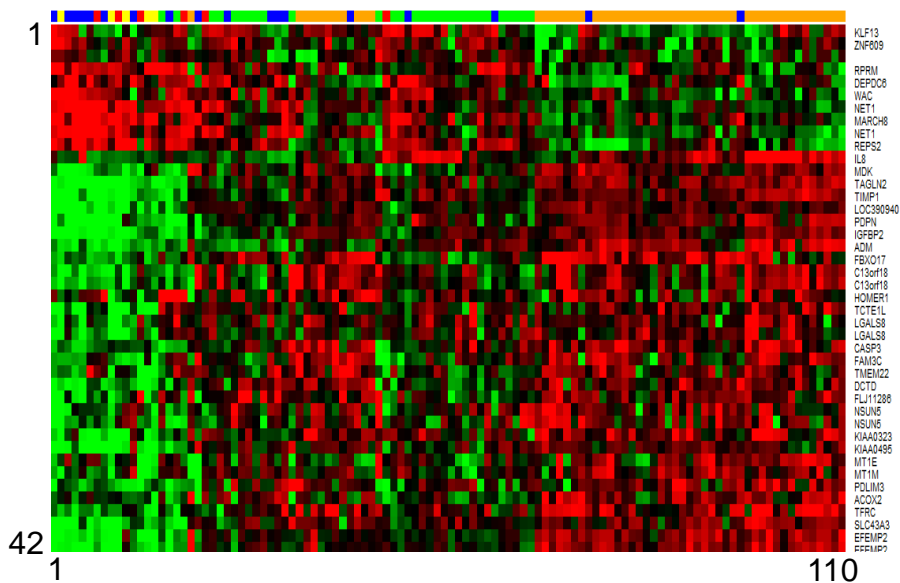


**Fig. S1**

### GSE4290 dataset

- Grade III Astrocytoma
- Grade II Astrocytoma
- Glioblastoma (Group 3)
- Glioblastoma (Group 2)
- Glioblastoma (Group 1)



### GSE4290 dataset

- Grade III Oligodendroglioma
- Grade II Oligodendroglioma
- Glioblastoma (Group 3)
- Glioblastoma (Group 2)
- Glioblastoma (Group 1)

