



A Prognostic Signature for Lower Grade Gliomas Based on Expression of Long Non-Coding RNAs

Manjari Kiran¹ · Ajay Chatrath¹ · Xiwei Tang² · Daniel Macrae Keenan² · Anindya Dutta¹

Received: 5 July 2018 / Accepted: 25 October 2018 / Published online: 3 November 2018
© Springer Science+Business Media, LLC, part of Springer Nature 2018

Abstract

Diffuse low-grade and intermediate-grade gliomas (together known as lower grade gliomas, WHO grade II and III) develop in the supporting glial cells of brain and are the most common types of primary brain tumor. Despite a better prognosis for lower grade gliomas, 70% of patients undergo high-grade transformation within 10 years, stressing the importance of better prognosis. Long non-coding RNAs (lncRNAs) are gaining attention as potential biomarkers for cancer diagnosis and prognosis. We have developed a computational model, UVA8, for prognosis of lower grade gliomas by combining lncRNA expression, Cox regression, and L1-LASSO penalization. The model was trained on a subset of patients in TCGA. Patients in TCGA, as well as a completely independent validation set (CGGA) could be dichotomized based on their risk score, a linear combination of the level of each prognostic lncRNA weighted by its multivariable Cox regression coefficient. UVA8 is an independent predictor of survival and outperforms standard epidemiological approaches and previous published lncRNA-based predictors as a survival model. Guilt-by-association studies of the lncRNAs in UVA8, all of which predict good outcome, suggest they have a role in suppressing interferon-stimulated response and epithelial to mesenchymal transition. The expression levels of eight lncRNAs can be combined to produce a prognostic tool applicable to diverse populations of glioma patients. The 8 lncRNA (UVA8) based score can identify grade II and grade III glioma patients with poor outcome, and thus identify patients who should receive more aggressive therapy at the outset.

Keywords Long non-coding RNAs · Gliomas · Gene expression profiling · Prognosis

Abbreviations

lncRNA	Long non-coding RNAs	IFNG	Interferon gamma
WHO	World Health Organization	Cindex	Concordance index
LGG	Lower grade gliomas	AUC	Area under curve
GBM	Glioblastoma multiforme	ROC	Receiver operating characteristics
CNS	Central nervous system	UVA8	University of Virginia 8
TCGA	The Cancer Genome Atlas	L1-LASSO	L1 least absolute shrinkage and selection operator
CGGA	Chinese Glioma Genome Atlas	MGMT	O6-methylguanine DNA methyltransferase
HR	Hazard ratio	FPKM	Fragment per kilobase per million
PFS	Progression-free survival	GTF	Gene transfer format

Electronic supplementary material The online version of this article (<https://doi.org/10.1007/s12035-018-1416-y>) contains supplementary material, which is available to authorized users.

✉ Anindya Dutta
ad8q@virginia.edu

¹ Department of Biochemistry and Molecular Genetics, University of Virginia School of Medicine, Pinn Hall 1232, Charlottesville, VA 22908, USA

² Department of Statistics, University of Virginia, Charlottesville, VA 22904, USA

Introduction

Over the past decade, high-throughput RNA-seq technology discovered many novel transcriptional units, which were otherwise missed by probe design based transcriptome profiling. Among these transcriptional units were many long non-coding RNAs (lncRNA), which are transcripts longer than 200 bases with almost no protein-coding potential or open reading frames of < 50 amino acids. These lncRNAs are

numerous in cells [1], are highly regulated, and are more cell-type specific than protein-coding genes [2]. LncRNAs are involved in a broad spectrum of function and recent studies suggest they have specific roles in different diseases like cancer (reviewed in [3, 4]).

Gliomas are the most common form of primary malignant brain tumor, which originate in the supporting glial cells in the brain, including astrocytes, oligodendrocytes, and ependymal cells. Based on WHO 2016 grading system, gliomas are classified into lower grade and much aggressive high-grade gliomas. Grade I is mostly benign, whereas diffuse low-grade and intermediate-grade gliomas make up the WHO grade II and III lesions. Grade IV gliomas include secondary glioblastomas (derived from lower grade gliomas) and primary glioblastoma multiforme (GBM). Surgical resection of tumor is the most common initial treatment for gliomas followed by radiation therapy and chemotherapy, which can increase survival to 12 months [5, 6]. Molecular markers like 1p/19q co-deletion, MGMT promoter methylation, and mutation in IDH1 gene are strong predictors of survival for gliomas [7]. Lower grade gliomas have a better prognosis than high-grade gliomas. Despite a better prognosis for lower grade gliomas than the grade IV tumors, 70% of patients from the former group undergo high-grade transformation within 10 years.

LncRNAs are widely expressed in the central nervous system (CNS) and are involved in several pathways related to CNS development [8–13]. LncRNA BRN1B is one of the critical lncRNAs for brain development [13]. LncRNA Sox2OT plays an important role in determining neural fate [14]. Dysregulation of many lncRNAs like DGCR5, NRON, H19, and DISC2 have been associated with different CNS diseases [15–18]. Previous studies have shown that specific lncRNA expression patterns are also associated with different histological subtypes and grade in gliomas [19, 20]. For example, expression of MALAT1, POU3F3, and H19 are highly correlated with glioma malignancy. More recently, lncRNAs are also found to be of prognostic significance suggesting their role in glioma malignancies and as a potential therapeutic target and biomarker [19, 20]. Li et al. revealed three molecular subtypes of gliomas based on lncRNAs expression that has a strong correlation with patient's survival [21]. Furthermore, analysis on previously published microarray data has explored lncRNA-based signature as a prognostic marker in gliomas ([20, 22–25]).

Many studies have highlighted the power of gene expression profiles to predict tumor classification, patient outcome, and tumor response to therapy. Differentially expressed genes in cancer patients versus normal individuals are often the starting set to predict prognostic signature associated with survival [26–28]. This strategy suffers from false negatives and from the fact that differentially expressed genes might not be associated with differences in survival at all [29]. Another limitation of this method is the requirement of perfect

matched normal to identify differentially expressed genes. This creates a major hurdle in case of brain cancer where getting a perfect matched normal tissue is not trivial. While high-throughput technologies have facilitated the search of biomarkers through multivariate data analyses, there still remain challenges with respect to meaningful statistical and biological information. Firstly, most of the biological datasets suffer with multicollinearity: the influence of one gene on expression of other genes. Secondly, there are more features (genes) than observations (patients), which lead to overfitting by most of existing learning algorithms and results in poor performance of the model in prediction in an unseen testing dataset. Thus, a more robust approach is required to find genes as prognostic signature from a multi-dimensional multivariate gene expression data. Regression models like lasso, ridge, and elastic net are some widely used approaches to penalize the effect of multicollinearity and are well suited for constructing models when there are large numbers of features.

In the present study, we develop an lncRNA-based prognostic signature in combination with Cox regression and L1-LASSO regularization to model survival of grade II and grade III glioma patients. This is the first study that combined Cox and lasso regularization to select lncRNAs that can predict survival in glioma patients. After controlling for covariates associated with glioma survival (age, grade, IDH1 mutation status), we selected 8 lncRNAs UVA8, to calculate a risk score, which successfully divides patients into high-risk and low-risk groups in both TCGA (461 patients) and CGGA (274 patients) dataset. The risk score calculated by these eight lncRNAs is an independent and better prognostic marker for grade II and grade III glioma patient survival. The guilt-by-association analysis of lncRNAs in UVA8 indicated their role in suppressing interferon signaling pathway and epithelial to mesenchymal transition. Besides their use as a biomarker, these lncRNAs need to be studied in detail to determine how they affect patient outcome.

Materials and Methods

Patients and Samples

Aligned bam files and clinical information for 512 LGG patients (grade II and III) were retrieved from The Cancer Genome Atlas (TCGA) data portal <https://portal.gdc.cancer.gov/>. The study is performed on 461 patients for which both RNA-seq and survival information were available. Most samples in TCGA are collected from patients from the USA and also from other countries, including Canada, Russia, and Italy. This dataset being the largest and most updated glioma dataset is used as training dataset in the present study. The raw sequencing data for 274 glioma patients (175 grade II and III) from Chinese Glioma Genome Atlas (CGGA) as independent

cohort was downloaded using accession no. SRP027383 [30]. The survival information for these Chinese patients was downloaded from CGGA <http://www.cgga.org.cn/>. IDH1 mutation data for all the LGG patients were retrieved from Tier 3 TCGA data accessed from the Broad GDAC Firehose; <https://gdac.broadinstitute.org>.

