range by using a linear transformation that subtracted the minimum and divided by the maximum. The mapping was done to assist with interpretation as well as integration with the stemness indices derived from other data platforms (i.e., DNA methylation and mRNA expression).

Additionally, we downloaded independent, non-TCGA datasets of gliomas [(Sturm et al., 2012) (GEO: GSE36245, GEO: GSE36278) and (Turcan et al., 2012) (GEO: GSE30339)] and BRCA samples (Reyngold et al., 2014) (GSE59000) and applied our metrics to measure the stemness in the validation data. For mRNA expression, the preprocessing consisted of mapping the Illumina probe IDs (Illumina HumanHT-12 V3.0 platform) to HUGO symbols, and then reducing the signature and the external dataset to a common set of genes. We then computed the Spearman correlation between the signature and the external samples. For DNA methylation, we applied the linear model.

### DNA Methylation Stemness Signatures

Due to the magnitude of the available DNA methylation platform Infinium HumanMethylation450 (HM450), we defined DNA methylation-based stemness signatures as a reduced input to the OCLR machine learning algorithm. For the DNA methylation-based stemness indices, three signatures were utilized, each defining a distinct, biologically relevant, molecular phenotype of stemness. First, we performed a supervised analysis between human pluripotent stem cells (ESC and iPSC) and stem cell-derived progenitors (embryoid bodies [EB], mesoderm [MESO], ectoderm [ECTO], and definitive endoderm [DE]) ( $\beta$  value mean difference  $< -0.4$  and false discovery rate [FDR]  $< 10e-22$ ;  $\beta$  value mean difference  $> 0.3$  and false discovery rate [FDR]  $< 10e-17$ ). All 'rs' and 'ch' probes were removed prior to analyses. To eliminate somatic tissue-specific probes, we removed probes that were consistently methylated (standard deviation  $\beta$  value  $> 0.05$ ) in non-tumor adult tissues available through TCGA. This resulted in a set of 62 pluripotent cell-specific and differentially methylated regions, which was then used as input for the OCLR to determine the stemness index for each TCGA tumor sample, named "differentially methylated probes-based stemness index" (DMPsi). Interestingly, most of these probes (85%) were positioned within intergenic regions known as open seas (Figure S2A).

Second, we defined a stem cell signature associated with genomic enhancer elements. Enhancers have been shown to be a critically relevant functional element for defining gene target expression and chromatin organization. For this, we downloaded Chromatin State data (ChromHMM) from the NIH Roadmap Epigenomics Consortium (<http://www.roadmapepigenomics.org>), which defined 18 chromatin states (based on 6 different histone marks: H3K4me3, H3K27ac, H3K4me1, H3K36me3, H3K9me3, and H3K27me3) across 98 different cell types (Kundaje et al., 2015). Briefly, by using ChromHMM data we mapped the HM450 probes to the chromatin states in each individual cell type; then we identified genomic regions corresponding to active enhancers that are specific to pluripotent stem cell states (ESC and iPSC), meaning that each region was defined as active enhancers (according to their states: 9-EnhA1 and 10-EnhA2 (Kundaje et al., 2015)) in all pluripotent stem cells ( $n = 9$ ) whereas not enhancer (enhancer in less than 25% of non-pluripotent stem cells ( $n = 89$ )) in non-pluripotent stem cells. We identified 82 DNA methylation probes of the HM450 platform that mapped to enhancer elements and considered them to be a DNA methylation-based pluripotent stem cell enhancer signature, which was then used as input for the OCLR to evaluate stemness signatures for TCGA samples, named "enhancer-based stemness index" (ENHsi) (Figure S2A).

Third, we applied ELMER (Enhancer Linking by Methylation/Expression Relationships), an R/Bioconductor package (Yao et al., 2015) that uses DNA methylation to identify enhancer elements and correlates enhancer state with expression of nearby genes to identify putative transcriptional targets. Using ELMER, we compared pluripotent stem cells (ESC and iPSC) to stem cell-derived progenitors (EB, MESO, ECTO, DE) from PCBC and identified 87 CpGs that were hypomethylated in the pluripotent state (ESC and iPSC) compared to stem cell-derived progenitors and that potentially regulate 103 genes. We confirmed the importance of these probe-gene pair targets by identifying that the SOX2-OCT4 transcription factor binding motif was among the most highly enriched signatures within these elements ( $\pm 250$  bp from the center). The SOX2-OCT4 complex is an important master regulator of pluripotency and stemness. We then derived a new set of signatures using the OCLR and defined TCGA samples' stemness as "epigenetically regulated stemness indices" for each molecular feature (RNA expression-based Epigenetically regulated-mRNAsi [EREG-mRNAsi] and DNA methylation-based [EREG-mDNAsi]).

Because there was high concordance among the three DNA methylation-based indices (DMPsi, ENHsi, and EREG-mDNAsi) (not shown) and each contributes important and complementary biological relevance to stemness, we combined the three stemness signatures (total of 219 probes) and derived a comprehensive DNA methylation index, named mDNAsi (Figure S2A). The lists of probes and genes used to derive the stemness indices are provided on the publication portal accompanying this publication (<https://gdc.cancer.gov/about-data/publications/PanCanStemness-2018>).