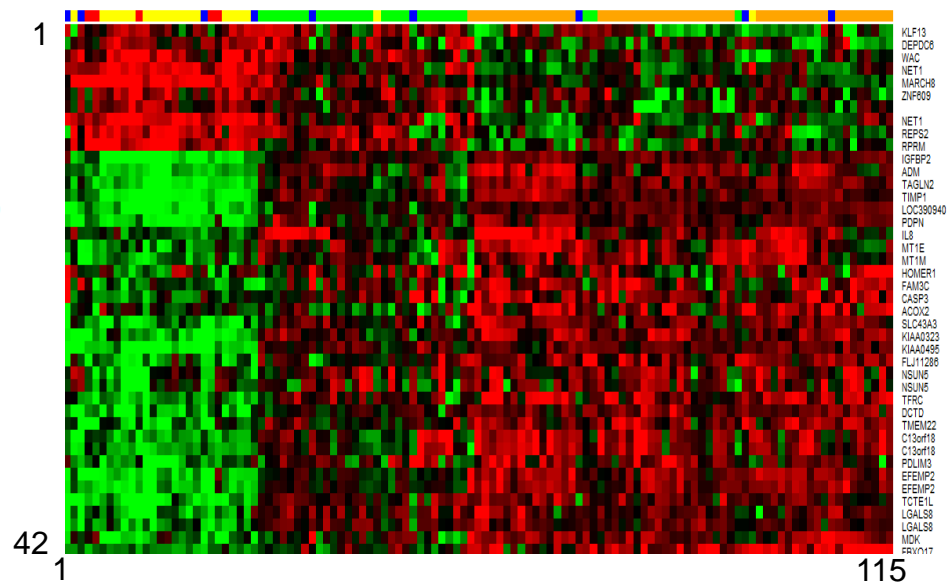
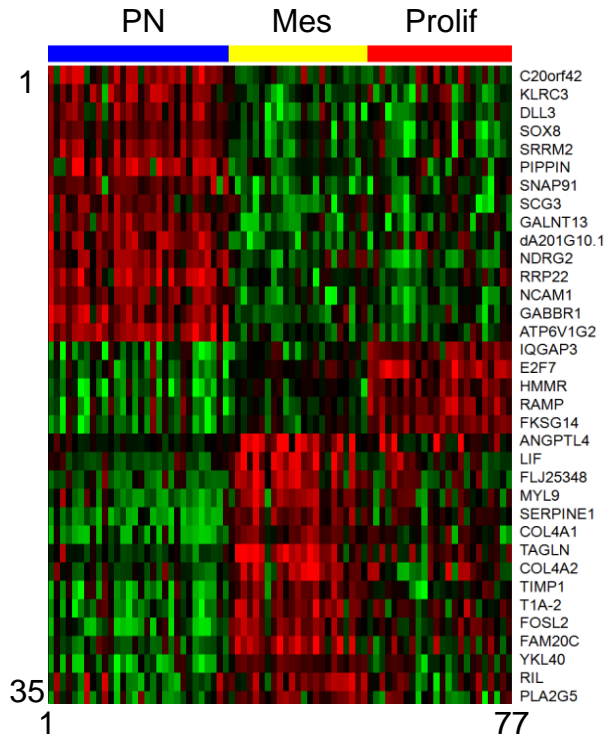
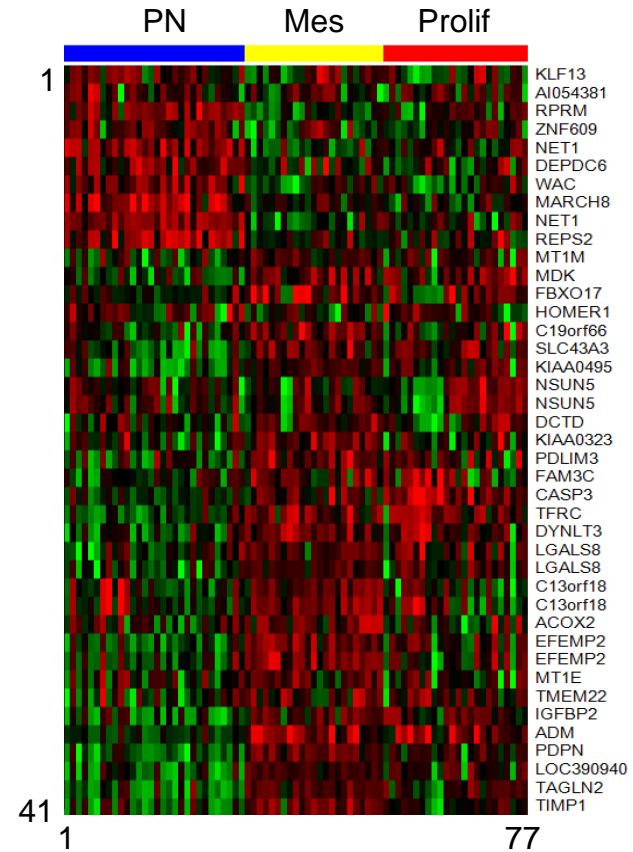