RNA-Seq Data Quantification and Analysis

The most recent version of Gencode (GENCODE v 26) GTF file available at the time of this study was used for the gene quantification [31]. Gene abundance in FPKM (fragment per kilobase per million) was obtained for 58,219 genes with 15,787 genes annotated as lncRNA in GENCODE v26 using Stringtie v1.3.3 [32]. Out of 15,787 lncRNAs, 1289 lncRNAs with a median expression of 1 FPKM in 512 LGG patients were finally considered for the survival model.

Survival Model Selection Process

The gene expression data for lncRNAs was Z-score transformed to avoid systematic error across different experiments. We first randomly selected 60% of TCGA patients for training set and remaining 40% of TCGA patients for testing set. Since, clinical information like age, gender, tumor grade, or IDH mutation status can have an effect on survival (Fig. S1), we assessed the prognostic potential of each lncRNA by multivariate Cox regression controlling the effects from these other variables. We used FDR corrected *p* value cutoff of 0.05 obtained after log-likelihood test comparing restricted (age, gender, tumor grade, and IDH mutation status) with unrestricted (lncRNA expression, age, gender, tumor grade, and IDH mutation status) model to identify the significant association of an lncRNA with survival. We used Cox-proportional hazards model based on L1-penalized (LASSO) estimation to select the best model comprising a subset of prognostic lncRNA [33–35]. We used LASSO because it is suited for constructing models when there is a large number of correlated covariates [34].

Risk Score Calculation

Risk score for each patient was established by including each of the selected genes weighted by their estimated regression coefficients in the multivariable Cox regression analysis as discussed in previous studies [36, 37].

UVA8 risk score = $(-0.378 \times \text{expression value of RP11-266 K4.14}) + (-0.301 \times \text{expression value of FLJ37035}) + (-0.280 \times \text{expression value of LINC01561}) + (-0.368 \times \text{expression value of RP11-118 K6.3}) + (-0.369 \times \text{expression value of DGCR9}) + (-0.299 \times \text{expression value of RP11-142A22.3}) + (-0.434 \times \text{expression value of LINC00641}) + (-0.543 \times \text{expression value of RP11-96H19.1})$.

Coefficients are median Cox coefficient (after lasso selection and multivariate Cox regression) for each of the eight lncRNAs from the successful models (models which can stratify patients in testing set).

Statistical Analysis

R package glmnet was used to perform L1-penalized cox regression (L1-least absolute shrinkage and selection operator) [38]. R package survival and survminer were used for survival data analysis and generating Kaplan–Meier plots. Different survival models were compared by time-dependent concordance index (Cindex) [39]. Cindex is the most commonly used performance measure for survival models, which calculates the fraction of pairs whose predicted survival time is correctly ordered. R package pec::cindex is used to calculate time-dependent cindex [40].

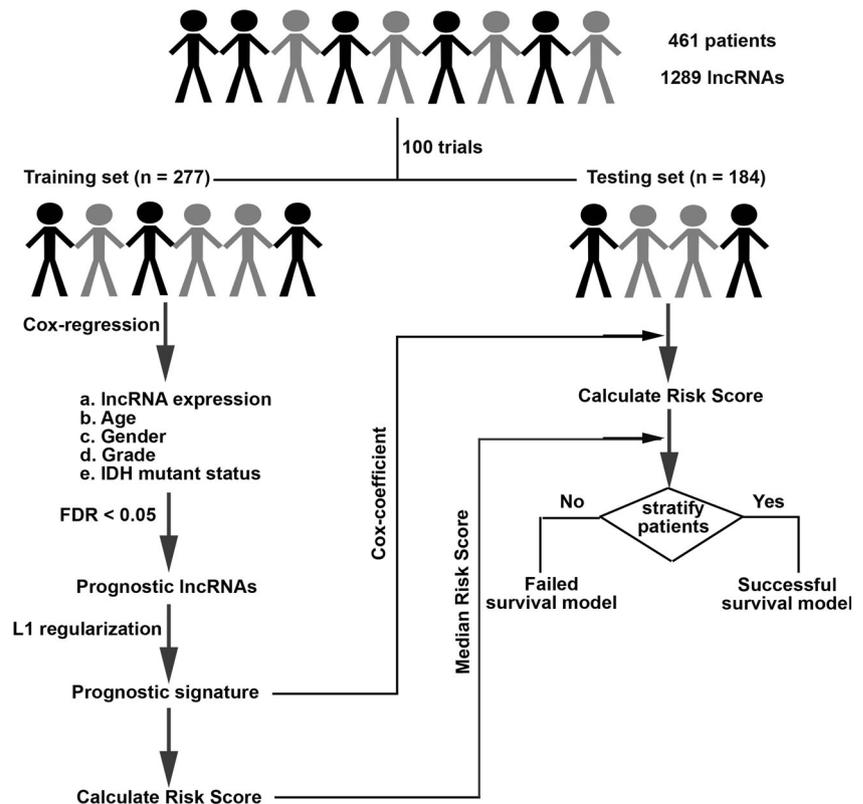
Results

Building the lncRNA-based Survival Model

We developed an lncRNA-based survival model for gliomas through the following steps (Fig. 1).

- 1) We first randomly selected 60% ($n = 277$) of the patients from TCGA as training set and reserved the remaining 40% ($n = 184$) of patients as testing set. The results remain similar with 70% patients in training and 30% in testing set (Fig. S3 A).
- 2) Cox multivariate regression was carried out in the training set on 1289 lncRNA controlling for effects from other covariates like age, gender, tumor grade, and IDH1 mutation status.
- 3) lncRNAs significantly associated with survival after likelihood ratio test (FDR, $p < 0.05$) were retained for selecting lncRNAs by lasso regularization.
- 4) After lasso regularization and lncRNA selection, a risk score formula was established by including selected lncRNAs weighted by their estimated regression coefficients in the multivariable Cox regression analysis. Risk score = $\sum_{i=1}^n \beta_i * x_i$ (where, β is coefficient and x is expression level of lncRNA i)
- 5) Patients were classified into high-risk and low-risk group by using the median risk score as the cutoff in the training set. The coefficient for each lncRNA and cutoff of risk score obtained from training set was used to calculate risk score and stratify patients into two groups in testing set.

Fig. 1 Flowchart showing steps involved in identification of lncRNA-based prognostic signature



- 6) Survival differences between the low-risk and high-risk groups in the training and testing sets were assessed by the Kaplan–Meier estimate and compared using the log-rank test.

Steps 1–6 were repeated 100 times to obtain up to 100 different lncRNA subsets (models). Only those models that separated patients in the testing set such that those with low-risk score had significantly better survival than those with high-risk score were considered as successful models and retained.

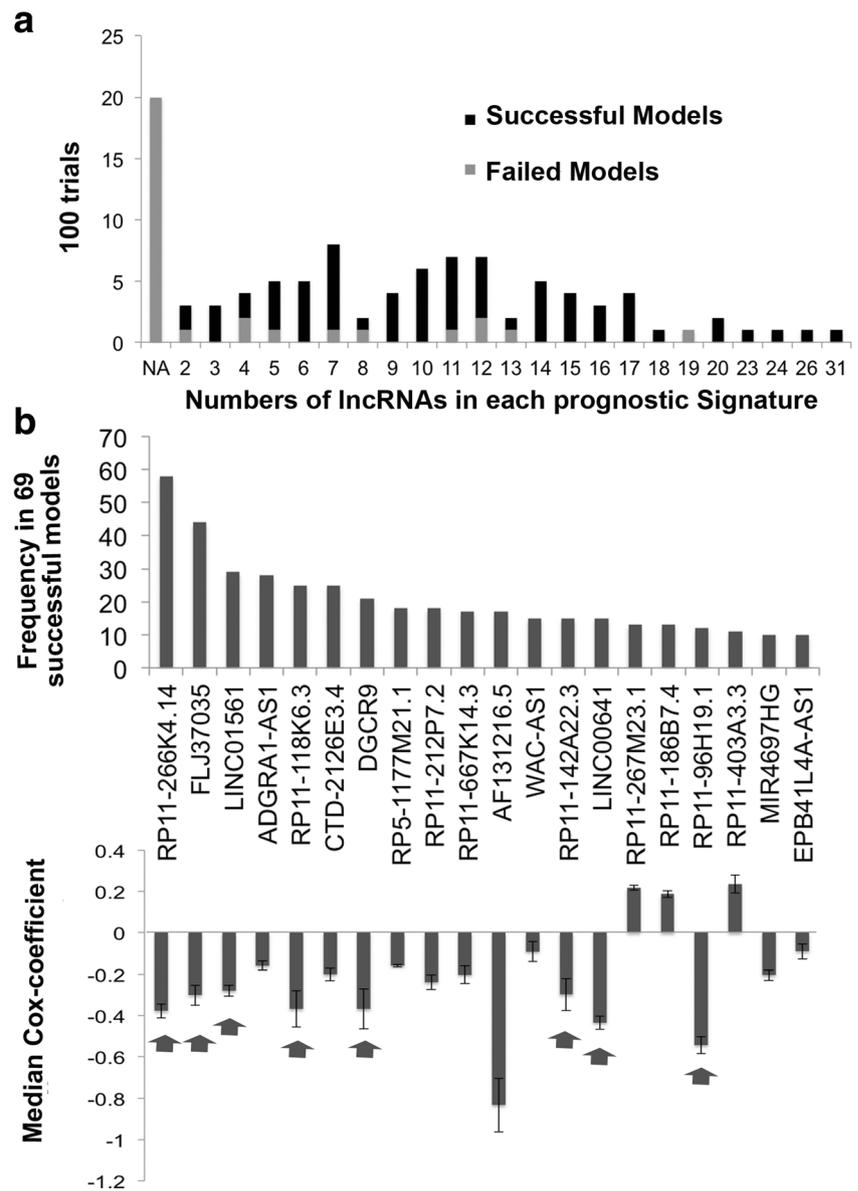
The result obtained from one such survival model is shown in Fig. S2. In ~20% of the trials the multivariate cox-regression and lasso regularization in the training set did not select any lncRNAs significantly associated with survival (NA in Fig. 2a). The remaining 80% of the survival models contained different numbers of lncRNAs (*x*-axis of Fig. 2a) that significantly stratify patients into low- and high-risk groups in training set (Fig. 2a). Among these 80% of survival models, 86% also significantly separated patients into high risk and low risk in the testing set and are referred to as successful survival models. In order to create a robust survival model we sorted the lncRNAs based on the number of times an lncRNA was selected by successful survival models (Fig. 2b). Out of 167 total prognostic lncRNA in 69 successful survival models, we first ranked lncRNAs based on number of