### Stemness versus Molecular and Clinical Features

To evaluate the performance of our stemness indices across the entire TCGA cohort, we performed an enrichment analysis by sorting TCGA samples by stemness index for each tumor type and looked for associations with all available genomic features (by using comprehensive mutation data [MC3]), molecular features (previously published TCGA molecular subtypes available at TCGAbiolinks

package (<http://bioconductor.org/packages/release/bioc/html/TCGAbiolinks.html>) (Colaprico et al., 2016; Silva et al., 2017), through the function “PanCancerAtlas\_subtypes(),” which provides full access to the curated matrix used for this study), and clinical features (more than 10,000 features). We used the fgsea R/Bioconductor package to compute the enrichment scores (Sergushichev, 2016). Briefly, for each tumor type we ranked the TCGA samples according to their stemness index (from -low to -high stemness index) and tested if any particular genomic/molecular/clinical feature was associated with either -low or -high stemness index in a non-random behavior. We performed 10,000 permutations for each parameter analyzed to calculate our enrichment score. We then normalized the enrichment scores to mean enrichment of random samples of the same size (NES - normalized enrichment score). Tables containing all the results can be accessed at <https://gdc.cancer.gov/about-data/publications/PanCanStemness-2018>. In addition, an interactive portal with the results across all tumor samples/types versus mDNAsi and mRNAsi can be accessed at [https://bioinformaticsfmrp.github.io/PanCanStem\\_Web/](https://bioinformaticsfmrp.github.io/PanCanStem_Web/) where the user can search for any gene or molecular/clinical feature of interest.

### Stemness versus Clinical Predictors

The associations between the three stemness indices and overall survival (OS) and progression free survival (PFS) in different tumors were evaluated in two stages. First, the proportional hazard (PH) model with the index as a single continuous covariate was used to test whether there was a statistically significant effect on OS or PFS. Given that, for each outcome, the effects of the three indices were tested for 33 cancer types. The significance level of the tests was adjusted for multiple testing to control the overall type I error probability at 5%. In the next stage, the cancer types for which at least one index showed a statistically significant association with either OS or PFS were analyzed in more detail by using a multivariable PH model that included relevant clinical factors. Moreover, the model included a functional form of the index obtained by using degree-2 fractional polynomials (Royston and Altman, 1994). The plausibility of the PH assumption was checked by using the test based on the scaled Schoenfeld residuals (Therneau and Grambsch, 2000). The analyses were conducted using STATA v13 software.

To select the clinical factors for inclusion in the PH model used in the second stage of the OS/PFS analysis for selected cancer types, a detailed analysis of the association between the stemness indices and demographic and clinical features (such as sex, age, race, stage, grade, etc.) was carried out by using linear models. mRNAsi and EREG-mRNAsi were analyzed on the original scale, while mDNAsi was transformed logarithmically to make its distribution more symmetric. The fit of the constructed models was assessed by using residual plots. The analyses were conducted using STATA v13 software.

The screening of the association between the stemness indices and OS (Figure 4E) by using univariable proportional hazard (PH) models indicated a statistically significant (using p values adjusted for multiple testing) effect of mRNAsi on OS for LGG ( $p < 0.0001$ ) and STAD ( $p = 0.005$ ) and on PFS for GBM ( $p = 0.04$ ), LGG ( $p < 0.0001$ ), LIHC ( $p = 0.05$ ), STAD ( $p = 0.04$ ), and UCEC ( $p = 0.03$ ). For mDNAsi, an effect on OS was found for LGG ( $p < 0.0001$ ) and on PFS for kidney renal papillary cell carcinoma (KIRP) ( $p = 0.04$ ) and LGG ( $p < 0.0001$ ). Finally, for EREG-mRNAsi, a statistically significant effect on OS was found for ACC ( $p = 0.005$ ), KIRC ( $p = 0.008$ ), and LGG ( $p = 0.03$ ), and on PFS for ACC ( $p = 0.03$ ), LGG ( $p = 0.03$ ), and UCEC ( $p = 0.04$ ). In these selected cases, multivariable analyses were conducted (using STATA v13 software), which took into account the effect of clinical factors. The analyses confirmed (by using unadjusted p values) the effect of mRNAsi on OS for STAD ( $p = 0.0001$ ) and for GBM/LGG ( $p = 0.002$ ) and the effect on PFS for GBM/LGG ( $p = 0.008$ ) and LIHC ( $p = 0.002$ ). For mDNAsi, the effect on PFS in KIRP was confirmed ( $p = 0.0001$ ), while for EREG-mRNAsi, the effect on PFS in UCEC was confirmed ( $p = 0.05$ ). These confirmed results indicate that the indices have a potential role as novel, independent prognostic factors for the indicated tumor types.

