

1 Title: Normal tissue content impact on the GBM molecular  
2 classification

3 AUTHORS: Rodrigo Madurga<sup>1,2,3</sup>, Noemí García-Romero<sup>1,2</sup>, Beatriz Jiménez<sup>1,2</sup>, Ana  
4 Collazo<sup>1,2</sup>, Francisco Pérez-Rodríguez<sup>1,2</sup>, Aurelio Hernández-Lain<sup>4</sup>, Carlos Fernández-  
5 Carballal<sup>5</sup>, Ricardo Prat-Acín<sup>6</sup>, Massimiliano Zanin<sup>7,8</sup>, Ernestina Menasalvas<sup>8,9</sup>, Ángel  
6 Ayuso-Sacido<sup>10, 11, 12\*</sup>.

7 1. Fundación de Investigación HM Hospitales, HM Hospitales, Madrid, Spain.

8 2. Instituto de Investigación Sanitaria HM Hospitales (IISHM), HM Hospitales, Madrid,  
9 Spain.

10 3. Faculty of Experimental Sciences, Universidad Francisco de Vitoria, Madrid, Spain

11 4. Unidad Multidisciplinar de Neurooncología, Hospital Universitario 12 de Octubre,  
12 Madrid, Spain.

13 5. Servicio de Neurocirugía, Hospital General Universitario Gregorio Marañón, Madrid,  
14 Spain.

15 6. Departamento de Neurocirugía, Hospital Universitario la Fe, Valencia, Spain.

16 7. Instituto de Física Interdisciplinar y Sistemas Complejos IFISC (CSIC-UIB), Campus UIB,  
17 07122 Palma de Mallorca, Spain.

18 8. Centro de Tecnología Biomédica, Universidad Politécnica de Madrid, Madrid, Spain.

19 9. ETS de Ingenieros Informáticos, Universidad Politécnica de Madrid, Madrid, Spain

20 10. Fundación Vithas, Vithas Hospitals, Madrid, Spain.

21 11. Formerly: Fundación de Investigación HM Hospitales, HM Hospitales, Madrid, Spain.

22 12. Facultad de Medicina (IMMA), Universidad San Pablo-CEU, Madrid, Spain.

23 \*. Corresponding author: +34 686966904, [ayusosacido@gmail.com](mailto:ayusosacido@gmail.com)

24 Authors description:

25 Rodrigo Madurga is a postdoc researcher in Biostatistics and Bioinformatics at  
26 Fundación de Investigación HM Hospitales and Professor at the Faculty of Experimental  
27 Sciences of the Universidad Francisco de Vitoria. His research focuses on computational  
28 and statistical modeling in cancer genomics.

29 Noemí García-Romero is a postdoc researcher in Molecular Biology at Fundación de  
30 Investigación HM Hospitales. Her research focuses on brain tumors, particularly on  
31 Glioblastoma.

32 Beatriz Jiménez is an oncologist physician specialized in brain tumors at HM Hospitales.

33 Ana Collazo is an oncologist physician specialized in brain tumors at HM Hospitales.

34 Francisco Pérez-Rodríguez is a pathologist at HM Hospitales.

35 Aurelio Hernández-Lain is a pathologist at Hospital Universitario 12 de octubre.

36 Carlos Fernández-Carballal is an oncologist physician specialized in brain tumors at  
37 Hospital Universitario Gregorio Marañón.

38 Ricardo Prat-Acín is an oncologist physician specialized in brain tumors at Hospital  
39 Universitario la Fe.

40 Massimiliano Zanin is a researcher at the Institute for Cross-Disciplinary Physics and  
41 Complex Systems, Spain. His main topics of interest are Complex Networks and Data  
42 Science.

43 Ernestina Menasalvas is a full professor at the Department of Computer Systems  
44 Languages and Software Engineering at the Faculty of Computer Science of Universidad  
45 Politecnica de Madrid. Her subject area is Data Mining.

46 Ángel Ayuso-Sacido is The Head of the Brain Tumour Laboratory, Scientific Director at

47 Fundación de Investigación Vithas Hospitales and Professor at the Medial School of CEU  
48 San Pablo University

## 49 Abstract

50 Molecular classification of glioblastoma has enabled a deeper understanding of the  
51 disease. The four-subtype model (including Proneural, Classical, Mesenchymal and  
52 Neural) has been replaced by a model that discards the Neural subtype, found to be  
53 associated with samples with a high content of normal tissue. These samples can be  
54 misclassified preventing biological and clinical insights into the different tumor subtypes  
55 from coming to light.

56 In this work, we present a model that tackles both the molecular classification of  
57 samples and discrimination of those with a high content of normal cells.

58 We performed a transcriptomic in silico analysis on GBM samples (n = 810) and tested  
59 different criteria to optimize the number of genes needed for molecular classification.

60 We used gene expression of normal brain samples (n = 555) to design an additional gene  
61 signature to detect samples with a high normal tissue content. Microdissection samples  
62 of different structures within GBM (n = 122) have been used to validate the final model.

63 Finally, the model was tested in a cohort of 43 patients and confirmed by histology.

64 Based on the expression of 20 genes, our model is able to discriminate samples with a  
65 high content of normal tissue and to classify the remaining ones. We have shown that  
66 taking into consideration normal cells can prevent errors in the classification and the  
67 subsequent misinterpretation of the results. Moreover, considering only samples with a  
68 low content of normal cells, we found an association between the complexity of the  
69 samples and survival for the three molecular subtypes.

## 70 Introduction

71 Glioblastoma (GBM) is the most lethal brain tumor with a median overall survival (OS)  
72 of 15 months and incidence rate of 3-4 new diagnosed cases per 100,000 population  
73 [1,2].

74 In the last decade, there has been increasing interest in the molecular classification of  
75 GBM [3]. In 2010, a four-subtype classification model (Proneural (PN), Classical (CL),  
76 Mesenchymal (ME) and Neural (NE)) was proposed, based on the expression levels of  
77 840 genes [4]. This classification has widely been used to analyze differences in  
78 treatment response patterns of different GBM subgroups [5,6].

79 However, when an unsupervised clustering of the samples was performed against  
80 tumoral-related genes, a three cluster GBM classification was obtained by analyzing the  
81 expression of 150 genes [7]. Notably, each cluster was strongly associated with one  
82 group of the four-subtype model, except for the NE subtype. One possibility is that this  
83 subtype includes samples with a high content of normal cells. In fact, at the infiltrative  
84 margins of GBM, normal cells have been found to far outnumber tumoral ones and the  
85 NE subtype is associated with this region [8].

86 This abundance of normal cells is a problem for transcriptional classification of samples,  
87 as long as RNA expression levels are affected by tumor purity [9]. The content of normal  
88 cells in a sample affects transcriptional classification, complexity of the sample  
89 measured by the simplicity score [7] and intratumoral heterogeneity, among others.

90 Different algorithms can calculate a purity score of tumor samples from CNV data [10]  
91 or from gene expression signatures [11]. However, few models are able to tackle both  
92 the molecular classification of GBM and the tumor purity of the sample.

93 In the present work, we develop a model that classifies GBM samples as PN, CL or ME  
94 and provides us with information about the abundance of normal cells in the samples  
95 based on the expression of 20 genes. This model not only integrates the molecular and  
96 purity classification of the samples but does so in a cost-effective way.

## 97 Materials and methods

98 Gene expression data processing and normalization

99 GBM IDH wt cohort (IDHWT)

100 Gene expression data from GBM patients with known IDH status were collected from  
101 TCGA [12], GlioVis [13] and from the Gene Expression Omnibus. Affymetrix data sets  
102 were normalized using robust multi-array average normalization (RMA) followed by  
103 quantile normalization as implemented in the 'affy' package for R/Bioconductor [14].  
104 Affymetrix data sets consisted in TCGA (n = 528) [12], GSE4271 (n = 76) [15] and  
105 GSE36245 (n = 46) [16]. Additionally, RNA-seq data were downloaded from GSE48865  
106 (n = 100) [17] and GSE121720 (n = 60). Collected RNA-seq had been mapped to the hg19  
107 human genome and log transformed. Because data sets were generated on different  
108 platforms and by different labs, we used ComBat to address the strong batch effects  
109 expected from such variable data sources [18]. We applied aggregation workflow, as  
110 described elsewhere [19], to select the probe set that represents each gene in each data  
111 set. Once different data sets had been aggregated, IDH wt and CIMP<sup>-</sup> samples were  
112 filtered obtaining a final cohort of 551 samples. Where the CIMP status was unknown,  
113 this was determined by the support vector machine, using the TCGA cohort as training  
114 data set [13]. The final cohort was divided in a training cohort (n = 367) and a validation  
115 cohort (n = 184).

116 Histology cohort (HIS)

117 Solid surgical tissue samples were obtained from patients operated in HM Hospitales,  
118 Madrid, Spain; Hospital General Universitario Gregorio Marañón, Madrid, Spain and  
119 Hospital Universitario la Fe, Valencia, Spain. A total of 43 GBM IDH wt patients were  
120 analyzed by qRT-PCR.

121 NormalBrain cohort (NB)

122 Normalized microarray (custom-designed Agilent 8x60K) gene expression from 6 human  
123 brains were downloaded from the ALLEN Human Brain Atlas ([http://human.brain-](http://human.brain-map.org/static/download)  
124 [map.org/static/download](http://human.brain-map.org/static/download)). One hundred samples of cerebral cortex were randomly  
125 selected from each brain. Outliers were removed using principal component analysis  
126 (PCA), which consisted in considering the first two principal components and marking all  
127 the samples with a distance greater than 2.5 as outliers. After outlier removal, to avoid  
128 batch effects we applied ComBat [18] to the NTB cohort using the TCGA cohort  
129 (including 10 normal samples) as a reference. Finally, the NTB cohort was divided into  
130 NTB-training (n = 370) and NTB-validation (n = 185).

131 Ivy GAP cohort (IVYGAP)

132 The Ivy Glioblastoma Atlas Project (Ivy GAP) analyzed the transcriptome of different  
133 anatomical structures from 10 different tumors [20]. The normalized read counts of  
134 these 122 samples were downloaded from the Ivy GAP portal  
135 (<http://glioblastoma.alleninstitute.org/static/download.html>) and log transformed.

136 Differentially expressed genes

137 Differentially expressed genes were identified using the R/Bioconductor package  
138 'multtest' [21] with 5000 bootstrap iterations and using FDR as method to control type  
139 one error rate. The significance level was set to 0.05.

140 Reduction of the number of genes in the gene signature  
141 The 50-gene signatures proposed elsewhere [7] were used for the subtype classification  
142 of the IDHWT training set in PN, CL or ME. These results were used as a reference in the  
143 following analysis. Afterwards, the gene signature of each subtype was reduced, one  
144 gene at a time, to a 2-gene signature. For each step, 1000 different combinations of  
145 signature genes were randomly generated and used for the subtype classification of the  
146 training cohort. The overlap of these classifications with the 50-gene signature  
147 classification was estimated, as well as the accuracy and precision achieved by the  
148 reduced gene signatures for each subtype.