times a given RNA was selected by successful models and then from the top 20 selected 8 lncRNAs with the highest median Cox coefficient (absolute value > 0.2) and least variance in the successful models in the testing set (absolute value < 0.10). Seven out of these 8 lncRNAs were also selected after 70–30% split of training and testing patients (Fig. S3A), after 1000 trials instead of 100 (Fig. S3B), and all 8 lncRNAs were selected when we used elastic net, instead of lasso, for regularization and lncRNA selection (Fig. S3C) suggesting the prognostic importance of these 8 lncRNAs in gliomas. For brevity, this set of eight lncRNAs as a prognostic signature of gliomas will be referred to as UVA8 (University of Virginia 8) in the manuscript. AF131216.5 is associated with high median Cox coefficient (−0.83) but with high variance (0.35, greater than defined cutoff of 0.10). Despite its high variance, we tested whether including AF131216.5 in the final model will improve the performance of UVA9 on our training (TCGA) and testing dataset (CGGA). The result is described below.

UVA8 is Predictive of Survival in Training and Independent Validation Set

We assessed the predictive power of UVA8 by comparing overall survival of low- and high-risk patients in the entire TCGA dataset stratified based on median risk score obtained by UVA8 (risk score calculation discussed in methods).

Fig. 2 Selection of lncRNAs with best predictors of outcome. **a** Barplot showing number of lncRNAs that predicted outcome in the training set in 100 trials. The successful models were those that also predicted outcome in the testing set. NA: no lncRNA predicted outcome in training set. **b** Barplot showing number of times each of the top 20 lncRNAs (out of 167) were present in successful survival models (significant in testing set). The lower panel shows median Cox coefficient (after lasso penalization and multivariate Cox regression) and the variance of the Cox coefficient for each of the above 20 lncRNAs from the successful models where they were selected. The arrow points towards lncRNAs selected for UVA8



Patients in the low-risk group showed longer overall survival than the high-risk group in TCGA dataset (Fig. 3a, median OS 741.5 vs. 639 days; $P = 3.1 \times 10^{-15}$, HR = 5.8). The risk scores of the patients in the TCGA dataset range from -4 to 4 with median risk score of -0.023 (Fig. 3b, top panel). Moreover, there are more patients alive in the low-risk group than in the high-risk group (Fig. 3b, middle panel). Interestingly, expression levels of all lncRNA in UVA8 are high- in low-risk patients than in high-risk patients indicating these lncRNAs as favorable prognostic genes (Fig. 3b, bottom panel). These findings were further validated in an independent validation dataset comprising of 274 patients obtained from CGGA. Using the same median coefficient of UVA8 obtained from the successful survival models in TCGA, patients showed longer overall survival in low-risk than in high-risk group in CGGA (Fig. 3c, median OS = 1120.5 vs. 587 days; $P =$

0.0017, HR = 1.68). Moreover, low-risk group in CGGA has also longer progression-free survival (PFS) than the high-risk group (Fig. 3d, median PFS 597.5 vs. 411.5 days; $P = 0.00088$, HR = 1.70). Thus, UVA8 can predict survival in both training and independent validation set.

Since, 32% of patients in CGGA are in grade IV, the difference in overall survival could be due to over-representation of grade IV patients in high-risk group. However, even when only lower grade gliomas (grade II and III) were separately examined we found significantly longer survival for low-risk versus high-risk patients (Fig. S4A). UVA8 fails to cluster grade IV patients from CGGA into two distinct groups highlighting the specificity of signature for lower grade gliomas (Fig. S4B). We also assessed the predictive capability of UVA9 (UVA8+ AF131216.5) on TCGA and CGGA and noticed almost no improvement in TCGA ($P = 3 \times 10^{-15}$, HR =

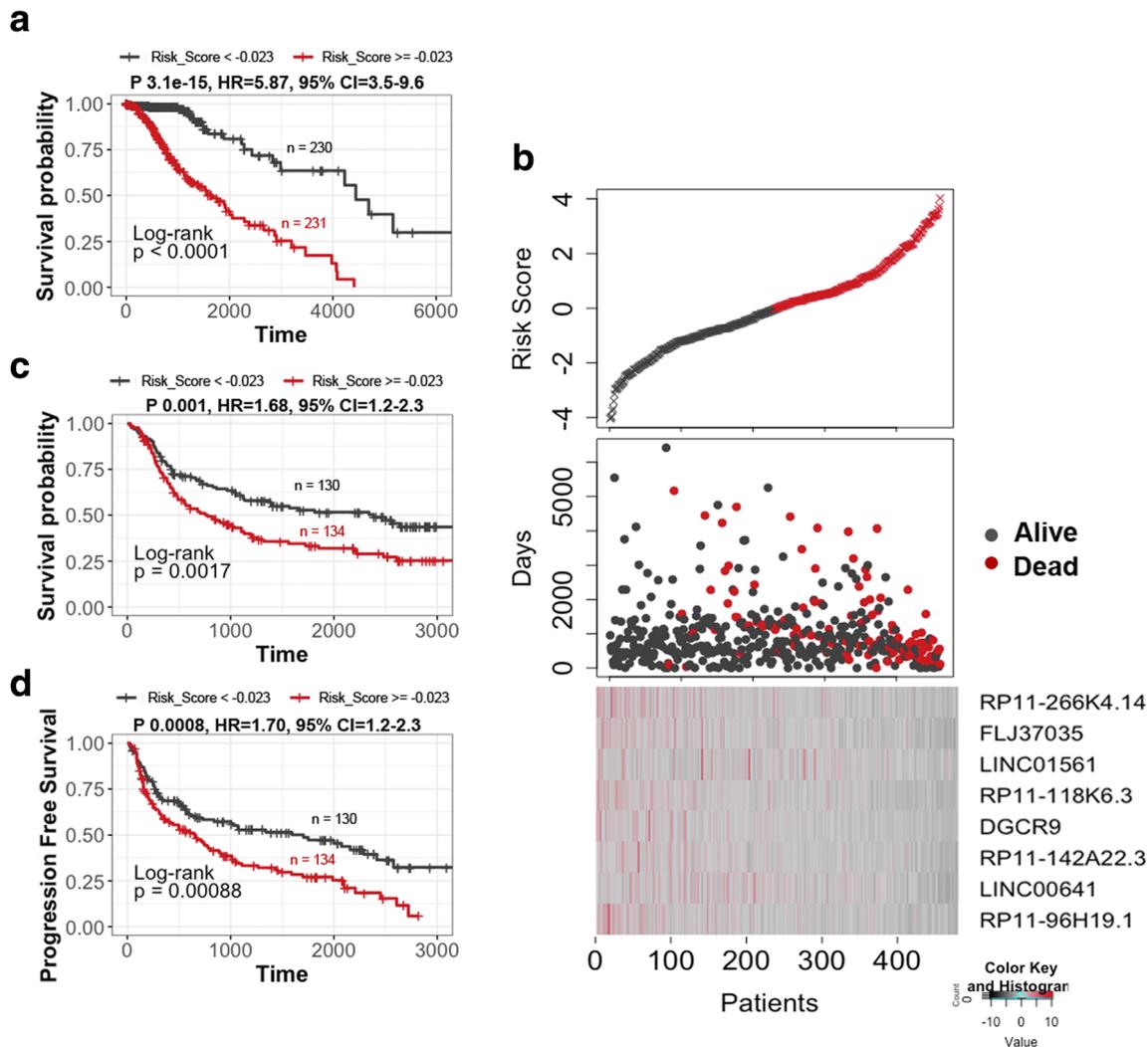


Fig. 3 Survival analysis of the patients divided by the prognostic lncRNAs in two data sets. **a** Patients in the entire TCGA dataset with risk score greater than median score of -0.023 show poor survival compared with patients with risk score less than median risk score. **b** Upper panel: plot showing patients sorted based on UVA8 risk score with black representing patient with risk score below median and red showing those with risk score above median. Middle panel: Number of days of survival indicated on *Y*-axis of patients sorted on the *X*-axis based