Fig. S2

**A** Phillips et al. 2006: using 35 signature genes

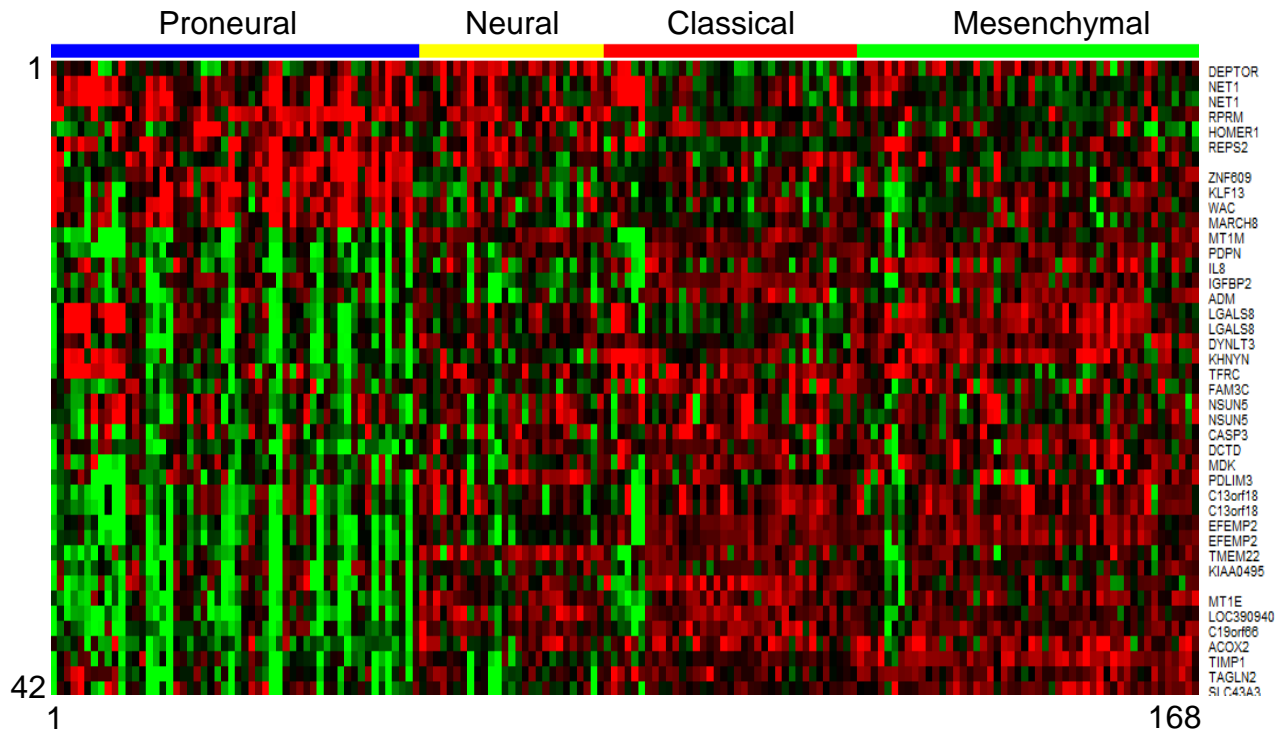


**B** Phillips et al. 2006: using 41 probesets (genes)



**Fig. S3**

**A** Verhaak et al. 2010: using 42 probesets in TCGA Core samples



**B**

Row Labels	Proneural	Neural	Classical	Mesenchymal	Grand Total
Group 1	13	13	37	41	104
Group 2	17	7	11	13	48
Group 3	14	2	0	0	16
Grand Total	44	22	48	54	168