149 Along with the randomly selected genes, three different criteria were used to rank genes  
150 inside each of the three gene-signatures, keeping the top ranked genes in each step of  
151 gene removal. These criteria were difference and relative difference in gene expression  
152 and statistical significance, measured by the value of the statistic, obtained when  
153 comparing the expression levels between subtypes. The overlap of the classification  
154 obtained from these criteria was compared to that obtained from the randomly selected  
155 gene signatures for each size of gene signature.

156 Molecular classification

157 Two methods for the classification of GBM samples have been used through this work:

158 Single sample gene set enrichment analysis (ssGSEA)

159 Single sample gene set enrichment analysis defines an enrichment score for a gene  
160 signature, in this case representative of a subtype, for each sample within a dataset. The  
161 process starts with the rank-normalization and rank-order of gene expression for a given  
162 sample. A statistic is then calculated from the difference between the cumulative  
163 empirical distribution functions (ECDF) of the gene signature and the remaining genes

164 [7]. A null distribution of the enrichment scores for each signature is obtained as follows.  
165 A large number of virtual samples (> 10,000) are generated assigning to each gene the  
166 expression level of the same gene in a randomly selected sample in the dataset. Null  
167 distributions are used to give an empirical p value to the enrichment scores obtained for  
168 each sample in the dataset [7]. A given sample is classified with the subtype with the  
169 lowest empirical p value.

170 Centroid-based classification

171 Verhaak's classification is based on a 210 gene signature for each of the four subtypes  
172 [4]. ClaNC software [22] was used to assign a category to the training samples based on  
173 Verhaak's gene signature.

174 CNA and mutations from the TCGA cohort

175 CNA and mutational information about samples from the TCGA cohort were  
176 downloaded from the TCGA repository using TCGAbiolinks package from Bioconductor  
177 [23].

178 Survival analysis

179 The optimal cutoff of the simplicity score for the survival analysis was determined by the  
180 *get.cutoff()* function described elsewhere [24]. The method used for the cutoff  
181 optimization was survival significance, where a Cox proportional hazard model is fitted to  
182 the dichotomized simplicity score and the survival variable. The point with the most  
183 significant split (measured by the log rank test) is defined as the optimal cutoff.

184 Histology

185 Six tissue microarrays (TMAs) were constructed from 32 Formalin Fixed Paraffin  
186 Embedded (FFPE) tissues using an arraying instrument (GALILEO CK 3500). From each  
187 tissue block a total of three tissue cores were made with a diameter of 0.6-1 mm. Then,



188 TMA blocks were cut at 4  $\mu\text{m}$  and stained with hematoxylin & eosin (H&E). The H&E  
189 stained tissue blocks were evaluated by a pathologist for tumoral cells density,  
190 abundance of pathogenic blood vessels and presence of necrosis.

191 RNA isolation and qRT-PCR

192 Total RNA was isolated using RNeasy Mini or Micro kit (QIAGEN) following the  
193 manufacturer's recommendations. cDNA synthesis (High-Capacity cDNA Reverse  
194 Transcription Kit; Applied Biosystems) was performed from one  $\mu\text{g}$  of RNA. An optical  
195 384-well plate equipped with an ABI PRISM 7900 HT sequence detection system  
196 (Applied Biosystems) was used for the performance of qRT-PCR reactions using SYBR  
197 Green. Two housekeeping genes were used to normalize data, the primers used for each  
198 gene can be seen in Supplementary Table S1. Because the ssGSEA classification is based  
199 on the ranking of genes rather than on absolute expression, we used the Ct value of the  
200 genes, which were scaled within the same sample:

201 
$$Z = \frac{r_i - 1}{n_{genes} - 1}$$

202 Statistical analysis

203 All the statistical analyses have been performed using R software. The statistical tests  
204 applied are indicated in the text. When relevant, p values were adjusted using Benjamin  
205 and Hochberg method.

## 206 Results

207 Transcriptomic data aggregation

208 As explained in the Materials and Methods, the initial cohort, comprised of 810 patients,  
209 was filtered to discard IDH mt or GCIMP<sup>+</sup> patients, which are already known to have  
210 favorable clinical outcomes [25,26]. A total of 86 patients were IDH mutant and 163

211 were IDH unknown. It is worth mentioning that IDH status of Phillips' cohort  
212 (GSE121720) were obtained elsewhere [27]. After filtering, a final cohort of 551 IDH wt  
213 / GCIMP<sup>-</sup> patients was obtained and divided, randomly, into two cohorts: IDHWT training  
214 cohort (n = 367) and IDHWT validation cohort (n = 184), as shown in Figure 1A.

215 5-Genes signature for molecular Glioma subtypes.

216 The IDHWT training cohort was classified into PN, CL and ME as proposed elsewhere [7].

217 The simplicity score proposed in the cited work was estimated. This gives a value

218 between 0 and 1, with high values corresponding to samples activating a single subtype

219 and low values for samples activating multiple subtypes. Afterwards, to performed a

220 clean comparison between subtypes, samples with a simplicity score higher than 0.95

221 (n = 121) were selected and a differential expression analysis was performed between

222 subtypes for the 150 genes involved in the classification process. Because the genes for

223 the original classification were selected on the basis of differences in gene expression

224 between groups [7], we used the results of the previous analysis to rank the gene

225 signatures according to differences and relative difference in mean expression and

226 statistical significance between subtypes.

227 The three rank criteria were used separately to remove the last gene from the gene

228 signature of each subtype at a time. The reduced gene signatures were used to classify

229 the IDHWT training cohort and the overlap with the original classification was estimated.

230 Additionally, 1,000 randomly ranked gene lists were generated and used to evaluate the

231 overlap with the original classification, in order to generate a null distribution for each

232 size of the gene signature.

233 To delucidate if the above mentioned rank criteria generate reduced gene signatures

234 with better performance than random, we used the overlap mean and standard

235 deviation obtained from the null distribution to scale the results obtained with the three  
236 different ranking criteria (see Figure 1B and SI Appendix Figure S1). Colored symbols  
237 represent the scaled overlap of the difference in mean ranking criterion and dashed line  
238 represent the mean Z-score obtained for each criterion. It can be appreciated that  
239 differences in means criterion reach higher Z-scores for different gene signature sizes  
240 than the other two criteria. It is worth mentioning that with just 5 genes per subtype an  
241 overlap higher than 90% is obtained and that when the number of genes is reduced to  
242 2 the overlap is still higher than 80%, with Z-scores higher than 2 in both cases (Figure  
243 1B). Comparing the classification by subtype, using 5 genes per subtype the true positive  
244 ratio (TPR) reached values higher than 90% for the CL and ME subtypes, and as high as  
245 83% for the PN subtype. On the other hand, the true negative ratio (TNR) reached values  
246 higher than 93% for the three subtypes (SI Appendix Figure S2A-C).

247 Considering these results, we decided to reduce the number of genes per subtype to the  
248 minimum for which an overlap with the original classification of at least 90% is achieved.  
249 Therefore, we used the difference in means expression rank criteria to reduce the gene  
250 signature of each subtype to 5 genes. At this point, we now have a 15 gene signature  
251 that classifies samples into 3 subtypes (5-gene signatures for each subtype). This gene  
252 signature seems to preserve an overlap close to 90% with the original signature of 150  
253 genes (50-gene signature for each subtype) and requires 1/10th of the molecular  
254 information used by the previous method.

255 Prediction of samples with high content of normal tissue

256 As mentioned above, the classification proposed by Wang et al. [7] stratifies samples  
257 into 3 subtypes with a high correspondence with 3 of the 4 subtypes previously  
258 proposed by Verhaak et al. [4]. The remaining subtype, called Neural (NE), was

259 hypothesized to include samples with a high content of normal tissue [7]. Under this  
260 assumption, we asked ourselves whether the presence of normal tissue could affect the  
261 classification of the samples. Therefore, we used the centroid-based classification,  
262 proposed by Verhaak et al. [4], to classify the training cohort and study the distribution  
263 of the simplicity score, obtained as proposed by Wang et al. [7], in the four subtypes.  
264 Figure 1C shows that the NE subtype is enriched for lower simplicity scores (Wilcoxon  
265 test,  $p \leq 0.001$ ). Because the simplicity score is based on the distances to the dominant  
266 subtype and between non-dominant subtypes [7], a low value can be obtained either if  
267 a non-dominant subtype is close to the dominant one, indicative of a complex sample;  
268 or if the dominant subtype is weakly activated, which may occur when the expression  
269 levels of the tumoral cells are masked by a high content of normal cells [9].

270 We started to search for genes of the NE signature that were overexpressed in the NE  
271 subtype compared with the non-neural samples, with a simplicity score higher than 0.95.  
272 This analysis gave 46 overexpressed genes (see Supplementary Table S2). We then used  
273 the NTB-training cohort (see Material and Methods) and analyzed the differential  
274 expression of the 46 overexpressed genes between normal tissue and non-neural  
275 samples with high simplicity scores ( $ss > 0.95$ ). The analysis resulted in 35 overexpressed  
276 genes in normal tissue (Supplementary Table S2). We then selected the top 5 of these  
277 genes in relation to differences in mean expression and formed a fourth gene signature  
278 (*CCK*, *CRYM*, *SERPINI1*, *KCNK1* and *GPR22*). We used Enrichr [28,29] to analyze which  
279 tissues were enriched in this new gene signature. From the ARCHS4 Tissues library six  
280 different structures, all of them from brain, were significantly enriched ( $p$  value  $< 1e-4$ ,  
281  $q$  value  $< 0.001$ ) (see Supplementary Table S3). This new gene signature was added to

282 the previously reduced gene signatures (Figure 1A) with the intention of detecting  
283 samples with a high content of normal tissue.

284 We used our 20-gene based signature to classify the samples into the different subtypes,  
285 or as normal tissue abundant (NT). Figure 1D shows the distribution of the simplicity  
286 scores from Wang's classification, for the different subtypes obtained with the reduced  
287 gene signature. It can be observed that the NT is enriched for lower simplicity scores  
288 (Wilcoxon test,  $p < 1e-4$ ) as was also found for the NE subtype (Figure 1C). Although the  
289 result is similar to that obtained for the NE subtype, we still need to prove that samples  
290 classified as NT have a high content of normal brain cells. In parallel, it is worth proving  
291 that the addition of a fourth subtype did not alter the overlap when compared with the  
292 50-gene signature classification.

293 Validation of the model

294 We performed splitting iteration 1000 times to generate different validation cohorts.  
295 Each of these cohorts were used to study the overlap between the classification  
296 performed by the reduced gene signature and the original gene signature (Figure 2A). If  
297 the samples classified as NT are not considered, a mean overlap of 89 % and 2% standard  
298 deviation are achieved.

299 When the results were analyzed by subtype, we found that our classification model, with  
300 5-gene signatures, is an excellent predictor of the results that would be obtained by the  
301 50-gene signature model as can be appreciated in the ROC space shown in Figure 2B.

