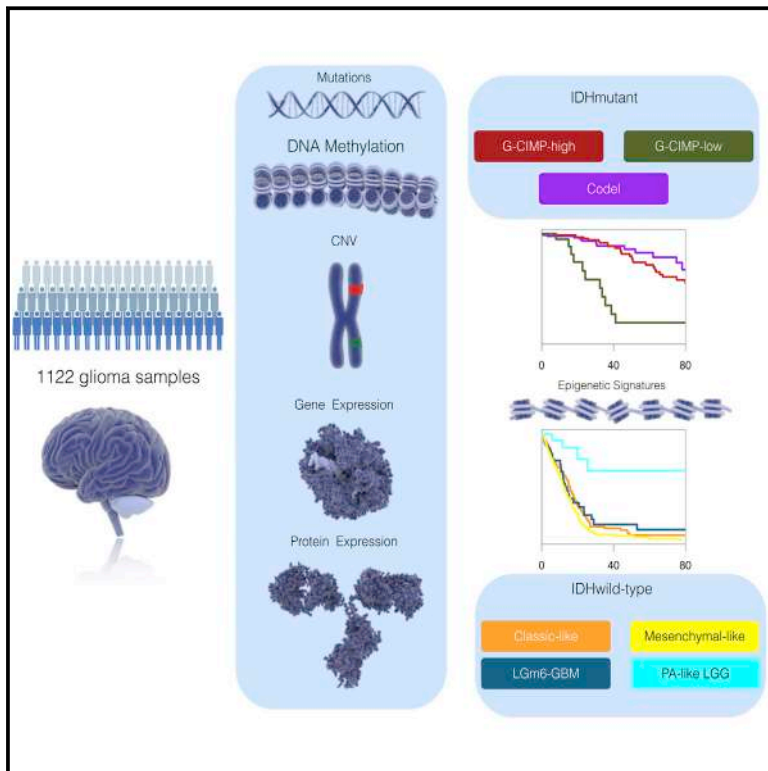


# Molecular Profiling Reveals Biologically Discrete Subsets and Pathways of Progression in Diffuse Glioma

## Graphical Abstract



## Authors

Michele Ceccarelli, Floris P. Barthel, Tathiane M. Malta, ..., Houtan Noushmehr, Antonio Iavarone, Roel G.W. Verhaak

## Correspondence

houtan@usp.br (H.N.),  
ai2102@columbia.edu (A.I.),  
rverhaak@mdanderson.org (R.G.W.V.)

## In Brief

Integration of a large sample size of glioma tumors with multidimensional ‘omic characterization and clinical annotation provides insights into molecular classification, telomere maintenance mechanisms, progression from low to high grade disease, driver mutations, and therapeutic options.

## Highlights

- Comprehensive molecular profiling of 1,122 adult diffuse grade II, III, and IV gliomas
- Telomere length and telomere maintenance defined by somatic alterations
- DNA methylation profiling reveals subtypes of IDH mutant and IDH-wild-type glioma
- Integrated molecular analysis of progression from low-grade to high-grade disease



# Molecular Profiling Reveals Biologically Discrete Subsets and Pathways of Progression in Diffuse Glioma

Michele Ceccarelli,<sup>1,2,24</sup> Floris P. Barthel,<sup>3,4,24</sup> Tathiane M. Malta,<sup>5,6,24</sup> Thais S. Sabedot,<sup>5,6,24</sup> Sofie R. Salama,<sup>7</sup> Bradley A. Murray,<sup>8</sup> Olena Morozova,<sup>7</sup> Yulia Newton,<sup>7</sup> Amie Radenbaugh,<sup>7</sup> Stefano M. Pagnotta,<sup>2,9</sup> Samreen Anjum,<sup>1</sup> Jiguang Wang,<sup>10</sup> Ganiraju Manyam,<sup>3</sup> Pietro Zoppoli,<sup>10</sup> Shiyun Ling,<sup>3</sup> Arjun A. Rao,<sup>7</sup> Mia Grifford,<sup>7</sup> Andrew D. Cherniack,<sup>8</sup> Hailei Zhang,<sup>8</sup> Laila Poisson,<sup>11</sup> Carlos Gilberto Carlotti, Jr.,<sup>5,6</sup> Daniela Pretti da Cunha Tirapelli,<sup>5,6</sup> Arvind Rao,<sup>3</sup> Tom Mikkelsen,<sup>11</sup> Ching C. Lau,<sup>12,13</sup> W.K. Alfred Yung,<sup>3</sup> Raul Rabadan,<sup>10</sup> Jason Huse,<sup>14</sup> Daniel J. Brat,<sup>15</sup> Norman L. Lehman,<sup>16</sup> Jill S. Barnholtz-Sloan,<sup>17</sup> Siyuan Zheng,<sup>3</sup> Kenneth Hess,<sup>3</sup> Ganesh Rao,<sup>3</sup> Matthew Meyerson,<sup>8,18</sup> Rameen Beroukhi,<sup>8,18,19</sup> Lee Cooper,<sup>15</sup> Rehan Akbani,<sup>3</sup> Margaret Wrensch,<sup>20</sup> David Haussler,<sup>7</sup> Kenneth D. Aldape,<sup>21</sup> Peter W. Laird,<sup>22</sup> David H. Gutmann,<sup>23</sup> TCGA Research Network, Houtan Noushmehr,<sup>5,6,25,\*</sup> Antonio Iavarone,<sup>10,25,\*</sup> and Roel G.W. Verhaak<sup>3,25,\*</sup>

<sup>1</sup>Qatar Computing Research Institute, Hamad Bin Khalifa University, Doha P.O. box 5825, Qatar

<sup>2</sup>Department of Science and Technology, University of Sannio, Benevento 82100, Italy

<sup>3</sup>Department of Genomic Medicine, Department of Bioinformatics and Computational Biology, Department of Biostatistics, Department of Neuro-Oncology, Department of Neurosurgery, Department of Pathology, University of Texas MD Anderson Cancer Center, Houston, TX 77030, USA

<sup>4</sup>Oncology Graduate School Amsterdam, Department of Pathology, VU University Medical Center, 1081 HV Amsterdam, the Netherlands

<sup>5</sup>Department of Genetics (CISBi/NAP), Department of Surgery and Anatomy, Ribeirão Preto Medical School, University of São Paulo, Monte Alegre, Ribeirão Preto-SP CEP: 14049-900, Brazil

<sup>6</sup>Center for Integrative Systems Biology (CISBi, NAP/USP), Ribeirão Preto Medical School, University of São Paulo, Ribeirão Preto, São Paulo 14049-900, Brazil

<sup>7</sup>UC Santa Cruz Genomics Institute, University of California, Santa Cruz, Santa Cruz, CA 95064, USA

<sup>8</sup>The Eli and Edythe L. Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, MA 02142, USA

<sup>9</sup>BIOGEM Istituto di Ricerche Genetiche "G. Salvatore," Campo Reale, 83031 Ariano Irpino, Italy

<sup>10</sup>Department of Neurology, Department of Pathology, Institute for Cancer Genetics, Department of Systems Biology and Biomedical Informatics, Columbia University Medical Center, New York, NY 10032, USA

<sup>11</sup>Henry Ford Hospital, Detroit, MI 48202, USA

<sup>12</sup>Texas Children's Hospital, Houston, TX 77030, USA

<sup>13</sup>Baylor College of Medicine, Houston, TX 77030, USA

<sup>14</sup>Memorial Sloan Kettering Cancer Center, New York, NY 10065, USA

<sup>15</sup>Winship Cancer Institute, Emory University, Atlanta, GA 30322, USA

<sup>16</sup>Department of Pathology, The Ohio State University, Columbus, OH 43210, USA

<sup>17</sup>Case Comprehensive Cancer Center, Case Western Reserve University, Cleveland, OH 44106, USA

<sup>18</sup>Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA 02215, USA

<sup>19</sup>Department of Medicine, Harvard Medical School, Boston, MA 02215, USA

<sup>20</sup>Department of Neurological Surgery, University of California, San Francisco, San Francisco, CA 94158, USA

<sup>21</sup>Princess Margaret Cancer Centre, Toronto, ON M5G 2M9, Canada

<sup>22</sup>Van Andel Research Institute, Grand Rapids, MI 49503, USA

<sup>23</sup>School of Medicine, Washington University, St. Louis, MO 63110, USA

<sup>24</sup>Co-first author

<sup>25</sup>Co-senior author

\*Correspondence: [houtan@usp.br](mailto:houtan@usp.br) (H.N.), [ai2102@columbia.edu](mailto:ai2102@columbia.edu) (A.I.), [rverhaak@mdanderson.org](mailto:rverhaak@mdanderson.org) (R.G.W.V.)  
<http://dx.doi.org/10.1016/j.cell.2015.12.028>

## SUMMARY

Therapy development for adult diffuse glioma is hindered by incomplete knowledge of somatic glioma driving alterations and suboptimal disease classification. We defined the complete set of genes associated with 1,122 diffuse grade II-III-IV gliomas from The Cancer Genome Atlas and used molecular profiles to improve disease classification, identify molecular correlations, and provide insights into

the progression from low- to high-grade disease. Whole-genome sequencing data analysis determined that *ATRX* but not *TERT* promoter mutations are associated with increased telomere length. Recent advances in glioma classification based on *IDH* mutation and 1p/19q co-deletion status were recapitulated through analysis of DNA methylation profiles, which identified clinically relevant molecular subsets. A subtype of *IDH* mutant glioma was associated with DNA demethylation and poor outcome;

a group of IDH-wild-type diffuse glioma showed molecular similarity to pilocytic astrocytoma and relatively favorable survival. Understanding of cohesive disease groups may aid improved clinical outcomes.

