

1 **Title:** A prognostic signature for lower-grade gliomas based on expression of long noncoding RNAs
2 **Authors and affiliations:** **Manjari Kiran**, Department of Biochemistry and Molecular Genetics, University of
3 Virginia School of Medicine, Pinn Hall 1232, Charlottesville, Virginia 22908
4 **Ajay Chatrath**, Department of Biochemistry and Molecular Genetics, University of Virginia School of Medicine,
5 Pinn Hall 1232, Charlottesville, Virginia 22908
6 **Xiwei Tang**, Department of Statistics, University of Virginia, Charlottesville, VA 22904
7 **Daniel Macrae Keenan**, Department of Statistics, University of Virginia, Charlottesville, VA 22904
8 **Anindya Dutta**, Department of Biochemistry and Molecular Genetics, University of Virginia School of Medicine,
9 Pinn Hall 1232, Charlottesville, Virginia 22908
10 **Corresponding Author:** Anindya Dutta, Department of Biochemistry and Molecular Genetics, University of
11 Virginia School of Medicine, Pinn Hall 1232, Charlottesville, Virginia 22908, Email: ad8q@virginia.edu
12 **Running title:** A long non-coding RNAs based prognostic signature for glioma
13 **Abbreviation:** lncRNA: Long Non-Coding RNAs, WHO: World Health Organization, LGG: Lower Grade
14 Gliomas, GBM: Glioblastoma Multiforme, CNS: Central Nervous System, TCGA: The Cancer Genome Atlas,
15 CGGA: Chinese Glioma Genome Atlas, HR: Hazard Ratio, PFS: Progression Free Survival, IFNG: Interferon
16 Gamma, Cindex: Concordance Index, AUC: Area Under Curve, ROC: Receiver Operating Characteristics
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

33 **Abstract**

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

49
50
51
52
53
54
55
56
57
58
59
60

Keywords: long non-coding RNAs, gliomas, gene expression profiling, prognosis

61 **Introduction**

62 Over the past decade, high-throughput RNA-seq technology discovered many novel transcriptional units,
63 which were otherwise missed by probe design based transcriptome profiling. Among these transcriptional units
64 were many long non-coding RNAs (lncRNA), which are transcripts longer than 200 bases with almost no protein-
65 coding potential or open reading frames of <50 amino acids. These lncRNAs are numerous in cells [1], are highly
66 regulated and are more cell-type specific than protein-coding genes [2]. LncRNAs are involved in a broad spectrum
67 of function and recent studies suggest they have specific roles in different diseases like cancer (reviewed in [3, 4]).

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

79 LncRNAs are widely expressed in the central nervous system (CNS) and are involved in several pathways
80 related to CNS development [8–13]. LncRNA BRN1B is one of the critical lncRNAs for brain development [13].
81 LncRNA Sox2OT plays an important role in determining neural fate [14]. Dysregulation of many lncRNAs like
82 DGCR5, NRON, H19, DISC2 have been associated with different CNS diseases [15–18]. Previous studies have
83 shown that specific lncRNA expression patterns are also associated with different histological subtypes and grade in
84 gliomas [19, 20]. For example, expression of MALAT1, POU3F3 and H19 are highly correlated with glioma
85 malignancy. More recently, lncRNAs are also found to be of prognostic significance suggesting their role in glioma
86 malignancies and as a potential therapeutic target and biomarker [19, 20]. Li et al, 2014 revealed three molecular

87 subtypes of gliomas based on lncRNAs expression that has a strong correlation with patient's survival [21].
88 Furthermore, analysis on previously published microarray data has explored lncRNA-based signature as a
89 prognostic marker in gliomas ([20, 22–25]).

90 Many studies have highlighted the power of gene expression profiles to predict tumor classification, patient
91 outcome and tumor response to therapy. Differentially expressed genes in cancer patients versus normal individuals
92 are often the starting set to predict prognostic signature associated with survival. This strategy suffers from false
93 negatives and from the fact that differentially expressed genes might not be associated with differences in survival at
94 all. Another limitation of this method is the requirement of perfect matched normal to identify differentially
95 expressed genes. This creates a major hurdle in case of brain cancer where getting a perfect matched normal tissue is
96 not trivial. While high-throughput technologies have facilitated the search of biomarkers through multivariate data
97 analyses, there still remain challenges with respect to meaningful statistical and biological information. Firstly, most
98 of the biological datasets suffer with multicollinearity: the influence of one gene on expression of other genes.
99 Secondly, there are more features (genes) than observations (patients), which leads to overfitting by most of existing
100 learning algorithms and results in poor performance of the model in prediction in an unseen testing dataset. Thus, a
101 more robust machine learning approach is required to find genes as prognostic signature from a multi-dimensional
102 multivariate gene-expression data. Regression models like lasso, ridge and elastic net are some widely used
103 approaches to penalize the effect of multicollinearity and are well suited for constructing models when there are
104 large numbers of features.

105 In the present study, we develop an lncRNA-based prognostic signature in combination with Cox
106 regression and L1-LASSO regularization to model survival of grade II and grade III glioma patients. This is the first
107 study that combined Cox and Lasso regularization to select lncRNAs that can predict survival in glioma patients.
108 After controlling for covariates associated with glioma survival (age, grade, IDH1 mutation status), we selected 8
109 lncRNAs UVA8, to calculate a risk-score, which successfully divides patients into high-risk and low-risk groups in
110 both TCGA (461 patients) and CGGA (274 patients) dataset. The risk score calculated by these 8 lncRNAs is an
111 independent and better prognostic marker for grade II and grade III glioma patient survival. The guilt-by-association
112 analysis of lncRNAs in UVA8 indicated their role in suppressing interferon signaling pathway and epithelial to

113 mesenchymal transition. Besides their use as a biomarker, these lncRNAs need to be studied in detail to determine
114 how they affect patient outcome.

115 **Materials and Methods**

116 **Patients and samples**

117 Aligned bam files and clinical information for 512 LGG patients (grade II and III) were retrieved from The Cancer
118 Genome Atlas (TCGA) data portal <https://portal.gdc.cancer.gov/>. The study is performed on 461 patients for which
119 both RNAseq and survival information were available. Most samples in TCGA are collected from patients from the
120 US and also from other countries including Canada, Russia, and Italy. This dataset being the largest and most
121 updated glioma dataset is used as training dataset in the present study. The raw sequencing data for 274 glioma
122 patients from Chinese Glioma Genome Atlas (CGGA) as independent cohort was downloaded using accession no.
123 SRP027383 [26]. The survival information for these Chinese patients was downloaded from CGGA
124 <http://www.cgga.org.cn/>. IDH1 mutation data for all the LGG patients were retrieved from Tier 3 TCGA data
125 accessed from the Broad GDAC Firehose; <https://gdac.broadinstitute.org>.