302 For all the subtypes, the model reaches a high sensitivity:  $0.94 \pm 0.04$ ,  $0.89 \pm 0.04$  and  
303  $0.87 \pm 0.03$  for PN, CL and ME respectively, and high specificity:  $0.95 \pm 0.02$ ,  $0.94 \pm 0.02$   
304 and  $0.95 \pm 0.02$  for PN, CL and ME respectively. As a measurement of the accuracy for  
305 each subtype we estimate the F1-score, with values of  $0.90 \pm 0.03$ ,  $0.89 \pm 0.03$  and  $0.89$

306  $\pm 0.02$  for PN, CL and ME respectively. In synthesis, the results presented here show that  
307 a significant reduction in the number of genes used for the classification does not  
308 dramatically affect the performance of the classification.

309 To confirm that the reduced gene signature gives valuable information we used the  
310 TCGA cohort, for which genomic data is available, to study the incidence of genomic  
311 alterations in the different subtypes. Figure 2C shows that the characteristic genomic  
312 alterations for each subtype are still found when samples are classified using the  
313 reduced gene signature.

314 Wang et al. [7] found that the ME subtype shows a reduced OS for single subtype  
315 activated samples, that is, samples with a simplicity score higher than 0.99 (~20% of the  
316 samples). We observed that the simplicity score calculated from the empirical p values  
317 obtained from our model was lower than that obtained from Wang's model. Only four  
318 samples had a simplicity score higher than 0.99. However, if we select the top 20% of  
319 the simplicity scores of the samples the same result is obtained as shown in Figure 2D  
320 (log rank test,  $p = 0.03$ ).

321 NT associates with abundance in normal brain cells

322 Once the reduced gene signatures have been proven to be a good predictor of the  
323 original classification, it is time to address whether or not NT is identifying samples with  
324 a high content of normal cells. We used the ABSOLUTE method [10], which gives a  
325 tumoral purity score based on copy number variation data to study tumoral purity of  
326 the samples in the different subtypes. The TCGA cohort was used for this analysis as it is  
327 the only one for which CNV data were available. Additionally, we used the ESTIMATE  
328 method [11], which gives a tumoral purity score based on the enrichment scores  
329 obtained for an immunological and a stromal gene signature, for the same purpose.

330 Figure 3A shows the ABSOLUTE and ESTIMATE purity scores obtained for the TCGA  
331 cohort. As reported previously [7], PN and CL subtypes showed higher scores than ME  
332 for both algorithms (Wilcoxon test,  $p < 1e-15$ ). This result was shown to be due to the  
333 higher infiltration of immunological cells occurring in ME tumors [7]. Our results showed  
334 that PN and CL subtypes also obtained a higher ABSOLUTE purity score than NT  
335 (Wilcoxon test,  $p = 1e-5$ ) indicating a higher content of normal brain cells in NT samples.  
336 However, no significant differences were observed for the ESTIMATE purity score.  
337 Therefore, there was no increase in the amount of immunological or stromal cells. This  
338 result was confirmed by repeating the analysis for the IDHWT validation cohort. The  
339 results can be seen in Figure 3B, with no significant differences between PN and CL  
340 subtypes compared with NT, although the ME subtype shows significantly lower values  
341 (Wilcoxon test,  $p = 1.9e-14$ ).

342 These results are in line with the hypothesis that NT are samples with a high content of  
343 normal cells. This hypothesis is also supported by the classification performed on the  
344 NTB-validation cohort. Of the 185 normal cortex samples, 184 (99.5 %) were classified  
345 as NT and only 1 sample was classified as tumoral, in this case PN.

346 To test NT in a cancer context we used the IVYGAP cohort, which consists of different  
347 tumoral structures obtained from GBM biopsies subjected to laser microdissection.  
348 Briefly, the IVYGAP cohort is composed of five different structures as defined elsewhere  
349 (<http://help.brain-map.org/display/glioblastoma/Documentation>): cellular tumor (CT)  
350 has the most core tumor with a tumor cell to non-tumor cell ratio between 100/1 and  
351 500/1, microvascular proliferation (MVP) are regions characterized by two or more  
352 vessels sharing common vessel walls, pseudopalisading cells around necrosis (PAN),  
353 which are generally found in the core tumor, infiltrating tumor (IT) that corresponds to

354 the intermediate region between the cellular tumor and the leading edge and has a  
355 tumor cell to non-tumor cell ratio of 10-20/100; and the leading edge (LE) that is the  
356 boundary of the tumor with a tumor cell to non-tumor cell ratio of 1-3/100. We classified  
357 the 122 samples of the IVYGAP cohort using the original 50-gene signatures and the  
358 reduced 5-gene signatures. The fraction of each subtype by structure is shown in Figure  
359 3C. Interestingly, all the LE structures were classified as NT as well as almost 75% of the  
360 IT samples. Most of these structures were classified as PN or CL by the 50-gene  
361 signatures. Besides, only one CT sample was classified as NT showing that the gene-  
362 signature detects samples with a high content of normal cells with high precision. It is  
363 worth mentioning that for CT, MVP and PAN structures, the results were highly  
364 coincident between the 5-gene and 50-gene signatures, with CT corresponding mainly  
365 to CL or PN and PAN and MVP mostly to ME, as reported elsewhere [20].

366 We used the results of the 50-gene signatures to obtain the simplicity scores of each  
367 sample. Figure 3D shows that LE is significantly enriched for lower simplicity scores in  
368 comparison to the other structures (Wilcoxon test,  $p < 0.02$ ). This result confirms that  
369 the presence of a high content of normal cells in a sample can affect interpretation of  
370 the results.

371 The ssGSEA classification system performed a random permutation of the experimental  
372 data to generate a null distribution, obtaining a p value for the association of a sample  
373 to each subtype. Figure 3E shows the p value for NT obtained from samples of different  
374 structures. LE, the structure with the lowest content of tumoral cells, has significantly  
375 lower p values (Wilcoxon test,  $p < 1e-7$ ). The p values increase slightly for IT and show  
376 median values close to 0.8 for structures found in the core tumor. Therefore, the



377 association with the NT subtype increases with the content of normal brain cells in the  
378 sample.

379 Two pathologists independently classified the HIS cohort according to: tumoral cell  
380 density, abundance of pathogenic blood vessels and presence of necrosis.  
381 Simultaneously, we used the expression levels obtained by qRT-PCR to classify the same  
382 cohort (Supplementary Table S4). Comparing the results, we found that of the 7 samples  
383 classified as infiltrating tumor by the pathologists, 5 (71%) were now classified as NT.  
384 The specificity of the NT class was 82%. We also observed that samples classified as CL  
385 and PN were indistinguishable in relation to the histological parameters. However,  
386 samples with an absence of pathogenic blood vessels and necrosis were mostly not  
387 classified as ME (80%). Figure 3F shows representative H&E stained histological images  
388 for each subtype.

389 Taken together, these results show that samples classified as NT have a lower tumoral  
390 cell density which is not due to immune cell infiltration, and that NT associates with  
391 samples or tumor regions with low cellularity.

#### 392 Survival analysis

393 To study the clinical relevance of the simplicity score obtained by our model, we  
394 classified all the samples from the IDHWT cohort and discarded those that fell into the  
395 NT group. We, then, evaluated the optimal cutoff for the simplicity score considering  
396 the hazard ratio for each subtype using the *get.cutoff()* function described elsewhere  
397 [24]. Results can be seen in Supplementary Figure S3. Using the corresponding cutoff to  
398 divide samples into the PN and CL subtypes we observed a significantly higher survival  
399 for the simpler samples. The difference in median survival was found to be around 8  
400 months in the PN subtype (HR = 0.54, 95% CI 0.32 – 0.93, log rank p-value = 0.02, Figure

401 4A) and 15.6 months for the CL subtype (HR = 0.38, 95%CI 0.21 – 0.7, log rank p-value =  
402 0.001, Figure 4B). The opposite was observed for the ME subtype in that the simpler  
403 samples presented a poorer survival, with a difference in median survival of 2 months  
404 (HR = 2.18, 95% CI 1.23 – 3.88, log rank p-value = 0.007, Figure 4C).

## 405 Discussion

406 We developed a model based on the expression of 20 genes for the molecular  
407 classification of GBM samples. This model can detect samples with a high content of  
408 normal tissue, classifying them as NT, and also classifies the samples into PN, CL or ME.  
409 Although it uses only 5 genes per subtype, our model showed an overlap of 87% with  
410 the 50 gene per subtype model proposed elsewhere [7]. It also detects the main  
411 characteristic genetic alterations of the different subtypes [4,25] and the difference in  
412 OS compared with subtypes for simpler samples [7]. These results show that molecular  
413 classification of GBM can be performed in a cost-effective way and we hope that this  
414 model will encourage researchers and physicians to use this classification method more  
415 frequently in the future.

416 The NT gene signature shows 71% sensitivity and 82% specificity in the HIS cohort. On  
417 the other hand, for microdissected samples of specific tumoral regions, where the  
418 above-mentioned variability is absent, we found that NT classification is strongly  
419 associated with samples from the boundary region of the tumor (83.7% sensitivity,  
420 98.7% specificity). This result is consistent with the association found between this  
421 region and the NE subtype [8,20]. It is noteworthy that the strength of the association  
422 between the different tumoral structures and NT, measured as an empirical p value,  
423 increases with the percent of normal cells. It remains to be tested whether or not the

424 empirical p value of NT can be used as an estimator of the percentage of normal brain  
425 cells in the sample. Further analyses are required to establish this.

426 We observed that NT samples are associated with low simplicity scores. Simplicity score  
427 was proposed as an estimate of the complexity of the GBM sample [7], where a low  
428 simplicity score indicates that the samples do not present a unique predominant  
429 subtype. Our results show that samples with a high content of normal cells can lead to  
430 an erroneous classification of a sample, considering it as one with high tumoral  
431 complexity, when actually it corresponds to a sample with low tumor purity. In fact,  
432 single cell RNA-seq analyses reveal that the subtype of a bulk tissue sample coincides  
433 with the subtype of the dominant cell population in the sample [30,31]. Moreover, it is  
434 important to know if the sample used for RNA extraction comes from the boundary of  
435 the tumor, because a low tumoral purity can alter gene expression measurements [9].

436 The incorporation of NT as a quality parameter of samples that are going to be classified  
437 brings important clinical and biological advantages, as discussed below. From a clinical  
438 point of view, Gill et al. suggested that the boundary of the tumor has to be classified as  
439 this is the tumoral region which cannot usually be resected [8]. The model we propose  
440 here not only classifies the infiltrating tumor mainly as NT, but also indicates the  
441 molecular subtype of the tumor as the second dominant group in that sample, which in  
442 88% of the cases matches the molecular subtype of the cellular tumor. Therefore, when  
443 the piece of tissue used for molecular classification comes from the infiltrating tumor,  
444 our model can detect the molecular subtype of the tumor at that time. However, if the  
445 sample has a high content in normal cells and is classified as NT, parameters like the  
446 simplicity score should not be taken into consideration and special care should be taken  
447 in the interpretation of experiments like gene expression measurements.