## INTRODUCTION

Diffuse gliomas represent 80% of malignant brain tumors (Schwartzbaum et al., 2006). Adult diffuse gliomas are classified and graded according to histological criteria (oligodendroglioma, oligoastrocytoma, astrocytoma, and glioblastoma; grade II to IV). Although histopathologic classification is well established and is the basis of the World Health Organization (WHO) classification of CNS tumors (Louis et al., 2007), it suffers from high intra- and inter-observer variability, particularly among grade II-III tumors (van den Bent, 2010). Recent molecular characterization studies have benefited from the availability of the datasets generated by The Cancer Genome Atlas (TCGA) (Brennan et al., 2013; Eckel-Passow et al., 2015; Frattini et al., 2013; Kim et al., 2015; Suzuki et al., 2015; Cancer Genome Atlas Research Network et al., 2015) and have related genetic, gene expression, and DNA methylation signatures with prognosis (Noushmehr et al., 2010; Sturm et al., 2012; Verhaak et al., 2010). For example, mutations in the isocitrate dehydrogenase genes 1 and 2 (*IDH1/IDH2*) define a distinct subset of glioblastoma (GBM) with a hypermethylation phenotype (G-CIMP) with favorable outcome (Noushmehr et al., 2010; Yan et al., 2009). Conversely, the absence of *IDH* mutations in LGG marks a distinct IDH-wild-type subgroup characterized by poor, GBM-like prognosis (Eckel-Passow et al., 2015; Cancer Genome Atlas Research Network et al., 2015). Recent work by us and others has proposed classification of glioma into *IDH* wild-type cases, *IDH* mutant group additionally carrying codeletion of chromosome arm 1p and 19q (*IDH* mutant-codel) and samples with euploid 1p/19q (*IDH* mutant-non-codel), regardless of grade and histology (Eckel-Passow et al., 2015; Cancer Genome Atlas Research Network et al., 2015). Mutation of the *TERT* promoter, which has been reported with high frequency across glioma, may be an additional defining feature. Current analyses have not yet clarified the relationships between LGGs and GBMs that share common genetic hallmarks like *IDH* mutation or *TERT* promoter mutation status. An improved understanding of these relationships will be necessary as we evolve toward an objective genome-based clinical classification.

To address the above issues, we assembled a dataset comprising all TCGA newly diagnosed diffuse glioma consisting of 1,122 patients and comprehensively analyzed using sequencing and array-based molecular profiling approaches. We have addressed crucial technical challenges in analyzing this comprehensive dataset, including the integration of multiple platforms and data sources (e.g., multiple methylation and gene expression platforms). We identified new diffuse glioma subgroups with distinct molecular and clinical features and shed light on the mechanisms driving progression of lower grade glioma (LGG) (WHO grades II and III) into full-blown GBM (WHO grade IV).

## RESULTS

### Patient Cohort Characteristics

The TCGA LGG and GBM cohorts consist of 516 and 606 patients, respectively. Independent analysis of the GBM dataset was previously described, as was analysis of 290 LGG samples (Brennan et al., 2013; Cancer Genome Atlas Research Network et al., 2015). 226 LGG samples were added to our current cohort (Table 1). Clinical data, including age, tumor grade, tumor histology, and survival, were available for 93% (1,046/1,122) of cases (Table S1). The majority of samples were grade IV tumors (n = 590, 56%), whereas 216 (21%) and 241 (23%) were grade II and III tumors, respectively. Similarly, 590 (56%) samples were classified as GBM, 174 (17%) as oligodendroglioma, 169 (16%) as astrocytoma, and 114 (11%) as oligoastrocytoma.

Among the data sources considered in our analysis were gene expression (n = 1,045), DNA copy number (n = 1,084), DNA methylation (n = 932), exome sequencing (n = 820), and protein expression (n = 473). Multiple and overlapping characterization assays were employed (Table S1). All data files that were used in our analysis can be found at [https://tcga-data.nci.nih.gov/docs/publications/lgggbm\\_2015/](https://tcga-data.nci.nih.gov/docs/publications/lgggbm_2015/).

### Identification of Novel Glioma-Associated Genomic Alterations

To establish the set of genomic alterations that drive gliomagenesis, we called point mutations and indels on the exomes of 513 LGG and 307 GBM using the Mutect, Indelocator, Varscan2, and RADIA algorithms and considered all mutations identified by at least two callers. Significantly mutated genes (SMGs) were determined using MutSigCV. This led to the identification of 75 SMGs, 10 of which had been previously reported in GBM (Brennan et al., 2013), 12 of which had been reported in LGG (Cancer Genome Atlas Research Network et al., 2015), and 8 of which had been identified in both GBM and LGG studies. 45 SMGs have not been previously associated with glioma and ranged in mutation frequency from 0.5% to 2.6% (Table S2A). We used GISTIC2 to analyze the DNA copy number profiles of 1,084 samples, including 513 LGG and 571 GBM, and identified 162 significantly altered DNA copy number segments (Table S2B). We employed PRADA and deFuse to detect 1,144 gene fusion events in the RNA-seq profiles available for 154 GBM and 513 LGG samples, of which 37 in-frame fusions involved receptor tyrosine kinases (Table S2C). Collectively, these analyses recovered all known glioma driving events, including in *IDH1* (n = 457), *TP53* (n = 328), *ATRX* (n = 220), *EGFR* (n = 314), *PTEN* (n = 168), *CIC* (n = 80), and *FUBP1* (n = 45). Notable newly predicted glioma drivers relative to the earlier TCGA analyses were genes associated with chromatin organization such as *SETD2* (n = 24), *ARID2* (n = 20), *DNMT3A* (n = 11), and the *KRAS/NRAS* oncogenes (n = 25 and n = 5, respectively).

We overlapped copy number, mutation (n = 793), and fusion transcript (n = 649) profiles and confirmed the convergence of genetic drivers of glioma into pathways, including the Ras-Raf-MEK-ERK, p53/apoptosis, PI3K/AKT/mTOR, chromatin modification, and cell cycle pathways. The Ras-Raf-MEK-ERK signaling cascade showed alterations in 106 of 119 members

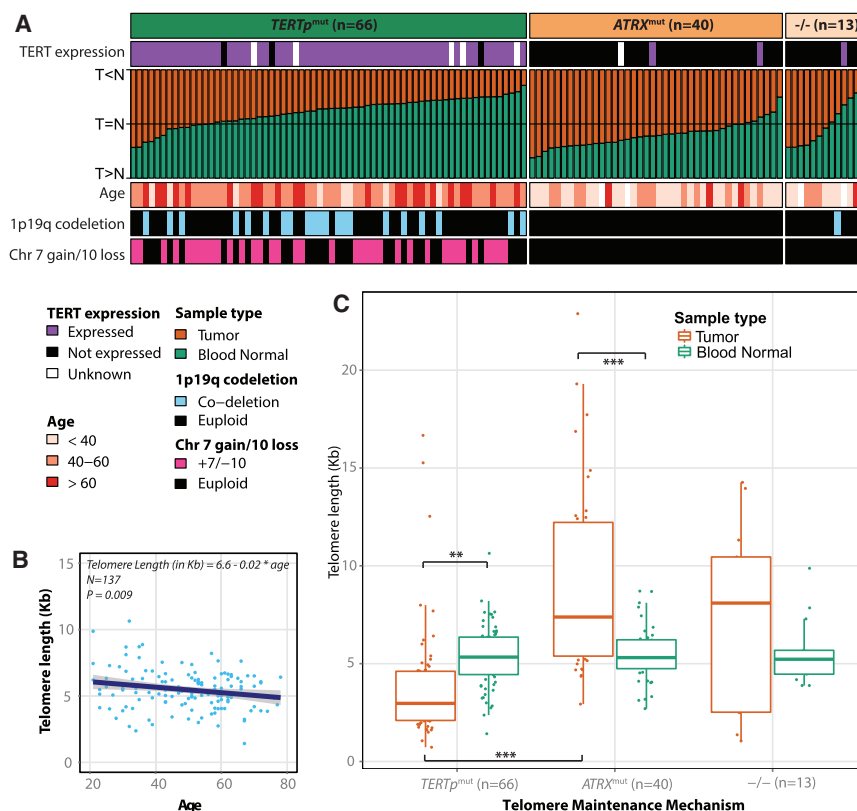
**Table 1. Clinical Characteristics of the Sample Set Arranged by IDH and 1p/19q Co-deletion Status**