**Fig. S4**



TCGA dataset

P = 0.05

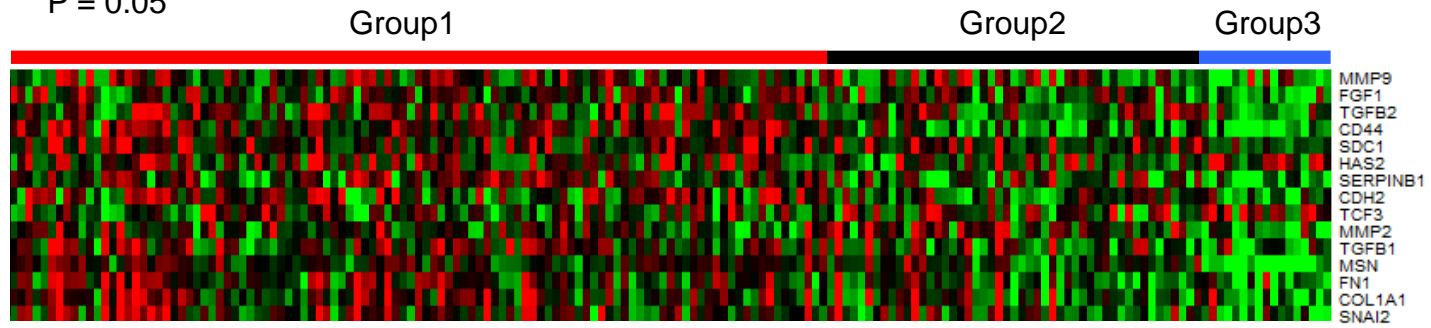
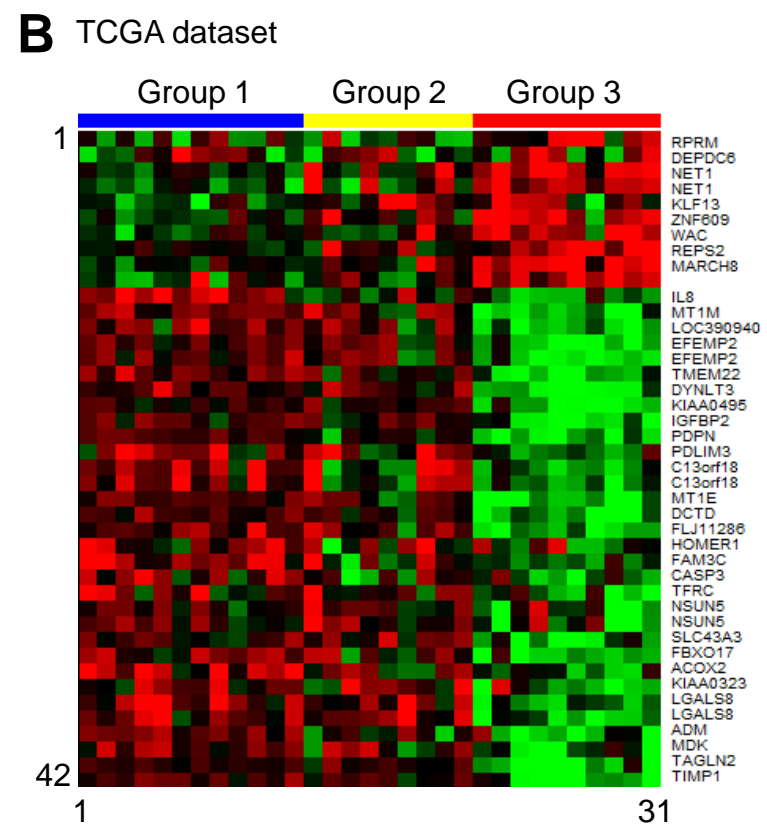
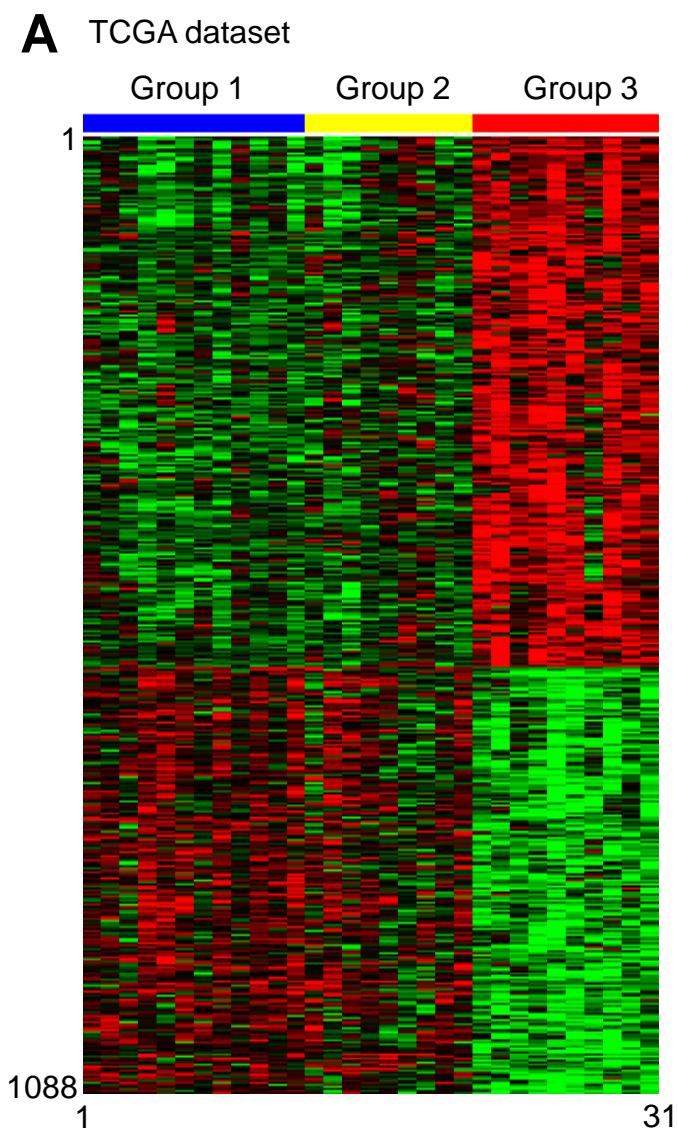