448 It is also relevant that when different cellular tumor sections of the same patients were  
449 analyzed we found that they were either PN or CL for all sections from the same patient.  
450 Puchalski et al. maintained that CT were PN, CL or NE [20]. Their results agree with ours  
451 if we consider NE to not be a real subtype. As reported in the cited work, we found that  
452 microvascular proliferation regions and regions around necrosis were mainly ME. In the  
453 same line, the ME subtype was reported to express markers of hypoxia and  
454 microvasculature [32].

455 We can, therefore, regard GBM as a PN or CL tumor that evolves to ME in response to  
456 different inputs, i.e. hypoxia [32]. Initially, the mesenchymal transition occurs in small  
457 regions of the tumor, but these grow and eventually become the predominant subtype  
458 (see SI Appendix Figure S4). This hypothesis is in line with our survival results. We  
459 observed that complex samples, those with no clear dominant subtype, showed a worse  
460 survival for PN and CL tumors; that is, when the ME regions of the tumor begin to grow  
461 the complexity of the tumor increases with the corresponding survival consequence.  
462 However, complex samples showed longer survivals for ME tumors, that is, a complex  
463 ME sample is the continuation of the evolution of complex PN or CL samples, but when  
464 the tumor becomes mainly ME it shifts to a low complexity ME sample with a worse  
465 survival. Further analyses are needed to confirm this hypothesis.

466 Sottoriva et al. reported the intratumoral heterogeneity of GBM after observing that 6  
467 out of 10 patients present regions of the same tumor with 2 or 3 different molecular  
468 subtypes [33] Nevertheless, recovering the idea that NE is not actually a tumor subtype,  
469 their results will be transformed into 5 out of 10 patients with 2 different molecular  
470 subtypes: 2 cases with ME and PN and 3 cases with ME and CL. The proposed hypothesis  
471 is in line with the shift of the ME subtype upon glioma recurrence [7,15]. Several works

472 have studied the PN-ME transition [34,35] but there is no evidence for a CL-ME  
473 transition.

474 We believe that our classification model, which takes into consideration tumor samples  
475 with a high content of normal tissue, will help to provide clinical insight into the different  
476 molecular subtypes of GBM and to better understand their biology.

## 477 Conclusions

478 In summary, we have developed a model which tackles both the classification of GBM  
479 samples into PN, CL or ME, and the detection of a high content of normal cells in a  
480 sample. The model shows an overlap of over 85% with the one proposed by Wang et al.  
481 and only requires the expression levels of 20 genes, making it a cost-effective alternative  
482 to other molecular classification models. The ability of our model to detect samples with  
483 high content of normal cells has been tested on microdissected regions of different GBM  
484 biopsies as well as on bulk tumor samples, contrasting the model results with the  
485 histological examination by two experts. We show the importance of determining the  
486 content of normal cells in GBM samples. Otherwise, normal tissue expression patterns  
487 can mask the expression patterns of other tumor types in the samples. This can lead to  
488 a misinterpretation of the results as we show with the simplicity score but can also affect  
489 the conclusions of tumor heterogeneity studies, among others.

## 490 Key points

- 491 • In this work, we present a cost-effective model based on the expression of 20  
492 genes, which can classify GBM samples into Proneural, Classical and  
493 Mesenchymal subtypes.

- 494 • The model incorporates a quality parameter that detects samples with a high  
495 content of normal tissue, preventing errors in the classification and  
496 interpretation of the results in clinical practice.
- 497 • Our results show that considering the abundance of normal cells in a sample can  
498 shed light on the interpretation of survival, tumor evolution or tumor  
499 heterogeneity.
- 500 • As the expression of 20 genes can be measured by qRT-PCR we believe that a  
501 greater volume of GBM samples will be classified and reported in the future.

## 502 Funding

503 This work was supported by grants from the ‘Fondo de Investigaciones Sanitarias’ (FIS)  
504 [PI17-01489], the Miguel Servet Program [CP11/00147] del Instituto de Salud Carlos III  
505 (AAS), and the Ministerio de Economía y Competitividad–FEDERER [RTC-2016-4990-1].  
506 RM was supported by “Convocatoria de ayudas para la contratación de investigadores  
507 predoctorales e investigadores postdoctorales cofinanciadas por Fondo Social Europeo  
508 a través del Programa Operativo de Empleo Juvenil y la Iniciativa de Empleo Juvenil  
509 (YEI)”.

## 510 Author contributions

511 R.M. participated in experimental design, performed the in-silico analysis and wrote the  
512 manuscript. N.G.R participated in experiments and revised the manuscript. A.H.L. and  
513 F.P.R. performed the histological examination of the samples. M.Z. and E.M. participated  
514 in the computational analysis, design and revision. B.J., A.C., C.F.C. and R.P provided  
515 patients samples. A.A.S participated in experimental design and was responsible for the  
516 financial support, and edited, revised and approved the final version of the manuscript.

## 517 Acknowledgments

518 We thank Victor González Rumayor for helping with the Tissue Microarray logistic and  
519 Marcos García Lorenzo from the Modeling and Virtual Reality Group of the Rey Juan  
520 Carlos University for helping with the computational equipment.

521 We want to particularly acknowledge the patients in this study for their participation  
522 and to the HGM BioBank integrated in RETICS, National Network Biobanks, and  
523 collaborating Centers for the generous gifts of clinical samples used in this work. The  
524 HGM BioBank, integrated in National Network Biobanks, is supported by Instituto de  
525 Salud Carlos III, Spanish Science and Innovation Ministry (Grant nº PT17/0015/0042)".

## 526 References

- 527 1. Louis DN, Perry A, Reifenberger G, et al. The 2016 World Health Organization  
528 Classification of Tumors of the Central Nervous System: a summary. *Acta Neuropathol*  
529 2016; 131:803–820
- 530 2. Ostrom QT, Gittleman H, Fulop J, et al. CBTRUS Statistical Report: Primary Brain and  
531 Central Nervous System Tumors Diagnosed in the United States in 2008-2012. *Neuro-*  
532 *oncology* 2015; 17 Suppl 4:iv1–iv62
- 533 3. Lee E, Yong RL, Paddison P, et al. Comparison of glioblastoma (GBM) molecular  
534 classification methods. *Seminars in Cancer Biology* 2018;
- 535 4. Verhaak R, Hoadley KA, Purdom E, et al. Integrated Genomic Analysis Identifies  
536 Clinically Relevant Subtypes of Glioblastoma Characterized by Abnormalities in PDGFRA,  
537 IDH1, EGFR, and NF1. *Cancer Cell* 2010; 17:98–110
- 538 5. Chen R, Smith-Cohn M, Cohen AL, et al. Glioma Subclassifications and Their Clinical  
539 Significance. *Neurotherapeutics* 2017; 14:284–297

- 540 6. Erdem-Eraslan L, Bent MJ van den, Hoogstrate Y, et al. Identification of Patients with  
541 Recurrent Glioblastoma Who May Benefit from Combined Bevacizumab and CCNU  
542 Therapy: A Report from the BELOB Trial. *Cancer Res* 2016; 76:525–534
- 543 7. Wang Q, Hu B, Hu X, et al. Tumor Evolution of Glioma-Intrinsic Gene Expression  
544 Subtypes Associates with Immunological Changes in the Microenvironment. *Cancer Cell*  
545 2017; 32:42-56.e6
- 546 8. Gill BJ, Pisapia DJ, Malone HR, et al. MRI-localized biopsies reveal subtype-specific  
547 differences in molecular and cellular composition at the margins of glioblastoma. *Proc*  
548 *National Acad Sci* 2014; 111:12550–12555
- 549 9. Ridder D de, Linden CE van der, Schonewille T, et al. Purity for clarity: the need for  
550 purification of tumor cells in DNA microarray studies. *Leukemia* 2005; 19:618–627
- 551 10. Carter SL, Cibulskis K, Helman E, et al. Absolute quantification of somatic DNA  
552 alterations in human cancer. *Nature Biotechnology* 2012; 30:413
- 553 11. Yoshihara K, Shahmoradgoli M, Martínez E, et al. Inferring tumour purity and stromal  
554 and immune cell admixture from expression data. *Nature Communications* 2013; 4:2612
- 555 12. Network TCGAR, Weinstein JN, Collisson EA, et al. The Cancer Genome Atlas Pan-  
556 Cancer analysis project. *Nat Genet* 2013; 45:1113
- 557 13. Bowman RL, Wang Q, Carro A, et al. GlioVis data portal for visualization and analysis  
558 of brain tumor expression datasets. *Neuro-oncology* 2016; 19:139–141
- 559 14. Gautier L, Cope L, Bolstad BM, et al. affy—analysis of Affymetrix GeneChip data at  
560 the probe level. *Bioinformatics* 2004; 20:307–315
- 561 15. Phillips HS, Kharbanda S, Chen R, et al. Molecular subclasses of high-grade glioma  
562 predict prognosis, delineate a pattern of disease progression, and resemble stages in  
563 neurogenesis. *Cancer Cell* 2006; 9:157–173



- 564 16. Sturm D, Witt H, Hovestadt V, et al. Hotspot Mutations in H3F3A and IDH1 Define  
565 Distinct Epigenetic and Biological Subgroups of Glioblastoma. *Cancer Cell* 2012; 22:425–  
566 437
- 567 17. Bao Z-S, Chen H-M, Yang M-Y, et al. RNA-seq of 272 gliomas revealed a novel,  
568 recurrent PTPRZ1-MET fusion transcript in secondary glioblastomas. *Genome Res* 2014;  
569 24:1765–1773
- 570 18. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data  
571 using empirical Bayes methods. *Biostatistics* 2007; 8:118–127
- 572 19. Guinney J, Dienstmann R, Wang X, et al. The consensus molecular subtypes of  
573 colorectal cancer. *Nat Med* 2015; 21:nm.3967
- 574 20. Puchalski RB, Shah N, Miller J, et al. An anatomic transcriptional atlas of human  
575 glioblastoma. *Science* 2018; 360:660–663
- 576 21. Pollard KS, Dudoit S, Laan MJ van der. *Multiple Testing Procedures: the multtest*  
577 *Package and Applications to Genomics*. Springer 2005;
- 578 22. Dabney AR. ClaNC: point-and-click software for classifying microarrays to nearest  
579 centroids. *Bioinformatics* 2006; 22:122–123
- 580 23. Colaprico A, Silva TC, Olsen C, et al. TCGAbiolinks: an R/Bioconductor package for  
581 integrative analysis of TCGA data. *Nucleic Acids Research* 2016; 44:e71–e71
- 582 24. Budczies J, Klauschen F, Sinn BV, et al. Cutoff Finder: A Comprehensive and  
583 Straightforward Web Application Enabling Rapid Biomarker Cutoff Optimization. *Plos*  
584 *One* 2012; 7:e51862
- 585 25. Brennan CW, Verhaak R, McKenna A, et al. The Somatic Genomic Landscape of  
586 Glioblastoma. *Cell* 2013; 155:462–477

