

Gene expression analysis of glioblastomas identifies the major molecular basis for the prognostic benefit of younger age.

Yohan Lee¹, Adrienne C Scheck², Timothy F Cloughesy^{3,6}, Albert Lai³, Jun Dong⁵, Haumith K Farooqi⁴, Linda M Liaw^{4,6}, Steve Horvath⁵, Paul S Mischel^{6,7}, Stanley F Nelson^{1,6*}

1. Department of Human Genetics, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, California 90095-7088
2. The Barrow Neurological Institute, St. Joseph's Hospital and Medical Center, Phoenix, Arizona 85013
3. Department of Neurology, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, California 90095-1769
4. Department of Neurosurgery, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, California 90095-6901
5. Department of Biostatistics, David Geffen School of Medicine, University of California Los Angeles, Los Angeles, California 90095-1772
6. Jonsson Comprehensive Cancer Center, University of California Los Angeles, Los Angeles, California 90095-1781

*Corresponding author

Email addresses:

YL: yohanlee@ucla.edu

ACS: Adrienne.Scheck@chw.edu

TFC: tcloughe@ucla.edu

AL: albertlai@mednet.ucla.edu

Continued ...

JD: jundong@ucla.edu

HKF: HKFarooqi@mednet.ucla.edu

LML: lliau@mednet.ucla.edu

SH: shorvath@mednet.ucla.edu

PSM: shorvath@mednet.ucla.edu

SFN: snelson@ucla.edu

Abstract

Background

Glioblastomas are the most common primary brain tumor in adults. While the prognosis for patients is poor, gene expression profiling has detected signatures that can sub-classify GBMs relative to histopathology and clinical variables. One category of GBM defined by a gene expression signature is termed ProNeural (PN), and has substantially longer patient survival relative to other gene expression-based subtypes of GBMs. Age of onset is a major predictor of the length of patient survival where younger patients survive longer than older patients. The reason for this survival advantage has not been clear.

Results

Here we collected 273 glioblastomas and explore the relationship between gene expression subtype, age at diagnosis, and survival. This meta-analysis of published data in addition to new data confirms the existence of four distinct GBM expression-signatures. Further, patients with PN subtype GBMs had longer survival, as expected. However, the age of the patient at diagnosis is not predictive of survival time when controlled for the PN subtype.

Conclusions

The survival benefit of younger age is nullified when patients are stratified by gene expression group. Thus, the main cause of the age effect in GBMs is the more frequent occurrence of PN GBMs in younger patients relative to older patients.

Background

Glioblastoma Multiforme (GBM) persists as one of the most lethal forms of human cancer with a median survival of 12 to 14.6 months for patients who receive the latest surgical, radiation, and chemotherapy treatment [1]. While increased efficacy of some chemotherapy drugs is evident, the prognosis remains dismal [2]. Despite the number of patients diagnosed with primary malignant brain and central nervous system tumors are small relative to other cancer types (1.35% of all primary malignant cancers, U.S.), the morbidity and mortality of these patients remains severe [3].

Gliomas are primarily identified by histological features established by the grading scale of the World Health Organization (WHO). Prominent features include necrosis, nuclear mitotic activity, vascular proliferation, and cellular atypia [4]. In clinical oncology, the primary method for identifying survival factors in high-grade gliomas has been through Cox proportional hazards models. Covariates such as histological classification, performance status, and patient age at the time of surgery are particularly useful to predict patient outcomes [5,6]. Numerous studies have indicated that lower tumor WHO grade and younger patient age are the most prominent indicators of longer survival. Of patients diagnosed with a malignant glioma, younger patients have a disproportionately higher likelihood of being diagnosed with a lower grade glioma relative to older patients. However, a frequently observed phenomenon is that even among patients with the same high grade of glioma, younger patients tend to have a longer survival time. The reason for this survival advantage is not clear [7].

Our group in 2004 used a heterogeneous group of high-grade gliomas that were categorized as WHO grade III and grade IV to identify molecularly classifiable

groups irrespective of histologic typing and correlated with survival [8]. Despite this histological heterogeneity, these gliomas could be successfully stratified into at least three molecularly-defined glioma subgroups whose classification predicted patient survival. In particular, the HC1A tumors were identified as a good survival group, which were categorized by the high expression of genes typically expressed during neuronal development. The remaining groups HC2A and HC2B were characterized by mitotic and extra-cellular matrix related genes, respectively, and both demonstrated poor survival prognoses. Recently, the increased generation of whole genome microarray expression data has supported the gene expression signatures of the three previously identified GBM subgroups [9, 10]. Phillips and colleagues suggested the following terms to be more descriptive of the cellular characteristics: ‘ProNeural’ (PN) for HC1A, ‘Proliferative’ (Pro) for HC2A, and ‘Mesenchymal’ (Mes) for HC2B, which we adopt here [9].

In order to investigate the robustness and generality of the gene expression-based groupings of GBMs, we have aggregated the publicly available genome scale expression data of histologically defined GBMs. Only samples performed on the Affymetrix platform, were included for joint analysis and analytical simplicity. In total, 181 GBMs were identified from the published literature for which Affymetrix microarrays were performed and CEL files were available. For many of these samples, survival and patient age were available [8, 9, 10, 11, 12, 13]. In order to have sufficient numbers of samples to explore the correlation of expression signatures with age in GBMs, we have also performed gene expression analysis of an additional 92 GBMs in addition to this publicly available dataset. From this combined analysis, we demonstrate that about 86% of all tumor biopsies classify strongly into one of the three molecular subgroups as defined by the available gene list from Freije *et al.*, and

for about 11%, there is strong evidence of a novel group of GBMs that share expression features of Pro and Mes subtypes from the bulk biopsy analysis. Further, we investigated if there was a correlation between the molecular subtypes in relation to age of the patient at time of diagnosis and survival. We find that PN tumors are substantially more commonly diagnosed in younger patients, and that there is no survival advantage of age independent of gene expression classification. In other words, the beneficial affect of younger age in patients diagnosed with GBM is entirely due to the observation that for yet undefined reasons, younger adult patients develop PN type GBMs more commonly than older patients.

Results

Agglomeration of publicly available high-grade glioma microarray data confirms the presence of four molecularly distinct prognostic groups

In our previous work in characterizing and exploring unrecognized subtypes of GBMs using genome-scale expression analysis, our laboratory was able to develop clear evidence of gene expression signatures, and further demonstrated that the HC1A subtype had prolonged patient survival relative to the HC2A and HC2B types [8]. Several groups have now published large gene expression analyses on GBMs that permit a more robust meta-analysis and exploration of the genomic landscape of GBMs. Through the efforts of several groups and the sharing of raw microarray data, the dataset available for addressing questions of gene expression status of individual genes and signatures has expanded to 273 glioblastomas (Table 1) and are organized and made available as a co-normalized dataset here. This larger scale data permits exploration and testing of hypotheses generated in the initial microarray studies. We originally hypothesized that there were at least three subtypes of molecularly-distinct gliomas. With the accumulated genome-scale gene expression data, we initially

thought that additional gene expression signature subtypes would become more evident. However, over 97% of our additional tumor samples (from multiple institutions) continue to bin clearly into one of the initially defined subtypes (Fig. 1A). A small set of tumors (11%) show evidence of a new category which has expression features of HC2A and HC2B. The existence of the HC1A, HC2A and HC2B subgroups was corroborated by Phillips *et al.* independently [9]. The Phillips *et al.* group suggested descriptive names for the gene expression signature based GBMs with HC1A named ProNeural (PN), HC2A named Proliferative (Pro), and HC2B named Mesenchymal (Mes). Here, figure 1A indicates that a group of GBMs with both Pro and Mes gene expression signatures exists which we name “ProMes”. The patient survival durations of each gene expression group confirms the prognostic utility of the gene signature-based predictor (Fig. 1B). PN GBM patients (n = 69) have a mean survival of 2.15 years (median = 519 days) while Pro (n = 62), Mes (n = 111), and ProMes (n = 31) GBM patients have a mean survival of 1.10 years (median = 302 days), 1.30 years (median = 359 days), and 1.37 years (median = 360 days), respectively. Moreover, only 36% of PN patients survive less than 1 year from the date of surgical resection, while about 60% of Pro, 52% of ProMes, and 53% of Mes patients succumb to their malignancies within the first year after resection. Part of the reason for the low survival is likely due to the lack of homogenous treatment across institutions over a period of several years. However, we note that even patients who received modern standard of care including combined high dose radiotherapy, aggressive surgery, and Temozolomide (N=37), only 46% of these patients were alive at 12 months post diagnosis and only 40% (n=15/37) survive up to or beyond 14.6 months. Thus, differences in therapeutic efficacy from the various studies included here are unlikely to alter the fundamental conclusions and that therapy is a minor

contributor to the effects measured in terms of patient survival. The findings from this expanded dataset firmly corroborate the importance of gene signature based classification in categorizing glioblastomas into prognostically meaningful molecular groups. While there is little difference in survival between Pro, Mes, and ProMes types, all have dramatically shorter patient survival times than the PN GBM patients. For the subset of the patients where treatment data was available, there were no differences in treatments attempted in the classification groups to account for the difference in survival length of the patients.

