

A pan-cancer analysis of prognostic genes

Jordan Anaya^{*1,4}, Brian J. Reon¹, Wei-Min Chen^{2,3}, Stefan Bekiranov¹, Anindya Dutta^{*1}

¹Department of Biochemistry and Molecular Genetics; ²Center for Public Health Genomics.

³Department of Public Health Sciences, Biostatistics Section; University of Virginia, Charlottesville, Virginia, USA.

⁴Current address: omnesres.com, Charlottesville, Virginia, USA

*Corresponding Authors:

Jordan Anaya, omnesresnetwork@gmail.com; Anindya Dutta, ad8q@virginia.edu

Abstract

Numerous studies have identified prognostic genes in individual cancers, but a thorough pan-cancer analysis has not been performed. In addition, previous studies have mostly used microarray data instead of RNA-SEQ, and have not published comprehensive lists of associations with survival. Using recently available RNA-SEQ and clinical data from the The Cancer Genome Atlas for 6,495 patients, we have investigated every annotated and expressed gene's association with survival across 16 cancer types. The most statistically significant harmful and protective genes were not shared across cancers, but were enriched in distinct gene sets which were shared across certain groups of cancers. These groups of cancers were independently reconstructed by unsupervised clustering of Cox coefficients (a measure of association with survival) for individual genes or for gene programs. This analysis has revealed unappreciated commonalities among cancers which may provide insights into cancer pathogenesis and rationales for co-opting treatments between cancers.

Main article text

Introduction

Led by The Cancer Genome Atlas, unprecedented efforts have been made to understand the molecular basis of cancer (<http://cancergenome.nih.gov>). Using standardized procedures, the TCGA Research Network has used whole genome sequencing, exome sequencing, RNA-SEQ, small RNA-SEQ, bisulfite-SEQ, and reverse phase arrays to identify the pathways commonly altered in different cancers (Brennan et al. 2013; Cancer Genome Atlas 2012a; Cancer Genome Atlas 2012b; Cancer Genome Atlas Research 2011; Cancer Genome Atlas Research 2012; Cancer Genome Atlas Research 2013a; Cancer Genome Atlas Research 2013b; Cancer Genome Atlas Research 2014a; Cancer Genome Atlas Research 2014b; Cancer Genome Atlas Research 2014c; Cancer Genome Atlas Research 2014d; Cancer Genome Atlas Research et al. 2013a). As a result, we now know the most commonly mutated genes in dozens of cancers and can use this information to give patients targeted therapeutics.

Whereas well established statistical techniques exist for identifying mutations which are drivers instead of simply passengers (mut-drivers), identifying copy number aberrations, methylation changes, or non-coding mutations that alter expression of a gene and result in a growth advantage (epi-drivers) are more difficult to identify and represent a "dark matter" of cancer (Vogelstein et al. 2013). Although it currently is challenging to identify epi-drivers which lead to development of a cancer (tumorigenesis), by correlating these changes to survival it is possible to detect their role in disease progression (pathogenesis), which is one of main goals of cancer research.

Of the possible genomic measures that can be correlated with survival, gene expression has been shown to be the strongest predictor of survival (Zhao et al. 2015), which is intuitive given that gene levels together with protein levels and posttranslational modifications are the final readout of the different possible alterations in a cell and are the final effectors of phenotype. To date many attempts have been made to identify genes whose expression is associated with survival to either identify markers that can predict patient survival or to identify mechanisms of pathogenesis (Chen et al. 2007; Valk et al. 2004; van de Vijver et al. 2002). One of the success stories of this approach is the identification of HER2 in breast cancer patients and the development of herceptin (Bange et al. 2001). This story also highlights the complications treatment regimens can have on interpreting survival data. Whereas HER2 overexpression used to predict poor survival for breast cancer patients, because of the progress of personalized medicine these patients now do well and HER2 would not show up as a prognostic marker in a data set with HER2 positive patients on herceptin. While treatments may introduce a confounding variable in understanding a disease, the ultimate goal of cancer studies is to improve patient outcome, and adding treatment to the equation adds more information and provides an opportunity to study genes in the context of the current standard of care.

The vast majority of studies to identify prognostic genes have focused on a single disease and have utilized microarrays instead of RNA-SEQ. In addition, these studies often only publish a small set of genes that together most significantly stratify patients. Even the TCGA Research Network publications do not provide lists of genes associated with survival. cBioPortal does allow users to make Kaplan Meier plots for most of the cancers which contain survival information, but users still have to input one gene at a time, leaving one to wonder where researchers should go to find the genes which are most highly correlated with survival for their disease of interest.

Through the TCGA Network, RNA-SEQ has only very recently become available for thousands of human cancer samples. RNA-SEQ has multiple advantages over microarray data, including having a higher dynamic range, no probe affinity effects, ability to identify novel transcripts, and lower and consistently falling cost. We took advantage of the availability of this data to 1) investigate the ability of RNA-SEQ to associate expression with clinical outcome in a range of cancers, 2) perform the largest analysis of prognostic genes to date, and 3) provide every gene's correlation with survival for hypothesis testing and further investigations by the scientific community. In addition, attempts are now being made to identify commonalities between cancers with the hope that this type of analysis may be able to identify treatments that can be co-opted for a molecularly similar cancer. Given that only survival correlations integrate treatment with the genomic data, prognostic genes represent an exclusive window for understanding how different cancers in the context of their individualized treatments relate to one another. The analysis identified reproducible groupings of cancers based on prognostic genes. This study serves as a starting point for better understanding how survival data can be used to understand the commonalities and differences of cancers.

Materials and methods

Code and files

All of the Python and R code used to generate the figures and tables in this study, including intermediate and final files, tables, and figures, is available at https://github.com/OmnesRes/pan_cancer. All scripts were run on a HP dv7t laptop with an i7-3820QM processor and 16GB of RAM running Windows 7, Python 2.7.5, and R 3.0.1.