587 26. Ceccarelli M, Barthel FP, Malta TM, et al. Molecular Profiling Reveals Biologically  
588 Discrete Subsets and Pathways of Progression in Diffuse Glioma. *Cell* 2016; 164:550–563

589 27. Lai A, Kharbanda S, Pope WB, et al. Evidence for Sequenced Molecular Evolution of  
590 IDH1 Mutant Glioblastoma From a Distinct Cell of Origin. *J Clin Oncol* 2011; 29:4482–  
591 4490

592 28. Kuleshov MV, Jones MR, Rouillard AD, et al. Enrichr: a comprehensive gene set  
593 enrichment analysis web server 2016 update. *Nucleic Acids Res* 2016; 44:W90–W97

594 29. Chen EY, Tan CM, Kou Y, et al. Enrichr: interactive and collaborative HTML5 gene list  
595 enrichment analysis tool. *Bmc Bioinformatics* 2013; 14:128

596 30. Darmanis S, Sloan SA, Croote D, et al. Single-Cell RNA-Seq Analysis of Infiltrating  
597 Neoplastic Cells at the Migrating Front of Human Glioblastoma. *Cell Reports* 2017;  
598 21:1399–1410

599 31. Patel AP, Tirosh I, Trombetta JJ, et al. Single-cell RNA-seq highlights intratumoral  
600 heterogeneity in primary glioblastoma. *Science* 2014; 344:1396–1401

601 32. Jin X, Kim LJY, Wu Q, et al. Targeting glioma stem cells through combined BMI1 and  
602 EZH2 inhibition. *Nat Med* 2017; 23:1352–1361

603 33. Sottoriva A, Spiteri I, Piccirillo SGM, et al. Intratumor heterogeneity in human  
604 glioblastoma reflects cancer evolutionary dynamics. *Proc National Acad Sci* 2013;  
605 110:4009–4014

606 34. Segerman A, Niklasson M, Haglund C, et al. Clonal Variation in Drug and Radiation  
607 Response among Glioma-Initiating Cells Is Linked to Proneural-Mesenchymal Transition.  
608 *Cell Reports* 2016; 17:2994–3009

609 35. Halliday J, Helmy K, Pattwell SS, et al. In vivo radiation response of proneural glioma  
610 characterized by protective p53 transcriptional program and proneural-mesenchymal  
611 shift. Proc National Acad Sci 2014; 111:5248–5253

## 612 Figure legends

613 Figure 1 (at Results/Transcriptomic data aggregation)

614 Figure1: A) Schematic representation of the process followed to generate the new  
615 model including data aggregation, reduction in the number of genes per gene signature  
616 and the incorporation of a fourth gene signature to detect samples with a high content  
617 of normal tissue. B) Z-score and overlap achieved when reducing the number of genes  
618 per subtype for the rank criterion based on differences in mean expression (red). Mean  
619 and standard deviation (black) used to obtain the Z-score belong to the null distribution,  
620 obtained from the results of a thousand randomly ordered genes. C) Comparison of the  
621 distribution of the simplicity score for the four subtypes obtained from Verhaak's  
622 classification [4]. P values are obtained from the comparison to the neural (NE) subtype  
623 using Wilcoxon test. D) Comparison of the distribution of the simplicity score for the four  
624 subtypes obtained by our model. P values are obtained from the comparison to the high  
625 content in normal cells (NT) group using Wilcoxon test.

626 Figure 2 (at Results/ Validation of the model)

627 Figure 2: A) Comparison of the classification obtained by Wang's [7] model and from our  
628 model. B) ROC space showing the sensitivity and specificity obtained by our model for  
629 each subtype using Wang's model as gold standard. Error bars represent the standard  
630 deviation. Values in parenthesis correspond to the F1 score of each subtype. C)  
631 Frequency of somatic genomic alterations for each subtype. Significance values were

632 obtained applying the Chi-squared test. D) Survival curves between subtypes for  
633 samples with the top 20% simplicity score.

634 Figure 3 (at Results/ NT associates with abundance in normal cells)

635 Figure 3: A) Tumor purity of TCGA-IDH-WT samples determined by ABSOLUTE and  
636 ESTIMATE, respectively. The difference in tumor purity between subtypes was evaluated  
637 by the Wilcoxon test. B) Tumor purity of the validation cohort determined by ESTIMATE.  
638 The difference in tumor purity between subtypes was evaluated by the Wilcoxon test.  
639 C) Frequency of the molecular classification for samples obtained by microdissection  
640 from different tumor structures. The results from Wang's classification are shown in gray  
641 and the results from our model in red. D) Simplicity score obtained by Wang's model for  
642 the different tumor structures. P values are obtained for the comparison of each  
643 structure against the leading edge evaluated by Wilcoxon test. E) Empirical p values  
644 associated with the NT group obtained for the different tumor structures. P values are  
645 obtained for the comparison of each structure against the leading edge evaluated by  
646 Wilcoxon test. F) Representative H&E stained histological images for samples that  
647 belong to each molecular subgroup as classified by our model. Scale bars represent 200  
648  $\mu\text{m}$ .

649 Figure 4 (at Results/Survival analysis)

650 Figure 4: Survival curves between optimal cut-off of the simplicity score within proneural  
651 (A), classical (B) and mesenchymal (C) samples.

## 652 Supplementary material legends

653 Table S1

654 Table S1: Primers used for qRT-PCR.

655 Table S2

656 Table S2: Differentially expression analysis of the NE genes for the NE vs non-NE  
657 subtypes and for tumor vs non-tumor tissue.

658 Table S3

659 Table S3: Results obtained from the enrichment analysis of the 5 gene signature for NT  
660 against tissue databases using Enrichr.

661 Table S4

662 Table S4: Histological examination of the TMAs from tissue samples and the  
663 classification obtained by the model based on qRT-PCR measurements.

664 Figure S1

665 Figure S1: Z-score and overlap achieved when reducing the number of genes per subtype  
666 for the rank criterion based on relative differences in mean expression (A) and on the  
667 statistic obtained from the differential expression analysis (B). Mean and standard  
668 deviation (black) used to obtain the Z-score belong to the null distribution, obtained  
669 from the results of a thousand randomly ordered genes.

670 Figure S2

671 Figure S2: True positive rate and true negative rate obtained from different gene  
672 signature sizes for proneural (A), classical (B) and mesenchymal (C) subtypes using Wang  
673 et al. model as gold standard of the molecular classification. Black dots represent the  
674 mean and standard deviation of 1,000 randomly ranked gene lists. Three different  
675 criteria to rank the genes were used: mean differences (red dots), fold change (blue  
676 dots) and statistical significance of the mean differences (green dots).

677 Figure S3

678 Figure S3: Hazard ratio optimal cutoff selection of the simplicity score for proneural (A),  
679 classical (B) and mesenchymal (C). Continuous line represents the hazard ratio and  
680 dashed lines the 95% confidence interval.

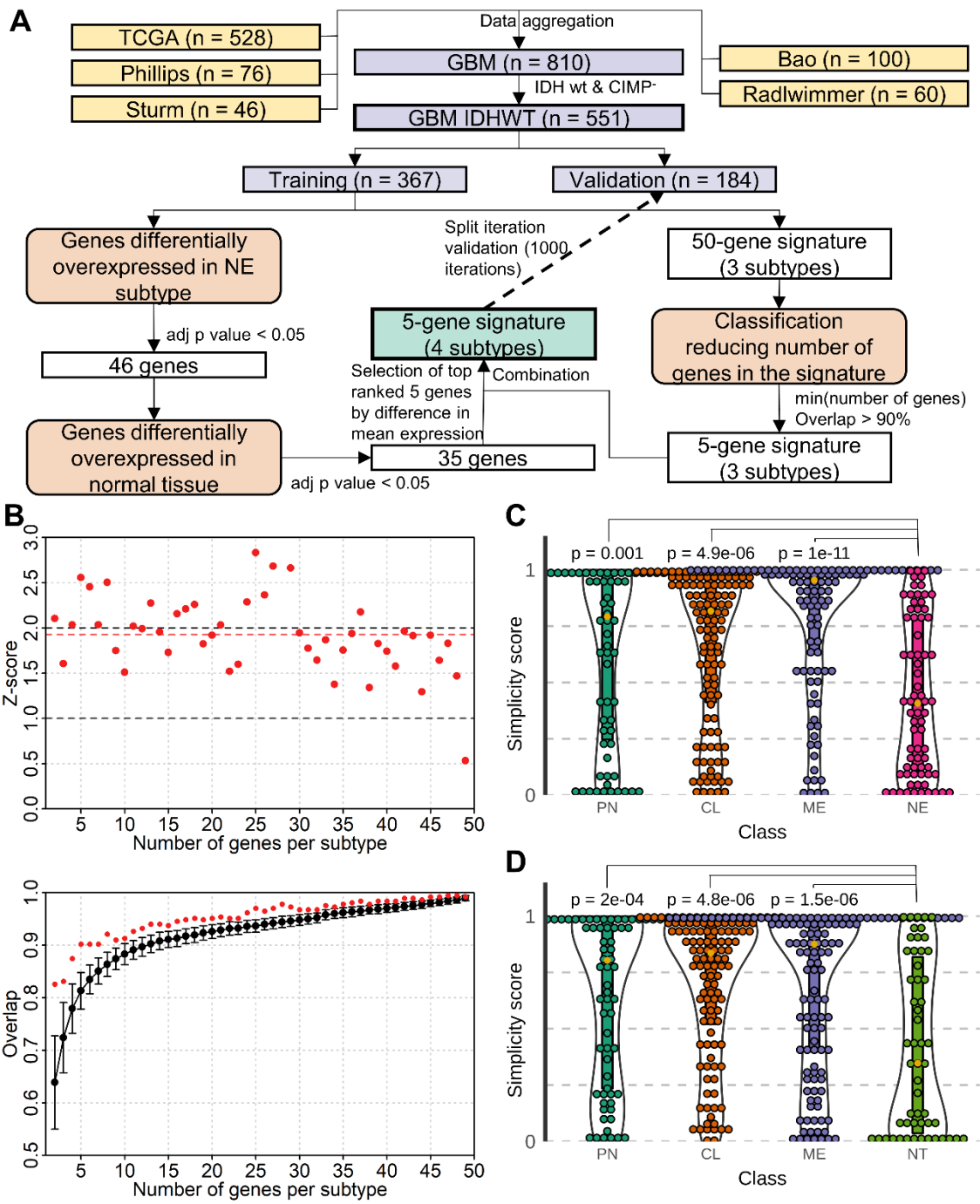
681 Figure S4

682 Figure S4: Schematic model of the evolution of a GBM tumor from a cellular tumor  
683 (proneural or classical) to a mesenchymal tumor.