Feature	IDH Wt (n = 520)	IDH mut - non-codel (n = 283)	IDH mut - codel (n = 171)	Unknown (n = 148)
<b>Clinical</b>				
<b>Histology (n)</b>				
Astrocytoma	52 (10.0%)	112 (39.6%)	4 (2.3%)	1 (0.7%)
Glioblastoma	419 (80.6%)	32 (11.3%)	2 (1.2%)	137 (92.6%)
Oligoastrocytoma	15 (2.9%)	69 (24.4%)	30 (17.5%)	0 (0%)
Oligodendroglioma	19 (3.7%)	37 (13.1%)	117 (68.4%)	1 (0.7%)
Unknown	15 (2.9%)	33 (11.7%)	18 (10.5%)	9 (6.1%)
<b>Grade (n)</b>				
G2	19 (3.7%)	114 (40.3%)	81 (47.4%)	2 (1.4%)
G3	67 (12.9%)	104 (36.7%)	70 (40.9%)	0 (0%)
G4	419 (80.6%)	32 (11.3%)	2 (1.2%)	137 (92.6%)
Unknown	15 (2.9%)	33 (11.7%)	18 (10.5%)	9 (6.1%)
<b>Age</b>				
Median (LQ-UQ)	59 (51–68)	38 (30–44)	46 (35–54)	55 (48–68)
Unknown (n)	16	33	18	9
<b>Survival</b>				
Median (CI)	14.0 (12.6–15.3)	75.1 (62.1–94.5)	115.8 (90.5–Inf)	12.6 (11.3–14.9)
Unknown (n)	14	32	18	12
<b>KPS</b>				
<70	85 (16.3%)	8 (2.8%)	5 (2.9%)	21 (14.2%)
70–80	196 (37.7%)	41 (14.5%)	18 (10.5%)	60 (40.5%)
90	29 (5.6%)	60 (21.2%)	32 (18.7%)	2 (1.4%)
100	51 (9.8%)	44 (15.9%)	30 (17.5%)	14 (9.5%)
Unknown	159 (30.6%)	129 (45.6%)	86 (50.3%)	51 (34.5%)
<b>Molecular</b>				
<b>MGMT promoter</b>				
Methylated	170 (32.7%)	242 (85.5%)	169 (98.8%)	32 (21.6%)
Unmethylated	248 (47.7%)	36 (12.7%)	1 (0.6%)	34 (23.0%)
Unknown	102 (19.6%)	5 (1.8%)	1 (0.6%)	82 (55.4%)
<b>TERT promoter</b>				
Mutant	67 (12.9%)	8 (2.8%)	86 (50.3%)	1 (0.7%)
Wild-type	19 (9.8%)	146 (51.6%)	2 (1.2%)	0 (0%)
Unknown	434 (83.5%)	129 (45.6%)	83 (48.5%)	135 (99.3%)
<b>TERT expression</b>				
Expressed	178 (34.2%)	14 (4.9%)	153 (89.5%)	6 (4.1%)
Not expressed	51 (9.8%)	242 (85.5%)	16 (9.4%)	7 (4.7%)
Unknown	291 (56.0%)	27 (9.5%)	2 (1.2%)	135 (91.2%)

detected across 578 cases (73%), mostly occurring in IDH-wild-type samples (n = 327 of 357, 92%). Conversely, we found that a set of 36 genes involved in chromatin modification was targeted by genetic alterations in 423 tumors (54%, n = 36 genes), most of which belonged to the IDH mutant-non-codel group (n = 230, 87%).

In order to identify new somatically altered glioma genes, we used MutComFocal to nominate candidates altered by mutation, as well as copy number alteration. Prominent among these genes was *NIPBL*, a crucial adherin subunit that is essential for loading cohesins on chromatin (Table S2D) (Peters and Nish-

iyama, 2012). The cohesin complex is responsible for the adhesion of sister chromatids following DNA replication and is essential to prevent premature chromatid separation and faithful chromosome segregation during mitosis (Peters and Nishiyama, 2012). Alterations in the cohesin pathway have been reported in 12% of acute myeloid leukemias (Kon et al., 2013). Mutations of the cohesin complex gene *STAG2* had been previously reported in GBM (Brennan et al., 2013). Taken together, 16% of the LGG/GBM showed mutations and/or CNAs in multiple genes involved in the cohesin complex, thus nominating this process as a prominent pathway involved in gliomagenesis.



**Figure 1. Telomere Length Associations in Glioma**

(A) Heatmap of relative tumor/normal telomere lengths of 119 gliomas, grouped by *TERTp* and *ATRX* mutation status.

(B) Telomere length decreases with increasing age (measured in years at diagnosis) in blood normal control samples (n = 137).

(C) Quantitative telomere length estimates of tumors and blood normal, grouped by *TERTp* mutant (n = 67, 56%), *ATRX* mutant (n = 40, 33%), and double negative (n = 13, 11%) status. \*\*\* = p < 0.0001; \*\* = p < 0.001.

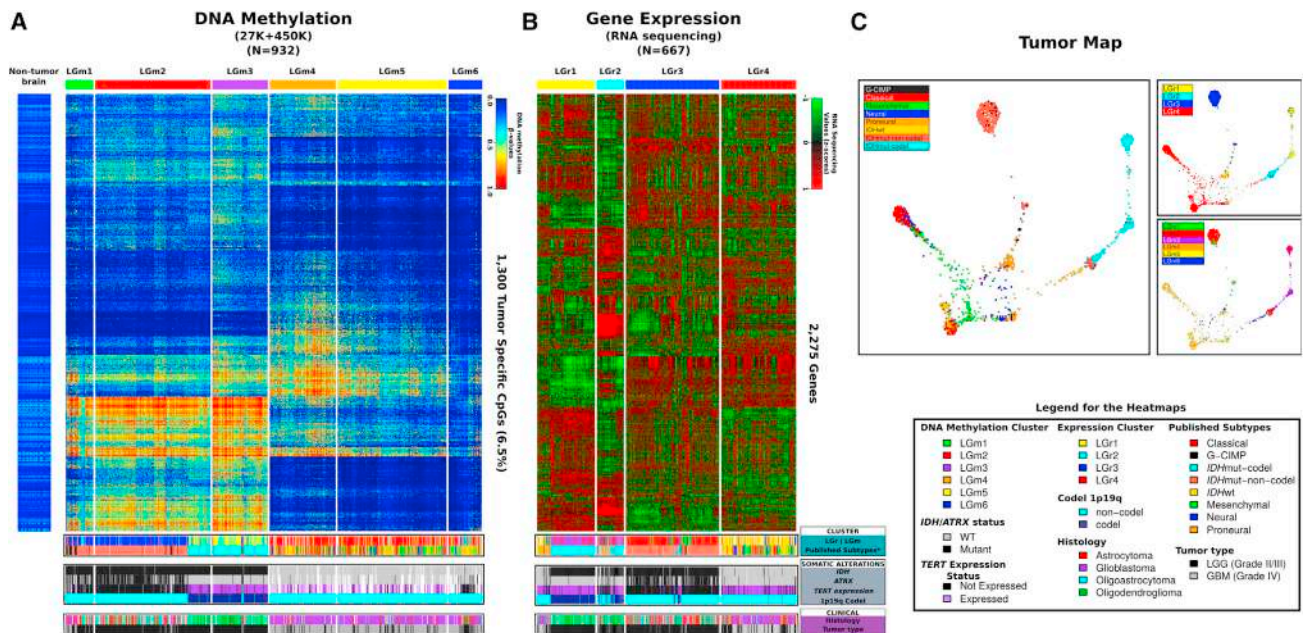
### Telomere Length Is Positively Correlated with *ATRX*, but Not *TERT* Promoter Mutations

Mutations in the *TERT* promoter (*TERTp*) have been reported in 80% of GBM (Killela et al., 2013). We used *TERTp* mutation calls from targeted sequencing (n = 287) and complemented them with *TERTp* mutations inferred from whole-genome sequencing (WGS) data (n = 42). *TERTp* mutations are nearly mutually exclusive with mutations in *ATRX* (Eckel-Passow et al., 2015), which was confirmed in our cohort. Overall, 85% of diffuse gliomas harbored mutations of *TERTp* (n = 157, 48%) or *ATRX* (n = 120, 37%). *TERTp* mutations activate *TERT* mRNA expression through the creation of a de novo E26 transformation-specific (ETS) transcription factor-binding site (Horn et al., 2013), and we observed significant *TERT* upregulation in *TERTp* mutant cases (p value < 0.0001, Figure S1A). *TERT* expression measured by RNA-seq was a highly sensitive (91%) and specific (95%) surrogate for the presence of *TERTp* mutation (Figure S1B). We correlated *TERTp* status with glioma driving alterations and observed that nearly all IDH-wild-type cases with chromosome 7 gain and chromosome 10 loss harbored *TERTp* mutations or upregulated *TERT* expression (n = 52/53 and n = 134/147, respectively; Figure 1A). Conversely, only 45% of IDH-wild-type samples lacking chromosome 7/chromosome 10 events showed *TERTp* mutations or elevated *TERT* expression (n = 15/33 and n = 43/82, respectively). Thus, *TERTp* mutations may precede the chr 7/chr 10 alterations that have been implicated in glioma initiation (Ozawa et al., 2014).

To correlate *TERTp* mutations to telomere length, we used whole-genome sequencing and low pass whole-genome sequencing data to estimate telomere length in 141 pairs of matched tumor and normal samples. As expected, we observed an inverse correlation of telomere length with age at diagnosis in matching blood normal samples (Figure 1B) and tumor samples (Figure S1C). Glioma samples harboring *ATRX* mutations showed significantly longer telomeres compared to *TERTp* mutant samples (t test p value < 0.0001; Figure 1C). Among *TERTp* mutation gliomas, there

was no difference in telomere length between samples with and without additional *IDH1/IDH2* mutations, despite a difference in age. *ATRX* forms a complex with DAXX and H3.3, and the genes encoding these proteins are frequently mutated in pediatric gliomas (Sturm et al., 2012). Mutations in DAXX and H3F3A were identified in only two samples in our WGS dataset. The *ATRX*-DAXX-H3.3 complex is associated with the alternative lengthening of telomeres (ALT) and our observations confirm previously hypothesized fundamental differences between the telomere control exerted by telomerase and ALT (Sturm et al., 2014).