on the risk scores in the top panel and alive/dead status indicated by color. Bottom panel: *z*-score transformed expression value of lncRNAs in UVA8 show higher expression in patients with low risk score. **c** Kaplan–Meier plot of overall survival of patients in CGGA dataset with risk score greater than (red) or less than (black) median risk score of TCGA dataset. **d** Kaplan–Meier plot for progression-free survival in CGGA dataset showed poor survival for patients with high-risk score. Rest as in (c)

5.623.95% CI=3.47–9.11) and marginal improvement in CGGA (P 0.02, HR = 1.74, 95% CI = 1.10–2.77) compared to UVA8 (TCGA P 3.1e-15, HR = 5.87, 95% CI = 3.5–9.6, CGGA P 0.03, HR = 1.66, 95% CI = 1.05–2.64) (Fig. S5).

8 lncRNA-based Risk Score is an Independent Predictor of Survival

Lower grade gliomas have poorer outcomes in older patients, in tumors of higher grade and tumors with wild-type IDH1 status (Fig. S1). Interestingly, the risk score derived from UVA8 is higher in patients older than 40 years, patients in grade III vs. grade II and patients harboring wild-type IDH1

gene (Fig. S6). It was therefore important to determine whether UVA8-derived risk score is an independent predictor of survival. We divided the patients into younger (age < 40) and older (age \geq 40) groups and found that risk score can still stratify the patients into low risk and high risk in both groups (Fig. 4a). Similarly, UVA8-based risk score can still separate the patients into low and high-risk groups in grade II or grade III gliomas (Fig. 4b). Although, IDH mutation status is a widely used prognostic and predictive biomarker, the UVA8-based risk score can also separate patients into two risk groups in patients presorted based on IDH mutation status (Fig. 4c). UVA8-derived risk score can also stratify patients into two risk groups among male and female patients (Fig. 4d).

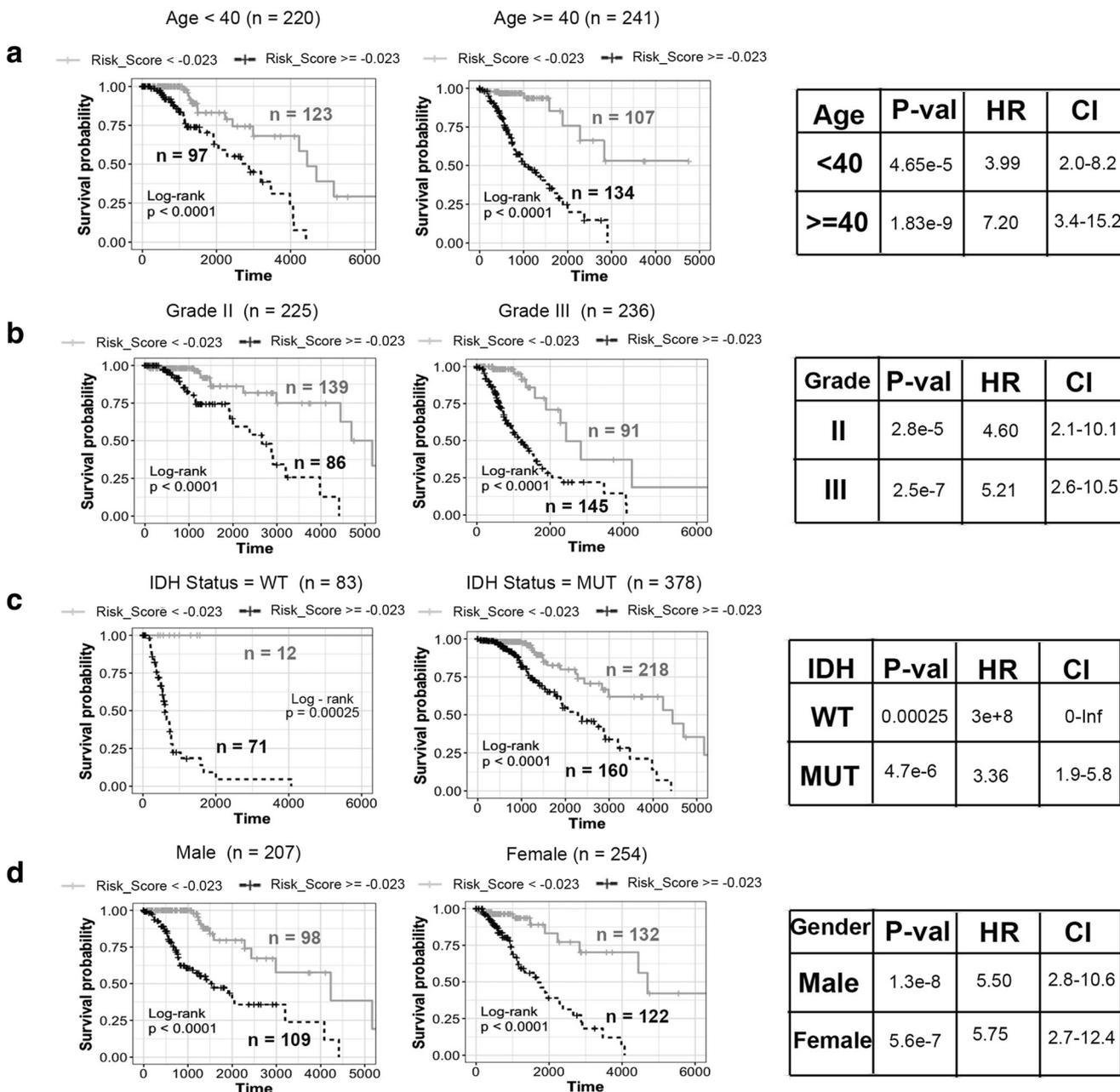


Fig. 4 Stratification analysis by different clinical variables. Kaplan–Meier curve analysis of overall survival in high- and low-risk groups for **a** younger (age < 40) and older patients (age ≥ 40). **b** Grade II and grade III patients **c** IDH mutation status as WT and mutation (MUT)

patients **d** male and female patients. Black-dashed line: patients with high-risk score, gray solid line: patients with low-risk score. The tables on the right show log-rank, *p* value, hazard ratio, and 95% confidence interval for each Kaplan–Meier plot

Conversely, we tested whether these standard clinically used parameters, age, gender, grade, and IDH mutation status, continue to independently stratify patients even after they have been presorted into two groups by UVA8 risk score (Fig. S7). In patients with high UVA8 risk score, age, grade, and IDH mutations status can further separate the patients into two groups of better or worse outcome. In contrast, in patients with low UVA8 risk scores, none of the clinical factors could further stratify patients into two different survival groups with a *p* value < 0.05 (Fig. S7). Consistent with the previous

observation (Fig. S1), gender is ineffective in stratifying patients into two categories within patients with high- or low-risk score.