126 **RNASeq data quantification and analysis**

127 The most recent version of Gencode (GENCODE v 26) GTF file available at the time of this study was used for the
128 gene quantification [27]. Gene abundance in FPKM was obtained for 58219 genes with 15787 genes annotated as
129 lncRNA in GENCODE v26 using Stringtie v1.3.3 [28]. Out of 15787 lncRNAs, 1289 lncRNAs with a median
130 expression of 1 FPKM in 512 LGG patients were finally considered for the survival model.

131 **Survival model selection process**

132 The gene-expression data for lncRNAs was Z-score transformed to avoid systematic error across different
133 experiments. We first randomly selected 60% of TCGA patients for training set and remaining 40% of TCGA
134 patients for testing set. Since, clinical information like age, gender, tumor grade or IDH mutation status can have an
135 effect on survival (**Figure S1**), we assessed the prognostic potential of each lncRNA by multivariate Cox-regression
136 controlling the effects from these other variables. We used FDR corrected p-value cutoff of 0.05 obtained after log-

137 likelihood test comparing restricted (Age, Gender, tumor grade and IDH mutation status) with unrestricted (lncRNA
138 expression, Age, Gender, tumor grade and IDH mutation status) model to identify the significant association of an
139 lncRNA with survival. We used Cox-Proportional Hazards model based on L1 – penalized (LASSO) estimation to
140 select the best model comprising a subset of prognostic lncRNA [29–31]. We used LASSO because it is suited for
141 constructing models when there is a large number of correlated covariates [30].

142 **Risk Score calculation**

143 Risk score for each patient was established by including each of the selected genes weighted by their estimated
144 regression coefficients in the multivariable Cox regression analysis as discussed in previous studies [32, 33].

145 UVA8 Risk score = (-0.378 x expression value of RP11-266K4.14) + (-0.301 x expression value of FLJ37035) + (-
146 0.280 x expression value of LINC01561) + (-0.368 x expression value of RP11-118K6.3) + (-0.369 x expression
147 value of DGCR9) + (-0.299 x expression value of RP11-142A22.3) + (-0.434 x expression value of LINC00641) +
148 (-0.543 x expression value of RP11-96H19.1).

149 Coefficients are median cox-coefficient (after lasso selection and multivariate cox-regression) for each of the 8
150 lncRNAs from the successful models (models which can stratify patients in testing set).

151 **Statistical Analysis**

152 R package glmnet was used to perform L1-penalized cox regression [34] . R package survival and survminer were
153 used for survival data analysis and generating Kaplan–Meier plots. Different survival models were compared by
154 time-dependent concordance index (Cindex) [35]. Cindex is the most commonly used performance measure for
155 survival models, which calculates the fraction of pairs whose predicted survival time is correctly ordered. R package
156 pec::cindex is used to calculate time dependent cindex [36].

157 **Results**

158 **Building the lncRNA based survival model**

159 We developed an lncRNA based survival model for gliomas through the following steps (**Figure 1**).

- 160 1) We first randomly selected 60% (n=277) of the patients from TCGA as training set and reserved the
161 remaining 40% (n=184) of patients as testing set. The results remain similar with 70% patients in training
162 and 30% in testing set (**Figure S3 A**).
- 163 2) Cox multivariate regression was carried out in the training set on 1289 lncRNA controlling for effects from
164 other covariates like age, gender, tumor grade and IDH1 mutation status.
- 165 3) lncRNAs significantly associated with survival after likelihood ratio test (FDR $p \leq 0.05$) were retained
166 for selecting lncRNAs by lasso regularization.
- 167 4) After lasso regularization and lncRNA selection, a risk score formula was established by including selected
168 lncRNAs weighted by their estimated regression coefficients in the multivariable Cox regression analysis.
169 Risk Score = $\sum_{i=1}^n \beta_i * x_i$ (where, β_i is coefficient and x_i is expression level of lncRNA i)
- 170 5) Patients were classified into high-risk and low-risk group by using the median risk score as the cutoff in the
171 training set. The coefficient for each lncRNA and cutoff of risk score obtained from training set was used to
172 calculate risk score and stratify patients into two groups in testing set.
- 173 6) Survival differences between the low-risk and high-risk groups in the training and testing sets were
174 assessed by the Kaplan–Meier estimate and compared using the log-rank test.

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

178 The result obtained from one such survival model is shown in **FigureS2**. In ~20% of the trials the multivariate cox-
179 regression and lasso regularization in the training set did not select any lncRNAs significantly associated with
180 survival (NA in **Figure 2A**). The remaining 80% of the survival models contained different numbers of lncRNAs (x-
181 axis of **Figure 2A**) that significantly stratify patients into low and high-risk groups in training set (**Figure 2A**).
182 Among these 80% of survival models, 86% also significantly separated patients into high-risk and low-risk in the
183 testing set and are referred to as successful survival models. In order to create a robust survival model we sorted the
184 lncRNAs based on the number of times an lncRNA was selected by successful survival models (**Figure 2B**). Out of
185 167 total prognostic lncRNA in 69 successful survival models, we first ranked lncRNAs based on number of times a
186 given RNA was selected by successful models and then from the top 20 selected 8 lncRNAs with the highest median

187 cox-coefficient (Absolute value > 0.2) and least variance in the successful models in the testing set (Absolute value
188 < 0.10). 7 out of these 8 lncRNAs were also selected after 70%-30% split of training and testing patients (**Figure**
189 **S3A**), after 1000 trials instead of 100 (**Figure S3B**) and all 8 lncRNAs were selected when we used Elastic net,
190 instead of Lasso, for regularization and lncRNA selection (**Figure S3C**) suggesting the prognostic importance of
191 these 8 lncRNAs in gliomas. For brevity, this set of 8 lncRNA as a prognostic signature of gliomas will be referred
192 to as UVA8 in the manuscript.