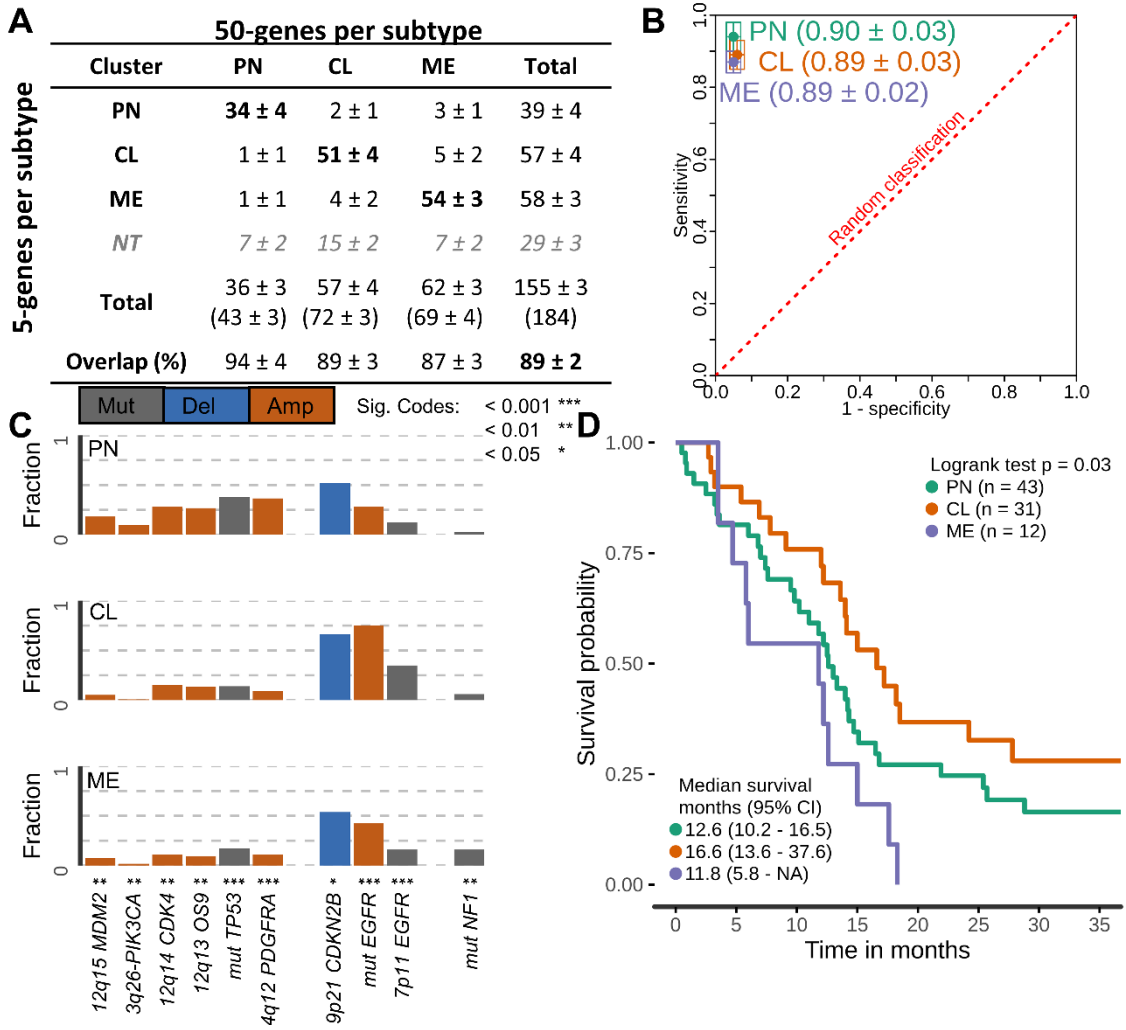
684 Methods S1

685 Methods S1: Blocks of code used for the development of the model.

686



689 Figure 2

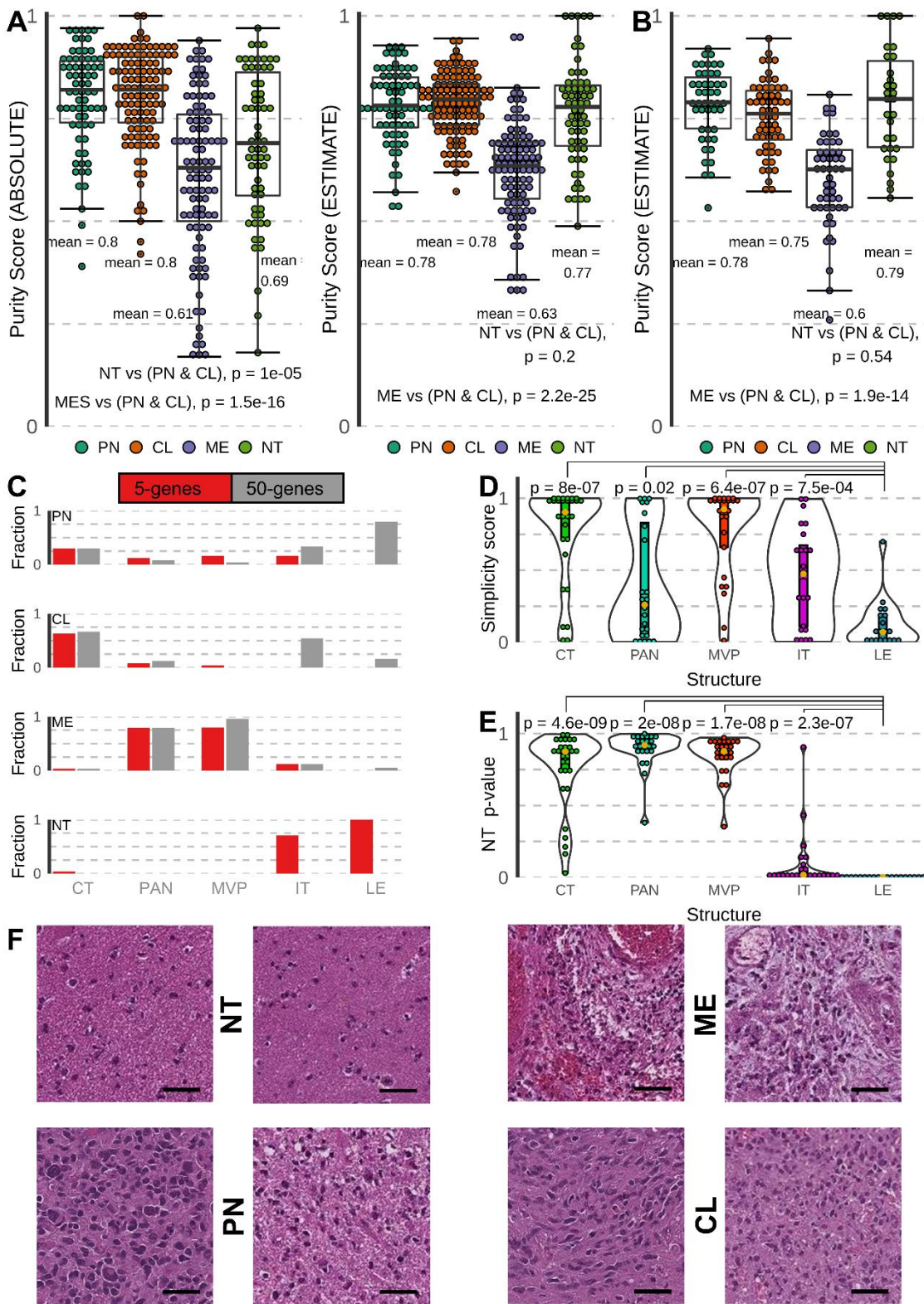


690

691

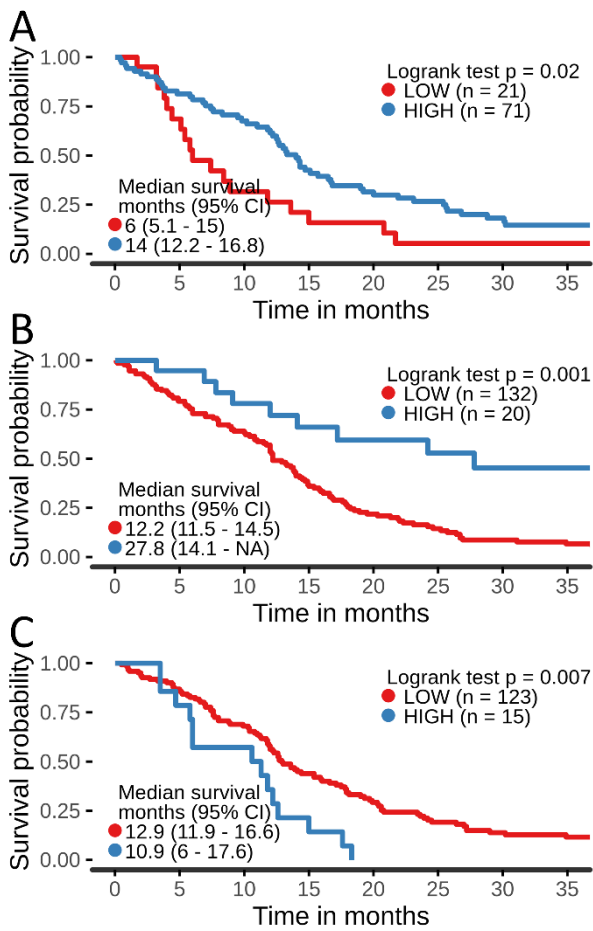


692 Figure 3



693

694 Figure 4



695

696

697 Table S1

Human primers	Forward 5'-----3'	Reverse 5'-----3'
<i>GPR17</i>	GGGCTTGTGATGGCTACAAT	CTGCCTTCAGGGTCTTTCTG
<i>CA10</i>	ATCCCACCTCAGTCAAATGC	TCATGAAGAAAGGGCCAATC
<i>UGT8</i>	TACTCTCCCACCAGGAGCTG	CCTTTTAACGGCAACATCGT
<i>HRASLS</i>	TCTTCTCATTCTGGGCTTG	TTCCTCCTCCCAAATTCCTT
<i>RAB33A</i>	GGAAGGTGCAGAACTGGAG	GGAAACAAGCAGGTGTCAG
<i>ELOVL2</i>	TCTTACCAAAGTGC GTTCCA	CTCCCTCCTTGCCATACAGA
<i>MLC1</i>	CGTAACAGCAGGAGCATGG	TCTGGTCAGGTCCAGAGAGC
<i>SLC4A4</i>	TCAAGACACAGACACGCACA	GGGACTCTGTCTGGAGGTCA
<i>CDH4</i>	GGACACCTGTCACCCTGAGT	GAGAGTGTCTGGGGTTTGA
<i>FGFR3</i>	TGCCCTCAGAGACTGAAAT	TCCGTTGTACCAGCCTTTTC
<i>LUM</i>	TGGAGCCAAATGTTATGCAG	GAAAGGCCGCTGTACCATAA
<i>PI3</i>	GCAAGAGCCAGTCAAAGGTC	TTCTTGATTCTGGGCAGTC
<i>SLPI</i>	CATATGGAGGAGGCTCTGGA	TCTTGAAAGCCTGCTGTGTG
<i>CYP1B1</i>	CTCCTGTGGAAGGCAGAGAA	TCCCAACTCTTGTACCTC
<i>NNMT</i>	ACCTTGCAGTGCCTCACTTT	CAAGCAATCTGTCTGCCTCA
<i>CCK</i>	TACATGGGCTGGATGGATTT	GTGAGGTGTGTGGTTGCACT
<i>CRYM</i>	GAATGGCAGTGAAGACACA	GGGACTGGACTCCCTCATT
<i>SERPINI1</i>	GACGAGTCATGCATCCTGAA	CCAGTTGCAAACATAATGTGC
<i>KCNK1</i>	CTGCAAACCATTGAGCGTAG	TGGGGTCACAGCTTCTTTGT
<i>GPR22</i>	CTCCCATTCTGGAAATCAACA	GCCAAGTCCCAACACAATTT