UVA8 is a Better Predictor of Glioma Patients’ Survival

We assessed the accuracy of UVA8 in prediction of survival by comparing its time-dependent AUC in a ROC curve with that of other clinical characteristics. For each prognostic factor (e.g., UVA8, IDH status, etc.), we varied the cutoff so as to

vary the false positive rate for 5-year survival prediction from 0 to 1. For each cutoff, the corresponding true positive rate for 5-year survival was calculated (Fig. 5a). Comparing the AUC for these ROC curves suggested that UVA8 performs best in predicting survival of the glioma patients compared to the other criteria. This calculation was extended to predict survival of other durations (1–16 years) and the AUC plotted for each predictor (Fig. 5b). UVA8 can predict survival better for all durations, particularly at the very early years after diagnosis when the prediction is worse for most of the predictors. Since, gender is not associated with glioma patients' survival (Fig. S1), the prediction of outcome was no better

than random guess (AUC = 0.5) (Fig. 5a, b). We employed Cox multivariable probability hazard model to identify the impact of UVA8 and different clinicopathological characteristics in estimating hazard (Fig. 5c). UVA8 is most significantly correlated with the survival information ($p = 1.4e-07$) and shows highest hazard ratio (HR = 4), indicating that the risk score performs better than any other currently used approaches for prognosis. Here, the hazard ratio of UVA8 is calculated by dichotomizing the risk score of > -0.023 (median risk score from TCGA) to 1 and < -0.023 to 0 to compare the hazard rates of high-risk versus low-risk patients. The hazard ratio of the eight lncRNAs individually and combined

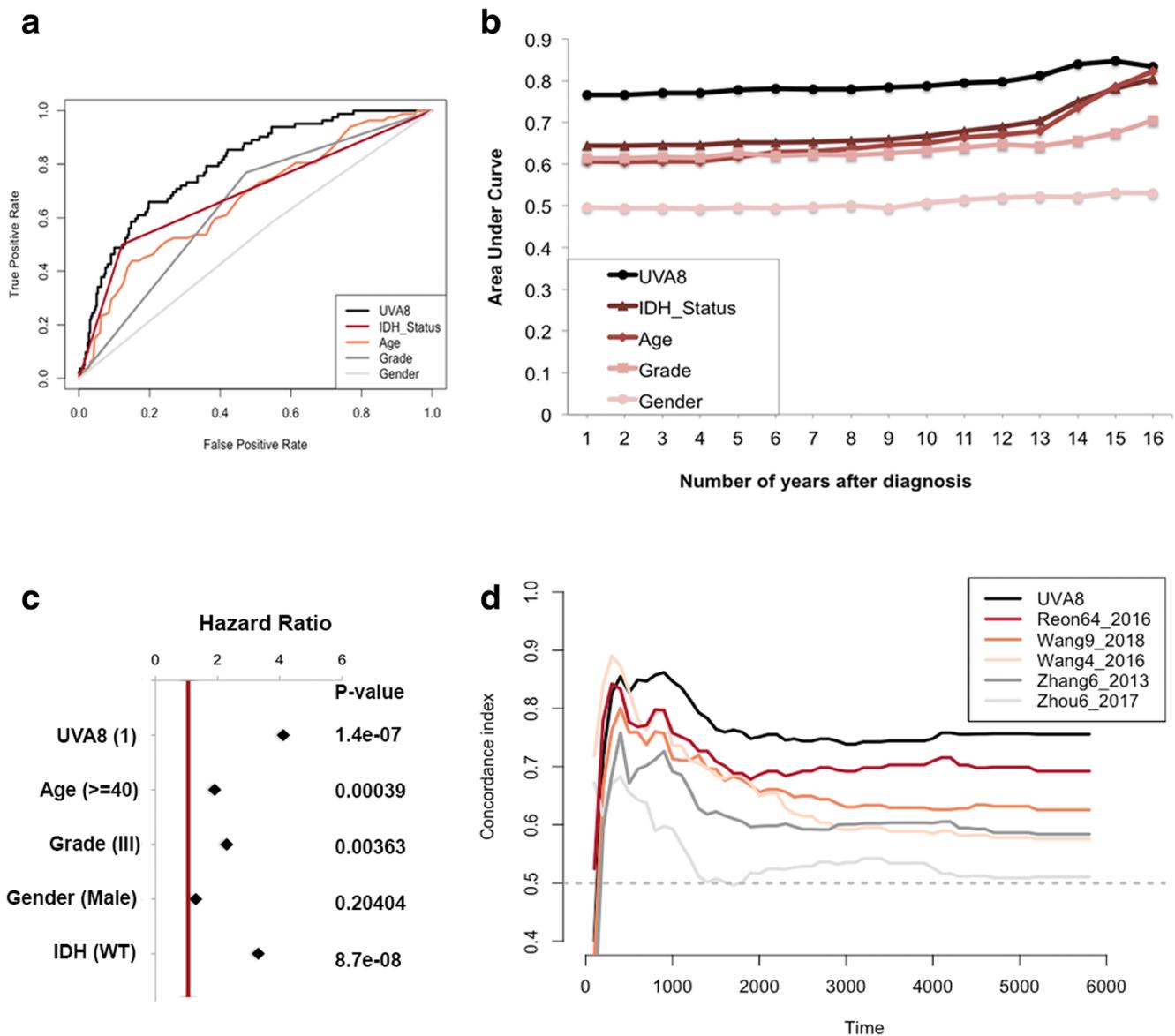


Fig. 5 Performance evaluation of the 8 lncRNA-based risk score. **a** Receiver operating characteristic curve for 5-year survival shows UVA8 has better area under curve compared with other predictors. **b** Area under curve plotted for different durations of survival for eight lncRNA-based risk score, tumor grade, age, IDH mutation status, and gender of patients

in TCGA cohort. **c** Cox multivariate regression with clinical information and risk score calculated from UVA8 for survival in TCGA cohort. **d** Concordance index showing measure of concordance of predictor with survival of patients in TCGA

as risk score is tabulated in Supplementary Table S1. The UVA8 risk score is associated with more hazard (HR = 2) than any of the individual lncRNA supporting the importance of a combinatorial signature than an individual RNA for predicting survival. The hazard ratio of UVA8 in Supplementary Table S1 is different from that in Fig. 5c because in the former the hazard ratio is calculated with the risk score as a continuous variable.

We then sought to compare the performance of UVA8-based survival model with published lncRNA-based survival models by calculating Cindex (as discussed in “Materials and Methods”) for TCGA dataset for each of the models. We first

calculated risk score for each patient by considering the expression level of the prognostic lncRNAs in each model weighted by their estimated regression coefficients retrieved from the respective studies (Supplementary Table S2). The patients were ordered based on their actual survival at a given time after diagnosis and based on their risk score in each model. The concordance of the two orders is measured in pairwise comparisons of the patients to calculate a single time-dependent concordance index for the model that is being evaluated. This is repeated for different survival times with an interval of 100 days. The concordance index for each survival time for UVA8 and all published models is tabulated as