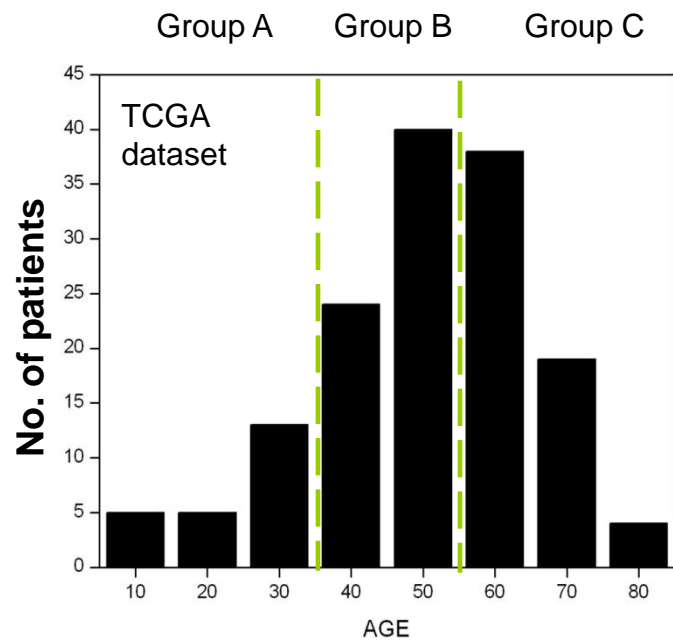
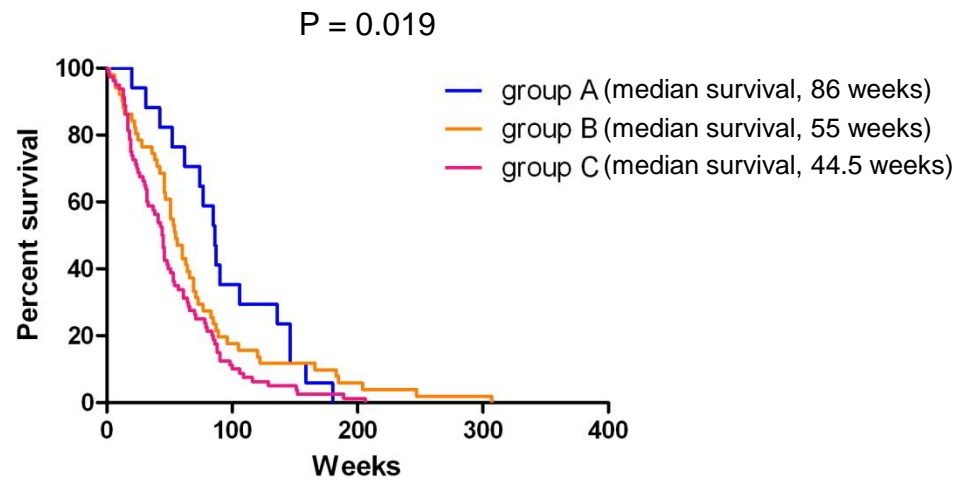
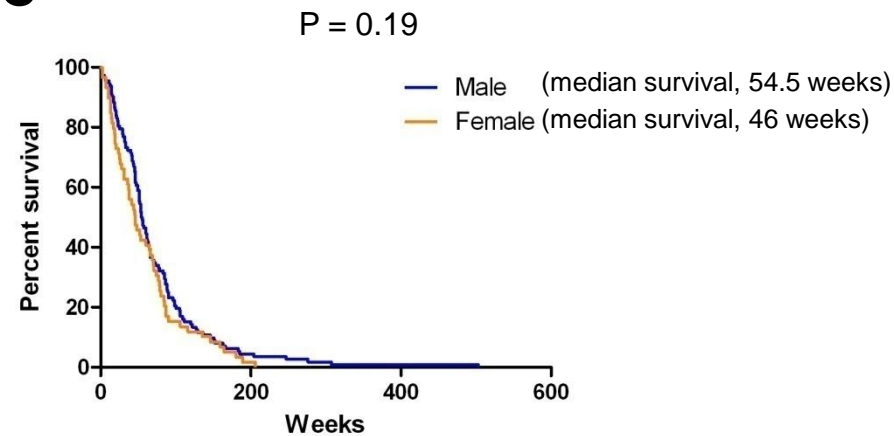