### Compounds Targeting with Cancer Stemness

To determine which target drugs might be useful against cancer stem cells, we used the Broad Institute’s Connectivity Map build 02 (CM) (Lamb et al., 2006), a public online tool (<https://portals.broadinstitute.org/cmap/>) (with registration) that allows users to predict compounds that can activate or inhibit based on a gene expression signature.

To further investigate about mechanism of actions (MoA) and drug-target we performed specific analysis within Connectivity Map tools (<https://clue.io/>) (Subramanian et al., 2017).

Using Connectivity Map (Query) in May 2017 having data available from a collection of cell lines (MCF, PC3, HL60 and SKMEL5) and 164 compounds as small molecules perturbagens. We obtained 33 mRNA expression signatures (one for each cancer type) by applying a differential expression analysis to samples with high mRNAsi and low mRNAsi, using the function TCGAanalyze\_DEA from the R/Bioconductor package TCGAbiolinks version 2.5.9 (Colaprico et al., 2016), carrying edgeR pipeline. The table with differentially expressed genes is reported as Table S3. Due to a limitation of the Connectivity Map tool that matches gene symbol and HG-U133A probe set (eg 200800\_s\_at) GPL96 platform ID, we had to remove duplicate IDs after sorting by decreasing  $|\logFC|$ . We selected the top 1000 genes (500 upregulated and 500 downregulated) where the number of differentially expressed genes was enough or considering the aggregation of upregulated or downregulated genes.

Connectivity MAP is a method similar to GSEA analysis and follows a 4 step approach: (i) looking for similarity between a query signature (diff.expr. genes) and expression profiles present in the dataset using pattern-matching strategy based on Kolmogorov-Smirnov test (ii) rank-ordering the list of genes according their diff.expr. relative to the control from the above expression profiles with significantly similarity (iii) comparison of each rank-ordered list with a query signature to specify when upregulated query genes appear in the proximity of the top of the list or near the bottom (“positive connectivity”) or vice versa (“negative connectivity”) producing an Enrichment Score (ES) from  $-1$  to  $1$ . (iv) All instances in the database are then ranked according to their connectivity

scores; those at the top are most strongly correlated to the query signature, and those at the bottom are most strongly anticorrelated.

For each cancer type we obtained two tables that applied the Connectivity Map's findings to stemness mRNA expression signatures, namely, "detailed results" and "permuted results." We used the permuted results and filter (with  $p < 0.05$ ), to identify an average of 74 compounds per tumor type that are predicted to repress or activate the stemness signature (Table S4A).

Connectivity Map (CMap) was recently updated (September 2017) (Subramanian et al., 2017), providing the end-users new functionalities and new graphical interface as web-server, previous registration (<https://clue.io/>) allowing easily the extraction of drug-interaction knowledge using as input a signature of genes or compounds.

The new interface (<https://clue.io/>), provided 7 different analysis (query, touchstone, proteomics query, command, data library, repurposing, morpheus).

In particular CMap Query it is a tool for perturbagens that give rise to similar (or opposing) expression signatures, for a technical limit, the CMap Query 2017 allows only to upload 150 genes max for upregulated genes and 150 genes for downregulated genes. For this reason we considered the results analyzed in May 2017 using 500 genes for up-downregulated genes.

### QUANTIFICATION AND STATISTICAL ANALYSIS

R version 3.3.1 was used for all statistical analyses, unless specified otherwise. The statistical details of all experiments are reported in the figure legends and figures, including statistical analysis performed, statistical significance and exact n values.

To identify differentially methylated DNA methylation probes, we used the Wilcoxon test followed by multiple testing using the Benjamini-Hochberg (BH) method to estimate false discovery rate (Benjamini and Hochberg, 1995).

To identify proteins and microRNAs differentially expressed between tumors with low versus high stemness index, we used a t test followed by multiple testing using BH.