Molecular Classification & Age Analysis

For adult patients who are diagnosed with GBM, younger age at diagnosis is a strong predictor of longer patient survival. The mechanism for this observation has not been clear, and in the Freije *et al.* work age and PN tumor type were clearly correlated [8]. With the larger group of patients and GBM tumor biopsies available here, we explore the relationship between the age of the patient survival, and gene expression signature-based subtype of GBM. We sought to determine the relative importance of the PN subtype as compared to age of the patient. For these analyses, 28 GBM samples were removed from the analysis as no available age data were provided from the published data.

In accordance with previous analyses, patients with PN GBMs (n=69) have substantially longer survival than patients with non-PN GBMs (Pro, Mes, or ProMes) (n=204) (P-value = $4.8e^{-5}$) (Fig. 2A). Patients with PN GBMs have an average survival of 2.2 years (median = 1.4 years) while the patients with non-PN GBMs survived on average 1.2 years (median = 0.9 years). We demonstrate that virtually all of the age effect observed in GBM patients is due to the increased likelihood of being diagnosed with the PN subtype in patients younger than 40. First, patients were

partitioned by age to determine if our dataset has the commonly observed beneficial young age affect. As expected, patients younger than 40 years (n=33) suffering from GBMs survived longer than those over 40 years of age (n=212) by about three-fold duration (P-value = $5.3e^{-4}$) (Fig. 2B). Patients diagnosed younger than age 40 survived on average 3.2 years (median = 3.0 years), and patients over age 40 survived on average 1.4 years (median = 1.0 years).

Next, we stratified the patients based on the subtype of GBM. Within the patients diagnosed with non-PN GBMs (Pro, Mes, and ProMes), there was no significant difference in survival between younger patients (n=18) and older patients (n=167) (p=0.91) (Fig. 3). Thus, when the tumor type was controlled by gene expression-based molecular type, there is no detectable beneficial effect of younger age. When patients were stratified based on age for all of the patients with PN subtypes, age was not a substantial predictor of survival (p=0.09, data not shown). If we ask whether gene expression based classification is a patient survival predictor within the different age groups, PN GBM subtype remains a strong predictor of longer patient survival. Of patients under 40 years of age at diagnosis, those diagnosed with PN GBMs (n = 13) demonstrate a significantly longer survival duration than those young patients who were diagnosed with non-PN GBMs (n = 18) (P-value = 0.01) (Fig. 4A).

Similarly, the beneficial gene expression based classifier was detected in the older group of patients. Patients diagnosed at over the age of 40 with PN GBMs (n = 41) survived significantly longer than equally age-matched patients over the age of 40 that suffered from non-PN GBMs (n = 167) (P-value = 0.046) (Fig 4B). Combined, these data indicate that the predominant cause of the beneficial age affect observed in GBM patients is due to the proportionally higher likelihood of patients younger than 40 to be diagnosed with the PN type GBM relative to the Pro, Mes or ProMes types. In our

sample set of the patients younger than 40 years of age, 42% (13/31) were PN, while only 20% (41/208) of the patients older than 40 years of age were diagnosed with PN type GBMs (Fishers exact p value = 0.011). GBMs increase in frequency with older age, and of those patients in our dataset over age 60, 25% have the PN type, which indicates a decreasing chance of developing this subtype of GBM with advancing age.

Multivariate and Univariate Cox proportional hazards analyses support that the molecular signature status of GBMs demonstrates greater hazard prediction than age.

To measure the robustness of the PN and non-PN molecular classifier amongst several additional covariates, multivariate and univariate Cox proportional hazards analyses were performed across the GBM patients where age data were available (N = 245: PN (n = 60), non-PN (n = 185) Tables 2 and 3). The model tested the following covariates: age, glioma gene signature, and expression status of MGMT, VEGF, and EGFR. MGMT, VEGF, and EGFR were added in as individual expression covariates due to prior reports of their possible individual contributions in relation to survival prediction [14, 15]. Our own analysis indicated that only increased VEGF expression significantly correlated with increasing age (data not shown), while MGMT and EGFR expression levels did not show significant correlations with age. While univariate analysis demonstrated as expected that the ‘Age 40 and above’ covariate represents a nominally significant marker for hazard (HR = 1.46, P-value = 0.054), age does not remain a significant survival predictor in the context of multivariate analysis (HR = 1.17, P-value = 0.44). Within the univariate model, the non-PN classification demonstrates a high hazard ratio (Table 3: HR = 1.85, P-value = 1.8e-4). Within the multivariate model, the most significant hazard covariate for patient

outcome was in fact the non-PN molecular subtype (Table 2: HR = 1.83, P-value = 7.5e-4).

Finally, several molecular studies have demonstrated that the aberrant expression of certain genes serve as useful prognostic indicators for patient survival. When the expression status of MGMT, VEGF, and EGFR were compared in the multivariate analysis with our molecular classifier, only MGMT expression demonstrated significant performance in the predictive models (Multivariate HR = 1.45, P-value = 7.8e-3). The higher expression of MGMT associated with poorer survival is consistent with other studies that report the expression of MGMT silenced by promoter hypermethylation enhances GBM patient survival [16]. MGMT silencing may also be responsible for conferring additional survival benefit when complemented by Temozolomide treatment, but this could not be studied in the current analysis [17]. In addition to MGMT, VEGF and EGFR expression levels were studied for possible relationships to poor survival outcomes as indicators of angiogenesis and constitutive signal transduction, respectively. However, from the multivariate results neither VEGF nor EGFR 'activation', as determined by the expression arrays, contributed significantly to hazard (VEGF HR = 1.21, P-value = 0.18; EGFR HR = 0.88, P-value = 0.37). Interestingly, with respect to univariate models, one published univariate Cox proportional hazards study reported that VEGF expression was a significant prognostic indicator for poor survival, while EGFR expression was not [15]. Our univariate analysis confirmed this finding as VEGF expression in the univariate model was the only gene whose expression indicated statistically significant hazard amongst these three genes (HR = 1.34, P-value = 0.029), and EGFR failed to display significant hazard implications in both multivariate and univariate models (Univariate HR = 0.95, P-value = 0.70). These

weaker effects from individual gene analyses are likely due to their partial contribution to the overall expression signatures. VEGF is part of the ProMes and Mes expression signatures while EGFR is part of the PN, Pro, and ProMes expression signatures. These data may highlight the complex nature of genetic effects within GBMs which are not derived from a single gene but rather a complex reprogramming of hundreds of genes being dysregulated. Thus, the relative contribution of individual, but important, genes is less than the aggregate set defined by the whole gene expression profile. The multivariate analysis supports this conclusion as the ProNeural versus non-ProNeural subtype of GBM within patients demonstrates itself as the primary favorable survival prognosis indicator.