Construction of multivariate Cox models

RNA-SEQ and clinical data were downloaded from the TCGA data portal, <https://tcga-data.nci.nih.gov/tcga/>. For each cancer, survival information was parsed from the "clinical_follow_up" files and "clinical_patient" file, and for each patient the most recent follow up information found in the multiple files was kept. Sex, age, and histological grade data were extracted from the "clinical_patient" file. For each cancer, only patients that had a follow up time greater than 0 days and had complete clinical information were included in the model. TCGA has used two different methods of reporting expression values, RSEM and RPKM. RPKM is simply the reads per kilobase per million mapped reads, while RSEM is a normalized value outputted by the RSEM software (Li & Dewey 2011). For each cancer, only genes which had a median RSEM value greater than 1 (for RNASeqV2), or median RPKM value greater than .1 (for RNASeq), and had 0 expression in less than one fourth of patients were included in the analysis. RNASeq uses a different gene annotation file from RNASeqV2, and because RNASeqV2 represents the most recent analysis, for RNASeq analyses only those genes present in the RNASeqV2 gene annotation file were included. Multivariate Cox models were run with the coxph function from the R survival library, and the equation for each model is shown in Table S1. Grade information was included in the model by separate terms, which were either 1 or 0, and model input gene expression values were inverse normal transformed. If a patient had replicates for their primary tumor, those expression values were averaged prior to inverse normal transformation. The scripts for performing Cox regression for each cancer are named "cox_regression.py".

Gene set analysis

For each cancer, the 250 most significant protective genes and 250 harmful genes were inputted separately into MSigDB with the "positional genes sets", "chemical and genetic perturbations", "canonical pathways", "KEGG gene sets", "microRNA targets", "transcription factor targets", "cancer modules", "GO biological process", and "oncogenic signatures" sets selected: <http://www.broadinstitute.org/gsea/msigdb>. The FDR q-value threshold was set at .05 and the top 100 enriched gene sets were saved, except for the 250 protective genes in BLCA, which only contained 27 overlaps below .05.

Normalization of Cox coefficients

In order to compare the Cox coefficients between cancers we robustly scaled the negative and positive coefficients, x , to their 5th and 95th percentile values, respectively, using the following sigmoidal normalization function:

$$z_{\pm} = \frac{2}{1 + e^{-\frac{2x}{|u_{\pm}|}}} - 1$$

where u_{-} and u_{+} are the 5th and 95th percentile values of the negative and positive Cox coefficients, respectively. The implementation of this code is present in the files named "normalizing_coeffs.py", and all the normalized coefficients are listed in Table S1.

Construction of gene programs

Gene programs from Table S4 of (Hoadley et al. 2014) were used. In general a nonredundant set of genes from gene sets which had a Pearson correlation of at least .9 (Hoadley et al. 2014) was generated for each program. The exact gene sets used are listed in Table S3. Lists of genes for the gene sets were obtained from <http://www.broadinstitute.org/gsea/msigdb> and (Fan et al. 2011).

Results

Cancers vary in number of prognostic genes

In order to perform the most comprehensive cancer analysis possible, we selected TCGA cancers that had sufficient numbers of patients with RNA-SEQ data and mature clinical follow up information, and did not contain any publication restrictions. This resulted in us studying a total of 16 cancers, 10 of which were present in the original pan-cancer initiative (Cancer Genome Atlas Research et al. 2013b): acute myeloid leukemia (LAML), bladder urothelial carcinoma (BLCA), breast invasive carcinoma (BRCA), colon adenocarcinoma (COAD), glioblastoma multiforme (GBM), head and neck squamous cell carcinoma (HNSC), kidney renal clear cell carcinoma (KIRC), lung adenocarcinoma (LUAD), lung squamous cell carcinoma (LUSC), and ovarian serous cystadenocarcinoma (OV), and 6 cancers which have been the focus of limited individual or pan-cancer studies: cervical squamous cell carcinoma and endocervical adenocarcinoma (CESC), brain lower grade glioma (LGG), kidney renal papillary cell carcinoma (KIRP), liver hepatocellular carcinoma (LIHC), skin cutaneous melanoma (SKCM), and stomach adenocarcinoma (STAD).

We were interested in the effect a gene has on prognosis independent of factors such as tumor grade and age of a patient. To achieve this we used a multivariate Cox proportions hazards model (Cox 1972), which is a standard regression method for studying survival data (Claus et al.

2015; Györfy et al. 2013; Wu & Stein 2012; Zhang et al. 2013). For every cancer, a model was generated separately for each gene, with the number of covariates depending on the cancer. Unlike microarray data, RNA-SEQ data has extreme values which may affect regression. To account for this we inverse normal transformed the expression values of each gene, which has been shown to increase the sensitivity and specificity for multivariate regression with RNA-SEQ data (Zwiener et al. 2014). Age and sex are also included in every model, and when a cancer contained strong histological grade information, grade was also included. If a patient was missing any of this information they were excluded from the analysis, and only primary tumors were considered, with the exception of SKCM, where metastatic tumors make up a large proportion of the patients.

A Cox model provides a p-value for each term in the model, indicating the significance of its association with the clinical outcome, and we recorded the p-values for every gene analyzed for the 16 different cancers. As can be seen in Table 1, there is a wide distribution among cancers in the number of genes that reached a Benjamini-Hochberg False Discovery Rate (FDR) adjusted p-value of less than or equal to .05. This can also be seen by looking at a distribution of the raw p-values for the different cancers (Fig. 1a and Fig. S1). This has important implications for understanding the significance of a gene being associated with survival in a specific cancer. For example, selecting a gene at random a researcher studying LGG has a 50% chance of being able to claim the gene is associated with survival, while a researcher studying STAD only has an 8% chance (using raw p-values).

Two factors that are known to be associated with power of a Cox model are sample size and number of events (deaths); however looking through Table 1 it is difficult to find a pattern that can explain why certain cancers have more significant expression level based prognostic genes (EPGs) than other cancers. For example, BRCA has around twice the number of patients of any other cancer, but only has 30 EPGs that meet a FDR cutoff. In contrast, KIRP has a fourth the number of patients of BRCA but has 2,415 EPGs. In addition, LUAD and LUSC have similar numbers of patients, median survivals, and events, yet have a large difference in number of EPGs. Interestingly, it has been shown that the number of prognostic genes for a cancer can be significantly different depending on whether microarray data or RNA-SEQ data is used (Yang et al. 2014), but that cannot be the explanation here. It is possible that the different numbers of EPGs between cancers are due to intra-disease heterogeneity and/or treatment differences that are not accounted for in the Cox model and are acting as confounding variables, or differences in the amount of transcriptional dysregulation between cancers.