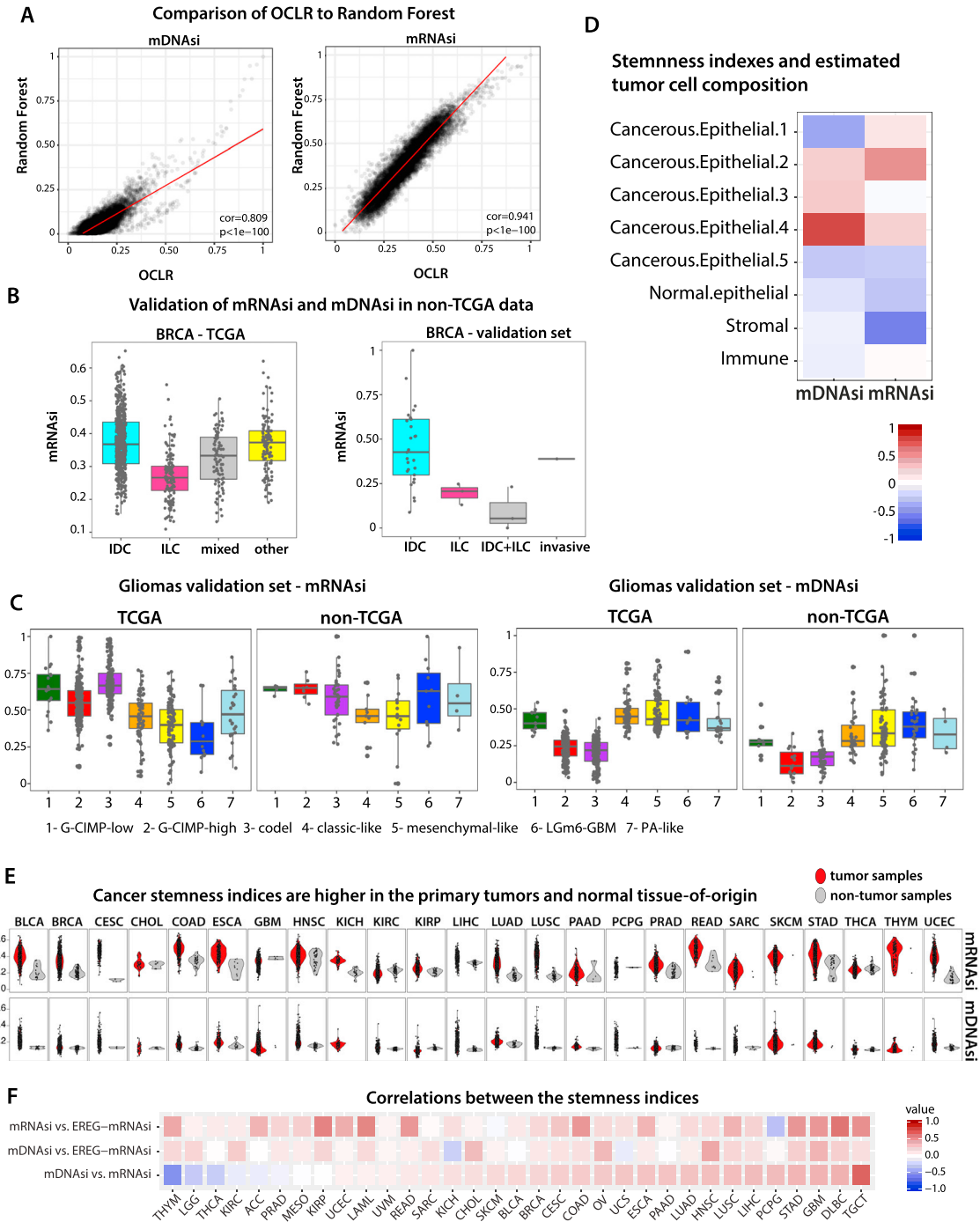
P values for the association between stemness index and continuous clinical data were also computed using a t test followed by multiple testing using BH.

### DATA AND SOFTWARE AVAILABILITY

All data are available on the NIH Genomic Data Commons (GDC), <https://gdc.cancer.gov/about-data/publications/PanCanStemness-2018>.

The workflow to reproduce the stemness index, including downloading PCBC and TCGA PanCan33 datasets, training a stemness signature, and applying it to score TCGA samples can be accessed at [https://bioinformaticsfmrp.github.io/PanCanStem\\_Web/](https://bioinformaticsfmrp.github.io/PanCanStem_Web/).

An interactive portal with the results for enrichment of molecular and clinical features and Stemness Indices across all tumor samples/types can be accessed at [https://bioinformaticsfmrp.github.io/PanCanStem\\_Web/](https://bioinformaticsfmrp.github.io/PanCanStem_Web/) where the user can search for any gene or molecular/clinical feature of interest.



**Figure S1. Validation of mRNAsi and mDNAsi, Related to Figures 1 and 3**

(A) Comparison of mRNAsi and mDNAsi values, as scored by signatures learned from PCBC data by OCLR (x axis) and Random Forest (y axis).  
 (B) Validation of mRNAsi in non-TCGA breast cancer (BRCA) samples (Reyngold et al., 2014) to define stemness status. Stratification of mRNAsi according to histology of TCGA samples (left) and of the validation cohort (right). IDC - invasive ductal carcinoma; ILC - invasive lobular carcinoma.  
 (C) Validation of stemness indices (mDNAsi and mRNAsi) in non-TCGA glioma samples (Sturm et al., 2012; Turcan et al., 2012) using our metric to define stemness status. The validation sets were previously classified into molecular subtypes of TCGA gliomas. Stratification of stemness indices according to molecular subtypes in the TCGA (left) and in the validation cohort (right).  
 (D) Correlation of stemness indices and estimated cell composition derived from deconvolution of tumor epigenomic features in TCGA breast cancer samples (Onuchic et al., 2016).

(legend continued on next page)



---

(E) Stemness indices in tumor versus non-tumor samples across tissue type. Top plots show mRNAsi and bottom plots show mDNAsi. See list of Abbreviations of the TCGA Tumor Types.