Discussion

This study uses available genome-scale gene expression based analysis across glioblastomas to further investigate observed age effects in patient survival and creates a unified dataset for further exploration of gene expression correlates in glioblastomas. We add to the literature 92 additional GBM biopsy gene expression profiles performed on the U133A and U133 Plus 2.0 platforms. We confirm that within histologically defined GBMs, there are robust and repeatedly observed gene expression signature based groups of GBMs, which produces a classification scheme within GBMs. Within this aggregate dataset of 273 GBMs, the three previously defined molecular glioma subgroups were robustly detected as defined by over-expression of a series of related genes: Neurogenesis / HC1A (aka ProNeural), Mitotic / HC2A (aka Proliferative), and Extra-Cellular Matrix related / HC2B (aka Mesenchymal). From the larger dataset, we observe evidence of a new subtype that highly co-expresses genes in both the Pro and Mes groups, which we term “ProMes”.

Of these four molecularly distinct groups identified, only PN portends a favorable prognosis relative to the other expression based groups.

In this study, we determine that the reason that younger patients diagnosed with GBMs have longer survival durations than older GBM patients is due to the observation that younger patients tend to develop the favorable PN GBM type more commonly relative to older patients. The reason for this observed age effect is not clear at this time, but may have to do with the precursor cell that develops into GBMs changing over time. One could hypothesize that the precursor cell that gives rise to the PN type diminishes in abundance in the CNS with advancing age. However, given that GBM incidence increases greatly with age, the absolute numbers of PN GBMs is actually numerically higher in older patient groups. Thus, we favor a model that the precursor cell type that gives rise to the Pro and Mes types of GBMs are increasingly likely to become neoplastic over time while this effect is not as pronounced within the PN type precursors.

PN type GBMs represent a unique tumor etiology whose idiopathic molecular mechanisms manifest in survival periods from two to ten years in contrast to ten or fifteen months for the Pro, Mes and ProMes types of GBMs. Historically, the percentage of GBM patients who have long term survival of 3 years or more has been reported to be approximately 5% [18]. This 5% incidence rate matches well with the observations reported here in which 5% (13/245) of our PN GBM patients are observed to survive 3 years or longer. The identification of the PN subgroup is important for patient management and stratification into small phase II clinical trials for experimental therapeutics as uneven representation of the PN GBM diagnoses would greatly alter observed survival times irrespective of potentially active agents

[19]. The identification of the gene expression subtype is clearly more important for patient stratification within clinical trials than age.

It is likely that the use of genome-wide expression based molecular classification will result in less variation in tumor diagnoses and provide more specific guidance to clinicians. The agglomeration of gene expression datasets permits meta-analyses that were insufficiently powered in the multiple individual publications. Resources for the sharing of genome-scale expression datasets have been set up at Array Express, Gene Expression Omnibus, and Celsius [20, 21]. Critical to the sharing of microarray data is providing raw microarray data as opposed to processed data. In order to facilitate this sharing, the NIH Neuroscience Microarray Consortium has established Celsius, which is a community resource of CEL (image) files performed on the Affymetrix platform for public distribution using programmatic tools. At the writing of this manuscript, the Celsius database contains CEL files from human experiments performed on U95Av2 arrays (n=5 006), U133A arrays (n=13 818), and U133 Plus 2.0 arrays (n=10 376). To fully capture and leverage the value of microarray expression data, a greater commitment must be made to capture and share clinical covariates and raw expression data. For instance, in this study of 273 glioblastomas, a total of 433 GBM CEL files were initially identified across the U95Av2, U133A, and U133 Plus 2.0 platforms. Of these, thirty-six percent (160/433) were immediately removed from this study due to a lack of any clinical data. Additionally, not all microarray CEL files or their matched clinical data points are systematically retrievable.

Amongst the samples gathered at UCLA, we simultaneously gathered additional covariates such as extent of surgical resection, Karnofsky Performance Scores (KPS), lesion locations, and MRI scans. We have begun to make all of these

data available through a web interface in order to promote data sharing and exploration of gene expression differences in gliomas. All of the data added here are deposited in Gene Expression Omnibus (GEO), and can be explored at our real-time survival-synchronized search engine “Probeset Analyzer” [22].

Clinical Decision Impact and Improved Public Disclosure.

In large academic hospitals, tumors come from a wide variety of patients from across different cities, states, or countries. In contrast, local hospitals treat their regional constituencies. The potential for demographically-biased patient populations and biased tumor subsets is a possibility. These trends can reinforce particular treatment strategies at local institutions over time. For example, if patients from a community highly populated by retirees (e.g. southern Florida) presented with a GBM, clinicians would be apt to predict that these older patients would likely succumb to their malignancies within one year. Current treatment for patients diagnosed with high grade gliomas consists of surgical resection followed by toxic and expensive therapy schedules that are minimally effective. But if these elderly patients were suffering from a PN tumor, they would have a high likelihood for surviving at least two to three years or longer. These patients could then be distinguished from patients who otherwise present identically under the microscope or according to their patient biographical sketch. This would permit time to enroll in potentially beneficial clinical trials. Thus, if grade and age alone were considered for prognosis, these factors would lead clinicians to prescribe unnecessary treatments due to trends reinforced by regional sampling biases.

Conclusions

We provide an explanation for why younger patients diagnosed with GBM patients have longer life expectancies than older patients using accumulated whole-genome

microarray expression data and clinical variables. We have discovered the reason that younger patients tend to survive longer is because they are more likely to present with PN gene signature tumors relative to the more common and aggressive Pro, Mes, and ProMes GBM types. The PN molecular classification predicts an enhanced survival performance by at least two to ten years irrespective of age when tested against age-matched, Pro, Mes, or ProMes molecularly-classified tumors. The application of Cox proportional hazards studies have also confirmed that having a non-PN tumor was the most statistically significant factor in predicting precipitous short survival over age by two orders of magnitude. This data lends more evidence to the clinical reality that high-grade glioma patients suffer from molecularly distinct tumors. Therefore, different tumors with distinct etiologies should be differentially segregated in terms of their treatment regimens and especially their clinical trial assignments. The benefits for clinicians would be a reduction in the heterogeneous admixture of genetic background for treatment cohorts and a potential reduction in the percentage of non-responders for molecularly-targeted investigational new drugs. Patients accurately classified and characterized for the biology of their tumor may potentially benefit from a variety of molecularly-targeted treatments as the inclusion criteria for clinical trials become simultaneously based on molecular signatures.

Methods

Microarray and Clinical Data Collection

Clinical data including histopathology, age, sex, and survival time from diagnosis were retrieved from 181 glioblastomas which have been reported within previous studies between 2003 and 2006 (Table 1) and for which CEL files (Affymetrix, Santa Clara, CA) were available from the authors. In addition, we collected 92 new patient-unique tumor biopsies from the UCLA Neuro-oncology Program (n = 61) and the

Barrow Neurological Institute (n = 31) for a grand total of 273 glioblastomas. Newly acquired tumors were collected through institutional review board approved protocols and assigned WHO grades at UCLA Neuropathology or Barrow Neuropathology by PSM. Time of survival (days), sex, and age were collected where available (See Supplementary Information Table S1). Patient age at the time of diagnosis was available for 245 patients and ranged from 18 to 86 years. Sex of the individual was available for those 245 patients (156 males and 89 females).

Microarray Experimentation

Total RNA was purified from fresh frozen tumor biopsies and visually inspected for tumor content using Qiagen RNeasy columns and standard manufacturer's protocols. Labeled one round cRNA was generated using kits (GeneChip One-Cycle Target Labeling and Control Reagent) from Affymetrix. cRNA was quantified and 15 micrograms were hybridized to U133A and U133 Plus 2.0 arrays at the UCLA DNA Microarray Facility using standard protocols recommended by the manufacturer (See Supplementary Information section S2 for full details). All newly generated CEL files were deposited into the Celsius microarray database and this system was used to normalize relative to other microarrays of the same Affymetrix platform using RMA with default settings from the Bioconductor R library [20, 21, 23, 24].

Combination of microarray data

The Celsius microarray database, which houses over 20 000 human CEL files on various array iterations, was used to quantify and normalize each dataset with a comparison group of 50 random samples selected from the database for RMA normalization and quantification using default parameters. Only probesets from the U133A portion of the Freije *et al.* paper that are retained or map to the same gene were analyzed from each tumor type.