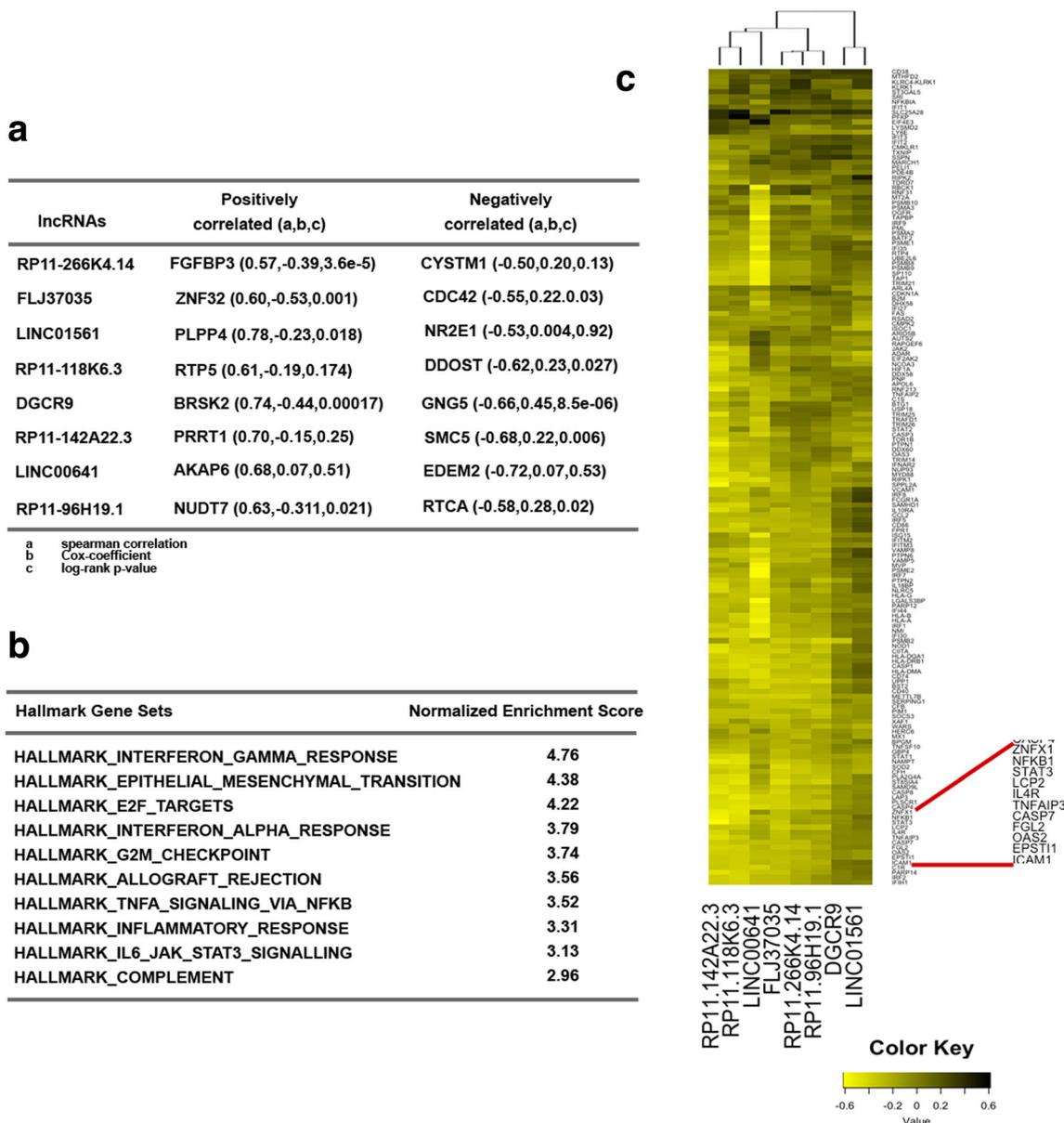


Fig. 6 Guilt-by-association analysis of the 8 lncRNAs in UVA8. **a** Correlation and Cox regression coefficient for the mRNAs that are most correlated (positive and negatively) with each of the lncRNAs in UVA8. a, b, and c defined below the table. **b** List of pathways that are most

enriched in protein-coding genes that are negatively correlated with the UVA8 lncRNAs. **c** Heatmap showing correlation of different genes in the interferon gamma response gene set (rows) to the lncRNAs in UVA8 (columns)

Supplementary Table S3. UVA8 outperforms all existing lncRNA-based survival models at different times after diagnosis (Fig. 5d). As expected, prognostic signatures that were specific to GBMs (Zhang6_2013 and Zhou6_2017) show poor concordance index when used to predict survival of lower grade glioma patients.

Interferon Signaling is the Most Enriched Pathway in Guilt-by-Association with UVA8

Although many lncRNAs have been identified there has been very little functional annotation of the RNAs. We therefore applied guilt-by-association to infer functions of the lncRNAs associated with survival in UVA8. First, we interrogated whether protein-coding genes most correlated with an lncRNA in TCGA glioma cohort are themselves predictive of outcome. All the lncRNAs in UVA8 are associated with a negative Cox coefficient (protective). Of the eight mRNAs most correlated positively with these eight lncRNAs, five also have a negative Cox coefficient with a significant p value. Conversely, of the eight mRNAs most anti-correlated with these lncRNAs, five have a positive Cox coefficient with a significant p value (Fig. 6a). This result is consistent with the expectation that the expression of these protective lncRNAs will be positively correlated with expression of protective mRNAs and negatively correlated with the expression of harmful mRNAs.

GSEA analysis on protein-coding genes pre-ranked from most positively correlated to most negatively correlated to the lncRNA revealed several common pathways co-regulated with each of the eight lncRNAs (Fig. 6b). Interestingly, among the mRNAs that are negatively correlated with the lncRNAs, genes involved in immune and inflammatory response (IFNG, IFNA, allograft rejection, NF κ B inflammatory response, and JAK-STAT pathway) are highly enriched. Similarly, genes involved in epithelial to mesenchymal transition and cell-cycle progressions are also most enriched. These gene set enrichments suggest a conventional tumor suppressor phenotype associated with these eight lncRNAs.

Many of the mRNAs are common in the IFNG, IFNA, allograft rejection, NF κ B inflammatory response, and JAK-STAT gene sets. The genes upregulated in response to IFNG are mostly negatively correlated to lncRNAs in UVA8. To visualize this, the correlation coefficients were plotted for each lncRNA (columns) with individual mRNAs in the IFNG response pathway (rows) (Fig. 6c). Out of eight, six lncRNAs (RP11-266K4.14, FLJ37035, RP11-118K6.3, RP11-142A22.3, LINC00641, and RP11-96H19.1) are clustered together because they are more negatively correlated with genes of interferon gamma response pathway (Fig. 6c).

We found both NF κ B and STAT3 genes as highly negatively correlated with the expression of the protective lncRNAs in UVA8. Genes involved in epithelial to mesenchymal transition and encoding cell cycle-related targets of E2F transcription factors and involved in G2/M checkpoints were also negatively correlated with UVA8 expression. On the other hand, genes that are downregulated upon activation of the oncogenes KRAS are positively correlated with the expression of the protective lncRNAs of UVA8.

eRNA (enhancer RNA) are another class of long non-coding RNAs which are 50–2000 bases long, unspliced, and non-polyA non-coding RNA expressed from enhancers involved in the activation of distantly located genes [41]. In order to check whether these lncRNAs can possibly act as eRNAs, we also checked the distance between lncRNAs and their correlated genes and found that these lncRNAs are correlated to several genes located in different location of genome suggesting a trans-regulation by these lncRNAs (data not shown). More experimental studies are required in future to decipher the role of these lncRNAs in regulating these genes and whether this regulation explains the effect of the lncRNAs on glioma tumor progression.

In order to investigate whether somatic mutations in these lncRNAs might account for the prognostic ability of the model, we used the somatic variant calls from The Cancer Genome Atlas (“Scalable Open Science Approach for Mutation Calling of Tumor Exomes Using Multiple Genomic Pipelines”). Based on this dataset, there seems to be somatic mutations in these lncRNAs in other cancers (most prominently DGCR9), but no somatic mutations were found in the eight lncRNAs in the TCGA lower grade glioma patients. Therefore, we feel that the predictive ability of our model is due to differences in the expression of these lncRNAs and not acquired somatic mutations in these lncRNAs. To see if copy number variation was the mechanism driving the differences in the expression of these eight lncRNAs, we tested whether the copy number of each gene was predictive of survival. We then correlated the expression of the gene to the copy number. Thus, copy number variation by itself may be predictive for two of the lncRNAs (FLJ37035 and LINC01561), but oddly for the second, there is no correlation between the CNV and level of expression. Thus, CNV is not the explanation for the expression differences of the lncRNAs, and is not a better predictor for prognosis (Supplementary Table S4).

Discussion

Gene expression profile reflects the underlying biological processes of disease. Cox regression is a widely used

approach to decipher correlation between gene expression profile and patient outcome. Previous analyses on microarray data explored protein-coding genes that could predict the prognosis of gliomas, particularly focusing on high-grade GBMs. lncRNAs are a class of RNA which can serve as a better prognostic marker than protein-coding mRNAs because they are numerous and cell-type specific [2, 3]. Additionally, since lncRNAs do not encode protein, they are the ultimate effectors, and their expression levels more accurately predict the levels of their activity. Recent studies have detected tumor-specific lncRNAs in exosomes, apoptotic bodies, and microparticles highlighting another advantage of considering lncRNAs in tumors, because they are expected to appear as fluid-based markers for the diagnosis of different cancers [42–44]. Among six published lncRNA-based prognostic signatures for gliomas, two are for predicting outcome in GBMs and one specifically for anaplastic gliomas. Wang et al. and Chen et al. have shown that a set of only four lncRNAs could predict survival in gliomas [23, 25]. However, the sequence of one of the lncRNAs in Chen et al., CR613436, was removed by the submitter on NCBI. Recently, the role of immune-related genes in glioma malignancies is gaining attention leading to the discovery of immune-related lncRNA-based prognostic markers for GBMs and anaplastic gliomas [22, 45]. Remarkably, there is no overlap between the prognostic lncRNAs identified in the aforementioned studies. Moreover, these studies are based on microarray data raising concerns particular to hybridization-based approaches, including reliance on current knowledge of expressed genes, problems of cross-hybridization, and cross-experiment comparison. Another issue is that association of lncRNAs with survival using Cox regression was sometimes carried out without controlling for any dependent variables and without penalizing for the effect of large number of variables.

In the present study, we have used an approach to screen lncRNAs from high-dimensional TCGA RNA-seq data, which is one of the largest and the most updated data for lower grade gliomas. After controlling for effects like age, grade, gender, and IDH mutation status, we applied regularization to penalize the effect of many dependent variables and select the lncRNAs based on 100 trials. We showed the robustness of eight lncRNA-based predictors in a completely independent cohort of Chinese glioma patients. The lncRNA prognostic signature identified in the present study, UVA8, is an independent predictor of survival in TCGA glioma patients. Since UVA8 is also a better predictor than the few patient and molecular characteristics currently used for prognosis in the clinic, a simple RNA quantification will aid the physician to decide whether to adopt more aggressive therapy at the outset.