As demonstrated by the identification of *TERTp* mutations, somatic variants affecting regulatory regions may play a role in gliomagenesis. Using 67 matched whole-genome and RNA-seq expression pairs, we similarly sought to identify mutations located within 2 kb upstream of transcription start sites and associated with a gene expression change. Using strict filtering methods, we identified 12 promoter regions with mutations in at least 6 samples. Three of 12 regions related to a significant difference in the expression of the associated gene expression, suggesting possible functional consequences. Other than *TERT* (n = 37), promoter mutations of the ubiquitin ligase *TRIM28* (n = 8) and the calcium channel gamma subunit *CACNG6* (n = 7) correlated with respectively upregulation and downregulation of these genes, respectively (Table S2E). *TRIM28* has been reported to mediate the ubiquitin-dependent degradation of AMP-activated protein kinase (AMPK) leading to activation of mTOR signaling and hypersensitization to AMPK agonists, such as metformin (Pineda et al., 2015).



**Figure 2. Pan-glioma DNA Methylation and Transcriptome Subtypes**

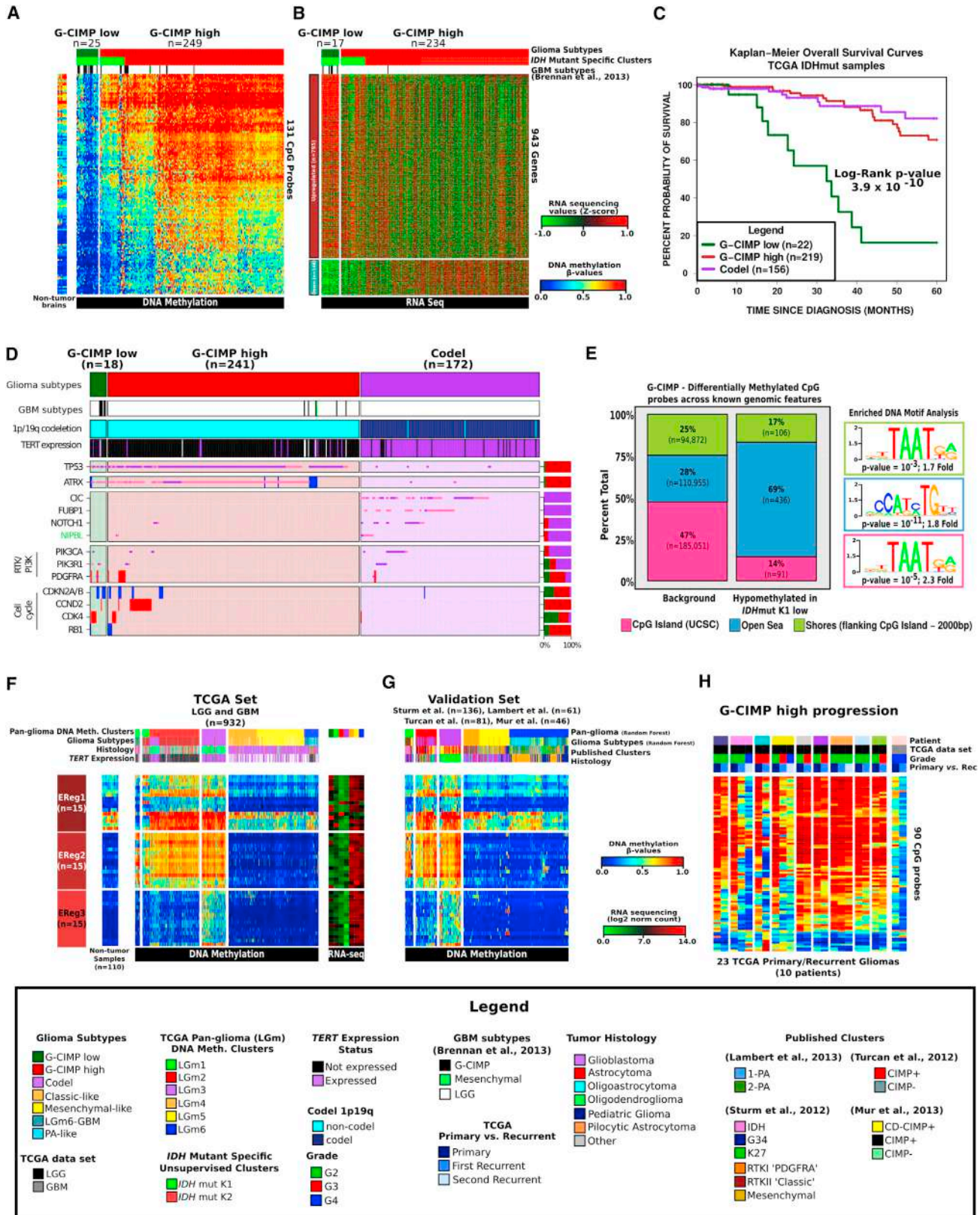
(A) Heatmap of DNA methylation data. Columns represent 932 TCGA glioma samples grouped according to unsupervised cluster analysis; rows represent DNA methylation probes sorted by hierarchical clustering. Non-neoplastic samples are represented on the left of the heatmap ( $n = 77$ ) (Guintivano et al., 2013). (B) Heatmap of RNA sequencing data. Unsupervised clustering analysis for 667 TCGA glioma samples profiled using RNA sequencing are plotted in the heatmap using 2,275 most variant genes. Previously published subtypes were derived from Brennan et al. (2013) and Cancer Genome Atlas Research Network et al., 2015. (C) Tumor Map based on mRNA expression and DNA methylation data. Each data point is a TCGA sample colored coded according to their identified status. A live interactive version of this map is available at <http://tumormap.ucsc.edu/?p=ynewton.gliomas-paper>.

### Unsupervised Clustering of Gliomas Identifies Six Methylation Groups and Four RNA Expression Groups Associated with IDH Status

To segregate the DNA methylation subtypes across the pan-glioma dataset, we analyzed 932 glioma samples profiled on the HumanMethylation450 platform (516 LGG and 129 GBM) and the HumanMethylation27 platform (287 GBM). In order to incorporate the maximum number of samples, we merged datasets from both methylation platforms yielding a core set of 25,978 CpG probes. To reduce computational requirements to cluster this large dataset, we eliminated sites that were methylated (mean  $\beta$  value  $\geq 0.3$ ) in non-tumor brain tissues and selected 1,300 tumor-specific methylated probes (1,300/25,978, 5%) to perform unsupervised k-means consensus clustering. This identified six distinct clusters, labeled LGM1–6 (Figure 2A and Tables S1 and S3A). Next, we sought to determine pan-glioma expression subtypes through unsupervised clustering analysis of 667 RNA-seq profiles (513 LGG and 154 GBM), which resulted in four main clusters labeled LGR1–4 (Figure 2B and Tables S1 and S3A). An additional 378 GBM samples with Affymetrix HT-HG-U133A profiles (but lacking RNA-seq data) were classified into the four clusters using a k-nearest neighbor classification procedure. *IDH* mutation status was the primary driver of methylation and transcriptome clustering and separated the cohort into two macro-groups. The LGM1/LGM2/LGM3 DNA methylation macro-group carried *IDH1* or *IDH2* mutations (449 of 450, 99%) and was enriched for LGG (421/454, 93%) while LGM4/

LGM5/LGM6 were *IDH*-wild-type (429/430, 99%) and enriched for GBM (383/478, 80%). LGM1–3 showed genome-wide hypermethylation compared to LGM4–6 clusters (Figure S2A), documenting the association between *IDH* mutation and increased DNA methylation (Noushmehr et al., 2010; Turcan et al., 2012). Principal component analysis using 19,520 probes yielded similar results, thus emphasizing that our probe selection method did not introduce unwanted bias (Figure S2B). The gene expression clusters LGR1–3 harbored *IDH1* or *IDH2* mutations (438 of 533, 82%) and were enriched for LGG (436/563, 77%), while the LGR4 was exclusively *IDH*-wild-type (376 of 387, 97%) and enriched for GBM (399/476, 84%).

We extended our analysis using Tumor Map (Supplemental Experimental Procedures) to perform integrated co-clustering analysis of the combined gene expression ( $n = 1,196$ ) and DNA methylation ( $n = 867$ ) profiles. An interactive Tumor Map version is publicly available at <http://tumormap.ucsc.edu/?p=ynewton.gliomas-paper>. Tumor Map assigns samples to a hexagon in a grid so that nearby samples are likely to have similar genomic profiles and allows visualizing complex relationships between heterogeneous genomic data samples and their clinical or phenotypical associations. Thus, clusters in the map indicate groups of samples with high similarity of integrated gene expression and DNA methylation profiles (Figure 2C). The map confirms clustering by *IDH* status and additionally shows islands of samples that share previously reported GBM cluster memberships (Noushmehr et al., 2010; Verhaak et al., 2010). To assess



(legend on next page)

clustering sensitivity to pre-processing, we tried complementary methods and obtained similar results (Figure S2C).

To identify genes whose copy number changes are associated with concordant changes in gene expression, we combined expression and copy number profiles from 659 samples to define a signature of 57 genes with strong functional copy number (fCN) change (Table S3B). The fCN signature clustered gliomas into three macro-clusters, LGfc1–3, strongly associated with IDH and 1p/19q status (Figure S2D). The fCN analysis revealed the functional activation of a cluster of *HOXA* genes in the IDH-wild-type LGfc2 cluster, which were previously associated with glioma stem cell maintenance (Kurscheid et al., 2015).

Finally, we clustered reverse phase protein array profiles, consisting of 196 antibodies on 473 samples. Two macro clusters were observed, and in contrast to the transcriptome/methylome/fCNV clustering, the primary discriminator was based on glioma grade (LGG versus GBM) rather than IDH status (Figure S2E). Compared to the LGG-like cluster, the GBM-like cluster had elevated expression of IGFBP2, fibronectin, PAI1, HSP70, EGFR, phosphoEGFR, phosphoAKT, Cyclin B1, Caveolin, Collagen VI, Annexin1, and ASNS, whereas the LGG class showed increased activity of PKC (alpha, beta, and delta), PTEN, BRAF, and phosphoP70S6K.