Protective and harmful genes display opposite expression patterns

The Cox model also provides a coefficient for each term, which is related to its contribution to the hazard ratio. A positive coefficient indicates that the gene increases the hazard ratio, i.e. high expression of the gene correlates with earlier patient death, while a negative coefficient indicates that expression of the gene is protective. Using the cancer with the highest number of EPGs

(LGG), we clustered patients with the 100 most significant genes which were harmful and the 100 most significant genes which were protective, and this revealed two broad clusters of patients: (1) those with high expression of harmful genes and low expression of protective genes, and (2) those with high expression of protective genes and low expression of harmful genes (Fig. 1b). As expected, a Kaplan Meier analysis with these two groups revealed that cluster 2 has a much higher survival than cluster 1 (Fig. 1c). This result has important implications for trying to find gene sets which can most accurately predict patient survival. The similar expression patterns indicate that there are numerous combinations of genes that would only differ slightly in their ability to predict survival, making the identification of a ‘best’ set of genes somewhat meaningless. In addition, given that each gene individually had a p-value less than or equal to $1.4E-8$, it is unlikely these patterns are due to chance but rather are explained either by common underlying gene regulatory pathways or by these genes being members of common cellular pathways.

Unlike LGG, some of the other cancers in the analysis yielded a much lower number of EPGs. While it might be tempting to disregard the results in these cancers, we checked if there are patterns of expression in the most significant good and bad genes like that observed for cancers with a high number of EPGs. Clustering of the patients of STAD, which has one of the lowest numbers of EPGs, with the 100 most significant harmful genes, and the 100 most significant protective genes, again divided patients into two broad clusters. Interestingly, a Kaplan Meier analysis on these two groups showed a very significant prognostic difference with a p-value of $2.73E-6$ (Fig. 1c). This indicates that despite the fact that none of the genes in STAD meet a 5% FDR cutoff, they still contain important biological information. As a result, further analyses included all the cancers regardless of their numbers of EPGs.

Cancers do not share prognostic genes, but do share gene sets

We next tested whether the most significantly prognostic genes were shared across cancers. However, there is very little overlap among the 100 most significant genes across the 16 cancers, consistent with previous results obtained from an analysis of four cancers (Yang et al. 2014) (Fig. 2a). Given the apparent co-regulation of the most significant genes in each cancer, we reasoned that although individual genes were not shared, maybe the genes were a part of gene sets which were shared between cancers. In addition, given that the harmful genes had an opposite pattern of expression from the protective genes, we hypothesized that they are regulated differently and would be enriched in different gene sets. To investigate this we took the 250 most significant harmful genes and 250 most significant protective genes in each cancer, and separately found the 100 most enriched gene sets through MSigDB (Subramanian et al. 2005). Consistent with the idea that harmful and protective genes are regulated differently, there was very little overlap between the 100 gene sets found with 250 harmful genes and the 100 gene sets found with 250 protective genes for a given cancer (Fig. 2b). In addition, the fact that even the protective and harmful gene sets from cancers with a low number of EPGs show almost no

overlap reinforces the idea that prognostic genes in these cancers still contain biologically significant information.

Next we assessed the extent to which these protective and harmful gene sets overlapped between the different cancers. The extent of overlap was investigated separately for the 100 harmful gene sets and 100 protective gene sets (Fig. 2c, 2d). Overall there was more overlap between the harmful gene sets, and there were three cancers which clearly shared a high number of harmful gene sets, LUAD, LIHC, and KIRP. Investigating these overlaps further showed that the three cancers shared 58 gene sets, and LUAD and KIRP shared 85 gene sets (Fig. 2e). Looking at the overlaps of the protective gene sets, the largest overlap was between COAD and LUSC, and these cancers also shared gene sets with GBM (Fig. 2f).

We next asked what are the most common harmful and protective gene sets across cancers. Table S2 shows frequency of every gene set, with gene sets that were shared between harmful and protective sets within a single cancer marked in bold as they may be nonspecific. As might be expected, the most common gene sets observed for harmful genes were associated with poor differentiation and metastasis. In contrast, the protective gene sets were enriched for apoptosis and good differentiation. Although when possible the grade of the tumor was included in the Cox model, and therefore should not be a confounding variable, it is possible that histological grade does not completely account for the differentiation of a tumor, indicating the importance of genomics for accurate profiling.

Cancers can be clustered by gene and gene program Cox coefficients

To date different cancers have been compared to each other through mRNA levels, miRNA levels, protein levels, networks, copy number alterations, DNA methylation, somatic mutations or some combination of these (Akbani et al. 2014; Ciriello et al. 2013; Hamilton et al. 2013; Hoadley et al. 2014; Kandoth et al. 2013; Knaack et al. 2014). The Cox coefficients in my analysis contain a level of information not present in any of these data types, and consequently can potentially reveal similarities or differences between cancers that were not appreciated before. Therefore we sought to attempt to cluster cancers using Cox coefficients of genes instead of expression levels. Because the Cox models for the different cancers contain different numbers of covariates, and different strengths of gene expression correlation to survival, the range of values of the Cox coefficients vary between cancers. To correct for this, we normalized the coefficients for each cancer using a sigmoidal function which robustly scaled both negative and positive coefficients to their 95th percentile values (see methods). In addition, whereas every gene has an expression value, only significant prognostic genes have Cox coefficients appreciably above or below 0. Performing clustering with large numbers of nonsignificant genes which all have very similar values for every cancer will only add noise to the clustering. As a result, the clustering was limited to genes which had a FDR less than or equal to .05 in at least four of the sixteen cancers.

Hierarchical clustering of the 16 cancers was performed with the sigmoidal normalized Cox coefficients of this set of genes (Fig. 3). The clustering grouped LIHC, LUAD, and KIRP together, which were the same cancers that shared the highest number of harmful gene sets. In addition, GBM, COAD and LUSC clustered together, which were the cancers that had the highest number of protective gene sets overlap. The fact that two separate methods, using different sets of genes, were able to find similar groupings of cancers suggests that the similarities between the cancers in each group is robust and possibly biologically significant.