193 **UVA8 is predictive of survival in training and independent validation set**

194 We assessed the predictive power of UVA8 by comparing overall survival of low and high risk patients in the entire
195 TCGA dataset stratified based on median risk score obtained by UVA8 (risk score calculation discussed in
196 methods). Patients in the low-risk group showed longer overall survival than the high-risk group in TCGA dataset
197 (**Figure 3A**, median OS 741.5 vs 639 days; $P = 3.1e-15$, HR=5.8). The risk scores of the patients in the TCGA
198 dataset range from -4 to 4 with median risk score of -0.023 (**Figure 3B**, top panel). Moreover, there are more
199 patients alive in the low risk group than in the high-risk group (**Figure 3B**, middle panel). Interestingly, expression
200 levels of all lncRNA in UVA8 are high in low risk patients than in high-risk patients indicating these lncRNAs as
201 favorable prognostic genes ((**Figure 3B**, bottom panel)). These findings were further validated in an independent
202 validation dataset comprising of 274 patients obtained from CGGA. Using the same median coefficient of UVA8
203 obtained from the successful survival models in TCGA, patients showed longer overall survival in low-risk than in
204 high-risk group in CGGA (**Figure 3C**, median OS = 1120.5 vs 587 days; $P = 0.0017$, HR=1.68). Moreover, low-risk
205 group in CGGA has also longer progression free survival (PFS) than the high-risk group (**Figure 3D**, median PFS
206 597.5 vs 411.5 days; $P = 0.00088$, HR=1.70). Thus, UVA8 can predict survival in both training and independent
207 validation set.

208 Since, 32% of patients in CGGA are in grade IV, the difference in overall survival could be due to over-
209 representation of grade IV patients in high-risk group. However, even when only lower-grade gliomas (grade II and
210 III) were separately examined we found significantly longer survival for low-risk versus high-risk patients (**Figure**
211 **S4A**). UVA8 fails to cluster grade IV patients from CGGA into two distinct groups highlighting the specificity of
212 signature for lower-grade gliomas (**Figure S4B**).

213 **8-lncRNA based risk score is an independent predictor of survival**

214 Lower grade gliomas have poorer outcomes in older patients, in tumors of higher grade and tumors with wild type
215 IDH1 status (**Figure S1**). Interestingly, the risk score derived from UVA8 is higher in patients older than 40 years,
216 patients in grade III vs grade II and patients harboring wild-type IDH1 gene (**Figure S5**). It was therefore important
217 to determine whether UVA8 derived risk score is an independent predictor of survival. We divided the patients into
218 younger (Age < 40) and older (Age >= 40) groups and found that risk-score can still stratify the patients into low-
219 risk and high risk in both groups (**Figure 4A**). Similarly, UVA8 based risk score can still separate the patients into
220 low and high-risk groups in grade II or grade III gliomas (**Figure 4B**). Although, IDH mutation status is a widely
221 used prognostic and predictive biomarker, the UVA8 based risk score can also separate patients into two risk groups
222 in patients presorted based on IDH mutation status (**Figure 4C**). UVA8 derived risk score can also stratify patients
223 into two risk groups among male and female patients (**Figure 4D**).

224 Conversely we tested whether these standard clinically used parameters, age, gender, grade and IDH mutation status,
225 continue to independently stratify patients even after they have been presorted into two groups by UVA8 risk score
226 (**Figure S6**). In patients with high UVA8 risk score, age, grade and IDH mutations status can further separate the
227 patients into two groups of better or worse outcome. In contrast, in patients with low UVA8 risk scores, none of the
228 clinical factors could further stratify patients into two different survival groups with a pvalue<0.05 (**Figure S6**).
229 Consistent with the previous observation (**Figure S1**), gender is ineffective in stratifying patients into two categories
230 within patients with high- or low-risk score.

231 **UVA-8 is a better predictor of glioma patients' survival**

232 We assessed the accuracy of UVA8 in prediction of survival by comparing its time-dependent area under curve
233 (AUC) with other clinical characteristics. For each prognostic factor (e.g. UVA8, IDH status etc.) we varied the cut-
234 off so as to vary the false positive rate for five-year survival prediction from 0 to 1. For each cut-off the
235 corresponding true positive rate for five-year survival was calculated (**Figure 5A**). Comparing the Area-under the
236 curve (AUC) for these ROC curves suggested that UVA8 performs best in predicting survival of the glioma patients
237 compared to the other criteria. This calculation was extended to predict survival of other durations (1-16 years) and
238 the AUC plotted for each predictor (**Figure 5B**). UVA8 can predict survival better for all durations, particularly at

239 the very early years after diagnosis when the prediction is worse for most of the predictors. Since, gender is not
240 associated with glioma patients' survival (**Figure S1**), the prediction of outcome was no better than random guess
241 (AUC = 0.5) (**Figure 5A and 5B**). We employed Cox multivariable probability hazard model to identify the impact
242 of UVA8 and different clinicopathological characteristics in estimating hazard (**Figure 5C**). UVA8 is most
243 significantly correlated with the survival information ($p = 1.4e-07$) and shows highest hazard ratio ($HR = 4$),
244 indicating that the risk score performs better than any other currently used approaches for prognosis. Here, the
245 hazard ratio of UVA8 is calculated by dichotomizing the risk score of > -0.023 (median risk score from TCGA) to 1
246 and < -0.023 to 0 to compare the hazard rates of high risk versus low risk patients. The hazard ratio of the 8
247 lncRNAs individually and combined as risk score is tabulated in **Supplementary Table S1**. The UVA8 Risk score
248 is associated with more hazard ($HR=2$) than any of the individual lncRNA supporting the importance of a
249 combinatorial signature than an individual RNA for predicting survival. The hazard ratio of UVA8 in
250 **Supplementary Table S1** is different from that in **Figure 5C** because in the former the hazard ratio is calculated
251 with the risk score as a continuous variable.

252 We then sought to compare the performance of UVA8 based survival model with published lncRNA based survival
253 models by calculating Cindex (as discussed in Methods) for TCGA dataset for each of the models. We first
254 calculated risk score for each patient by considering the expression level of the prognostic lncRNAs in each model
255 weighted by their estimated regression coefficients retrieved from the respective studies (**Supplementary Table**
256 **S2**). The patients were ordered based on their actual survival at a given time after diagnosis and based on their risk
257 score in each model. The concordance of the two orders is measured in pairwise comparisons of the patients to
258 calculate a single time-dependent concordance index for the model that is being evaluated. UVA8 outperforms all
259 existing lncRNA based survival models at different times after diagnosis (**Figure 5D**). As expected, prognostic
260 signatures that were specific to GBMs (Zhang6_2013 and Zhou6_2017) show poor concordance index when used to
261 predict survival of lower-grade glioma patients.