698

699

Gene	NEURAL VS NON-NEURAL		BRAIN VS TUMOR	
	Difference in mean expression	Adjusted p value	Difference mean expression	Adjusted p-value
CCK	1,449	4,0E-04	6,128	< 1E-04
CRYM	1,306	2,8E-03	5,6	< 1E-04
SERPINI1	1,166	2,0E-03	4,521	< 1E-04
KCNK1	0,895	4,4E-02	4,43	< 1E-04
GPR22	0,461	6,4E-03	4,032	< 1E-04
HPCAL4	0,516	1,4E-02	3,797	< 1E-04
CPNE6	0,342	5,6E-03	3,13	< 1E-04
CA4	0,534	4,4E-03	2,421	< 1E-04
UROS	0,441	1,7E-02	2,067	< 1E-04
KCNJ3	0,184	4,0E-04	2	< 1E-04
DHRS9	1,183	< 1E-04	1,93	< 1E-04
MDH1	0,345	4,0E-04	1,771	< 1E-04
ANXA3	0,94	4,0E-04	1,583	< 1E-04
CRYZL1	0,518	< 1E-04	1,311	< 1E-04
MGST3	0,48	< 1E-04	1,234	< 1E-04
SNCG	0,184	4,0E-04	1,173	< 1E-04
ACYP2	0,669	< 1E-04	1,168	< 1E-04
YPEL5	0,451	< 1E-04	1,161	< 1E-04
CLCA4	0,422	4,6E-02	1,052	< 1E-04
PEX11B	0,248	1,2E-02	0,947	< 1E-04
ADD3	0,655	< 1E-04	0,918	< 1E-04
MYBPC1	1,431	4,0E-04	0,884	< 1E-04
CASQ1	0,4	6,4E-03	0,884	< 1E-04
SEPW1	0,39	1,6E-03	0,851	< 1E-04
CRBN	0,388	2,0E-03	0,666	< 1E-04
ANXA7	0,578	< 1E-04	0,654	< 1E-04
TMEM144	0,513	1,7E-02	0,59	< 1E-04
TCEAL1	0,394	4,0E-04	0,547	< 1E-04
COX5B	0,346	3,2E-03	0,547	< 1E-04
TTC1	0,3	4,0E-04	0,533	< 1E-04
GUK1	0,381	< 1E-04	0,528	< 1E-04
PEX19	0,245	1,9E-02	0,373	< 1E-04
IMPA1	0,369	2,2E-02	0,32	< 1E-04
RBKS	0,699	< 1E-04	0,258	7,2E-03
MAT2B	0,421	< 1E-04	0,204	7,2E-03
CRYL1	0,703	4,0E-04	NS	NS
SEPP1	0,701	< 1E-04	NS	NS
MRPL49	0,532	< 1E-04	NS	NS
LYRM1	0,503	2,4E-03	NS	NS
TSNAX	0,447	4,0E-04	NS	NS
ATP5L	0,366	< 1E-04	NS	NS
AKR7A3	0,345	6,4E-03	NS	NS
SNX11	0,336	4,0E-04	NS	NS
CCDC121	0,334	2,4E-03	NS	NS
ATP5F1	0,223	1,3E-02	NS	NS
NSL1	0,206	4,4E-02	NS	NS

701

702

703 Table S3

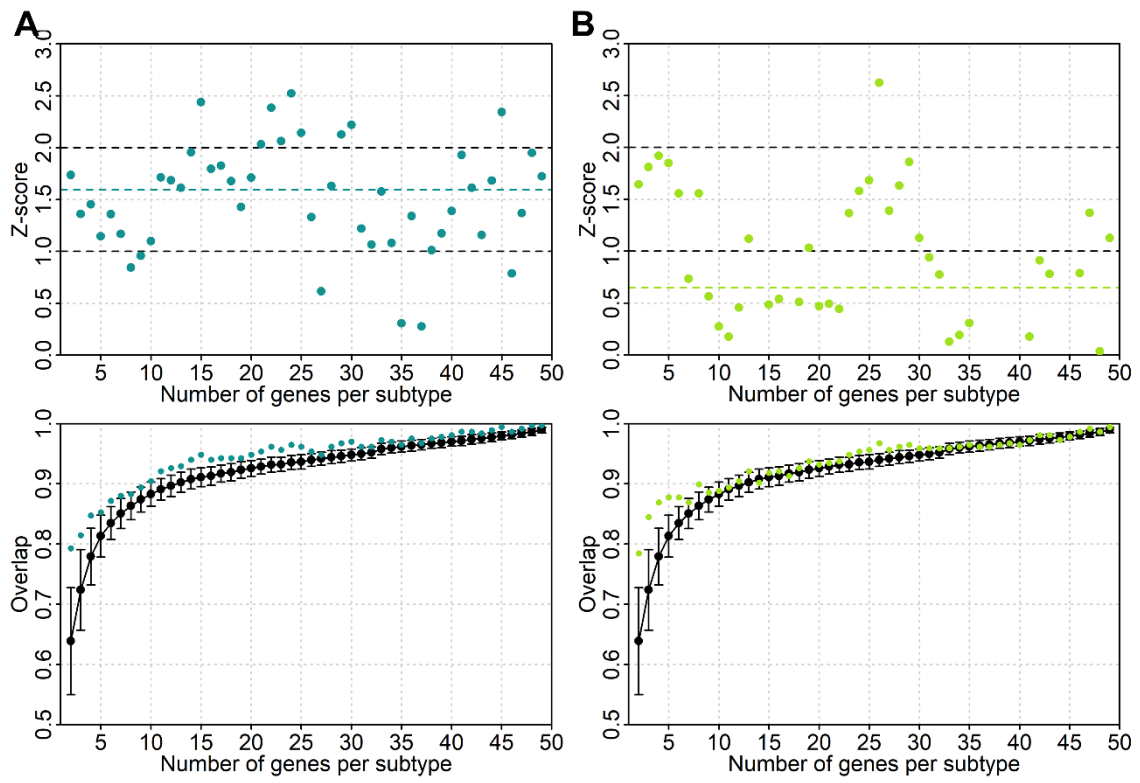
<b>Index</b>	<b>Name</b>	<b>P-value</b>	<b>Adjusted p-value</b>	<b>Odds Ratio</b>	<b>Combined score</b>
<b>1</b>	BRAIN (BULK)	0.00002074	0.002240	8.64	93.12
<b>2</b>	CEREBRAL CORTEX	0.00002074	0.001120	8.64	93.12
<b>3</b>	CINGULATE GYRUS	0.00002074	0.0007467	8.64	93.12
<b>4</b>	DENTATE GRANULE CELL	0.00002074	0.0005600	8.64	93.12
<b>5</b>	DORSAL STRIATUM	0.00002074	0.0004480	8.64	93.12
<b>6</b>	SUPERIOR FRONTAL GYRUS	0.00002074	0.0003734	8.64	93.12
<b>7</b>	ATRIUM	0.01294	0.1997	5.18	22.52
<b>8</b>	HEART (BULK TISSUE)	0.01294	0.1747	5.18	22.52
<b>9</b>	PREFRONTAL CORTEX	0.01294	0.1553	5.18	22.52
<b>10</b>	VENTRICLE	0.01294	0.1398	5.18	22.52