We next tested whether there were established pathways that distinguished the groupings of cancers from each other. Using a list of nonredundant gene programs that have been shown to distinguish cancers from one another on the basis of expression levels (Hoadley et al. 2014), we sought to distinguish cancers using Cox coefficients of pathways. For each pathway the average sigmoidal normalized Cox coefficient was calculated in each cancer. Because a Cox coefficient can be positive or negative, if a pathway has some genes which are protective and some genes which are harmful, the average Cox score will be near zero. In addition, if a pathway only contains genes which are not prognostic, all of those Cox scores will be near zero and the pathway score will be near 0. The only way for a pathway to have a positive or negative score is for it to contain prognostic genes which are either consistently protective or consistently harmful.

Hierarchical clustering was performed with the Cox scores for these 22 gene programs (Fig. 4). The values were column scaled to highlight which gene programs are most important for each cancer. Overall the same groupings that were seen with gene sets and individual genes were recapitulated from clustering the Cox scores of gene programs, with LUAD, KIRP, and LIHC again forming a cluster and COAD, LUSC, and GBM grouped together. In the LUAD/KIRP/LIHC group poor prognosis is associated with high proliferation rates and glycolysis, while good prognosis is associated with apoptosis and a dependence on oxidative phosphorylation. In contrast, for GBM/LUSC/COAD, proliferation is protective while genes associated with the EGF response predict poor survival.

The analysis also found cancer specific protective/harmful pathway enrichments that are consistent with known cancer biology. For example, in KIRC the highest intensity gene program is "fatty acid oxidation", and KIRC is a cancer that is known to depend on dysregulation of metabolism and is a classic example of the "Warburg effect" (Linehan et al. 2010). The results show that patients with high expression of genes utilizing oxygen survive longer, which underscores the importance of a metabolic shift in this cancer. As another example, EGFR is the most commonly mutated gene in GBM (Brennan et al. 2013), and in our analysis increased EGFR activity is associated with poorer outcomes. BLCA and SKCM, which are known for being responsive to immunotherapy, both benefit from increased interferon response and an immune cell signature which is likely a proxy for immune cell infiltration.

Discussion

Cancer researchers are increasingly looking to focus on factors which have clinical significance, and many different resources now allow researchers to identify if a protein of interest has clinical implications, including OMIM, dbSNP, ClinVar, cBioPortal, FINDbase, and others (Hamosh et al. 2005; Landrum et al. 2014; Papadopoulos et al. 2014; Smigielski et al. 2000). Despite this, it currently is not possible to find comprehensive lists of genes which are associated with survival in different cancers. Using recently available RNA-SEQ and clinical data from the TCGA for 6,495 patients, we correlated every expressed annotated gene to survival in 16 different cancers, providing the scientific community with thousands of highly significant genes for further study.

There is an unexpectedly large variation between cancers in the number of statistically significant prognostic genes, which should be used to inform our evaluation of prognostic genes from different cancers. For example, a significant p-value for a gene from a cancer such as LGG or KIRC should not be surprising, given the thousands of genes that survive a stringent p-value cutoff in these tumors (Table 1, Fig. 1a, Fig. S1). In contrast, weaker p-values for predicting prognosis in cancers such as STAD or COAD are still biologically important although they have no genes that pass a stringent p-value threshold for biological significance (Table 1, Fig. 1a,c).

RNA-SEQ is a relatively new technology, and its ability to identify prognostic genes in many cancers has not been explored. Although the number of expression level based prognostic genes (EPGs) varied among cancers, regardless of the cancer we identified expression profiles which significantly separated patients into high risk and low risk groups. One of the main advantages of RNA-SEQ over microarrays is the ability to identify unannotated transcripts. In fact, recent studies have investigated the expressions of pseudogenes and long noncoding RNAs in large numbers of TCGA RNA-SEQ data sets (Han et al. 2014; Iyer et al. 2015). It would be interesting to see if these transcripts show the same trends as protein coding genes across these cancers.

This comprehensive analysis of prognostic genes allowed us to explore the ability of the prognostic genes themselves, enriched gene sets, and Cox coefficients (a measure of strength of correlation to better or worse survival) to identify similarities and differences among cancers. The most prognostically significant genes were not shared between cancers. However, protective genes and harmful genes are enriched in very different gene sets, and there were large overlaps of these gene sets for LUAD, LIHC, and KIRP, and for COAD, LUSC, and GBM. The groupings of these cancers were recapitulated by clustering with both Cox coefficients of individual genes, and average Cox coefficients of gene programs, suggesting that these findings are biologically significant and that this is an effective strategy for incorporating genomic and clinical data to compare cancers.

Although it is important not to mistake a correlation for causation, the analysis suggests intriguing insights into the pathogenesis of different cancers. For example, currently EGFR inhibitors are recommended for LUAD patients with EGFR mutations, but EGFR mutations are rare in LUSC and patients with mutations do not respond well to tyrosine kinase inhibitors (Chiu

et al. 2014). Despite this, response rates to EGFR inhibitors for LUSC studies are threefold higher than expected (Chiu et al. 2014), suggesting that although EGFR itself may not be mutated, responders may still have a cancer which is dependent on EGFR signaling. This is consistent with the gene program analysis in this paper, where EGFR response was most strongly associated with poor survival in LUSC, and LUSC was consistently associated with GBM, which is a cancer known for EGFR dysregulation. This suggests that using a measure of EGFR activity other than mutational status could be used to find LUSC patients that would benefit from a tyrosine kinase inhibitor. In addition, this type of analysis may be used to suggest treatments for cancers which are not well studied. For example, KIRP does not have successful treatments and there is a current search for drugs which may be of benefit (Schuller et al. 2015). This analysis suggests that the pathogenesis of KIRP is very similar to LIHC and LUAD, indicating that treatments currently used for those cancers may be co-opted for KIRP.