262 **Interferon signaling is the most enriched pathway in guilt by association with UVA8**

263 Although many lncRNAs have been identified there has been very little functional annotation of the RNAs. We
264 therefore applied guilt-by-association to infer functions of the lncRNAs associated with survival in UVA8. First we

265 interrogated whether protein-coding genes most correlated with an lncRNA in TCGA glioma cohort are themselves
266 predictive of outcome. All the lncRNAs in UVA8 are associated with a negative cox coefficient (protective). Of the
267 8 mRNAs most correlated positively with these 8 lncRNAs, 5 also have a negative cox-coefficient with a significant
268 p-value. Conversely, of the 8 mRNAs most anti-correlated with these lncRNAs, 5 have a positive cox coefficient
269 with a significant p-value (**Figure 6A**). This result is consistent with the expectation that the expression of these
270 protective lncRNAs will be positively correlated with expression of protective mRNAs and negatively correlated
271 with the expression of harmful mRNAs.

272 GSEA analysis on protein-coding genes pre-ranked from most positively correlated to most negatively correlated to
273 the lncRNA revealed several common pathways co-regulated with each of the 8 lncRNAs (**Figure 6B**).
274 Interestingly, among the mRNAs that are negatively correlated with the lncRNAs, genes involved in immune and
275 inflammatory response (IFNG, IFNA, allograft rejection, NFkB inflammatory response and JAK-STAT pathway)
276 are highly enriched. Similarly genes involved in epithelial to mesenchymal transition and cell-cycle progressions
277 are also most enriched. These gene-set enrichments suggest a conventional tumor suppressor phenotype associated
278 with these 8 lncRNAs.

279 Many of the mRNAs are common in the IFNG, IFNA, allograft rejection, NFkB inflammatory response and JAK-
280 STAT gene sets. The genes up-regulated in response to IFNG are mostly negatively correlated to lncRNAs in
281 UVA8. To visualize this, the correlation coefficients were plotted for each lncRNA (columns) with individual
282 mRNAs in the IFNG response pathway (rows) (**Figure 6C**). Out of 8, 6 lncRNAs (RP11-266K4.14, FLJ37035,
283 RP11-118K6.3, RP11-142A22.3, LINC00641 and RP11-96H19.1) are clustered together because they are more
284 negatively correlated with genes of interferon gamma response pathway (**Figure 6C**).

285 We found both NFkB and STAT3 genes as highly negatively correlated with the expression of the protective
286 lncRNAs in UVA8. Genes involved in epithelial to mesenchymal transition and encoding cell cycle related targets
287 of E2F transcription factors and involved in G2/M checkpoints were also negatively correlated with UVA8
288 expression. On the other hand, genes that are down regulated upon activation of the oncogenes KRAS are positively
289 correlated with the expression of the protective lncRNAs of UVA8.

290 In order to check whether these lncRNAs can possibly act as eRNAs, we also checked the distance between
291 lncRNAs and their correlated genes and found that these lncRNAs are correlated to several genes located in different
292 location of genome suggesting a trans-regulation by these lncRNAs (data not shown). More experimental studies are
293 required in future to decipher the role of these lncRNAs in regulating these genes and whether this regulation
294 explains the effect of the lncRNAs on glioma tumor progression.

295 **Discussion**

296 Gene expression profile reflects the underlying biological processes of disease. Cox regression is a widely used
297 approach to decipher correlation between gene expression profile and patient outcome. Previous analyses on
298 microarray data explored protein coding genes that could predict the prognosis of gliomas, particularly focusing on
299 high grade GBMs. lncRNAs are a class of RNA which can serve as a better prognostic marker than protein coding
300 mRNAs because they are numerous and cell-type specific [2, 3]. Additionally, since lncRNAs do not encode
301 protein, they are the ultimate effectors, and their expression levels more accurately predict the levels of their
302 activity. Recent studies have detected tumor-specific lncRNAs in exosomes, apoptotic bodies and microparticles
303 highlighting another advantage of considering lncRNAs in tumors, because they are expected to appear as fluid-
304 based markers for the diagnosis of different cancers [37–39]. Among six published lncRNA-based prognostic
305 signatures for gliomas two are for predicting outcome in GBMs and one specifically for anaplastic gliomas. Wang et
306 al, 2016 and Chen et al, 2017 have shown that a set of only four lncRNAs could predict survival in gliomas [23, 25].
307 However, the sequence of one of the lncRNAs in Chen et al., 2017, CR613436, was removed by the submitter on
308 NCBI. Recently, the role of immune-related genes in glioma malignancies is gaining attention leading to the
309 discovery of immune-related lncRNA-based prognostic markers for GBMs and anaplastic gliomas [40, 41].
310 Remarkably, there is no overlap between the prognostic lncRNAs identified in the aforementioned studies.
311 Moreover, these studies are based on microarray data raising concerns particular to hybridization-based approaches
312 including reliance on current knowledge of expressed genes, problems of cross-hybridization and cross-experiment
313 comparison. Another issue is that association of lncRNAs with survival using cox-regression was sometimes carried
314 out without controlling for any dependent variables and without penalizing for the effect of large number of
315 variables.

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

324 The protective lncRNAs that constitute UVA8 are negatively correlated with protein coding genes involved in
325 interferon gamma and inflammatory response highlighting the role of immune-response genes in glioma
326 progression. Except LINC01561, all 7 lncRNAs (RP11-266K4.14, FLJ37035, RP11-118K6.3, DGCR9, RP11-
327 142A22.3, LINC00641 and RP11-96H19.1) are negatively correlated to most of the protein-coding genes which are
328 up-regulated in response to interferon gamma/alpha, genes regulated by NF- κ B in response to TNF, inflammatory
329 response, and genes up-regulated by IL6 via STAT3. This suggests that an active immune reaction perhaps in
330 response to cytokines secreted from tumor and immune cells is predictive of poor outcome in gliomas. NF- κ B and
331 JAK/STAT pathways are known to be aberrantly up-regulated in GBMs. The level of NF- κ B increases as the tumors
332 progress in astrocytic tumors [42, 43] and STAT3 is constitutively active in GBMs [44, 45]. Immune related
333 pathways are also known to be involved in glioma tumor cell proliferation [46], survival [40], invasion [47] and
334 chemoresistance [48]. In addition, epithelial-mesenchymal transition (associated with invasion) and active cell
335 proliferation are suppressed if UVA8 lncRNAs are high, and this leads to better outcome, consistent with our
336 understanding of how invasion and cell proliferation negatively impact outcome. On the other hand, genes that were
337 positively correlated with the expression of UVA8 are enriched in genes that are down regulated by activation of the
338 oncogene KRAS.