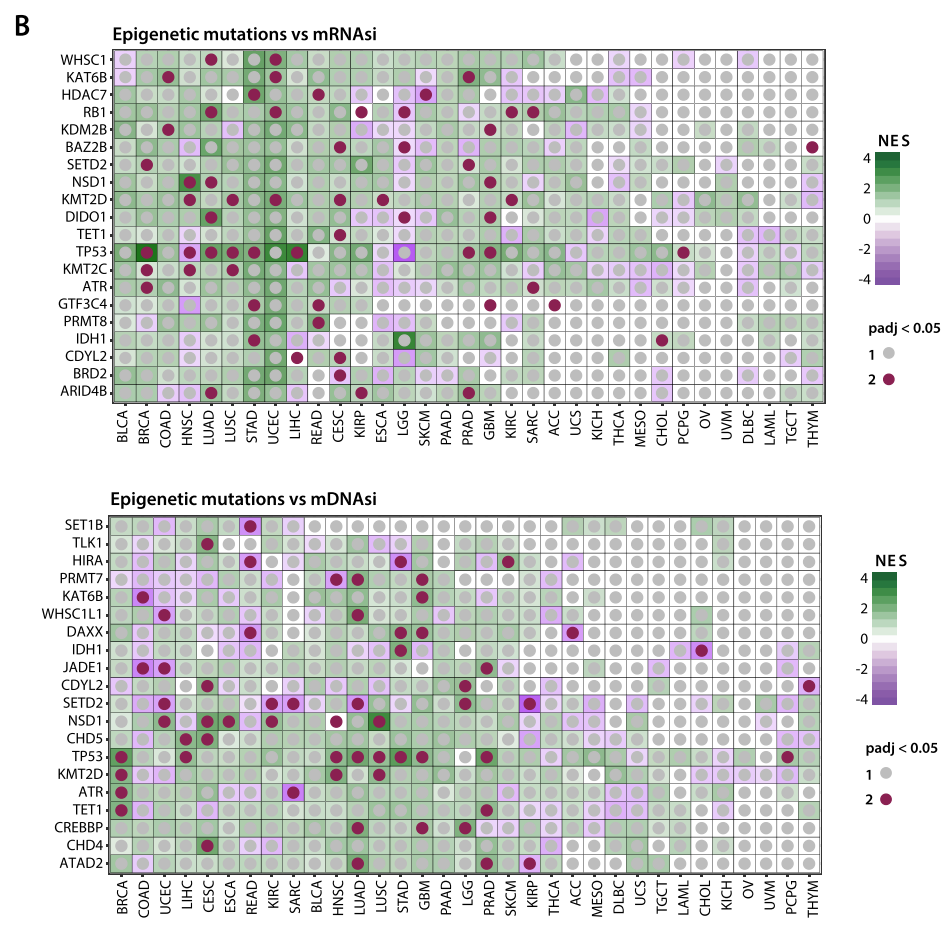
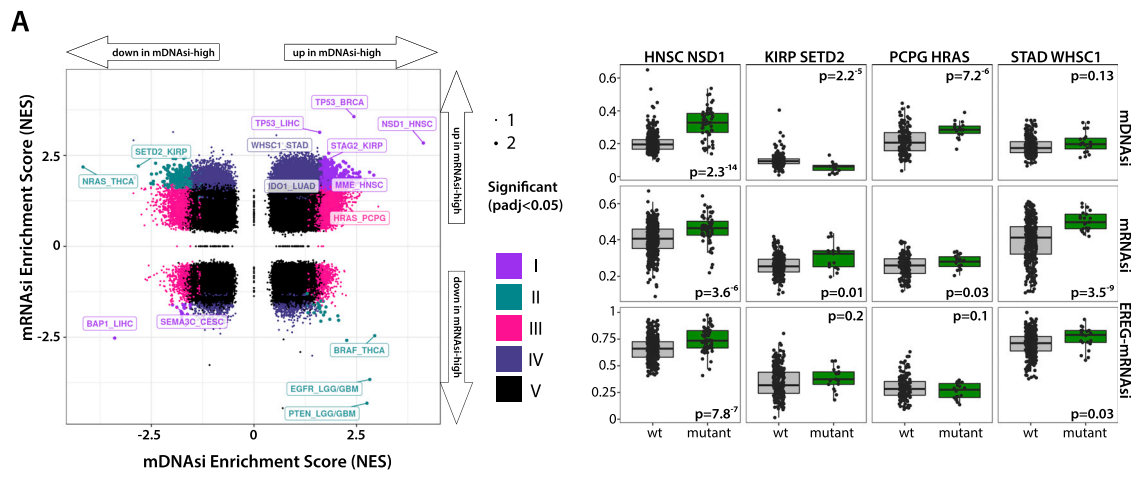
(F) Correlation of stemness indices across TCGA tumor types. Tumor types on the x axis are sorted by low-to-high correlation between mDNAsi and mRNAsi. See list of Abbreviations of the TCGA Tumor Types. mRNAsi, RNA-based stemness index. mDNAsi, DNA methylation-based stemness index. EREG-mRNAsi, epigenetically regulated-mRNAsi.