The above results confirm IDH status as the major determinant of the molecular footprints of diffuse glioma. To further elucidate the subtypes of diffuse glioma, we performed unsupervised clustering within each of the two IDH-driven macroclusters. We used 1,308 tumor-specific CpG probes defined among the IDH mutation cohort ( $n = 450$ ) and identified three IDH mutant-specific DNA methylation clusters (Figure S3A). Using 914 tumor-specific CpG probes in the IDH-wild-type cohort ( $n = 430$ ), we uncovered three IDH-wild-type-specific clusters (Figure S4A). The sets of CpG probes used to cluster each of the two IDH-driven datasets overlapped significantly with the 1,300 probes that defined the pan-glioma DNA methylation clustering (1162/1,300, 89% and 853/1,300, 66%, for IDH mutant and IDH-wild-type, respectively). The clusters identified by separating IDH mutant and IDH-wild-type gliomas showed strong overall concordance with pan-glioma DNA methylation subtypes (Table S3A). Similarly, unsupervised clustering of 426 IDH mutant RNA-seq profiles resulted in three subtypes (Figure S3A), and analysis of the 234 IDH-wild-type samples led to four mixed LGG/GBM clusters that showed enrichment for previously identified GBM expression subtypes (Figure S4C) (Verhaak et al., 2010).

### An Epigenetic Signature Associated with Activation of Cell Cycle Genes Segregates a Subgroup of IDH Mutant LGG and GBM with Unfavorable Clinical Outcome

The three epigenetic subtypes defined by clustering IDH mutant glioma separated samples harboring the 1p/19q co-deletion into a single cluster and non-codel glioma into two clusters (Figure S3A). Conversely, non-codel glioma grouped nearly exclusively into a single expression cluster, and codels were split in two separated expression clusters (Figure S3A). A distinct subgroup of samples within the IDH mutant-non-codel DNA methylation clusters manifested relatively reduced DNA methylation (Figure S3B). The unsupervised clustering of IDH mutant glioma was unable to segregate the lower methylated non-codel subgroup as the 1,308 probes selected for unsupervised clustering included only 19 of the 131 differentially methylated probes characteristic for this subgroup ( $FDR < 10^{-15}$ , difference in mean methylation beta value  $> 0.27$ ). The low-methylation subgroup consisted of both G-CIMP GBM (13/25) and LGGs (12/25) and was confirmed using a non-TCGA dataset (Figure S3C). The tumors with higher methylation in the split cluster were very similar to those grouped in the second non-codel cluster, and a supervised comparison identified only 12 probes as differentially DNA methylated (Figures 3A and 3B). We concluded that IDH mutant glioma is composed of three coherent subgroups: (1) the Codel group, consisting of IDH mutant-codel LGGs; (2) the G-CIMP-low group, including IDH mutant-non-codel glioma (LGG and GBM) manifesting relatively low genome-wide DNA methylation; and (3) the G-CIMP-high group, including IDH mutant-non-codel glioma (LGG and GBM) with higher global levels of DNA methylation. The newly identified G-CIMP-low group of glioma was associated with significantly worse survival as compared to the G-CIMP-high and Codel groups (Figure S3D). The clinical outcome of the tumors classified as G-CIMP-high was as favorable as that of Codel tumors, the subgroup generally thought to have the best prognosis among glioma patients (Figures 3C and S3D). We compared the frequencies of glioma driver gene alterations between the three types of IDH mutant glioma and found that 15 of 18 G-CIMP-low cases carried abnormalities in cell cycle pathway genes such as *CDK4* and *CDKN2A*, relative to 36/241 and 2/172 for G-CIMP-high and Codels, respectively (Figure 3D). Supervised analysis between gene expression of G-CIMP-low and G-CIMP-high resulted in 943 differentially expressed genes. We mapped the 943 deregulated genes to 767 nearest CpG probes (max distance 1 kb) and found the majority

#### Figure 3. Identification of a Distinct G-CIMP Subtype Defined by Epigenomics

- (A) Heatmap of probes differentially methylated between the two IDH mutant-non-codel DNA methylation clusters allowed the identification of a low-methylation subgroup named G-CIMP-low. Non-tumor brain samples ( $n = 12$ ) are represented on the left of the heatmap.
- (B) Heatmap of genes differentially expressed between the two IDH mutant-non-codel DNA methylation clusters.
- (C) Kaplan-Meier survival curves of IDH mutant methylation subtypes. Ticks represent censored values.
- (D) Distribution of genomic alterations in genes frequently altered in IDH mutant glioma.
- (E) Genomic distribution of 633 CpG probes differentially demethylated between co-clustered G-CIMP-low and G-CIMP-high. CpG probes are grouped by UCSC genome browser-defined CpG Islands, shores flanking CpG island  $\pm 2$  kb and open seas (regions not in CpG islands or shores).
- (F) DNA methylation heatmap of TCGA glioma samples ordered per Figure 2A and the epigenetically regulated (EReg) gene signatures defined for G-CIMP-low, G-CIMP-high, and Codel subtypes. The mean RNA sequencing counts for each gene matched to the promoter of the identified cgID across each cluster are plotted to the right.
- (G) Heatmap of the validation set classified using the random forest method applying the 1,300 probes defined in Figure 2A.
- (H) Heatmap of probes differentially methylated between G-CIMP-low and G-CIMP-high in longitudinally matched tumor samples.



of the CpG probes (486/767, 63%) to show a significant methylation difference (FDR < 0.05, difference in mean methylation beta value > 0.01) between G-CIMP-low and G-CIMP-high, suggesting a mechanistic relation between loss of methylation and increased transcript levels.

Recent analysis of epigenetic profiles derived from colon cancers showed that transcription factors may bind to regions of demethylated DNA (Berman et al., 2012). Therefore, we asked whether transcription factors may be recruited to the DNA regions differentially methylated between G-CIMP-low samples and G-CIMP-high samples from the same methylation cluster, using 450K methylation profiles (n = 39). Globally, we detected 643 differentially methylated probes between 27 G-CIMP-low and 12 G-CIMP-high samples (absolute diff-mean difference  $\geq 0.25$ , FDR  $\leq 5\%$ ). Most of these probes (69%) were located outside of any known CpG island but positioned within intergenic regions known as open seas (Figure 3E). This represents a 2.5-fold open sea enrichment compared to the expected genome-wide distribution of 450K CpG probes (chi-square p value <  $2.2 \times 10^{-16}$ ). We also observed a 3.4-fold depletion within CpG islands (chi-square p value <  $2.2 \times 10^{-16}$ ).

Using this set of intergenic CpG probes, we asked whether a DNA motif signature associated with distal regulatory elements. Such a pattern would point to candidate transcription factors involved in tumorigenesis of the G-CIMP-low group. A de novo motif scan and known motif scan identified a distinct motif signature TGTT (geometric test p value =  $10^{-11}$ , fold enrichment = 1.8), known to be associated with the OLIG2 and SOX transcription factor families (Figure 3E) (Lodato et al., 2013). This observation was corroborated by the higher expression levels of SOX2, as well as 17 out of 20 other known SOX family members in G-CIMP-low compared to G-CIMP-high (fold difference > 2). The primary function of SOX2 in the nervous system is to promote self-renewal of neural stem cells and, within brain tumors, the glioma stem cell state (Graham et al., 2003). Interestingly, SOX2 and OLIG2 have been described as neurodevelopmental transcription factors being essential for GBM propagation (Suvà et al., 2014). Supervised gene expression pathway analysis of the genes activated in the G-CIMP-low group as opposed to G-CIMP-high group revealed activation of genes involved in cell cycle and cell division consistent with the role of SOX in promoting cell proliferation (Figure S3E). The enrichment in cell cycle gene expression provides additional support to the notion that development of the G-CIMP-low subtype is associated with activation of cell cycle progression and may be mediated by a loss of CpG methylation and binding of SOX factors to candidate genomic enhancer elements.

To validate the G-CIMP-low, G-CIMP-high, and Codel IDH mutant subtypes, we compiled a validation cohort from published studies, including 324 adult and pediatric gliomas (Lambert et al., 2013; Mur et al., 2013; Sturm et al., 2012; Turcan et al., 2012). The CpG probe methylation signatures used to classify the validation set are provided on the publication portal accompanying this publication ([https://tcga-data.nci.nih.gov/docs/publications/lgggbm\\_2015/](https://tcga-data.nci.nih.gov/docs/publications/lgggbm_2015/)). Among them, 103 were identified as IDH mutant on the basis of their genome-wide DNA methylation profile. We classified samples in the validation set using the probes that defined the IDH mutant-specific DNA