339 There are reports of the same lncRNA being predictive of outcome in the same manner in multiple tumor types.
340 For example, DRAIC expression predicts good outcome in gliomas, melanomas, and cancers of the prostate,
341 stomach, liver, kidney and lung [49]. In contrast, expression of LINC00152/CYTOR is predictive of poor outcome
342 in gliomas, and cancers of the head & neck, lung, kidney, liver and pancreas (our unpublished work). Such
343 observations are particularly exciting because they imply that the lncRNA has an important role in tumor biology

344 that transcends tumor types, and these RNAs should be prioritized for cell- and molecular-biology studies to discern
345 their function. It will thus be very interesting to explore whether any of the lncRNAs of UVA8 will be protective in
346 other tumor types. Finally, future studies will address whether structural variation, copy number variations and
347 sequence polymorphism of these lncRNAs contribute to the prognostic outcome. We are excited that UVA8 was
348 also predictive of outcome in a completely different tumor cohort (CGGA) from a patient population that is from an
349 entirely different geographical location with attendant differences in environment and population genotypes. It will
350 be interesting to see if UVA8 is equally predictive of outcome in other patient populations from other parts of the
351 world.

352 **Acknowledgments**

353 We thank Dutta lab members for helpful discussions. M.K. is supported by a DOD award PC151085. The work was
354 supported by a V foundation award D2018-002 and R01 AR067712 from NIAMS.

355 **References**

- 356 1. Derrien T, Johnson R, Bussotti G, et al (2012) The GENCODE v7 catalog of human long noncoding RNAs:
357 Analysis of their gene structure, evolution, and expression. *Genome Res* 22:1775–1789 . doi:
358 10.1101/gr.132159.111
- 359 2. Cabili M, Trapnell C, Goff L, et al (2011) Integrative annotation of human large intergenic noncoding RNAs
360 reveals global properties and specific subclasses. *Genes Dev* 25:1915–1927 . doi: 10.1101/gad.17446611
- 361 3. Huarte M (2015) The emerging role of lncRNAs in cancer. *Nat Med* 21:1253–1261 . doi: 10.1038/nm.3981
- 362 4. Schmitt AM, Chang HY (2016) Long Noncoding RNAs in Cancer Pathways. *Cancer Cell* 29:452–463 . doi:
363 10.1016/j.ccell.2016.03.010
- 364 5. Stupp R, Mason WP, van den Bent MJ, et al (2005) Radiotherapy plus Concomitant and Adjuvant
365 Temozolomide for Glioblastoma. *N Engl J Med* 352:987–996 . doi: 10.1056/NEJMoa043330
- 366 6. Huang J, Samson P, Perkins SM, et al (2017) Impact of concurrent chemotherapy with radiation therapy for
367 elderly patients with newly diagnosed glioblastoma: a review of the National Cancer Data Base. *J*
368 *Neurooncol* 131:593–601 . doi: 10.1007/s11060-016-2331-6
- 369 7. Ducray F, Idbaih A, Wang X-W, et al (2011) Predictive and prognostic factors for gliomas. *Expert Rev*

- 370 Anticancer Ther 11:781–789 . doi: 10.1586/era.10.202
- 371 8. Carninci P, Kasukawa T, Katayama S, et al (2005) The transcriptional landscape of the mammalian genome.
372 Science 309:1559–63 . doi: 10.1126/science.1112014
- 373 9. Ravasi T, Suzuki H, Pang KC, et al (2006) Experimental validation of the regulated expression of large
374 numbers of non-coding RNAs from the mouse genome. Genome Res 16:11–19 . doi: 10.1101/gr.4200206
- 375 10. Mehler MF, Mattick JS (2007) Noncoding RNAs and RNA editing in brain development, functional
376 diversification, and neurological disease. Physiol Rev 87:799–823 . doi: 10.1152/physrev.00036.2006
- 377 11. Taft RJ, Pang KC, Mercer TR, et al (2010) Non-coding RNAs: Regulators of disease. J. Pathol. 220:126–
378 139
- 379 12. Qureshi IA, Mattick JS, Mehler MF (2010) Long non-coding RNAs in nervous system function and disease.
380 Brain Res. 1338:20–35
- 381 13. Mercer TR, Dinger ME, Sunken SM, et al (2008) Specific expression of long noncoding RNAs in the mouse
382 brain. Proc Natl Acad Sci U S A 105:712–716 . doi: 0706729105 [pii] 10.1073/pnas.0706729105
- 383 14. Amaral PP, Neyt C, Wilkins SJ, et al (2009) Complex architecture and regulated expression of the Sox2ot
384 locus during vertebrate development. RNA 15:2013–27 . doi: 10.1261/rna.1705309
- 385 15. Johnson R, Teh CH-L, Jia H, et al (2009) Regulation of neural macroRNAs by the transcriptional repressor
386 REST. RNA 15:85–96 . doi: 10.1261/rna.1127009
- 387 16. Arron JR, Winslow MM, Polleri A, et al (2006) NFAT dysregulation by increased dosage of DSCR1 and
388 DYRK1A on chromosome 21. Nature 441:595–600 . doi: 10.1038/nature04678
- 389 17. Wang J, Zhao H, Fan Z, et al (2017) Long Noncoding RNA H19 Promotes Neuroinflammation in Ischemic
390 Stroke by Driving Histone Deacetylase 1-Dependent M1 Microglial Polarization. Stroke 48:2211–2221 .
391 doi: 10.1161/STROKEAHA.117.017387
- 392 18. Chubb JE, Bradshaw NJ, Soares DC, et al (2008) The DISC locus in psychiatric illness. Mol Psychiatry
393 13:36–64 . doi: 10.1038/sj.mp.4002106
- 394 19. Zhang X, Sun S, Pu JKS, et al (2012) Long non-coding RNA expression profiles predict clinical phenotypes
395 in glioma. Neurobiol Dis 48:1–8 . doi: 10.1016/J.NBD.2012.06.004
- 396 20. Reon BJ, Anaya J, Zhang Y, et al (2016) Expression of lncRNAs in Low-Grade Gliomas and Glioblastoma
397 Multiforme: An In Silico Analysis. PLOS Med 13:e1002192 . doi: 10.1371/journal.pmed.1002192