The protective lncRNAs that constitute UVA8 are negatively correlated with protein-coding genes involved in interferon gamma and inflammatory response highlighting the role of immune-response genes in glioma progression. Except LINC01561, all seven lncRNAs (RP11-266K4.14, FLJ37035, RP11-118K6.3, DGCR9, RP11-142A22.3, LINC00641, and RP11-96H19.1) are negatively correlated to most of the protein-coding genes, which are upregulated in response to interferon gamma/alpha, genes regulated by NF κ B in response to TNF, inflammatory response, and genes upregulated by IL6 via STAT3. This suggests that an active immune reaction perhaps in response to cytokines secreted from tumor and immune cells is predictive of poor outcome in gliomas. NF κ B and JAK/STAT pathways are known to be aberrantly upregulated in GBMs. The level of NF κ B increases as the tumors progress in astrocytic tumors [46, 47] and STAT3 is constitutively active in GBMs [48, 49]. Immune-related pathways are also known to be involved in glioma tumor cell proliferation [50], survival [45], invasion [51], and chemoresistance [52]. In addition, epithelial-mesenchymal transition (associated with invasion) and active cell proliferation are suppressed if UVA8 lncRNAs are high, and this leads to better outcome, consistent with our understanding of how invasion and cell proliferation negatively impact outcome. On the other hand, genes that were positively correlated with the expression of UVA8 are enriched in genes that are down regulated by activation of the oncogene KRAS.

There are reports of the same lncRNA being predictive of outcome in the same manner in multiple tumor types. For example, DRAIC expression predicts good outcome in gliomas, melanomas, and cancers of the prostate, stomach, liver, kidney, and lung [53]. In contrast, expression of LINC00152/CYTOR is predictive of poor outcome in gliomas, and cancers of the head and neck, lung, kidney, liver, and pancreas (our unpublished work). Such observations are particularly exciting because they imply that the lncRNA has an important role in tumor biology that transcends tumor types, and these RNAs should be prioritized for cell- and molecular-biology studies to discern their function. It will thus be very interesting to explore whether any of the lncRNAs of UVA8 will be protective in other tumor types. Finally, future studies will address whether structural variation, copy number variations, and sequence polymorphism of these lncRNAs contribute to the prognostic outcome. We are excited that UVA8 was also predictive of outcome in a completely different tumor cohort (CGGA) from a patient population that is from an entirely different geographical location with attendant differences in environment and population genotypes. It will be interesting to see if UVA8 is equally predictive of outcome in other patient populations from other parts of the world.

Acknowledgments We thank Dr. Stefan Bekiranov, Dr. William Pearson, and Dutta lab members for helpful discussions. M.K. is supported by a DOD award PC151085.

Funding Information The work was supported by a V foundation award D2018-002 and R01 AR067712 from NIAMS.

References

- Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, Guernec G, Martin D et al (2012) The GENCODE v7 catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and expression. *Genome Res* 22:1775–1789. <https://doi.org/10.1101/gr.132159.111>
- Cabili M, Trapnell C, Goff L et al (2011) Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses. *Genes Dev* 25:1915–1927. <https://doi.org/10.1101/gad.17446611>
- Huarte M (2015) The emerging role of lncRNAs in cancer. *Nat Med* 21:1253–1261. <https://doi.org/10.1038/nm.3981>
- Schmitt AM, Chang HY (2016) Long noncoding RNAs in cancer pathways. *Cancer Cell* 29:452–463. <https://doi.org/10.1016/j.ccell.2016.03.010>
- Stupp R, Mason WP, van den Bent MJ, Weller M, Fisher B, Taphoorn MJ, Belanger K, Brandes AA et al (2005) Radiotherapy plus concomitant and adjuvant temozolomide for glioblastoma. *N Engl J Med* 352:987–996. <https://doi.org/10.1056/NEJMoa043330>
- Huang J, Samson P, Perkins SM, Anstas G, Chheda MG, DeWees TA, Tsien CI, Robinson CG et al (2017) Impact of concurrent chemotherapy with radiation therapy for elderly patients with newly diagnosed glioblastoma: a review of the National Cancer Data Base. *J Neuro-Oncol* 131:593–601. <https://doi.org/10.1007/s11060-016-2331-6>
- Ducray F, Idhah A, Wang X-W, Cheneau C, Labussiere M, Sanson M (2011) Predictive and prognostic factors for gliomas. *Expert Rev Anticancer Ther* 11:781–789. <https://doi.org/10.1586/era.10.202>
- Carninci P, Kasukawa T, Katayama S et al (2005) The transcriptional landscape of the mammalian genome. *Science* 309:1559–1563. <https://doi.org/10.1126/science.1112014>
- Ravasi T, Suzuki H, Pang KC, Katayama S, Furuno M, Okunishi R, Fukuda S, Ru K et al (2006) Experimental validation of the regulated expression of large numbers of non-coding RNAs from the mouse genome. *Genome Res* 16:11–19. <https://doi.org/10.1101/gr.4200206>
- Mehler MF, Mattick JS (2007) Noncoding RNAs and RNA editing in brain development, functional diversification, and neurological disease. *Physiol Rev* 87:799–823. <https://doi.org/10.1152/physrev.00036.2006>
- Taft RJ, Pang KC, Mercer TR, Dinger M, Mattick JS (2010) Non-coding RNAs: regulators of disease. *J Pathol* 220:126–139
- Qureshi IA, Mattick JS, Mehler MF (2010) Long non-coding RNAs in nervous system function and disease. *Brain Res* 1338:20–35
- Mercer TR, Dinger ME, Sunkin SM, Mehler MF, Mattick JS (2008) Specific expression of long noncoding RNAs in the mouse brain. *Proc Natl Acad Sci U S A* 105:712–716. <https://doi.org/10.1073/pnas.0706729105>
- Amaral PP, Neyt C, Wilkins SJ, Askarian-Amiri ME, Sunkin SM, Perkins AC, Mattick JS (2009) Complex architecture and regulated expression of the Sox2ot locus during vertebrate development. *RNA* 15:2013–2027. <https://doi.org/10.1261/rna.1705309>
- Johnson R, Teh CH-L, Jia H, Vanisri RR, Pandey T, Lu ZH, Buckley NJ, Stanton LW et al (2009) Regulation of neural macroRNAs by the transcriptional repressor REST. *RNA* 15:85–96. <https://doi.org/10.1261/ma.1127009>
- Arron JR, Winslow MM, Polleri A, Chang CP, Wu H, Gao X, Neilson JR, Chen L et al (2006) NFAT dysregulation by increased dosage of DSCR1 and DYRK1A on chromosome 21. *Nature* 441:595–600. <https://doi.org/10.1038/nature04678>
- Wang J, Zhao H, Fan Z, Li G, Ma Q, Tao Z, Wang R, Feng J et al (2017) Long noncoding RNA H19 promotes neuroinflammation in ischemic stroke by driving histone deacetylase 1-dependent M1 microglial polarization. *Stroke* 48:2211–2221. <https://doi.org/10.1161/STROKEAHA.117.017387>
- Chubb JE, Bradshaw NJ, Soares DC, Porteous DJ, Millar JK (2008) The DISC locus in psychiatric illness. *Mol Psychiatry* 13:36–64. <https://doi.org/10.1038/sj.mp.4002106>
- Zhang X, Sun S, Pu JKS, Tsang ACO, Lee D, Man VOY, Lui WM, Wong STS et al (2012) Long non-coding RNA expression profiles predict clinical phenotypes in glioma. *Neurobiol Dis* 48:1–8. <https://doi.org/10.1016/j.nbd.2012.06.004>
- Reon BJ, Anaya J, Zhang Y, Mandell J, Purow B, Abounader R, Dutta A (2016) Expression of lncRNAs in low-grade gliomas and glioblastoma multiforme: an in silico analysis. *PLoS Med* 13:e1002192. <https://doi.org/10.1371/journal.pmed.1002192>
- Li R, Qian J, Wang Y-Y, Zhang JX, You YP (2014) Long noncoding RNA profiles reveal three molecular subtypes in glioma. *CNS Neurosci Ther* 20:339–343. <https://doi.org/10.1111/cns.12220>
- Wang W, Zhao Z, Yang F, Wang H, Wu F, Liang T, Yan X, Li J et al (2018) An immune-related lncRNA signature for patients with anaplastic gliomas. *J Neuro-Oncol* 136:263–271. <https://doi.org/10.1007/s11060-017-2667-6>
- Wang W, Yang F, Zhang L et al (2016) LncRNA profile study reveals four-lncRNA signature associated with the prognosis of patients with anaplastic gliomas. *Oncotarget* 7:77225–77236. <https://doi.org/10.18632/oncotarget.12624>
- Zhang X-Q, Sun S, Lam K-F, Kiang KMY, Pu JKS, Ho ASW, Lui WM, Fung CF et al (2013) A long non-coding RNA signature in glioblastoma multiforme predicts survival. *Neurobiol Dis* 58:123–131. <https://doi.org/10.1016/j.nbd.2013.05.011>
- Chen G, Cao Y, Zhang L et al (2017) Analysis of long non-coding RNA expression profiles identifies novel lncRNA biomarkers in the tumorigenesis and malignant progression of gliomas. *Oncotarget* 8:67744–67753. <https://doi.org/10.18632/oncotarget.18832>
- van de Vijver MJ, He YD, van't Veer LJ et al (2002) A gene-expression signature as a predictor of survival in breast cancer. *N Engl J Med* 347:1999–2009. <https://doi.org/10.1056/NEJMoa021967>
- Spentzos D, Levine D, Ramoni M et al (2004) Gene expression signature with independent prognostic significance in epithelial ovarian cancer. *J Clin Oncol* 22:4700–4710. <https://doi.org/10.1200/jco.2004.04.070>
- Bullinger L, Döhner K, Bair E, Fröhling S, Schlenk RF, Tibshirani R, Döhner H, Pollack JR (2004) Use of gene-expression profiling to identify prognostic subclasses in adult acute myeloid leukemia. *N Engl J Med* 350:1605–1616. <https://doi.org/10.1056/NEJMoa031046>
- Chibon F (2013) Cancer gene expression signatures—the rise and fall? *Eur J Cancer* 49:2000–2009. <https://doi.org/10.1016/j.ejca.2013.02.021>
- Bao ZS, Chen HM, Yang MY, Zhang CB, Yu K, Ye WL, Hu BQ, Yan W et al (2014) RNA-seq of 272 gliomas revealed a novel, recurrent PTPRZ1-MET fusion transcript in secondary glioblastomas. *Genome Res* 24:1765–1773. <https://doi.org/10.1101/gr.165126.113>
- Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL, Barrell D et al (2012) GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Res* 22:1760–1774. <https://doi.org/10.1101/gr.135350.111>