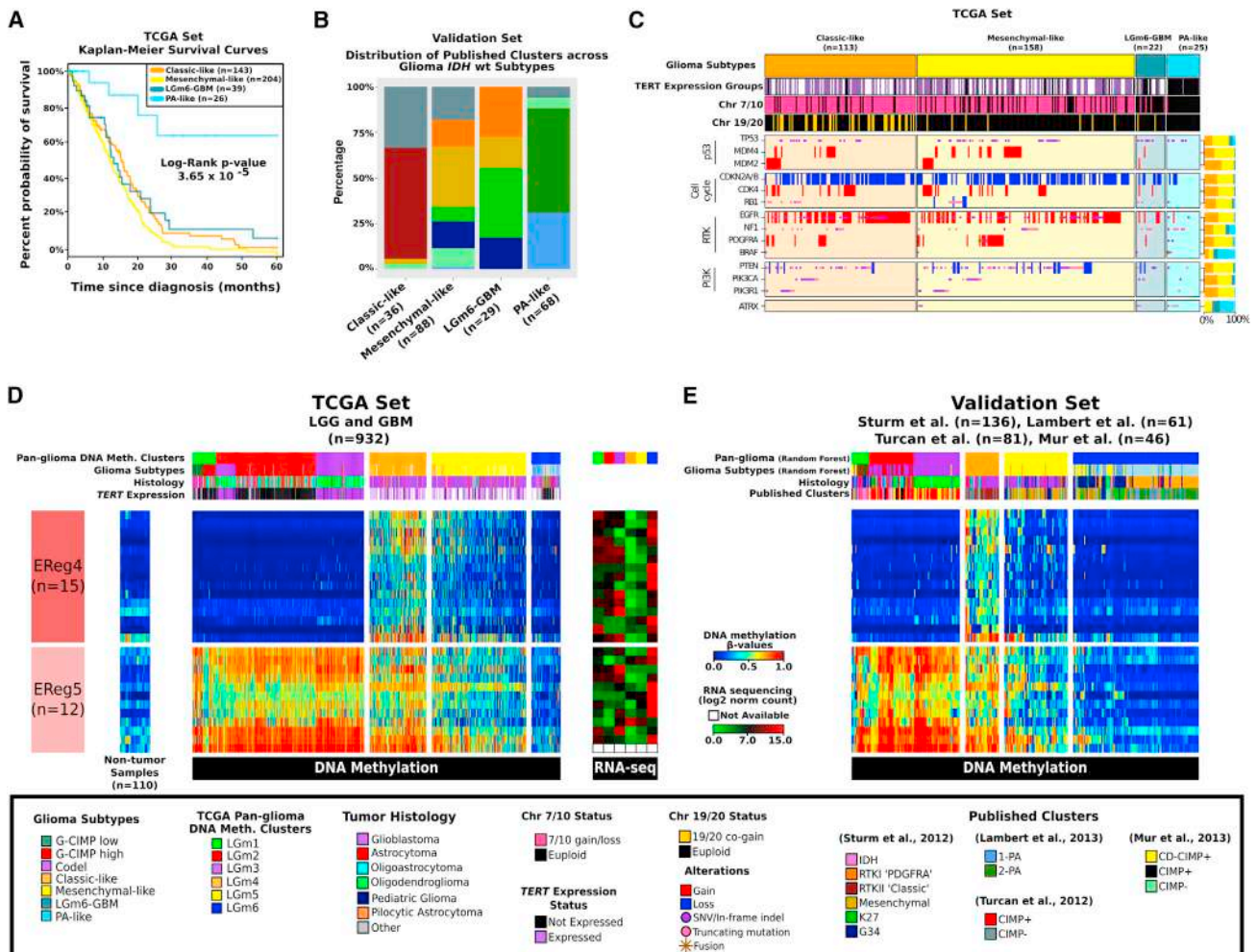
methylation cluster analysis integrated in a supervised random forest method. The analysis recapitulated the clusters generated from the TCGA collection (Figure S3C). In order to determine epigenetically regulated (EReg) genes that may be characteristic of the biology of the IDH mutant diffuse glioma subtypes, we compared 450k methylation DNA methylation profiles and gene expression levels between 636 IDH mutant and IDH-wild-type gliomas and 110 non-tumor samples from 11 different tissue types. From the list of epigenetically regulated genes, we extracted 263 genes that were grouped into EReg gene signatures, which showed differential signals among the three IDH mutant subtypes (Figure 3F). These trends were confirmed in the validation set (Figure 3G).

We investigated the possibility that the G-CIMP-high group is a predecessor to the G-CIMP-low group by comparing the DNA methylation profiles from ten IDH mutant-non-codel LGG and GBM primary-recurrent cases with the TCGA cohort. We evaluated the DNA methylation status of probes identified as differentially methylated (n = 90) between G-CIMP-low and G-CIMP-high (FDR <  $10^{-13}$ , difference in mean methylation beta-value > 0.3 and < -0.4). Four out of ten IDHmut-non-codel cases showed a demethylation pattern after disease recurrence, while partial demethylation was demonstrated in the remaining six recurrences, supporting the notion of a progression from G-CIMP-high to G-CIMP-low phenotype (Figure 3H).

### An IDH-Wild-Type Subgroup of Histologically Defined Diffuse Glioma Is Associated with Favorable Survival and Shares Epigenomic and Genomic Features with Pilocytic Astrocytoma

IDH-wild-type gliomas segregated into three DNA methylation clusters (Figure S4A). The first is enriched with tumors belonging to the classical gene expression signature and was labeled Classic-like, whereas the second group, enriched with mesenchymal subtype tumors, was labeled Mesenchymal-like (Table S1) (Verhaak et al., 2010). The third cluster contained a larger fraction of LGG in comparison to the other IDH-wild-type clusters. We observed that the IDH-wild-type LGGs but not the IDH-wild-type GBM in this cluster displayed markedly longer survival (log-rank p value =  $3.6 \times 10^{-5}$ ; Figure 4A) and occurred in younger patients (mean 37.6 years versus 50.8 years, t test p value = 0.002). Supervised analysis of differential methylation between LGG and GBM in the third DNA methylation cluster did not reveal any significant probes despite significant differences in stromal content (p value < 0.005; Figure S4D), suggesting that this group cannot be further separated using CpG methylation markers.

Next, we sought to validate the methylation-based classification of IDH-wild-type glioma in an independent cohort of 221 predicted IDH-wild-type glioma samples, including 61 grade I pilocytic astrocytomas (PAs). Toward this aim, we used a supervised random forest model built with the probes that defined the IDH-wild-type clusters. Samples classified as Mesenchymal-like showed enrichment for the Sturm et al. (2012) Mesenchymal subtype (29/88), and gliomas predicted as Classic-like were all RTK II "Classic" (22/22), per the Sturm et al. (2012) classification (Figures 4B and S4B). We observed that PA tumors were unanimously classified as the third,



**Figure 4. A Distinct Subgroup of IDH-Wild-Type Diffuse Glioma with Molecular Features of Pilocytic Astrocytoma**

(A) Kaplan-Meier survival curves for the IDH-wild-type glioma subtypes. Ticks represent censorship. (B) Distribution of previous published DNA methylation subtypes in the validation set, across the TCGA IDH-wild-type-specific DNA methylation clusters. (C) Distribution of genomic alterations in genes frequently altered in IDH-wild-type glioma. (D) Heatmap of TCGA glioma samples ordered according to Figure 2A and two EReg gene signatures defined for the IDH-wild-type DNA methylation clusters. Mean RNA sequencing counts for each gene matched to the promoter of the identified cgID across each cluster are plotted to the right. (E) Heatmap of the validation set classified using the random forest method using the 1,300 probes defined in Figure 2A.

LGG-enriched group (Figure S4B). Based on the molecular similarity with PA, we labeled the LGGs in the third methylation cluster of IDH-wild-type tumors as PA-like. The GBMs in this group were best described as LGM6-GBM for their original pan-glioma methylation cluster assignment and tumor grade.

Pilocytic astrocytomas are characterized by frequent alterations in the MAPK pathway, such as *FGFR1* mutations, *KIAA1549-BRAF*, and *NTRK2* fusions (Jones et al., 2013). The frequency of mutations, fusions, and amplifications in eight PA-associated genes (*BRAF*, *NF1*, *NTRK1*, *NTRK2*, *FGFR1*, and *FGFR2*) rated from 11% (n = 12/113) of Classic-like, 13% (n = 21/158) of Mesenchymal-like IDH-wild-type tumors to 32% (n = 7/22) of LGM6-GBM and 52% (n = 13/25) of PA-like LGG (Fisher's exact test [FET] p value < 0.0001; Figure 4C). Conversely, only 2 of 25 (8%) PA-like LGG tumors showed

*TERT* expression, compared to 5 of 12 LGM6-GBM (43%), 60 of 65 Classic-like (92%), and 82 of 98 Mesenchymal-like (84%, FET p value < 0.0001). The PA-like group was characterized by relatively low frequency of typical GBM alterations in genes such as *EGFR*, *CDKN2A/B*, and *PTEN* and displayed euploid DNA copy number profiles (Figure S4E). To ascertain that the histologies of the PA-like subgroup had been appropriately classified, we conducted an independent re-review. This analysis confirmed the presence of the histologic features of diffuse glioma (grade II or grade III) in 23 of the 26 cases in the cluster. The remaining three cases were re-named as PA (grade I). An independent review of the magnetic resonance diagnostic images from 13 cases showed a similar pattern, with the majority of tumors showing behavior consistent with grade II or grade III glioma. Taken together, the epigenetic analysis of the

**Table 2. DNA Methylation Subtypes Are Prognostically Relevant in Multivariable Analysis and in External Validation Data**

		Discovery (n = 809)			Validation (n = 183)		
		C-Index: 0.835 ± 0.019			C-Index: 0.745 ± 0.032		
Predictor	Levels	n	HR (95% CI)	Signif.	n	HR (95% CI)	Signif.
Age at diagnosis	per year	809	1.05 (1.03–1.06)	***	183	1.02 (1–1.04)	*
WHO Grade	II	214	1.0 (ref)		41	1.0 (ref)	
	III	241	1.96 (1.15–3.33)	*	51	1.24 (0.55–2.76)	
	IV	354	2.38 (1.3–4.34)	*	91	2.6 (1.08–6.3)	*
Subgroup	IDHmut-codel	156	1.0 (ref)		57	1.0 (ref)	
	G-CIMP-low	22	5.6 (2.49–12.62)	***	2	0 (0–Inf)	
	G-CIMP-high	219	1.92 (1.05–3.51)	*	15	1.25 (0.43–3.66)	
	classic-like	143	5.4 (2.79–10.44)	***	22	4.55 (1.8–11.49)	*
	mesenchymal-like I	204	8.71 (4.59–16.53)	***	61	5.55 (2.52–12.21)	***
	LGM6-GBM	39	5.79 (2.78–12.1)	***	22	6.8 (2.58–17.91)	**
PA-like	26	2.02 (0.71–5.71)		4	3.64 (0.79–16.78)	.	