- 398 21. Li R, Qian J, Wang Y-Y, et al (2014) Long Noncoding RNA Profiles Reveal Three Molecular Subtypes in
399 Glioma. *CNS Neurosci Ther* 20:339–343 . doi: 10.1111/cns.12220
- 400 22. Wang W, Zhao Z, Yang F, et al (2018) An immune-related lncRNA signature for patients with anaplastic
401 gliomas. *J Neurooncol* 136:263–271 . doi: 10.1007/s11060-017-2667-6
- 402 23. Wang W, Yang F, Zhang L, et al (2016) LncRNA profile study reveals four-lncRNA signature associated
403 with the prognosis of patients with anaplastic gliomas. *Oncotarget* 7:77225–77236 . doi:
404 10.18632/oncotarget.12624
- 405 24. Zhang X-Q, Sun S, Lam K-F, et al (2013) A long non-coding RNA signature in glioblastoma multiforme
406 predicts survival. *Neurobiol Dis* 58:123–131 . doi: 10.1016/J.NBD.2013.05.011
- 407 25. Chen G, Cao Y, Zhang L, et al (2017) Analysis of long non-coding RNA expression profiles identifies novel
408 lncRNA biomarkers in the tumorigenesis and malignant progression of gliomas. *Oncotarget* 8:67744–67753
409 . doi: 10.18632/oncotarget.18832
- 410 26. Bao ZS, Chen HM, Yang MY, et al (2014) RNA-seq of 272 gliomas revealed a novel, recurrent PTPRZ1-
411 MET fusion transcript in secondary glioblastomas. *Genome Res* 24:1765–1773 . doi:
412 10.1101/gr.165126.113
- 413 27. Harrow J, Frankish A, Gonzalez JM, et al (2012) GENCODE: the reference human genome annotation for
414 The ENCODE Project. *Genome Res* 22:1760–74 . doi: 10.1101/gr.135350.111
- 415 28. Pertea M, Pertea GM, Antonescu CM, et al (2015) StringTie enables improved reconstruction of a
416 transcriptome from RNA-seq reads. *Nat Biotechnol* 33:290–295 . doi: 10.1038/nbt.3122
- 417 29. Tibshirani R (1997) The lasso method for variable selection in the cox model. *Stat Med* 16:385–395 . doi:
418 10.1002/(SICI)1097-0258(19970228)16:4<385::AID-SIM380>3.0.CO;2-3
- 419 30. Tibshirani R (2011) Regression shrinkage and selection via the lasso: A retrospective. *J R Stat Soc Ser B*
420 *Stat Methodol* 73:273–282 . doi: 10.1111/j.1467-9868.2011.00771.x
- 421 31. Goeman JJ (2010) L1 penalized estimation in the Cox proportional hazards model. *Biometrical J* 52:70–84 .
422 doi: 10.1002/bimj.200900028
- 423 32. Alizadeh AA, Gentles AJ, Alencar AJ, et al (2011) Prediction of survival in diffuse large B-cell lymphoma
424 based on the expression of 2 genes reflecting tumor and microenvironment. *Blood* 118:1350–8 . doi:
425 10.1182/blood-2011-03-345272

- 426 33. Lossos IS, Czerwinski DK, Alizadeh AA, et al (2004) Prediction of Survival in Diffuse Large-B-Cell
427 Lymphoma Based on the Expression of Six Genes. *N Engl J Med* 350:1828–1837 . doi:
428 10.1056/NEJMoa032520
- 429 34. Friedman J, Hastie T, Tibshirani R (2010) Regularization Paths for Generalized Linear Models via
430 Coordinate Descent. *J Stat Softw* 33: . doi: 10.18637/jss.v033.i01
- 431 35. Raykar VC, Steck H, Krishnapuram B, et al On Ranking in Survival Analysis: Bounds on the Concordance
432 Index
- 433 36. Gerds TA, Kattan MW, Schumacher M, Yu C (2013) Estimating a time-dependent concordance index for
434 survival prediction models with covariate dependent censoring. *Stat Med* 32:2173–2184 . doi:
435 10.1002/sim.5681
- 436 37. Panzitt K, Tschernatsch MMO, Guelly C, et al (2007) Characterization of HULC, a Novel Gene With
437 Striking Up-Regulation in Hepatocellular Carcinoma, as Noncoding RNA. *Gastroenterology* 132:330–342 .
438 doi: 10.1053/J.GASTRO.2006.08.026
- 439 38. Du Z, Fei T, Verhaak RGW, et al (2013) Integrative genomic analyses reveal clinically relevant long
440 noncoding RNAs in human cancer. *Nat Struct Mol Biol* 20:908–13 . doi: 10.1038/nsmb.2591
- 441 39. Mohankumar S, Patel T (2016) Extracellular vesicle long noncoding RNA as potential biomarkers of liver
442 cancer. *Brief Funct Genomics* 15:249–56 . doi: 10.1093/bfgp/elv058
- 443 40. Zhou M, Zhang Z, Zhao H, et al (2017) An Immune-Related Six-lncRNA Signature to Improve Prognosis
444 Prediction of Glioblastoma Multiforme. *Mol. Neurobiol.* 1–14
- 445 41. Wang W, Zhao Z, Yang F, et al (2018) An immune-related lncRNA signature for patients with anaplastic
446 gliomas. *J Neurooncol* 136:263–271 . doi: 10.1007/s11060-017-2667-6
- 447 42. Angileri FF, Aguenouz M, Conti A, et al (2008) Nuclear factor- κ B activation and differential expression of
448 survivin and Bcl-2 in human grade 2-4 astrocytomas. *Cancer* 112:2258–2266 . doi: 10.1002/cncr.23407
- 449 43. Korkolopoulou P, Levidou G, Saetta AA, et al (2008) Expression of nuclear factor- κ B in human
450 astrocytomas: relation to pIkBa, vascular endothelial growth factor, Cox-2, microvascular characteristics,
451 and survival. *Hum Pathol* 39:1143–1152 . doi: 10.1016/J.HUMPATH.2008.01.020
- 452 44. Schaefer LK, Ren Z, Fuller GN, Schaefer TS (2002) Constitutive activation of Stat3 α in brain tumors:
453 localization to tumor endothelial cells and activation by the endothelial tyrosine kinase receptor (VEGFR-2).