Fig. S6



**Fig. S7**

**A****B****C****Fig. S8**



## Supplementary Figure Legends

**Figure S1.** Hierarchical clustering analysis. **a**, Hierarchical clustering of the GSE4290 dataset of astrocytomas (II and III). Fourteen patients were classified as long-term survivors ( $\geq 5$  years) and 14 as short-term survivors ( $< 5$  years). The data are presented in matrix format in which rows represent individual genes and columns represent each tissue. Each cell in the matrix represents the expression level of a gene feature in an individual tissue. Red and green in cells reflect high and low expression levels, respectively. Each blue bar represents a more than 5-year survival in astrocytoma patients. **b**, Kaplan–Meier plot of overall survival of astrocytoma patients grouped on the basis of gene expression profiling. The difference between the groups was not significant ( $P = 0.18$ , log-rank test). **c**, Hierarchical clustering of the GSE4290 data set of oligodendrogliomas (II and III). Thirteen patients were classified as long-term survivors ( $\geq 5$  years) and 20 as short-term survivors ( $< 5$  years). Each blue bar represents a more than 5-year survival in oligodendroglioma patients. **d**, Kaplan–Meier plot of overall survival of oligodendrogliomas. The difference between the groups was not significant ( $P = 0.99$ , log-rank test).

**Figure S2.** Hierarchical clustering of astrocytomas (II and III) and oligodendrogliomas (II and III) with GBMs of GSE4290 dataset, respectively. The data are presented in matrix format in which rows represent individual genes and columns represent each tissue.

**Figure S3.** Hierarchical clustering of 77 gliomas observed in Phillips et al., 2006. **a**, Confirmation of hierarchical clustering of 35 signature genes as shown in Phillips et al., 2006. **b**, Hierarchical clustering of 77 gliomas using 41 probe sets. The data are presented in matrix format in which rows represent individual genes and columns represent each tissue.

**Figure S4.** Hierarchical clustering of 168 TCGA Core GBMs samples observed in Verhaak et al., 2010. **a**, Hierarchical clustering of 168 GBMs was used 42 probe sets. The data are presented in matrix format in which rows represent individual genes and columns represent each tissue. **b**, 168/173 samples in our set are also part of the Verhaak et al paper. We pulled the data from our data and the Verhaak dataset to make the table for the 168 samples.

**Figure S5.** Hierarchical clustering of 173 TCGA GBM tumors. Columns represent each tissue and rows represent outcomes of various prediction models as indicated. Each cell represents memberships of tissues when a particular prediction model was applied in the set. Blue represents a more than 2-year survival in GBM patients. LDA, linear discriminator analysis; SVM, support vector machines; NC, nearest centroid; NN, nearest neighbor; CCP, compound covariate prediction.

**Figure S6.** Hierarchical clustering of EMT-associated genes in TCGA dataset. EMT-associated 15 genes (29 probe sets) were generated using a two-sample t test ( $P = 0.05$ ) among the three groups. Among those, the maximally expressed 15 genes were measured by average intensity across arrays. The 15 genes are presented using supervised hierarchical clustering in matrix format, where rows represent individual genes and columns represent each tissue. Each cell in the matrix represents the expression level of a gene in an individual tissue. Red and green cells reflect high and low expression levels, respectively.

**Figure S7.** Hierarchical clustering of 31 TCGA GBM tumors. a, 31 patients were classified as long-term survivors ( $\geq 2$  years). 1088 probe sets were generated using a two-sample t test ( $P = 0.0005$ ) between the three groups. b, Hierarchical clustering of 31 GBMs using 42 probe sets. The data are presented in matrix format in which rows represent individual genes and columns represent each tissue.

**Figure S8.** a, Age distribution of 173 TCGA GBM patients. Patients ( $n = 173$ ) were classified as Group A ( $\leq 35$  years), Group B ( $> 35$  and  $\leq 55$  years) and Group C ( $> 55$  years). b, Kaplan-Meier plot of overall survival in GBM patients grouped on the basis of age. The difference between the groups was significant ( $P = 0.019$ , log-rank test). c, Kaplan-Meier plot of overall survival in GBM patients who were grouped on the basis of sex. The difference between the groups was not significant ( $P = 0.19$ , log-rank test).