---

(B) Scatterplots showing correlation of mDNasi (x axis) and mRNasi (y axis) by tumor type. Samples are colored according to TCGA classification in: primary, recurrent, metastatic, and non-tumors.

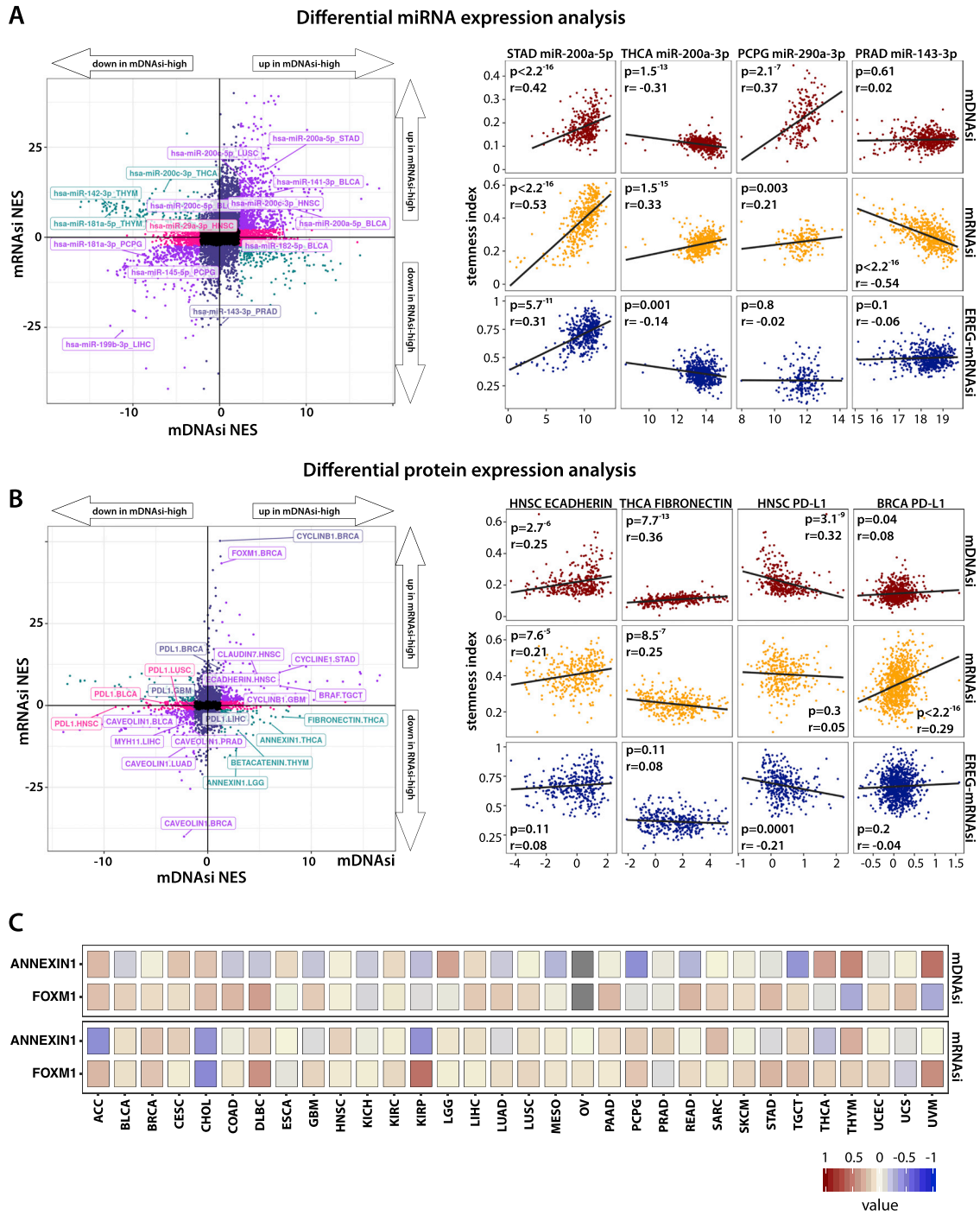
(C) Correlation between mRNasi and mRNA expression of epithelial-to-mesenchymal transition (EMT) markers. The markers are grouped according to whether they are associated with epithelial (pink) or mesenchymal (teal blue) phenotype.



**Figure S3. Mutation Enrichment Analysis in Correlation with the mRNAi and mDNAsi, Related to Figures 3 and 4**

(A) (Left Panel) Comparison of normalized enrichment scores (NES) in mDNAsi (x axis) and mRNAi (y axis) computed by tumor type in the mutation enrichment analysis. Positive and negative NES entail enrichment and depletion of a given mutation being associated with a high stemness index, respectively. The quadrants and colors represent genes with agreement (I) or disagreement (II) between the two distinct indices and being associated only with mDNAsi (III) or associated only with mRNAi (IV) between the two distinct indices (V, not significant). (Right Panel) Stemness indices stratified by representative mutation status and tumor types.