Survival regression analysis indicates that an optimal model of prognosis includes age, grade, and methylation subtype. These predictors are statistically significant in both our discovery dataset and an external validation dataset. Significance codes: 0 “\*\*\*\*”; 0.001 “\*\*\*”; 0.05 “\*\*”; 0.1 “.”

IDH-wild-type group of adult glioma revealed the existence of a novel subgroup sharing genetic and DNA methylation features with pediatric PA and favorable clinical outcome compared to diffuse IDH-wild-type glioma. This group may include but extends beyond *BRAF*-mutated grade II oligodendroglioma that were previously recognized as a unique clinical entity (Chi et al., 2013).

Through comparison of the methylation profiles of 636 glioma and 110 non-neoplastic normal samples from different tissue types, we defined EReg signatures consisting of 27 genes that showed differential signals among IDH-wild-type subtypes in the TCGA (Figure 4D) and the validation set (Figure 4E). EReg4 comprised a group of 15 genes hypermethylated and downregulated in particularly Classic-like. EReg5 was defined as a group of 12 genes associated with hypomethylation in LGM6/PA-like compared to all other LGM clusters. These ERegs aided in characterizing the biological importance of IDH-wild-type subtypes and were subsequently used to evaluate the prognostic importance of the IDH-wild-type clusters.

### The Epigenetic Classification of Glioma Provides Prognostic Value Independent of Age and Grade

In order to assess whether the DNA methylation-based subtypes we identified carry prognostically relevant information independent of known overall survival predictors, we constructed a series of survival regression models. To find the optimal model for survival prediction, we studied covariates individually and in combination with other covariates. Age at diagnosis, histology, IDH/codel subtype, *TERT* expression, and epigenetic subtype all contribute to survival in single-predictor analysis (log-rank p value < 0.05, Table S4). As expected, age was a highly significant predictor (p < 0.0001, C-Index 0.78) and was included in all subsequent multi-predictor models. We found that histology and grade are highly correlated. Histology provided only marginal improvement to a model that includes grade (likelihood ratio test [LRT] p value = 0.08) and was therefore not included in further analyses. Conversely, grade markedly impacted a histol-

ogy-based predictor model (LRT p value = 0.0005, Table S4) and was retained in the subsequent models. In contrast to previous reports (Eckel-Passow et al., 2015), we failed to observe a statistically significant and independent survival association with *TERT* expression (LRT p value = 0.82, Table S4) or *TERT*p mutations after accounting for age and grade (LRT p value = 0.85, data not shown). Thus, the optimal survival prediction model includes age, grade, and epigenetic subtype (LRT p value < 0.0001, C-Index 0.836; Table 2).

To confirm that the epigenetic subtypes provide independent prognostic information, we tested the survival model on the validation dataset. Epigenetic subtypes in these samples were determined as described above. The distinction between LGM6-GBM and PA-like gliomas was made on the basis of tumor grade and not by DNA methylation signature. Using a subset of 183 samples in the validation cohort with known survival, age, and grade, we found that epigenetic subtypes are significant independent predictors of survival in the multivariate analysis (LRT p value < 0.0001, C-Index 0.746, Table 2). This generalization of our model supports the epigenetic subtypes as a means to improve the prognostication of glioma.

### Activation of Cell Cycle/Proliferation and Invasion/Microenvironmental Changes Marks Progression of LGG to GBM

We observed that, in spite of morphological differences between LGG and GBM, such as high cell density and microvascular proliferation, clustering of gene expression profiles frequently grouped LGG and GBM together within the same subtype. Gene Set Enrichment Analysis of the genes activated in G-CIMP GBM as opposed to the IDH mutant-non-codel within LGr3 (Figure 2B) revealed four major groups, including cell cycle and hyperproliferation, DNA metabolic processes, response to stress, and angiogenesis (Figure S5A and Table S5). These biological functions are consistent with the criteria based on mitotic index used by pathologists to discriminate lower and high-grade glioma and the significance of activated microglia for tumor

aggressiveness (Roggendorf et al., 1996). Conversely, compared with the G-CIMP GBM, IDH mutant-non-codel LGG in LGr3 were characterized by enrichment of genes associated with neuro-glial functions such as ion transport and synaptic transmission, possibly suggesting a more differentiated nature. The comparison of co-clustered GBM and LGG in LGr3 by the PARADIGM algorithm that integrates DNA copy number and gene expression to infer pathway activity confirmed that GBMs express genes associated with cell cycle, proliferation, and aggressive phenotype through activation of a number of cell cycle, cell replication, and NOTCH signaling pathways whereas LGGs exhibit an enrichment of neuronal-differentiation-specific categories, including synaptic pathways (Figure S5C and Table S5).

The analysis of the genes activated in GBM versus the LGG component of LGr4, which grouped IDH-wild-type tumors, identified an inflammation and immunologic response signature characterized by the activation of several chemokines (*CCL18*, *CXCL13*, *CXCL2*, and *CXCL3*) and interleukins (*IL8* and *CXCR2*) enriching sets involved in inflammatory and immune response, negative regulation of apoptosis, cell cycle and proliferation, and the I $\kappa$ B/NF $\kappa$ B kinase cascade Map (Figure S5B and Table S5). These characteristics suggest differences in the relative amount of microglia. We used the ESTIMATE method to estimate the relative presence of stromal cells, which revealed significantly lower (p value  $10^{-6}$ ) stromal scores of LGG IDH-wild-type versus GBM IDH-wild-type (Figure S5F) (Yoshihara et al., 2013). Resembling the functional enrichment for LGG within LGr3, functional enrichment of LGG IDH-wild-type in comparison to GBM within LGr4 showed activation in LGG of special glial-neuronal functions involved in ion transport, synaptic transmission, and nervous system development.

Finally, we aimed to identify transcription factors that may exert control over prominent gene expression programs, known as master regulators. Master regulator analysis comparing the IDH-wild-type group to the IDH mutant group revealed transcription factors that were upregulated in IDH-wild-type gliomas and showed an increase in expression of target genes, including *NKX2-5*, *FOSL1*, *ETV4*, *ETV7*, *RUNX1*, *CEBPB*, *NFE2L3*, *ELF4*, *RUNX3*, *NR2F2*, *PAX8*, and *IRF1* (Table S5). No transcription factors (TFs) were found to be upregulated in IDH mutant gliomas relative to IDH-wild-type gliomas (at a log fold change > 1).

## DISCUSSION

This study represents the largest multi-platform genomic analysis performed to date of adult diffuse glioma (WHO grades II, III, and IV). A simplified graphical summary of the identified groups and their main clinical and biological characteristics is reported in Figure 5. The clustering of all diffuse glioma classes and grades within similarly shaped methylation-based and expression-based groups has allowed us to pinpoint specific molecular signatures with clinical relevance. The DNA methylation classification proposed should be considered as a basis and it is likely that future studies involving significantly larger cohorts and more refined profiling methods will be able to further reduce intra-subtype heterogeneity. The dissection of the IDH mutant non-codel G-CIMP LGG and GBM into two separate subgroups (G-CIMP-