704

705

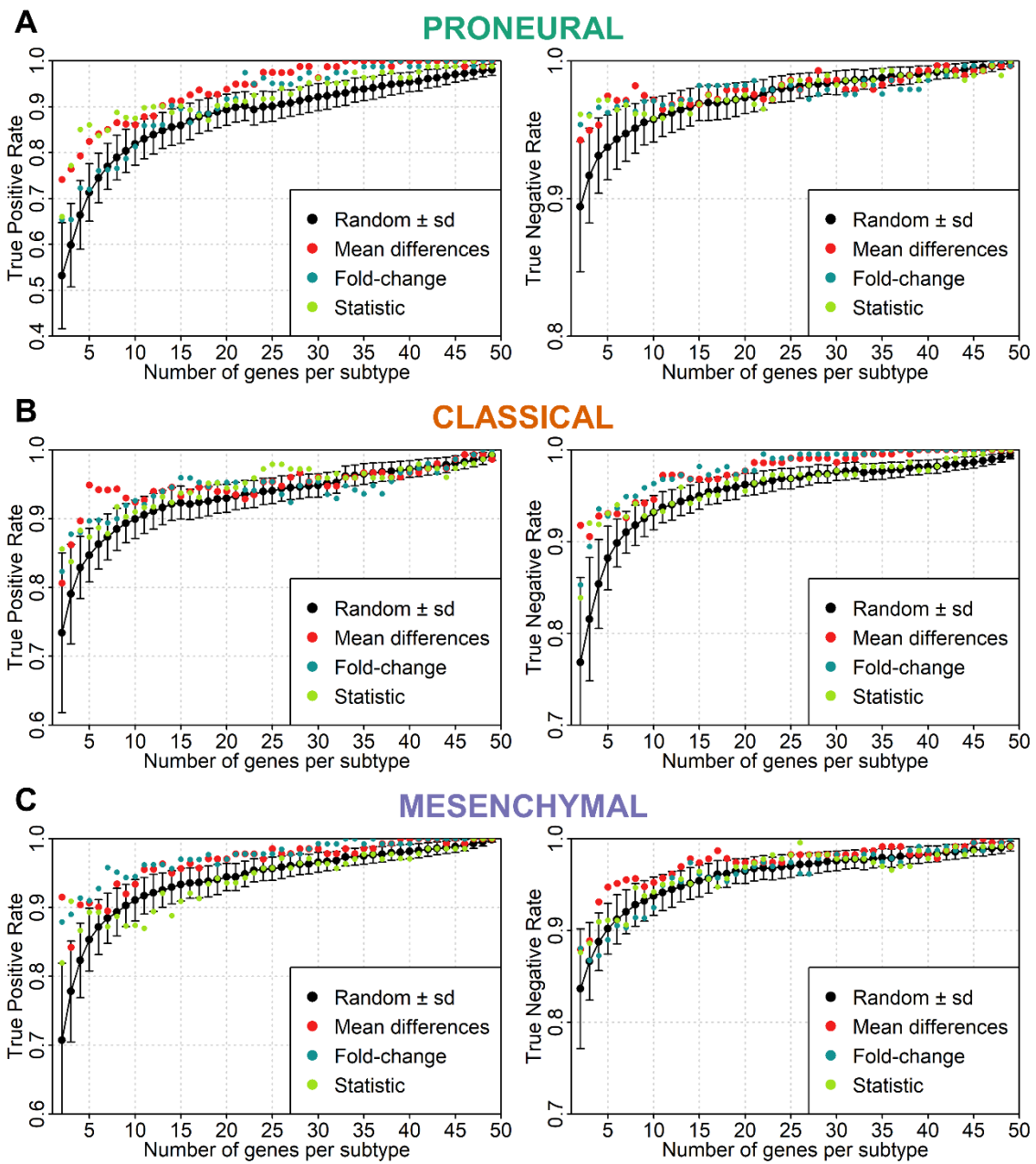
<b>Sample</b>	<b>Cellularity</b>	<b>Pathogenic blood vessels</b>	<b>Necrosis</b>	<b>Molecular classification</b>
<b>S1</b>	MODERATE / MODERATE	YES / YES	NO / NO	ME
<b>S2</b>	LOW / LOW	YES / YES	YES / NO	CL
<b>S3</b>	MODERATE / HIGH	YES / YES	NO / NO	ME
<b>S4</b>	LOW / LOW	NO / NO	NO / NO	NT
<b>S5</b>	MODERATE / MODERATE	YES / YES	YES / NO	ME
<b>S6</b>	LOW / LOW	YES / YES	YES / NO	ME
<b>S7</b>	HIGH / HIGH	YES / YES	NO / NO	ME
<b>S8</b>	MODERATE / MODERATE	YES / YES	YES / NO	ME
<b>S9</b>	LOW / LOW	NO / NO	NO / NO	PN
<b>S10</b>	MODERATE / MODERATE	YES / YES	YES / NO	ME
<b>S11</b>	LOW / LOW	YES / YES	NO / NO	NT
<b>S12</b>	LOW / LOW	YES / YES	YES / NO	CL
<b>S13</b>	MODERATE / MODERATE	YES / NO	YES / NO	CL
<b>S14</b>	MODERATE / MODERATE	YES / YES	YES / NO	ME
<b>S15</b>	LOW / LOW	NO / NO	NO / NO	PN
<b>S16</b>	HIGH / MODERATE	YES / YES	YES / NO	ME
<b>S17</b>	HIGH / HIGH	NO / NO	NO / NO	PN
<b>S18</b>	MODERATE / MODERATE	YES / NO	NO / NO	ME
<b>S19</b>	LOW / LOW	NO / NO	NO / NO	NT
<b>S20</b>	LOW / LOW	NO / NO	NO / NO	NT
<b>S21</b>	LOW / HIGH	YES / YES	YES / NO	NT
<b>S22</b>	MODERATE / MODERATE	NO / NO	NO / NO	ME
<b>S23</b>	HIGH / HIGH	YES / NO	NO / NO	NT
<b>S24</b>	HIGH / HIGH	YES / YES	NO / NO	ME
<b>S25</b>	HIGH / HIGH	YES / NO	NO / NO	ME
<b>S26</b>	LOW / LOW	NO / NO	NO / NO	NT
<b>S27</b>	HIGH / HIGH	NO / NO	NO / NO	CL
<b>S28</b>	MODERATE / MODERATE	YES / NO	YES / YES	NT
<b>S29</b>	HIGH / MODERATE	YES / NO	YES / NO	ME
<b>S30</b>	LOW / HIGH	NO / NO	YES / NO	PN
<b>S31</b>	HIGH / MODERATE	YES / YES	YES / YES	CL
<b>S32</b>	HIGH / MODERATE	YES / NO	YES / NO	CL
<b>S33</b>	LOW / MODERATE	NO / NO	NO / NO	NT
<b>S34</b>	HIGH / MODERATE	YES / YES	NO / YES	CL

708 Figure S1



709

710

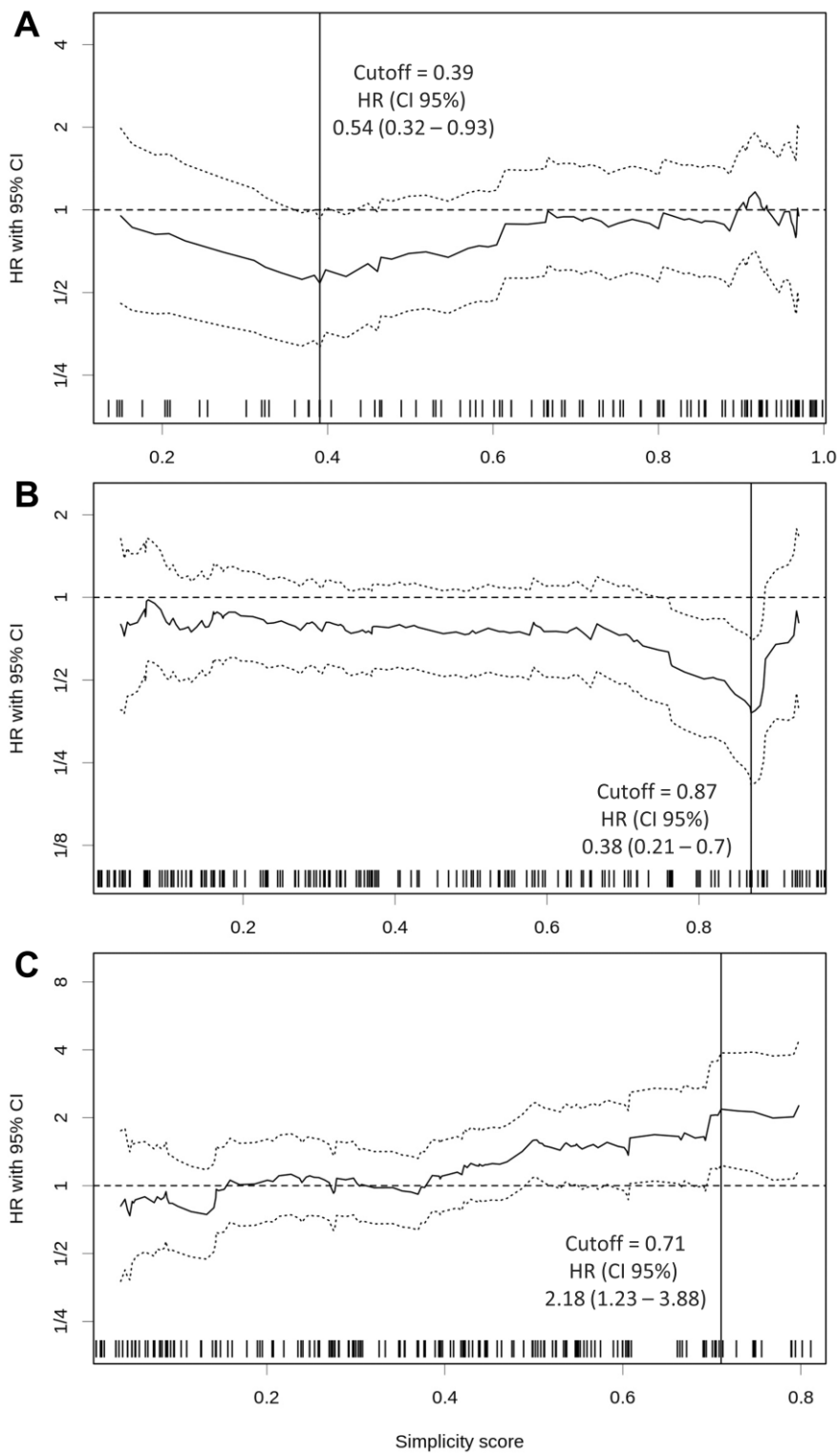


712

713



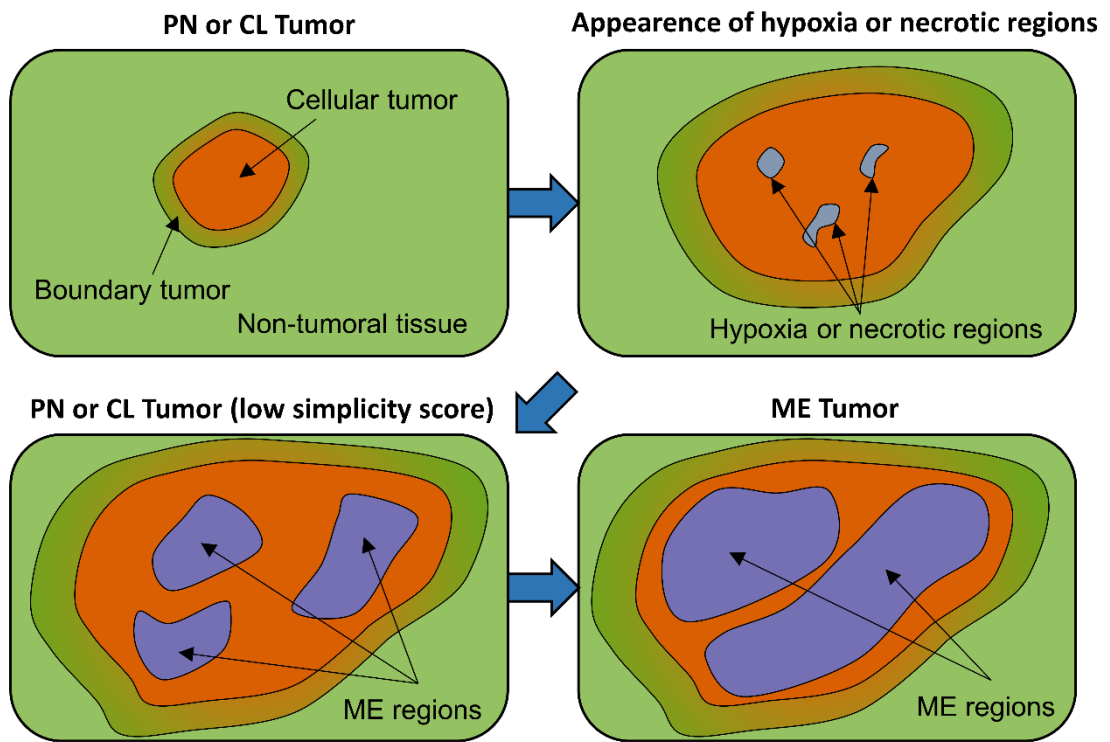
714 Figure S3



715

716

717 Figure S4



718

719