Sample membership by HC Classification Gene Voting

The Hierarchical Clustering (HC) classification for each glioma was determined by the gene voting strategy as described previously [8]. Briefly, the mean value of each probeset was evaluated from all samples within each of the three microarray platforms U95Av2, U133A, and U133 Plus 2.0 separately. Second, the probesets from each sample were assigned “yes” or “no” votes if that probeset’s value was above or below the aforementioned probeset mean of its platform. Third, the “yes” or “no” votes of each probeset from the 377 probesets contained within the U133A portion of the 595 HC probeset classifier (which was based on U133A and U133B data) were tallied and used to categorize every glioma into one of three HC molecular groups. All of the 377 probesets from the 377 used for classification were on the U133A and the U133 Plus 2.0 array types, and 200 of the probesets were able to be matched to the same genes from the older U95Av2 arrays. Lastly, each tumor was voted into one of three HC groups based on the highest vote tally: 1A vote = (tally of 1A probesets above the mean)/(total 1A probesets); 2A vote = (tally of 2A probesets above the mean)/(total 2A probesets); 2B vote = (tally of 2B probesets above the mean)/(total 2B probesets). A few gliomas appeared to vote almost equally well into both the 2A and 2B categories and are defined as “2A2B”, which is defined where the highest vote category must be a 2A or 2B and the second highest vote must be 2B or 2A , respectively with at least a 33% vote. The list of probesets used for voting each tumor across each of the Affymetrix platforms is available in the Supplementary Information section S3.

Kaplan-Meier and Cox Regression Analysis

Kaplan-Meier survival plots and Cox proportional hazard regression analyses were implemented in R version 2.5.0 using the “survival” library. Survival data and vital status for UCLA samples were based on time in days elapsed from surgical resection to the date of death up to May 1, 2006. The full R code is available and documented in the Supplementary Information section S4.

MGMT, VEGFA, EGFR Expression State Determination

Expression states for the MGMT, VEGFA and EGFR probesets were called “ON” or “OFF” based on whether their probesets’ expression levels were above or below their respective mean within each platform. The probesets employed for the genes from each platform are available in the Supplementary Table S5.

List of Abbreviations

GBM: Glioblastoma

WHO: World Health Organization

HC1A/PN: Hierarchical Cluster group 1A (Neurogenesis) / ProNeural

HC2A/Pro: Hierarchical Cluster group 2A (Mitotic) / Proliferative

HC2A/2B/ProMes: Hierarchical Cluster group 2A/2B / Proliferative Mesenchymal

HC2B/Mes: Hierarchical Cluster group 2B (Extra-Cellular Matrix) / Mesenchymal

MGMT: *O-6-methylguanine-DNA-methyltransferase* (GeneID:4255)

VEGFA: *vascular endothelial growth factor A* (GeneID: 7422)

EGFR: *epidermal growth factor receptor* (GeneID: 1956)

Temozolomide (Brand Name: Temodar®)

Authors' contributions

YL carried out the public and newly obtained data acquisition, analyzed the microarray expression data, performed the gene voting, identified the 2A2B / ProMes group, performed all statistical analyses (except the positive correlation between age and VEGF expression), and drafted the manuscript. ACS participated actively in revising the manuscript, obtained a substantial portion of the newly acquired glioma biopsies, collected the clinical covariates for those glioma biopsies, and performed patient follow-up analysis. TFC participated in the design of the study, collected the clinical covariates for the entire UCLA subset of the glioma biopsies, performed patient follow-up analysis, proposed the VEGF and EGFR expression studies, and revised the manuscript with the most important clinical intellectual content. AL participated in the design of the study, collected additional clinical covariates for the entire UCLA subset of the glioma biopsies, performed additional patient follow-up analysis, proposed the MGMT expression studies, and screened all of the UCLA clinical covariates. JD participated in the design of the statistical analyses, provided substantive assistance for interpretation of the statistical calculations, and performed the correlation analysis between age and positive expression between MGMT, VEGF, and EGFR activation. HKF performed all the tissue biopsy collection, micro-dissection, and electronic patient identification assignment and confirmation for all the newly acquired UCLA tumor biopsies. LML collected all publicly and newly acquired UCLA tumor biopsies, participated in the conception of the manuscript, and revised the manuscript with critical intellectual comments and suggestions. SH was involved in the initial conception, design, of the statistical analyses, participated and advised indispensably on the interpretation and revision of the statistical analysis and computer code. PSM participated in the design of the study, performed

histopathological grading on the public and newly acquired UCLA tumor biopsies, and provided substantial critical input and insightful revisions and comments. SFN conceived of the study and participated in its design, coordination, and helped draft the manuscript. All authors read and approved the final manuscript.

Acknowledgements

The authors wish to thank the UCLA and Barrow Institute subjects for their participation in this study. The work was supported by the NIH Neuroscience Microarray Consortium UCLA Site (U24NS052108), The UCLA Gene Expression Shared Resource and the Singleton Brain Tumor Program.

References

1. Stupp R MW vdBM, Weller M, Fisher B, Taphoorn MJ, Belanger K, Brandes AA, Marosi C, Bogdahn U, Curschmann J, Janzer RC, Ludwin SK, Gorlia T, Allgeier A, Lacombe D, Cairncross JG, Eisenhauer E, Mirimanoff RO;European Organisation for Research and Treatment of Cancer Brain Tumor and Radiotherapy Groups; National Cancer Institute of Canada Clinical Trials Group.: **Radiotherapy plus concomitant and adjuvant temozolomide for glioblastoma.** *N Engl J Med* 2005, **352**(10):987-996.
2. Omuro AM FS, Raymond E.: **Lessons learned in the development of targeted therapy for malignant gliomas.** *Mol Cancer Ther* 2007, **6**(7):1909-1919.
3. **CBTRUS. 2005-2006 Statistical Report: Primary Brain Tumors in the United States Statistical Report, 1998-2002 (Years Data Collected).** 2006.
4. Louis DN PJ, Jacobs T: **Report of the Brain Tumor Progress Review Group.** 2000, **National Institutes of Health: Bethesda, MD.**
5. Dong S NC, Betensky RA, Stemmer-Rachamimov AO, Denko NC, Ligon KL, Rowitch DH, Louis DN.: **Histology-based expression profiling yields novel prognostic markers in human glioblastoma.** *J Neuropathol Exp Neurol* 2005, **64**(11):948-955.

6. Krex D KB, Hartmann C, Deimling AV, Pietsch T, Simon M, Sabel M, Steinbach JP, Heese O, Reifenberger G, Weller M, Schackert G;: **Long-term survival with glioblastoma multiforme.** *Brain* 2007(2007 Sep 4):2596-2606.
7. Wrensch M FJ, Schwartzbaum JA, Bondy M, Berger M, Aldape KD.: **The molecular epidemiology of gliomas in adults.** *Neurosurg Focus* 2005, **19**(5):E5: pp1-11.
8. Freije WA C-VF, Fang Z, Horvath S, Cloughesy T, Liao LM, Mischel PS, Nelson SF.: **Gene expression profiling of gliomas strongly predicts survival.** *Cancer Res* 2004, **64**(18):6503-6510.
9. Phillips HS KS, Chen R, Forrest WF, Soriano RH, Wu TD, Misra A, Nigro JM, Colman H, Soroceanu L, Williams PM, Modrusan Z, Feuerstein BG, Aldape K.: **Molecular subclasses of high-grade glioma predict prognosis, delineate a pattern of disease progression, and resemble stages in neurogenesis.** *Cancer Cell* 2006, **9**(3):157-173.
10. Rich JN HC, Jones B, Iversen ES, McLendon RE, Rasheed BK, Dobra A, Dressman HK, Bigner DD, Nevins JR, West M.: **Gene expression profiling and genetic markers in glioblastoma survival.** *Cancer Res* 2005, **65**(10):4051-4058.
11. Mischel PS SR, Shi T, Horvath S, Lu KV, Choe G, Seligson D, Kremen TJ, Palotie A, Liao LM, Cloughesy TF, Nelson SF.: **Identification of molecular subtypes of glioblastoma by gene expression profiling.** *Oncogene* 2003, **22**(15):2361-2373.