This analysis is among the first attempts at using clinical correlations to compare cancers. Although we utilized the most up to date information possible, well established statistical techniques, and obtained robust findings, there are many ways this type of analysis can be improved. For example, it is now being recognized that cancer is not a single disease, but rather a group of molecularly and clinically distinct diseases which share a tissue of origin. Through a combination of genomic measurements, the TCGA Research Network has divided individual cancers into four or five subtypes, for example GBM has been divided into proneural, neural, classical, and mesenchymal subtypes (Brennan et al. 2013). Currently, clear subtypes have not been found for all 16 of the cancers in this study, and for many cancers dividing the cancers into the subtypes would result in a loss of power due to the limited number of patients. However, as these classifications are refined and the number of patient samples continues to grow, a natural extension of this study would be to repeat it for individual subtypes, which would potentially decrease the heterogeneity of the data. In addition, treatment is one the largest confounding variables in survival analyses, but the TCGA pharmacological data is currently incomplete making it impossible to incorporate this information into the model. Despite these current limitations, this study has shown that incorporating clinical information into pan-cancer analyses is capable of yielding insights into cancer pathogenesis that have thus far been unappreciated by other methods.

Additional Information and Declarations

Competing Interests

Jordan Anaya has started a company, Omnes Res, that may carry out analyses similar to that reported in this paper. The other authors have no competing interests to declare.

Author Contributions

This project was initiated and completed by JA while he was a Ph.D. student in the laboratories of Drs. Anindya Dutta and Stefan Bekiranov. The project was conceived out of a joint discussion between JA, BR and AD. JA performed all analyses, prepared the figures and wrote the first draft of the paper, which was then edited by SB and AD. WMC and SB provided analysis tools and advice on statistical analyses. AD provided material for this study. The GitHub repository will be maintained by JA and AD.

Funding

JA was supported by the Cell and Molecular Biology Training Grant T32-GM008136 and by an F31-CA189446 from the NIH. The work was supported by R01-CA60499 to AD from the NIH.

Acknowledgments

We acknowledge the contributions of the TCGA Research Network and its TCGA Pan-Cancer Analysis Working Group led by J.M. Stuart, C. Sander and I. Shmulevich. Without their efforts this type of analysis would not be possible.

References

- Akbani R, Ng PK, Werner HM, Shahmoradgoli M, Zhang F, Ju Z, Liu W, Yang JY, Yoshihara K, Li J, Ling S, Seviour EG, Ram PT, Minna JD, Diao L, Tong P, Heymach JV, Hill SM, Dondelinger F, Stadler N, Byers LA, Meric-Bernstam F, Weinstein JN, Broom BM, Verhaak RG, Liang H, Mukherjee S, Lu Y, and Mills GB. 2014. A pan-cancer proteomic perspective on The Cancer Genome Atlas. *Nat Commun* 5:3887. 10.1038/ncomms4887
- Bange B, Zwick E, and Ullrich A. 2001. Molecular targets for breast cancer therapy and prevention. *Nature Medicine* 7:548-552.
- Brennan CW, Verhaak RG, McKenna A, Campos B, Noushmehr H, Salama SR, Zheng S, Chakravarty D, Sanborn JZ, Berman SH, Beroukhi R, Bernard B, Wu CJ, Genovese G, Shmulevich I, Barnholtz-Sloan J, Zou L, Vegesna R, Shukla SA, Ciriello G, Yung WK, Zhang W, Sougnez C, Mikkelsen T, Aldape K, Bigner DD, Van Meir EG, Prados M, Sloan A, Black KL, Eschbacher J, Finocchiaro G, Friedman W, Andrews DW, Guha A, Iacocca M, O'Neill BP, Foltz G, Myers J, Weisenberger DJ, Penny R, Kucherlapati R, Perou CM, Hayes DN, Gibbs R, Marra M, Mills GB, Lander E, Spellman P, Wilson R, Sander C, Weinstein J, Meyerson M, Gabriel S, Laird PW, Haussler D, Getz G, Chin L, and Network TR. 2013. The somatic genomic landscape of glioblastoma. *Cell* 155:462-477. 10.1016/j.cell.2013.09.034
- Cancer Genome Atlas N. 2012a. Comprehensive molecular characterization of human colon and rectal cancer. *Nature* 487:330-337. 10.1038/nature11252
- Cancer Genome Atlas N. 2012b. Comprehensive molecular portraits of human breast tumours. *Nature* 490:61-70. 10.1038/nature11412
- Cancer Genome Atlas Research N. 2011. Integrated genomic analyses of ovarian carcinoma. *Nature* 474:609-615. 10.1038/nature10166
- Cancer Genome Atlas Research N. 2012. Comprehensive genomic characterization of squamous cell lung cancers. *Nature* 489:519-525. 10.1038/nature11404

Cancer Genome Atlas Research N. 2013a. Comprehensive molecular characterization of clear cell renal cell carcinoma. *Nature* 499:43-49. 10.1038/nature12222

Cancer Genome Atlas Research N. 2013b. Genomic and epigenomic landscapes of adult de novo acute myeloid leukemia. *N Engl J Med* 368:2059-2074. 10.1056/NEJMoa1301689

Cancer Genome Atlas Research N. 2014a. Comprehensive molecular characterization of gastric adenocarcinoma. *Nature* 513:202-209. 10.1038/nature13480

Cancer Genome Atlas Research N. 2014b. Comprehensive molecular characterization of urothelial bladder carcinoma. *Nature* 507:315-322. 10.1038/nature12965

Cancer Genome Atlas Research N. 2014c. Comprehensive molecular profiling of lung adenocarcinoma. *Nature* 511:543-550. 10.1038/nature13385

Cancer Genome Atlas Research N. 2014d. Integrated genomic characterization of papillary thyroid carcinoma. *Cell* 159:676-690. 10.1016/j.cell.2014.09.050

Cancer Genome Atlas Research N, Kandoth C, Schultz N, Cherniack AD, Akbani R, Liu Y, Shen H, Robertson AG, Pashtan I, Shen R, Benz CC, Yau C, Laird PW, Ding L, Zhang W, Mills GB, Kucherlapati R, Mardis ER, and Levine DA. 2013a. Integrated genomic characterization of endometrial carcinoma. *Nature* 497:67-73. 10.1038/nature12113

Cancer Genome Atlas Research N, Weinstein JN, Collisson EA, Mills GB, Shaw KR, Ozenberger BA, Ellrott K, Shmulevich I, Sander C, and Stuart JM. 2013b. The Cancer Genome Atlas Pan-Cancer analysis project. *Nat Genet* 45:1113-1120. 10.1038/ng.2764