- 454 Oncogene 21:2058–2065 . doi: 10.1038/sj.onc.1205263
- 455 45. Abou-Ghazal M, Yang DS, Qiao W, et al (2008) The incidence, correlation with tumor-infiltrating
456 inflammation, and prognosis of phosphorylated STAT3 expression in human gliomas. Clin Cancer Res
457 14:8228–35 . doi: 10.1158/1078-0432.CCR-08-1329
- 458 46. Puliappadamba VT, Hatanpaa KJ, Chakraborty S, Habib AA (2014) The role of NF- κ B in the pathogenesis
459 of glioma. Mol Cell Oncol 1:e963478 . doi: 10.4161/23723548.2014.963478
- 460 47. Kesanakurti D, Chetty C, Rajasekhar Maddirela D, et al (2013) Essential role of cooperative NF- κ B and
461 Stat3 recruitment to ICAM-1 intronic consensus elements in the regulation of radiation-induced invasion
462 and migration in glioma. Oncogene 32:5144–5155 . doi: 10.1038/onc.2012.546
- 463 48. Coupienne I, Bontems S, Dewaele M, et al (2011) NF-kappaB inhibition improves the sensitivity of human
464 glioblastoma cells to 5-aminolevulinic acid-based photodynamic therapy. Biochem Pharmacol 81:606–616 .
465 doi: 10.1016/J.BCP.2010.12.015
- 466 49. Sakurai K, Reon BJ, Anaya J, Dutta A (2015) The lncRNA DRAIC/PCAT29 Locus Constitutes a Tumor-
467 Suppressive Nexus. Mol Cancer Res 13:828–38 . doi: 10.1158/1541-7786.MCR-15-0016-T
- 468
- 469
- 470
- 471
- 472
- 473
- 474
- 475
- 476
- 477
- 478
- 479
- 480
- 481

482 **Figure Legend**

483 **Figure 1. Flowchart showing steps involved in identification of lncRNA based prognostic signature.**

484 **Figure 2. Selection of lncRNAs with best predictors of outcome.** A) Barplot showing number of lncRNAs that
485 predicted outcome in the training set in 100 trials. The successful models were those that also predicted outcome in
486 the testing set. NA: no lncRNA predicted outcome in training set. B) Barplot showing number of times each of the
487 top 20 lncRNAs (out of 167) were present in successful survival models (significant in testing set). The lower panel
488 shows median Cox-coefficient (after lasso penalization) and the variance of the cox-coefficient for each of the above
489 20 lncRNAs from the successful models where they were selected. The arrow points towards lncRNAs selected for
490 UVA8.

491 **Figure 3. Survival analysis of the patients divided by the prognostic lncRNAs in two data sets.** A) Patients in
492 the entire TCGA dataset with risk score greater than median score of -0.023 show poor survival compared with
493 patients with risk score less than median risk score. B) **Upper panel:** Plot showing patients sorted based on UVA8
494 risk score with black representing patient with risk score below median and red showing those with risk score above
495 median. **Middle panel:** Number of days of survival indicated on Y-axis of patients sorted on the X-axis based on
496 the risk scores in the top panel and alive/dead status indicated by color. **Bottom panel:** z-score transformed
497 expression value of lncRNAs in UVA8 show higher expression in patients with low risk score. C) Kaplan Meier plot
498 of overall survival of patients in CGGA dataset with risk score greater than (red) or less than (black) median risk
499 score of TCGA dataset. D) Kaplan-Meier plot for progression free survival in CGGA dataset showed poor survival
500 for patients with high-risk score. Rest as in C.

501 **Figure 4. Stratification analysis by different clinical variables.** Kaplan-Meier curve analysis of overall survival in
502 high- and low-risk groups for A) younger (Age < 40) and older patients (Age >= 40). B) Grade II and Grade III
503 patients C) IDH mutation status as WT and mutation (MUT) patients D) Male and Female patients. Black dashed
504 line: patients with high risk score, Gray solid line: patients with low risk score. The tables on the right show log-
505 rank p-value, hazard ratio and 95% confidence interval for each Kaplan-Meier plot.

506 **Figure 5. Performance evaluation of the 8-lncRNA based risk score.** A) Receiver operating characteristic curve
507 for 5-year survival shows UVA8 has better Area-Under-Curve compared with other predictors. B) Area-Under-
508 Curve plotted for different durations of survival for 8-lncRNA based risk score, tumor grade, Age, IDH mutation
509 status and gender of patients in TCGA cohort. C) Cox multivariate regression with clinical information and risk

510 score calculated from UVA8 for survival in TCGA cohort. **D)** Concordance-index showing measure of concordance
511 of predictor with survival of patients in TCGA.

512 **Figure 6. Guilt-by-association analysis of the 8lncRNAs in UVA8.** **A)** Correlation and Cox-regression coefficient
513 for the mRNAs that are most correlated (positive and negatively) with each of the lncRNAs in UVA8. a, b and c
514 defined below the table. **B)** List of pathways that are most enriched in protein-coding genes that are negatively
515 correlated with the UVA8 lncRNAs. **C)** Heatmap showing correlation of different genes in the interferon-gamma
516 response gene set (rows) to the lncRNAs in UVA8 (columns).

517

518

519

520

521

522

523

Figure 1

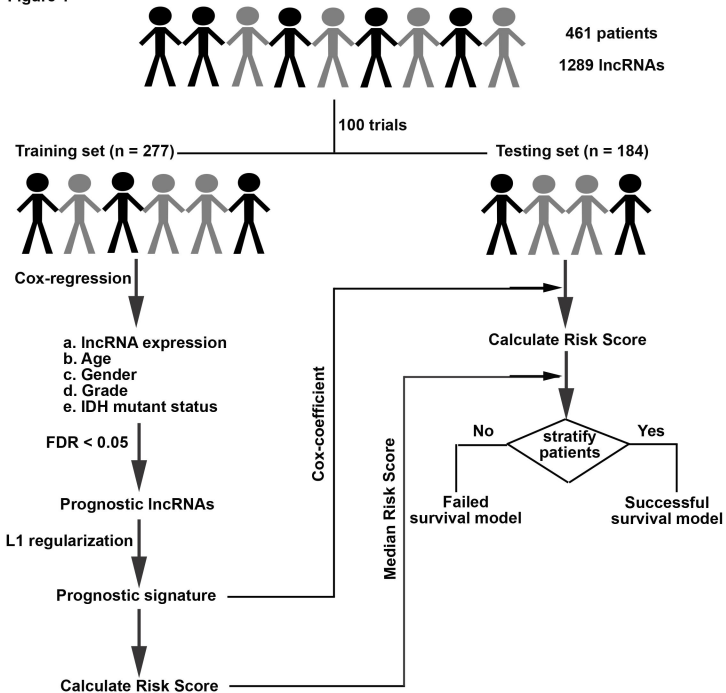
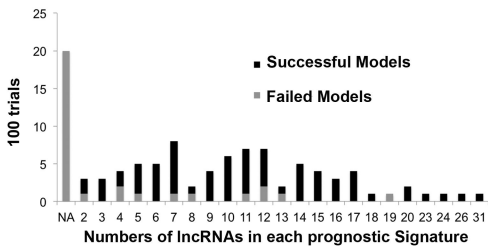