12. Shai R ST, Kremen TJ, Horvath S, Liau LM, Cloughesy TF, Mischel PS, Nelson SF.: **Gene expression profiling identifies molecular subtypes of gliomas.** *Oncogene* 2003, **22**(31):4918-4923.
13. Nutt CL MD, Betensky RA, Tamayo P, Cairncross JG, Ladd C, Pohl U, Hartmann C, McLaughlin ME, Batchelor TT, Black PM, von Deimling A, Pomeroy SL, Golub TR, Louis DN.: **Gene expression-based classification of malignant gliomas correlates better with survival than histological classification.** *Cancer Res* 2003, **63**(7):1602-1607.
14. Herrlinger U RJ, Koch D, Loeser S, Blaschke B, Kortmann RD, Steinbach, JP HT, Wick W, Meyermann R, Tan TC, Sommer C, Bamberg M, Reifenberger, G WM: **Phase II trial of lomustine plus temozolomide chemotherapy in addition to radiotherapy in newly diagnosed glioblastoma: UKT-03.** *J Clin Oncol* 2007, **24**(27):4412-4417.
15. Zhou YH TF, Hess KR, Yung WK.: **The expression of PAX6, PTEN, vascular endothelial growth factor, and epidermal growth factor receptor in gliomas: relationship to tumor grade and survival.** *Clin Cancer Res* 2003, **9**(9):3369-3375.
16. Martinez R SG, Yaya-Tur R, Rojas-Marcos I, Herman JG, Esteller M.: **Frequent hypermethylation of the DNA repair gene MGMT in long-term survivors of glioblastoma multiforme.** *J Neurooncol* 2007, **83**(1):91-93.

17. Ishii D NA, Wakabayashi T, Hatano H, Asano Y, Takeuchi H, Shimato S, Ito M, Fujii M, Yoshida J.: **Efficacy of temozolomide is correlated with 1p loss and methylation of the deoxyribonucleic acid repair gene MGMT in malignant gliomas.** *Neurol Med Chir (Tokyo)* 2007, **47(8):341-349** discussion 350.
18. Hegi ME DA, Gorlia T, Hamou MF, de Tribolet N, Weller M, Kros JM, Hainfellner JA, Mason W, Mariani L, Bromberg JE, Hau P, Mirimanoff RO, Cairncross JG, Janzer RC, Stupp R.: **MGMT gene silencing and benefit from temozolomide in glioblastoma.** *N Engl J Med* 2005, **352(10):997-1003.**
19. Mischel PS CT, Nelson SF.: **DNA-microarray analysis of brain cancer: molecular classification for therapy.** *Nat Rev Neurosci* 2004, **5(10):782-792.**
20. Day A CM, Dong J, O'connor BD, Nelson SF.: **Celsius: a community resource for Affymetrix microarray data.** *Genome Biol* 2007, **8(6):R112.**
21. **The Celsius Project**
[<http://genomics.ctrl.ucla.edu/wiki/index.php/Celsius>]
22. **Probeset Analyzer** [<http://www.probesetalyzer.com>]
23. Bolstad BM IR, Astrand M, Speed TP.: **A comparison of normalization methods for high density oligonucleotide array data based on variance and bias.** *Bioinformatics* 2003, **19(2):185-193.**

24. Gentleman RC CV, Bates DM, Bolstad B, Dettling M, Dudoit S, Ellis B, Gautier L, Ge Y, Gentry J, Hornik K, Hothorn T, Huber W, Iacus S, Irizarry R, Leisch F, Li C, Maechler M, Rossini AJ, Sawitzki G, Smith C, Smyth G, Tierney L, Yang JY, Zhang J.: **Bioconductor: open software development for computational biology and bioinformatics.** *Genome Biol* 2004, **5**(10):R80.

Figures

Figure 1: Expression-based high grade glioma profiles related to survival duration across GBM and AA histologic tumors.

A. 377 genes were used to classify the tumors available on the U133A and U133 Plus 2.0 arrays into the three glioma types: PN: ProNeural (yellow, n = 71), Pro: Proliferative (blue, n = 45), ProMes: Proliferative-Mesenchymal (purple n = 31), and Mes: Mesenchymal (red, n = 98).

B. Percent composition of Long (LTS), Medium (MTS), and Short (STS) term survival across the three glioma gene-signature tumor types. Over one-third of PN-classified tumor patients survive over 2 years, while over half of Pro and Mes classified tumor patients succumb in less than 1 year.

Figure 2: Kaplan-Meier analysis of GBM patient survival according to predictive gene expression signatures and age alone.

A. PN GBM patients vs. non-PN GBM patients.

Molecularly categorized GBM patient survival comparison shows PN GBMs survive longer than non-PN GBMs (Pro, Mes, and ProMes). PN GBM patient survival (mean = 2.2 years, median = 1.4 years) vs. non-PN GBM patient survival (mean = 1.2 years, median = 0.9 years).

B. GBM patient age difference for survival: younger than age 40 vs. older than age 40. GBM patients younger than age 40 survive longer than GBM patients older than age 40. GBM age < 40 patient survival (mean = 3.2 years, median = 3.0 years) vs. GBM age > 40 patient survival (mean = 1.4 years, median = 1.0 years).

Figure 3: Kaplan-Meier analysis of GBM patient survival for Pro, Mes, and ProMes GBM patients according to age.

GBM patients suffering from Pro, Mes, and ProMes tumors do not survive any better if younger than age 40 than identically gene expression-classified GBM patients older than the age of 40.

Figure 4: Kaplan-Meier analysis of GBM patient survival partitioned by age or predictive expression signature.

A. GBM patients younger than age 40 partitioned by PN status versus non-PN status.

Younger patients with PN GBMs survive longer than fellow younger patients with non-PN GBMs. PN GBM patients age < 40 survival (mean = 3.2 years, median = 3.0 years) vs. GBM age < 40 non-PN patient survival (mean = 1.5 years, median = 0.9 years).

B. GBM patients older than age 40 partitioned by PN status versus non-PN status.

Older patients with PN GBMs survive longer than fellow older patients with non-PN GBMs. Non-PN GBM patients age > 40 PN survival (mean = 1.8 years, median = 1.1 years) vs. GBM age > 40 non-PN patient survival (mean = 1.3 years, median = 1.0 years).

Tables

Table 1 - Combined Data Sources

Studies	No. GBM IV	Set	Array	References	PMID
Freije <i>et al.</i>	46	0	U133A	[8]	15374961
Phillips <i>et al.</i>	55	2,3	U133A	[9]	16530701
Rich <i>et al.</i>	31	4	U133A	[10]	15899794
Mischel <i>et al.</i>	2	7	U95Av2	[11]	12700671
Shai <i>et al.</i>	19	8	U95Av2	[12]	12894235
Nutt <i>et al.</i>	28	6	U95Av2	[13]	12670911
UCLA	28	1	U133A	New	
UCLA	33	1	U1332.0	New	
Barrow	31	5	U133A	New	
Total	273				

Table 2 - Multivariate Cox Proportional Hazards Ratio: Five Covariates

Model				
Covariates	Covariate	Hazard Ratio	95% CI	P-value
AGE	>= 40	1.17	0.78 – 1.76	0.44
PN status	non-PN	1.83	1.29 – 2.60	7.5e-4
MGMT	ON	1.45	1.10 – 1.90	7.8e-3
VEGF	ON	1.21	0.92 – 1.59	0.18
EGFR	ON	0.88	0.66 – 1.17	0.37

Table 3 - Univariate Cox Proportional Hazards Ratio: Five Covariates

Model				
Covariates	Covariate	Hazard Ratio	95% CI	P-value
AGE	>= 40	1.46	0.99 – 2.16	0.054
PN status	non-PN	1.85	1.34 – 2.56	1.8e-4
MGMT	ON	1.26	0.97 – 1.64	0.084
VEGF	ON	1.34	1.03 – 1.75	0.029
EGFR	ON	0.95	0.71 – 1.26	0.70

Additional files

Additional file 1 – Supplementary Information Table S1

File Name: Supplementary_Information_Table_S1

File Format: Microsoft Office Excel 2003 spreadsheet (.xls)

Title of Data: Supplementary Information Table S1

Description of Data: Tumor sample clinical covariates meta-data e.g. Time of survival (days), sex, and age where available.