Chen HY, Yu SL, Chen CH, Chang GC, Chen CY, Yuan A, Cheng CL, Wang CH, Terng HJ, Kao SF, Chan WK, Li HN, Liu CC, Singh S, Chen WJ, Chen JJ, and Yang PC. 2007. A five-gene signature and clinical outcome in non-small-cell lung cancer. *N Engl J Med* 356:11-20. 10.1056/NEJMoa060096

Chiu CH, Chou TY, Chiang CL, and Tsai CM. 2014. Should EGFR mutations be tested in advanced lung squamous cell carcinomas to guide frontline treatment? *Cancer Chemother Pharmacol* 74:661-665. 10.1007/s00280-014-2536-3

Ciriello G, Miller ML, Aksoy BA, Senbabaoglu Y, Schultz N, and Sander C. 2013. Emerging landscape of oncogenic signatures across human cancers. *Nat Genet* 45:1127-1133. 10.1038/ng.2762

Claus EB, Walsh KM, Wiencke JK, Molinaro AM, Wiemels JL, Schildkraut JM, Bondy ML, Berger M, Jenkins R, and Wrensch M. 2015. Survival and low-grade glioma: the emergence of genetic information. *Neurosurg Focus* 38:E6. 10.3171/2014.10.FOCUS12367

Cox DR. 1972. Regression Models and Life-Tables. *Journal of the Royal Statistical Society Series B (Methodological)* 34:187-220.

Fan C, Prat A, Parker JS, Liu Y, Carey LA, Troester MA, and Perou CM. 2011. Building prognostic models for breast cancer patients using clinical variables and hundreds of gene expression signatures. *BMC Med Genomics* 4:3. 10.1186/1755-8794-4-3

Gyorffy B, Surowiak P, Budczies J, and Lanczky A. 2013. Online survival analysis software to assess the prognostic value of biomarkers using transcriptomic data in non-small-cell lung cancer. *PLoS One* 8:e82241. 10.1371/journal.pone.0082241

Hamilton MP, Rajapakshe K, Hartig SM, Reva B, McLellan MD, Kandoth C, Ding L, Zack TI, Gunaratne PH, Wheeler DA, Coarfa C, and McGuire SE. 2013. Identification of a pan-cancer oncogenic microRNA superfamily anchored by a central core seed motif. *Nat Commun* 4:2730. 10.1038/ncomms3730

Hamosh A, Scott AF, Amberger JS, Bocchini CA, and McKusick VA. 2005. Online Mendelian Inheritance in Man (OMIM), a knowledgebase of human genes and genetic disorders. *Nucleic Acids Res* 33:D514-517. 10.1093/nar/gki033

Han L, Yuan Y, Zheng S, Yang Y, Li J, Edgerton ME, Diao L, Xu Y, Verhaak RG, and Liang H. 2014. The Pan-Cancer analysis of pseudogene expression reveals biologically and clinically relevant tumour subtypes. *Nat Commun* 5:3963. 10.1038/ncomms4963

Hoadley KA, Yau C, Wolf DM, Cherniack AD, Tamborero D, Ng S, Leiserson MD, Niu B, McLellan MD, Uzunangelov V, Zhang J, Kandoth C, Akbani R, Shen H, Omberg L, Chu A, Margolin AA, Van't Veer LJ, Lopez-Bigas N, Laird PW, Raphael BJ, Ding L, Robertson AG, Byers LA, Mills GB, Weinstein JN, Van Waes C, Chen Z, Collisson EA, Cancer Genome Atlas Research N, Benz CC, Perou CM, and Stuart JM. 2014. Multiplatform analysis of 12 cancer types reveals molecular classification within and across tissues of origin. *Cell* 158:929-944. 10.1016/j.cell.2014.06.049

Iyer MK, Niknafs YS, Malik R, Singhal U, Sahu A, Hosono Y, Barrette TR, Prensner JR, Evans JR, Zhao S, Poliakov A, Cao X, Dhanasekaran SM, Wu YM, Robinson DR, Beer DG, Feng FY, Iyer HK, and Chinnaiyan AM. 2015. The landscape of long noncoding RNAs in the human transcriptome. *Nat Genet* 47:199-208. 10.1038/ng.3192

Kandoth C, McLellan MD, Vandin F, Ye K, Niu B, Lu C, Xie M, Zhang Q, McMichael JF, Wyczalkowski MA, Leiserson MD, Miller CA, Welch JS, Walter MJ, Wendl MC, Ley TJ, Wilson RK, Raphael BJ, and Ding L. 2013. Mutational landscape and significance across 12 major cancer types. *Nature* 502:333-339. 10.1038/nature12634

Knaack SA, Siahpirani AF, and Roy S. 2014. A pan-cancer modular regulatory network analysis to identify common and cancer-specific network components. *Cancer Inform* 13:69-84. 10.4137/CIN.S14058

Landrum MJ, Lee JM, Riley GR, Jang W, Rubinstein WS, Church DM, and Maglott DR. 2014. ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res* 42:D980-985. 10.1093/nar/gkt1113

Li B, and Dewey CN. 2011. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics* 12:323. 10.1186/1471-2105-12-323

Linehan WM, Srinivasan R, and Schmidt LS. 2010. The genetic basis of kidney cancer: a metabolic disease. *Nat Rev Urol* 7:277-285. 10.1038/nrurol.2010.47

Papadopoulos P, Viennas E, Gkantouna V, Pavlidis C, Bartsakoulia M, Ioannou ZM, Ratbi I, Sefiani A, Tsaknakis J, Poulas K, Tzimas G, and Patrinos GP. 2014. Developments in FINdbase worldwide database for clinically relevant genomic variation allele frequencies. *Nucleic Acids Res* 42:D1020-1026. 10.1093/nar/gkt1125

Schuller A, Barry E, Jones RD, Henry R, Frigault MM, Beran G, Linsenmayer D, Mertens Hattersley M, Smith A, Wilson J, Cairo S, Deas O, Nicolle D, Adam A, Zinda M, Reimer C, Fawell S, Clark E, and D'Cruz C. 2015. The MET inhibitor AZD6094 (Savolitinib, HMPL-504) induces regression in papillary renal cell carcinoma patient derived xenograft models. *Clin Cancer Res*. 10.1158/1078-0432.ccr-14-2685

Smigielski EM, Sirotkin K, Ward M, and Sherry ST. 2000. dbSNP: a database of single nucleotide polymorphisms. *Nucleic Acids Res* 28:352-355.

Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, and Mesirov JP. 2005. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* 102:15545-15550. 10.1073/pnas.0506580102

Valk PJ, Verhaak RG, Beijen MA, Erpelinck CA, Barjesteh van Waalwijk van Doorn-Khosrovani S, Boer JM, Beverloo HB, Moorhouse MJ, van der Spek PJ, Lowenberg B, and Delwel R. 2004. Prognostically useful gene-expression profiles in acute myeloid leukemia. *N Engl J Med* 350:1617-1628. 10.1056/NEJMoa040465

van de Vijver MJ, He YD, van't Veer LJ, Dai H, Hart AA, Voskuil DW, Schreiber GJ, Peterse JL, Roberts C, Marton MJ, Parrish M, Atsma D, Witteveen A, Glas A, Delahaye L, van der Velde T, Bartelink H, Rodenhuis S, Rutgers ET, Friend SH, and Bernards R. 2002. A gene-expression signature as a predictor of survival in breast cancer. *N Engl J Med* 347:1999-2009. 10.1056/NEJMoa021967

Vogelstein B, Papadopoulos N, Velculescu VE, Zhou S, Diaz. LA, and Kinzler KW. 2013. Cancer Genome Landscapes. *Science* 339:1546-1558. 10.1126/science.1235122

Wu G, and Stein L. 2012. A network module-based method for identifying cancer prognostic signatures. *Genome Biol* 13:R112. 10.1186/gb-2012-13-12-r112

Yang Y, Han L, Yuan Y, Li J, Hei N, and Liang H. 2014. Gene co-expression network analysis reveals common system-level properties of prognostic genes across cancer types. *Nat Commun* 5:3231. 10.1038/ncomms4231

Zhang W, Ota T, Shridhar V, Chien J, Wu B, and Kuang R. 2013. Network-based survival analysis reveals subnetwork signatures for predicting outcomes of ovarian cancer treatment. *PLoS Comput Biol* 9:e1002975. 10.1371/journal.pcbi.1002975

Zhao Q, Shi X, Xie Y, Huang J, Shia B, and Ma S. 2015. Combining multidimensional genomic measurements for predicting cancer prognosis: observations from TCGA. *Brief Bioinform* 16:291-303. 10.1093/bib/bbu003

Zwiener I, Frisch B, and Binder H. 2014. Transforming RNA-Seq data to improve the performance of prognostic gene signatures. *PLoS One* 9:e85150. 10.1371/journal.pone.0085150

Legends

Figure 1

Distinct expression patterns of protective and harmful prognostic genes

(A) Raw gene p-value distributions from multivariate Cox models for a cancer with high number of expressionally prognostic genes (EPGs; LGG, left), and a cancer with low number of EPGs (STAD, right). Distributions for the other 14 cancers are displayed in Fig. S1. (B) Unsupervised hierarchical clustering (Pearson correlation distance metric) of patients using the inverse normal transformed expression values from the 100 most significant protective genes and 100 most significant harmful gene for LGG, left, and STAD, right. (C) Kaplan Meier plots comparing survival times for the two broad clusters of patients identified in B and logrank p-values for LGG, left, and STAD, right.

Figure 2

Overlaps of prognostic genes and gene sets

(A) Heatmap displaying the overlaps between cancers of the 100 most significant genes of each cancer. (B) Overlaps within cancers of the 100 most significantly enriched gene sets for protective genes, and the 100 most significantly enriched gene sets for harmful genes. (C,D) Overlaps between cancers of the 100 most significantly enriched gene sets for harmful genes (C) and protective genes (D). (E) Venn diagram showing the overlaps of the 100 harmful gene sets for LIHC, LUAD, and KIRP. (F) Venn diagram showing the overlaps of the 100 protective gene sets for COAD, GBM, and LUSC.

Figure 3

Clustering of cancers using gene Cox coefficients

Clustering of genes and cancers using the sigmoidal normalized Cox coefficients of a list of genes that had an FDR less than or equal to .05 for at least four cancers. Pearson correlation distance metric was used for both row and column clustering, and Cox coefficients were row scaled (z-score).

Figure 4

589 Clustering of cancers using gene programs

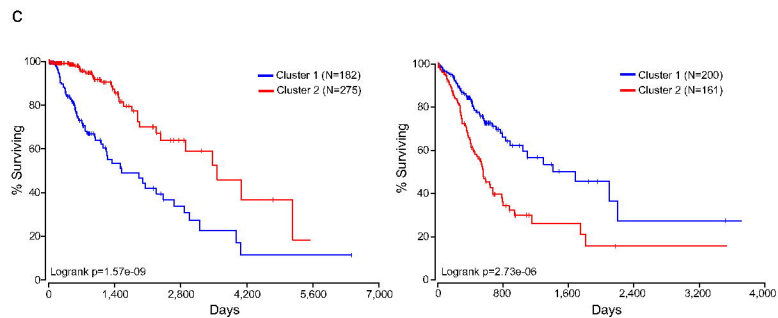
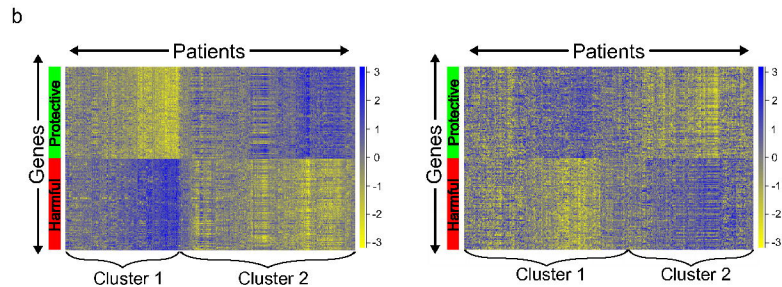
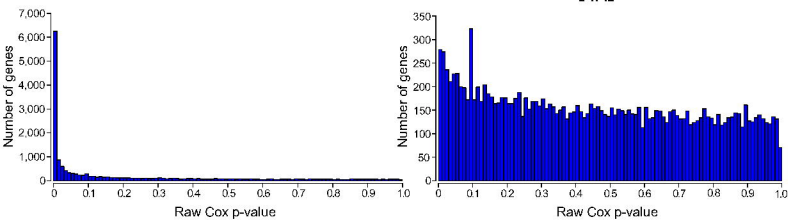
590 Using an established list of gene programs (see methods and Table S3), cancers and gene
 591 programs were clustered using the means of sigmoidal normalized Cox coefficients of the genes
 592 present in each program. Pearson correlation distance metric was used for both row and column
 593 clustering, and the average Cox coefficients were column scaled (z-score).