(B) Relationship of mutations in epigenetic modifier genes and stemness indices mRNAi (top) and mDNAsi (bottom). NES, normalized enrichment scores.



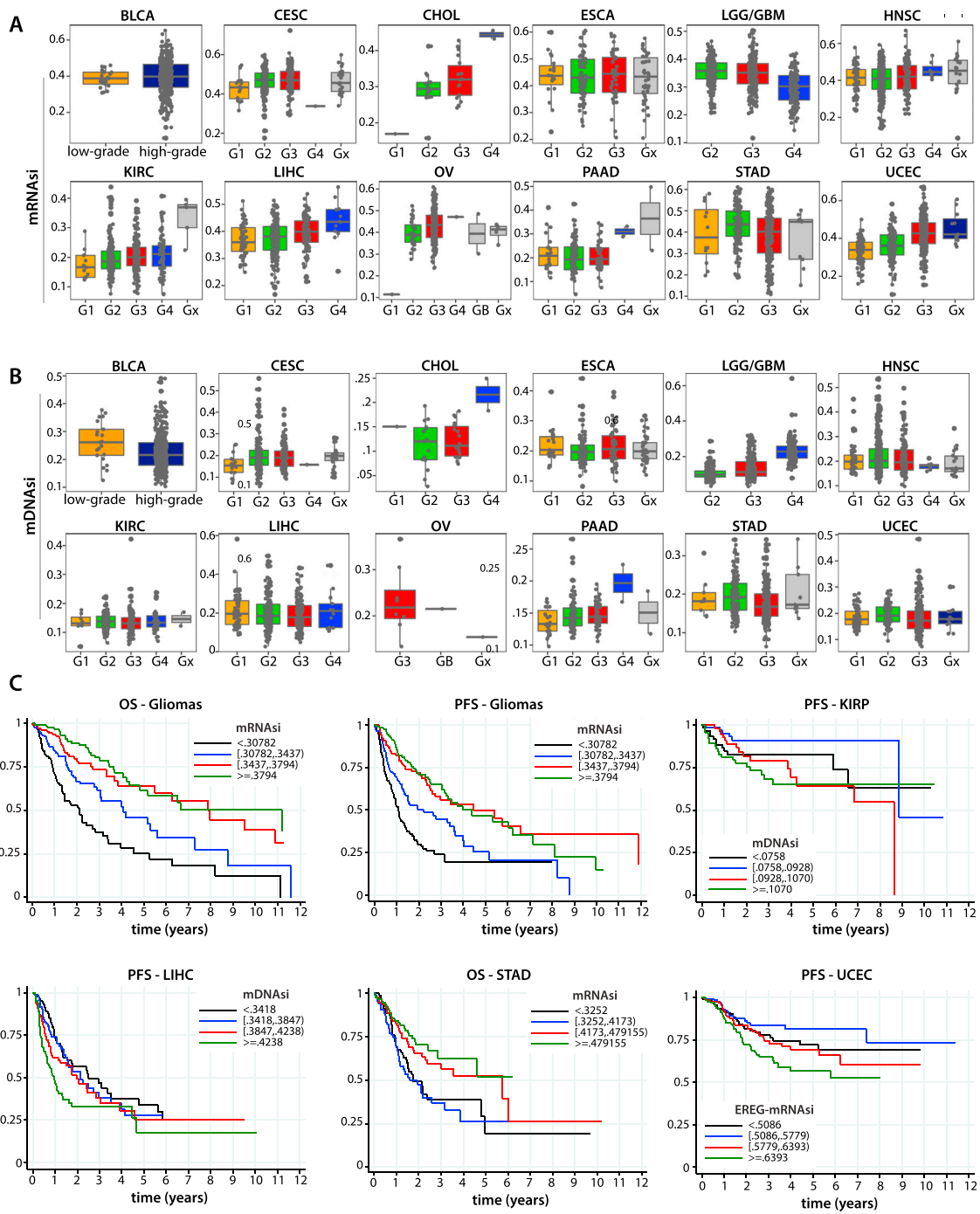
**Figure S4. Correlation of miRNA and Protein Expression against mRNAsi and mDNAsi, Related to Figures 3 and 4.**

(A) (Left Panel) Comparison of miRNA expression associated with mDNAsi and mRNAsi. Log<sub>10</sub>-normalized FDR-adjusted P value is plotted for mDNAsi (x axis) and mRNAsi (y axis) for each miRNA probe for each tumor type. If a particular miRNA is upregulated in samples with a high stemness index, the values are multiplied by  $-1$ . The quadrants and colors represent proteins with agreement or disagreement between the two distinct indices. (Right Panel) Correlation of stemness indices and representative miRNAs for selected tumor types.