Additional file 2 – Supplementary Information section S2

File Name: Supplementary_Information_section_S2

File Format: Microsoft Office Word 2003 text (.doc)

Title of Data: Supplementary Information section S2

Description of Data: Standard protocol recommended by Affymetrix used at UCLA DNA Microarray Facility.

Additional file 3 – Supplementary Information section S3

File Name: Supplementary_Information_section_S3

File Format: Microsoft Office Excel 2003 spreadsheet (.xls)

Title of Data: Supplementary Information section S3

Description of Data: List of Affymetrix probesets used for voting each tumor across each of the Affymetrix platforms.

Additional file 4 – Supplementary Information section S4

File Name: Supplementary_Information_section_S4

File Format: Microsoft Office Word 2003 text (.doc)

Title of Data: Supplementary Information section S4

Description of Data: The full R code used for Kaplan-Meier survival curves, Multivariate/Univariate Cox proportional Hazards models, and correlations between age and gene expression.

Genomic coordinates and gene names (e.g., SLC6A4, MAOA, MAO-B, MAO-A, MAO-A1, MAO-A2, MAO-A3, MAO-A4, MAO-A5, MAO-A6, MAO-A7, MAO-A8, MAO-A9, MAO-A10, MAO-A11, MAO-A12, MAO-A13, MAO-A14, MAO-A15, MAO-A16, MAO-A17, MAO-A18, MAO-A19, MAO-A20, MAO-A21, MAO-A22, MAO-A23, MAO-A24, MAO-A25, MAO-A26, MAO-A27, MAO-A28, MAO-A29, MAO-A30, MAO-A31, MAO-A32, MAO-A33, MAO-A34, MAO-A35, MAO-A36, MAO-A37, MAO-A38, MAO-A39, MAO-A40, MAO-A41, MAO-A42, MAO-A43, MAO-A44, MAO-A45, MAO-A46, MAO-A47, MAO-A48, MAO-A49, MAO-A50, MAO-A51, MAO-A52, MAO-A53, MAO-A54, MAO-A55, MAO-A56, MAO-A57, MAO-A58, MAO-A59, MAO-A60, MAO-A61, MAO-A62, MAO-A63, MAO-A64, MAO-A65, MAO-A66, MAO-A67, MAO-A68, MAO-A69, MAO-A70, MAO-A71, MAO-A72, MAO-A73, MAO-A74, MAO-A75, MAO-A76, MAO-A77, MAO-A78, MAO-A79, MAO-A80, MAO-A81, MAO-A82, MAO-A83, MAO-A84, MAO-A85, MAO-A86, MAO-A87, MAO-A88, MAO-A89, MAO-A90, MAO-A91, MAO-A92, MAO-A93, MAO-A94, MAO-A95, MAO-A96, MAO-A97, MAO-A98, MAO-A99, MAO-A100) are listed along the top of the heatmap.