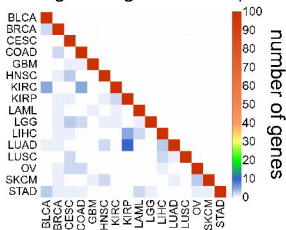
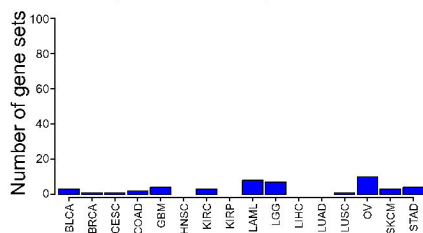
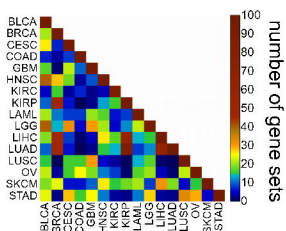
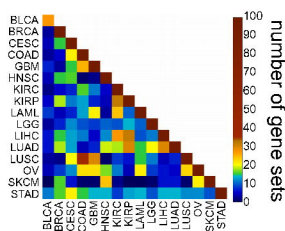
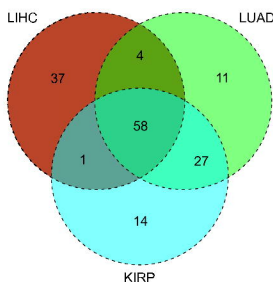
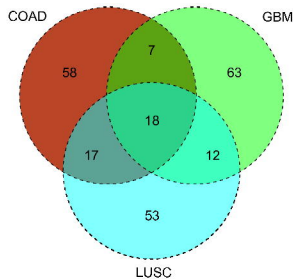
594 Table 1

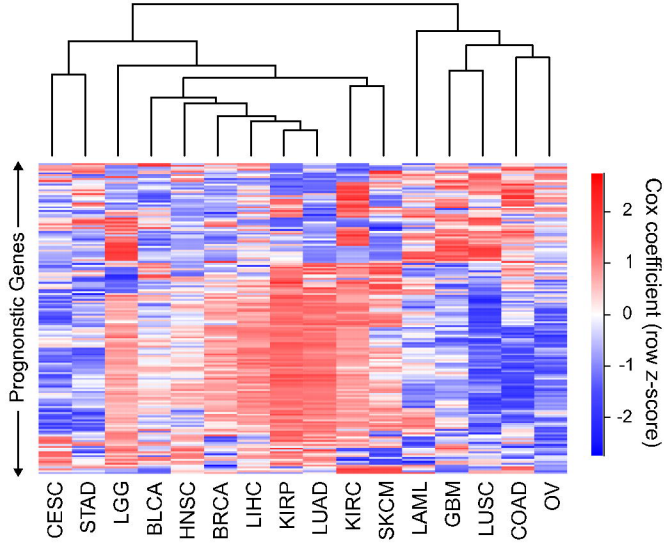
595 Characteristics of datasets and patients included in this study

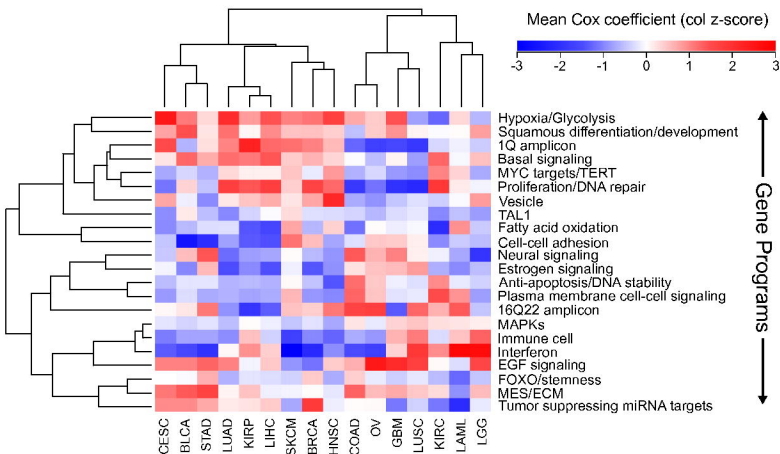
596 Events are the number of deaths in the data set. Age is the average age and is in years. Median
 597 survival is in days. The median survival for KIRP could not be calculated.

598



a Prognostic genes overlaps**b** Harmful:protective overlaps with cancers**c** Harmful gene set overlaps**d** Protective gene set overlaps**e****f**





Cancer	Patients	Median Survival	Events	Age at Diagnosis	Male/Female	RNASeqV2	Genes in study	FDR<=.05
BLCA	347	1064	139	67.9	254/93	YES	16385	532
BRCA	981	3669	116	58.4	10/971	YES	16649	30
CESC	259	3097	60	48.1	0/259	YES	16358	146
COAD	434	2475	89	66.6	232/202	YES	16414	0
GBM	152	406	119	59.8	98/54	YES	16833	0
HNSC	484	1671	190	61.2	353/131	YES	16652	45
KIRC	516	2386	167	60.7	335/181	YES	16677	5785
KIRP	247	NA	32	60.8	181/66	YES	16430	2415
LAML	149	577	92	54.7	80/69	YES	15255	4
LGG	457	2875	91	43.1	255/202	YES	16818	7186
LIHC	324	2116	105	60.1	217/107	YES	15855	2
LUAD	486	1379	146	65.3	224/262	YES	16784	1179
LUSC	471	1655	180	67.4	352/119	YES	16979	0
SKCM	427	2889	195	57.5	264/163	YES	16067	1548
OV	401	1321	224	59.6	0/401	NO	15748	0
STAD	361	874	130	65.2	238/123	NO	15560	0