(B) Comparison of protein expression associated with mDNAsi and mRNAsi (left panel). Log<sub>10</sub> (FDR-adjusted P value) is plotted for mDNAsi (x axis) and mRNAsi (y axis) for each protein for each tumor type. If the protein is upregulated in high stemness index,  $-1$  is multiplied to log<sub>10</sub> providing positive values. The quadrants and colors represent proteins with agreement or disagreement between the two distinct indices. Right, correlation scores of stemness indices and representative proteins and tumor types.

(C) Correlation of ANNEXIN1 and FOXM1 with mDNAsi (top) and mRNAsi (bottom) across tumor types.



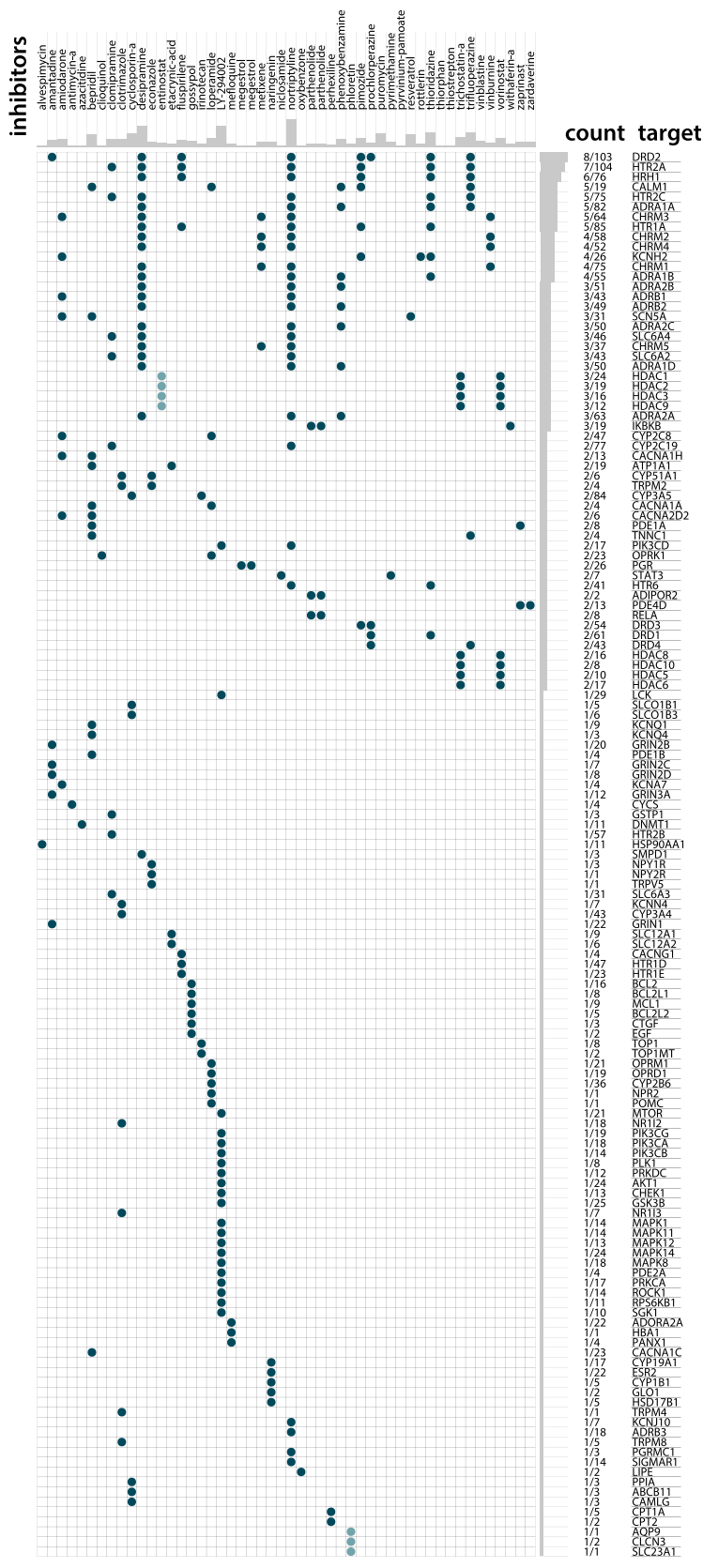


**Figure S5. Tumor Pathology and Clinical Outcome in Correlation with the mRNasi and mDNasi, Related to Figures 3 and 4.**

(A and B) Tumor pathology grade associations with mRNasi (A) and mDNasi (B) by TCGA tumors types.

(C) OS and PFS curves for four (quartile-based) categories of the relevant indices. The figure also includes the estimated (solid line) functional forms of the dependence of the logarithm of hazard ratio (log-HR) of OS or PFS on the indices. The forms in general correspond to the ordering of the survival curves. For selected indices and outcomes the estimated functional form suggests a non-monotonic relationship between log-HR and the index.





(legend on next page)

---

**Figure S7. Connectivity Map, Related to Figure 7**

Heatmap showing each compound (perturbagen) from the Connectivity Map that share gene targets in (rows). Sorted by descending number of targets. See also [Table S3](#).