HC

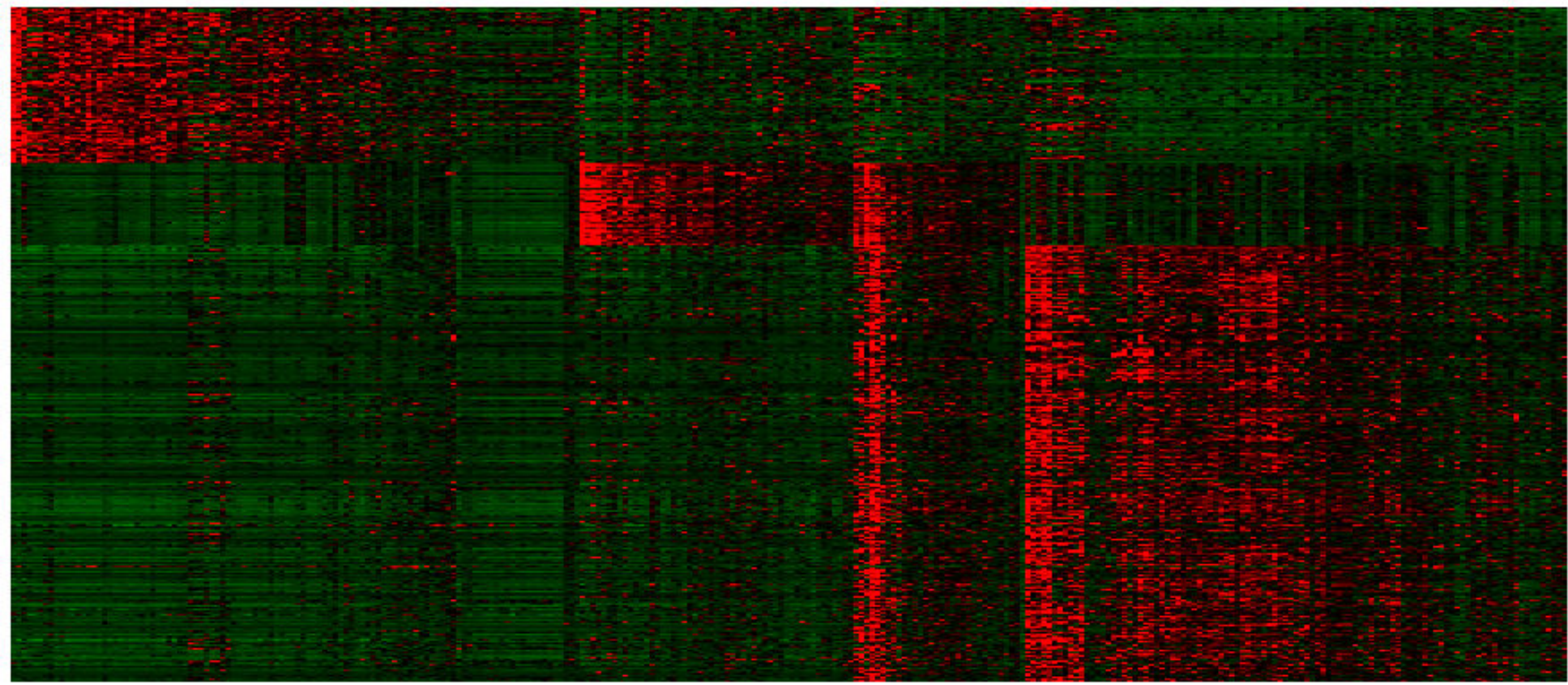


Figure 1

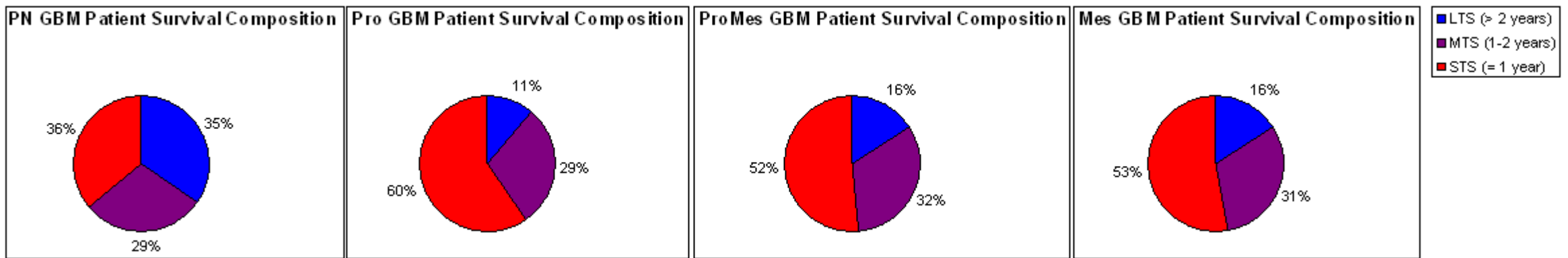


Figure 2

Survival Difference between PN vs. Pro+Mes+ProMes GBMs p-value= 4.8e-05

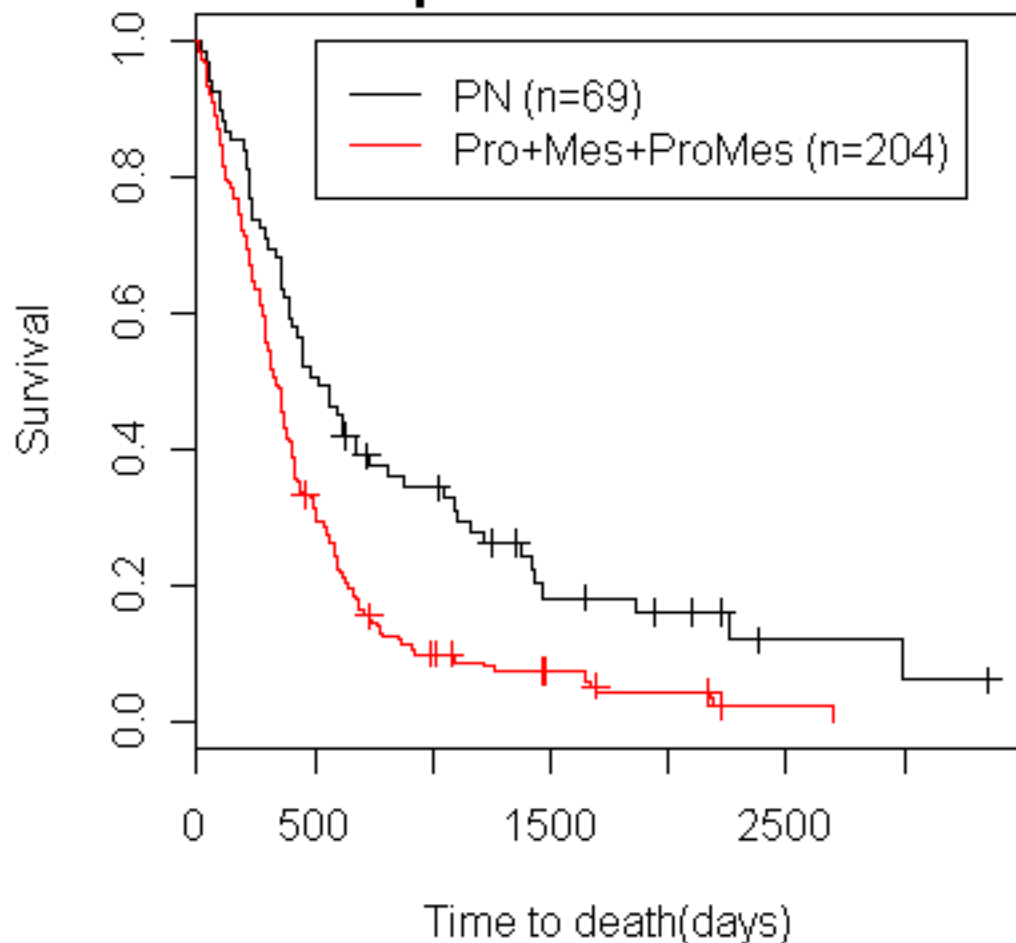


Figure 3

Age Difference in GBMs

p-value= 0.00053

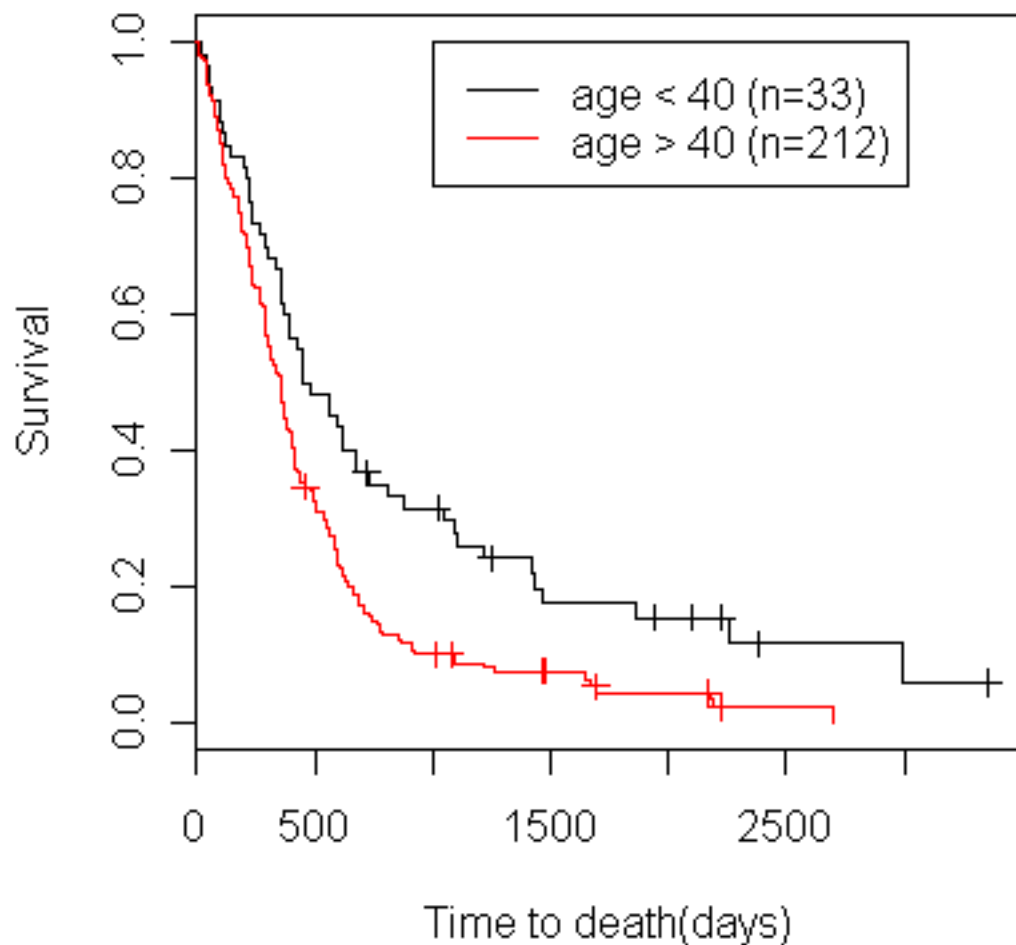


Figure 4

Age Difference in Pro+Mes+ProMes GBMs p-value= 0.91

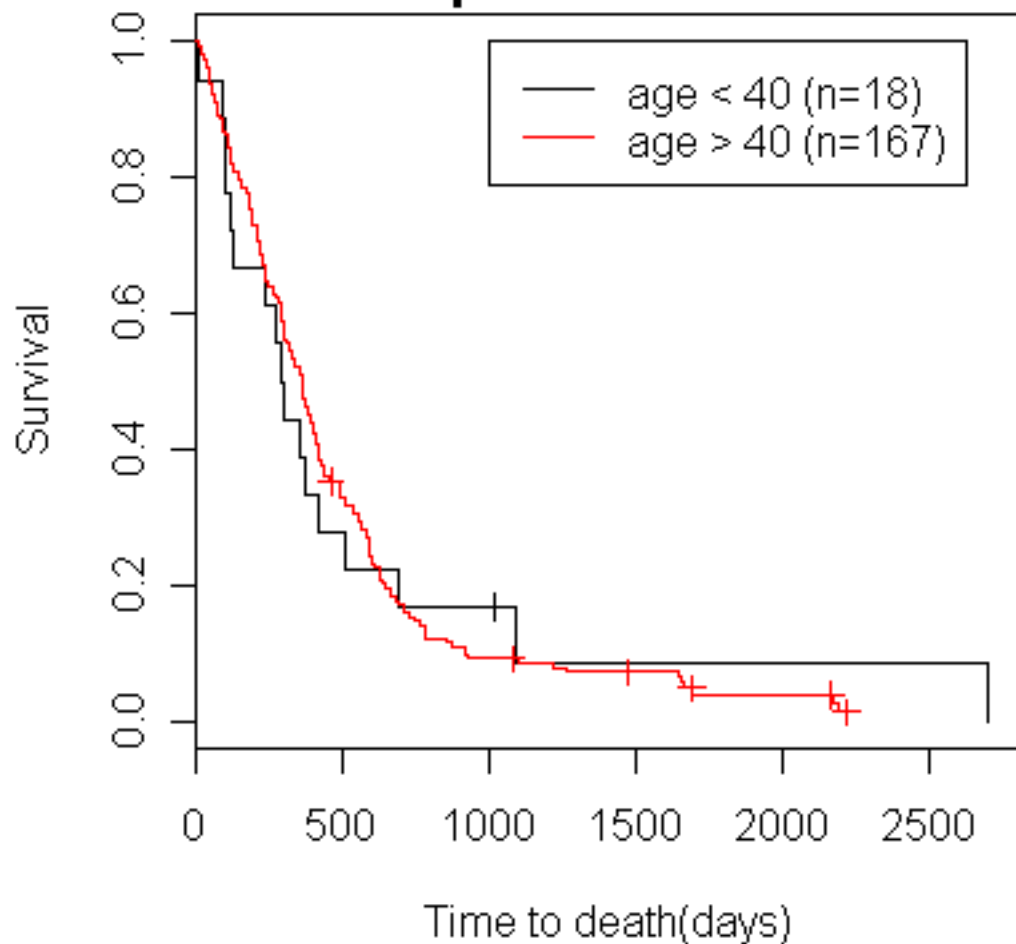


Figure 5

**Survival Difference between
PN vs. Pro+Mes+ProMes GBMs
in subjects age < 40
p-value= 0.01**

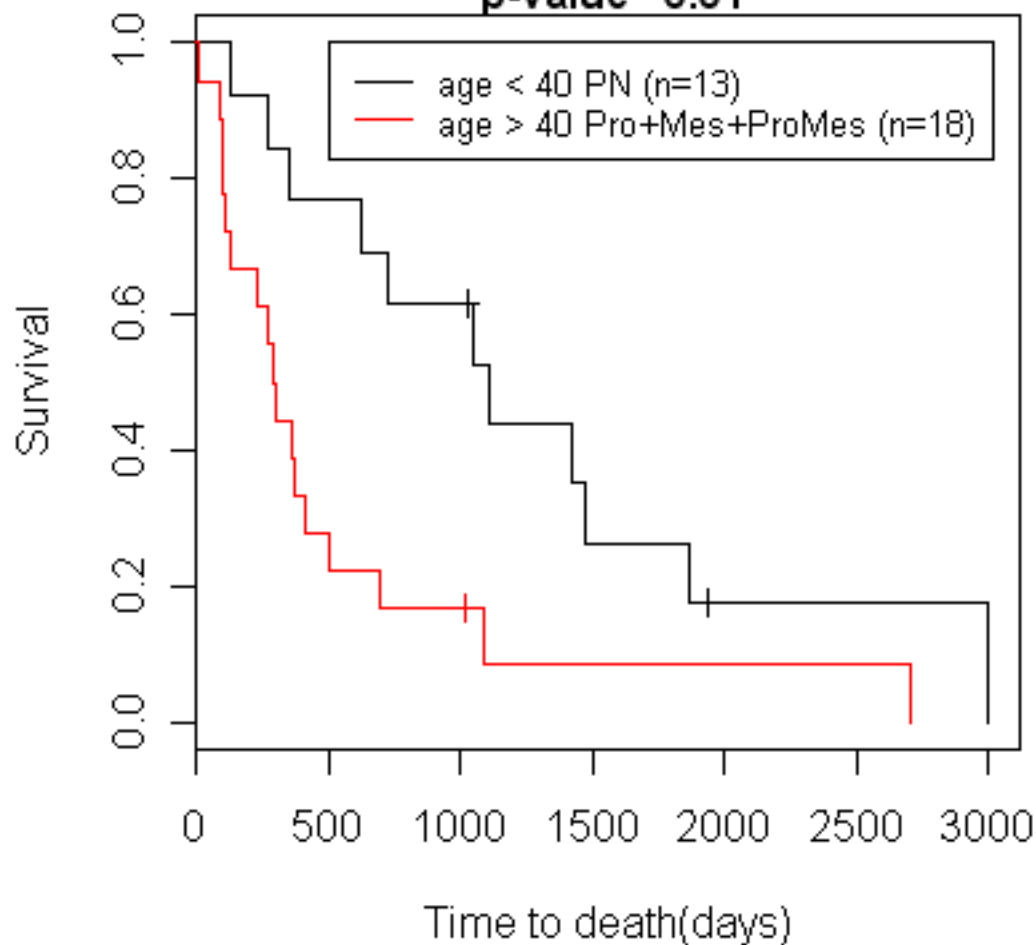


Figure 6

**Survival Difference between
PN vs. Pro+Mes+ProMes GBMs
in subjects age > 40