32. Perteu M, Perteu GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL (2015) StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol* 33:290–295. <https://doi.org/10.1038/nbt.3122>
33. Tibshirani R (1997) The lasso method for variable selection in the cox model. *Stat Med* 16:385–395. [https://doi.org/10.1002/\(SICI\)1097-0258\(19970228\)16:4<385::AID-SIM380>3.0.CO;2-3](https://doi.org/10.1002/(SICI)1097-0258(19970228)16:4<385::AID-SIM380>3.0.CO;2-3)
34. Tibshirani R (2011) Regression shrinkage and selection via the lasso: a retrospective. *J R Stat Soc Ser B Stat Methodol* 73:273–282. <https://doi.org/10.1111/j.1467-9868.2011.00771.x>
35. Goeman JJ (2010) L1 penalized estimation in the Cox proportional hazards model. *Biom J* 52:70–84. <https://doi.org/10.1002/bimj.200900028>
36. Alizadeh AA, Gentles AJ, Alencar AJ, Liu CL, Kohrt HE, Houot R, Goldstein MJ, Zhao S et al (2011) Prediction of survival in diffuse large B-cell lymphoma based on the expression of 2 genes reflecting tumor and microenvironment. *Blood* 118:1350–1358. <https://doi.org/10.1182/blood-2011-03-345272>
37. Lossos IS, Czerwinski DK, Alizadeh AA, Wechser MA, Tibshirani R, Botstein D, Levy R (2004) Prediction of survival in diffuse large-B-cell lymphoma based on the expression of six genes. *N Engl J Med* 350:1828–1837. <https://doi.org/10.1056/NEJMoa032520>
38. Friedman J, Hastie T, Tibshirani R (2010) Regularization paths for generalized linear models via coordinate descent. *J Stat Softw* 33: . doi: <https://doi.org/10.18637/jss.v033.i01>
39. Raykar VC, Steck H, Krishnapuram B, et al On ranking in survival analysis: bounds on the concordance index
40. Gerds TA, Kattan MW, Schumacher M, Yu C (2013) Estimating a time-dependent concordance index for survival prediction models with covariate dependent censoring. *Stat Med* 32:2173–2184. <https://doi.org/10.1002/sim.5681>
41. Kim TK, Hemberg M, Gray JM (2015) Enhancer RNAs: a class of long noncoding RNAs synthesized at enhancers. *Cold Spring Harb Perspect Biol* 7:a018622. <https://doi.org/10.1101/cshperspect.a018622>
42. Panzitt K, Tschernatsch MMO, Guelly C, Moustafa T, Stradner M, Strohmaier HM, Buck CR, Denk H et al (2007) Characterization of HULC, a novel gene with striking up-regulation in hepatocellular carcinoma, as noncoding RNA. *Gastroenterology* 132:330–342. <https://doi.org/10.1053/J.GASTRO.2006.08.026>
43. Du Z, Fei T, Verhaak RGW et al (2013) Integrative genomic analyses reveal clinically relevant long noncoding RNAs in human cancer. *Nat Struct Mol Biol* 20:908–913. <https://doi.org/10.1038/nsmb.2591>
44. Mohankumar S, Patel T (2016) Extracellular vesicle long noncoding RNA as potential biomarkers of liver cancer. *Brief Funct Genomics* 15:249–256. <https://doi.org/10.1093/bfpg/elv058>
45. Zhou M, Zhang Z, Zhao H, et al (2017) An immune-related six-lncRNA signature to improve prognosis prediction of glioblastoma multiforme. *Mol Neurobiol* 1–14
46. Angileri FF, Aguenouz M, Conti A et al (2008) Nuclear factor- κ B activation and differential expression of survivin and Bcl-2 in human grade 2–4 astrocytomas. *Cancer* 112:2258–2266. <https://doi.org/10.1002/cncr.23407>
47. Korkolopoulou P, Levidou G, Saetta AA, el-Habr E, Eftichiadis C, Demenagas P, Thymara I, Xiromeritis K et al (2008) Expression of nuclear factor- κ B in human astrocytomas: relation to pI κ B α , vascular endothelial growth factor, Cox-2, microvascular characteristics, and survival. *Hum Pathol* 39:1143–1152. <https://doi.org/10.1016/J.HUMPATH.2008.01.020>
48. Schaefer LK, Ren Z, Fuller GN, Schaefer TS (2002) Constitutive activation of Stat3 α in brain tumors: localization to tumor endothelial cells and activation by the endothelial tyrosine kinase receptor (VEGFR-2). *Oncogene* 21:2058–2065. <https://doi.org/10.1038/sj.onc.1205263>
49. Abou-Ghazal M, Yang DS, Qiao W, Reina-Ortiz C, Wei J, Kong LY, Fuller GN, Hiraoka N et al (2008) The incidence, correlation with tumor-infiltrating inflammation, and prognosis of phosphorylated STAT3 expression in human gliomas. *Clin Cancer Res* 14: 8228–8235. <https://doi.org/10.1158/1078-0432.CCR-08-1329>
50. Puliappadamba VT, Hatanpaa KJ, Chakraborty S, Habib AA (2014) The role of NF- κ B in the pathogenesis of glioma. *Mol Cell Oncol* 1:e963478. <https://doi.org/10.4161/23723548.2014.963478>
51. Kesanakurti D, Chetty C, Rajasekhar Maddirela D, Gujrati M, Rao JS (2013) Essential role of cooperative NF- κ B and Stat3 recruitment to ICAM-1 intronic consensus elements in the regulation of radiation-induced invasion and migration in glioma. *Oncogene* 32: 5144–5155. <https://doi.org/10.1038/ncr.2012.546>
52. Coupienne I, Bontems S, Dewaele M, Rubio N, Habraken Y, Fulda S, Agostinis P, Piette J (2011) NF- κ B inhibition improves the sensitivity of human glioblastoma cells to 5-aminolevulinic acid-based photodynamic therapy. *Biochem Pharmacol* 81:606–616. <https://doi.org/10.1016/J.BCP.2010.12.015>
53. Sakurai K, Reon BJ, Anaya J, Dutta A (2015) The lncRNA DRAIC/PCAT29 locus constitutes a tumor-suppressive nexus. *Mol Cancer Res* 13:828–838. <https://doi.org/10.1158/1541-7786.MCR-15-0016-T>