Figure 2

A



B

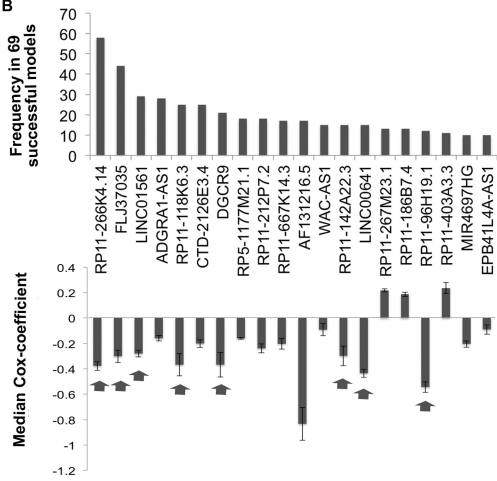
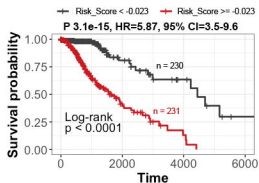
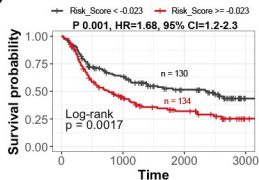


Figure 3

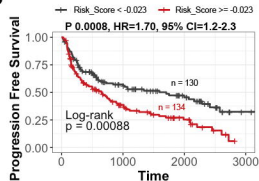
A



C



D



B

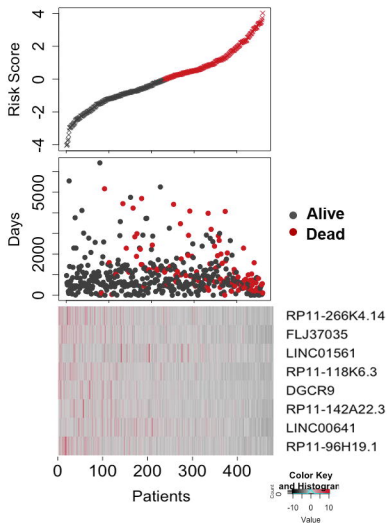
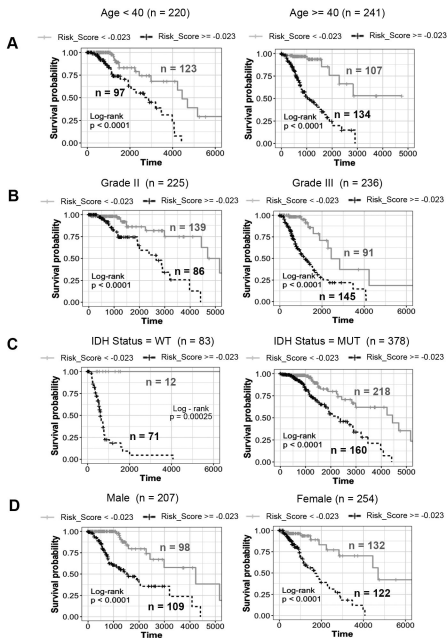


Figure 4



Age	P-val	HR	CI
<40	4.65e-5	3.99	2.0-8.2
>=40	1.83e-9	7.20	3.4-15.2

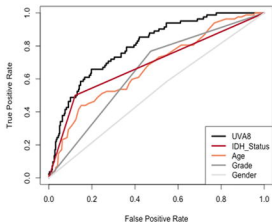
Grade	P-val	HR	CI
II	2.8e-5	4.60	2.1-10.1
III	2.5e-7	5.21	2.6-10.5

IDH	P-val	HR	CI
WT	0.00025	3e+8	0-Inf
MUT	4.7e-6	3.36	1.9-5.8

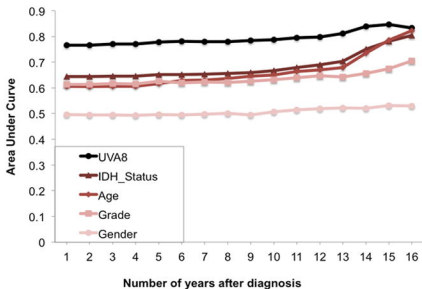
Gender	P-val	HR	CI
Male	1.3e-8	5.50	2.8-10.6
Female	5.6e-7	5.75	2.7-12.4

Figure 5

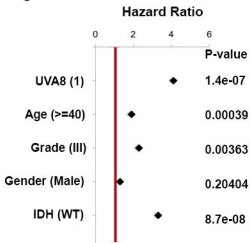
A



B



C



D

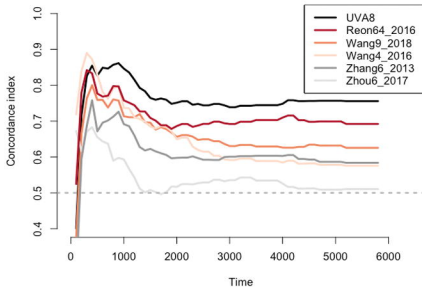


Figure 6

A

IncrNAs	Positively correlated (a,b,c)	Negatively correlated (a,b,c)
RP11-266K4.14	FGFBP3 (0.57,-0.39,3.6e-5)	CYSTM1 (-0.50,0.20,0.13)
FLJ37035	ZNF32 (0.60,-0.53,0.001)	CDC42 (-0.55,0.22,0.03)
LINC01561	PLPP4 (0.78,-0.23,0.018)	NR2E1 (-0.53,0.004,0.92)
RP11-118K6.3	RTP5 (0.61,-0.19,0.174)	DDOST (-0.62,0.23,0.027)
DGCR9	BRSK2 (0.74,-0.44,0.00017)	GN5 (-0.66,0.45,8.5e-06)
RP11-142A22.3	PRRT1 (0.70,-0.15,0.25)	SMC5 (-0.68,0.22,0.006)
LINC00641	AKAP6 (0.68,0.07,0.51)	EDEM2 (-0.72,0.07,0.53)
RP11-96H19.1	NUDT7 (0.63,-0.311,0.021)	RTCA (-0.58,0.28,0.02)

a spearman correlation
 b Cox-coefficient
 c log-rank p-value

B

Hallmark Gene Sets	Normalized Enrichment Score
HALLMARK_INTERFERON_GAMMA_RESPONSE	4.76
HALLMARK_EPITHELIAL_MESENCHYMAL_TRANSITION	4.38
HALLMARK_E2F_TARGETS	4.22
HALLMARK_INTERFERON_ALPHA_RESPONSE	3.79
HALLMARK_G2M_CHECKPOINT	3.74
HALLMARK_ALLOGRAFT_REJECTION	3.56
HALLMARK_TNFA_SIGNALING_VIA_NFKB	3.52
HALLMARK_INFLAMMATORY_RESPONSE	3.31
HALLMARK_IL6_JAK_STAT3_SIGNALLING	3.13
HALLMARK_COMPLEMENT	2.96

C