p-value= 0.046**

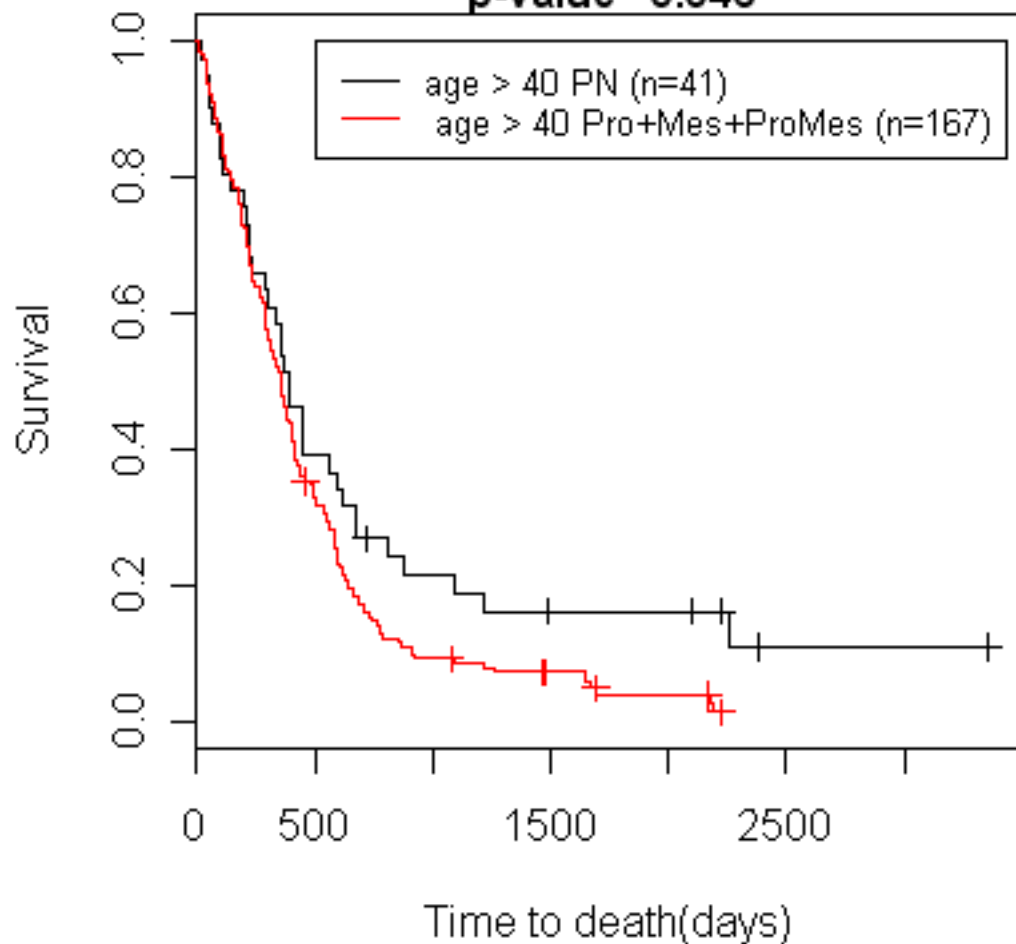


Figure 7

Additional files provided with this submission:

Additional file 1: supplementary_information_table_s1.xls, 318K

<http://www.biomedcentral.com/imedia/2096100982187190/supp1.xls>

Additional file 2: supplementary_information_section_s2.doc, 24K

<http://www.biomedcentral.com/imedia/2108698706187190/supp2.doc>

Additional file 3: supplementary_information_section_s3.xls, 214K

<http://www.biomedcentral.com/imedia/1478135415187190/supp3.xls>

Additional file 4: supplementary_information_section_s4.doc, 42K

<http://www.biomedcentral.com/imedia/1510882268187190/supp4.doc>