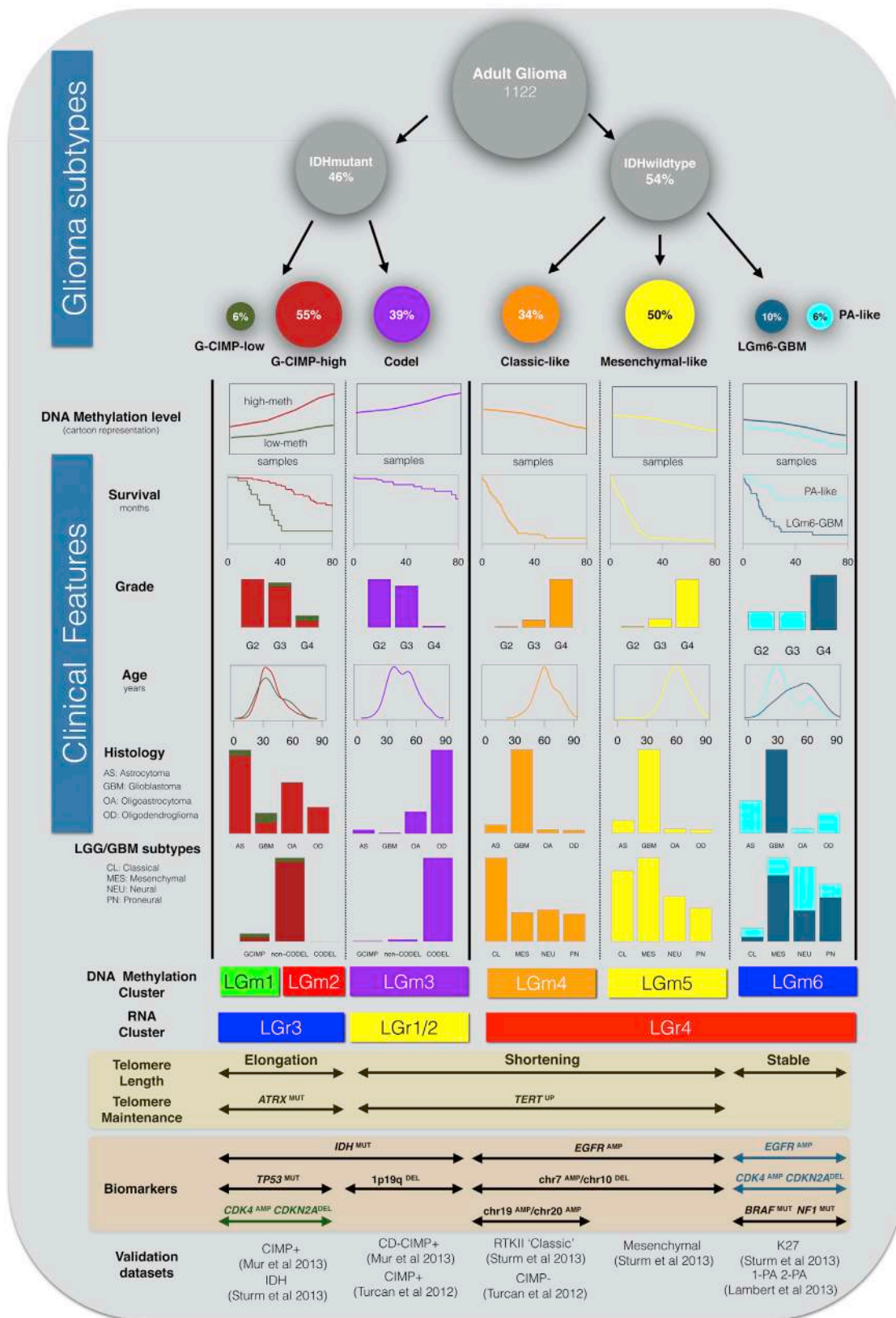
low and G-CIMP-high) based on the extent of genome-wide DNA methylation has crucial biological and clinical relevance. In particular, the identification of the G-CIMP-low subset, characterized by activation of cell cycle genes mediated by SOX binding at hypomethylated functional genomic elements and unfavorable clinical outcome, is an important finding that will guide more accurate segregation and therapeutic assessment in a group of patients in which correlations of conventional grading with outcome are modest (Olar et al., 2015; Reuss et al., 2015). The finding that G-CIMP-high tumors can emerge as G-CIMP-low glioma at recurrence identifies variations in DNA methylation as crucial determinants for glioma progression and provides a clue to the mechanisms driving evolution of glioma. Our results unify previous observations that linked the cell cycle pathway to malignant progression of low-grade glioma (Mazor et al., 2015). Future updates of the TCGA glioma clinical annotation and independent validation of our findings may be able to consider additionally important clinical confounders such as extent of resection and performance status to further optimize the weights of the currently known prognostic variables and their association to the molecular subtypes we identified.

Analysis of IDH-wild-type glioma revealed the PA-like LGG subset that harbors a silent genomic landscape, confers favorable prognosis relative to other IDH-wild-type diffuse glioma, and displays a molecular profile with high similarity to PA. Re-view by neuropathologists and neuroradiologists confirmed that the majority were correctly diagnosed as diffuse glioma, emphasizing the need for integration of molecular signatures into clinical classification (Chi et al., 2013) for this subgroup of patients that may be spared potentially unnecessary intensive treatments.

The large number of exomes in our dataset allowed identification of novel glioma-associated somatic alterations, including in the *KRAS* and *NRAS* genes, which were frequently used in genetically engineered glioma mouse models (Holland et al., 2000). Our analysis further nominates glial tumors to join an increasing number of tumor types characterized by a deactivated cohesin pathway (Kon et al., 2013; Solomon et al., 2011). Cohesin mutant tumors may infer increased sensitivity to DNA damage agents and PARP inhibitors (Bailey et al., 2014), suggesting that gliomas with genetic alterations of key cohesin regulatory factors may represent biomarkers and therapeutic opportunities.

Overexpression of *TERT* mRNA was found to be associated with increased telomere length in urothelial cancer (Borah et al., 2015). Our results revealed that, in gliomas, increased telomere length is associated with *ATRX* mutations, suggesting an alternative lengthening of telomeres (ALT) mechanism. ALT has been associated with sensitivity to inhibition of the protein kinase ATR (Flynn et al., 2015).

In summary, our pan-glioma analysis has expanded our knowledge of the glioma somatic alteration landscape, emphasized the relevance of DNA methylation profiles as a modality for clinical classification, and quantitatively linked somatic *TERT* pathway alterations to telomere maintenance. Combined, these findings are an important step forward in our understanding of glioma as discrete disease subsets and the mechanisms driving gliomagenesis.



(legend on next page)

## EXPERIMENTAL PROCEDURES

### Patient and Sample Characteristics

Specimens were obtained from patients with appropriate consent from institutional review boards. Details of sample preparation are described in the [Supplemental Experimental Procedures](#).

### Data Generation

In total, tumors from 1,132 patients were assayed on at least one molecular profiling platform, which platforms included: (1) whole-genome sequencing, including high coverage and low pass whole-genome sequencing; (2) exome sequencing; (3) RNA sequencing; (4) DNA copy-number and single-nucleotide polymorphism arrays, including Agilent CGH 244K, Affymetrix SNP6.0, and Illumina 550K Infinium HumanHap550 SNP Chip microarrays; (5) gene expression arrays, including Agilent 244K Custom Gene Expression, Affymetrix HT-HGU133A and Affymetrix Human Exon 1.0 ST arrays; (6) DNA methylation arrays, including Illumina GoldenGate Methylation, Illumina Infinium HumanMethylation27, and Illumina Infinium HumanMethylation450 Bead-Chips; (7) reverse phase protein arrays; (8) miRNA sequencing; and (9) miRNA Agilent 8 × 15K Human miRNA-specific microarrays. Details of data generation have been previously reported ([Brennan et al., 2013](#); [Cancer Genome Atlas Research Network et al., 2015](#)). To ensure cross-platform comparability, features from all array platforms were compared to a reference genome.

### Data Analysis

The data and analysis results can be explored through the Broad Institute FireBrowse portal (<http://firebrowse.org/?cohort=GBMLGG>), the cBioPortal for Cancer Genomics ([http://www.cbioportal.org/study.do?cancer\\_study\\_id=lgggbm\\_tcga\\_pub](http://www.cbioportal.org/study.do?cancer_study_id=lgggbm_tcga_pub)), in a Tumor Map (<http://tumormap.ucsc.edu/?p=newton.gliomas-paper>), the TCGA transcript fusion portal (<http://www.tumorfusions.org>), TCGA Batch Effects (<http://bioinformatics.mdanderson.org/tcgambatch/>), Regulome Explorer (<http://explorer.cancerregulome.org/>), Next-Generation Clustered Heat Maps (<http://bioinformatics.mdanderson.org/TCGA/NGCHMPortal/>). See also [Supplemental Information](#) and the TCGA publication page ([https://tcga-data.nci.nih.gov/docs/publications/lgggbm\\_2015/](https://tcga-data.nci.nih.gov/docs/publications/lgggbm_2015/)).

## SUPPLEMENTAL INFORMATION

Supplemental Information includes Supplemental Experimental Procedures, five figures, and five tables and can be found with this article online at <http://dx.doi.org/10.1016/j.cell.2015.12.028>.

## AUTHOR CONTRIBUTIONS

Conceptualization and project administration: R.G.W.V., A.I., and H.N.; supervision: S.R.S., K.D.A., P.W.L., M.G., D.H., D.J.B., D.H.G., R.R., C.C.L., J.S.B.-S., C.G.C., D.P.C.T., W.K.A.Y., J.H., L.C., M.M., and T.M.; formal analysis: R.G.W.V., A.I., H.N., M.C., F.P.B., T.M.M., T.S.S., O.M., Y.N., S.M.P., P.Z., L.P., A. Radenbaugh, G.R., R.A., J.W., G.M., S.L., S.A., A. Rao, B.A.M., A.D.C., and H.Z.; investigation: D.J.B., L.C., and L.P.; data curation: D.J.B., L.P., and F.P.B.; writing - original draft: R.G.W.V., A.I., H.N., M.C., F.P.B., T.M.M., and T.S.S.; manuscript review: D.J.B., K.A.D., S.R.S., M.W., N.L., and D.H.G.

## ACKNOWLEDGMENTS

This study was supported by NIH grants U24CA143883, U24CA143858, U24CA143840, U24CA143799, U24CA143835, U24CA143845,

U24CA143882, U24CA143867, U24CA143866, U24CA143848, U24CA144025, U54HG003067, U54HG003079, U54HG003273, U24CA126543, U24CA126544, U24CA126546, U24CA126551, U24CA126554, U24CA126561, U24CA126563, U24CA143731, U24CA143843, P30CA016672, P50 CA127001, U54CA193313, R01CA179044, R01CA185486, R01 CA190121, and P01 CA085878; Cancer Prevention & Research Institute of Texas (CPRI) R140606; and São Paulo Research Foundation (FAPESP) 2014/02245-3, 2015/07925-5, 2015/02844-7, and 2015/08321-3. D.J.W. is a consultant for Zymo Research Corporation. R.B. is a consultant for and received grant funding from Novartis. A.D.C. and M.M. received grant support from Bayer.

Received: July 17, 2015

Revised: October 20, 2015

Accepted: December 11, 2015

Published: January 28, 2016

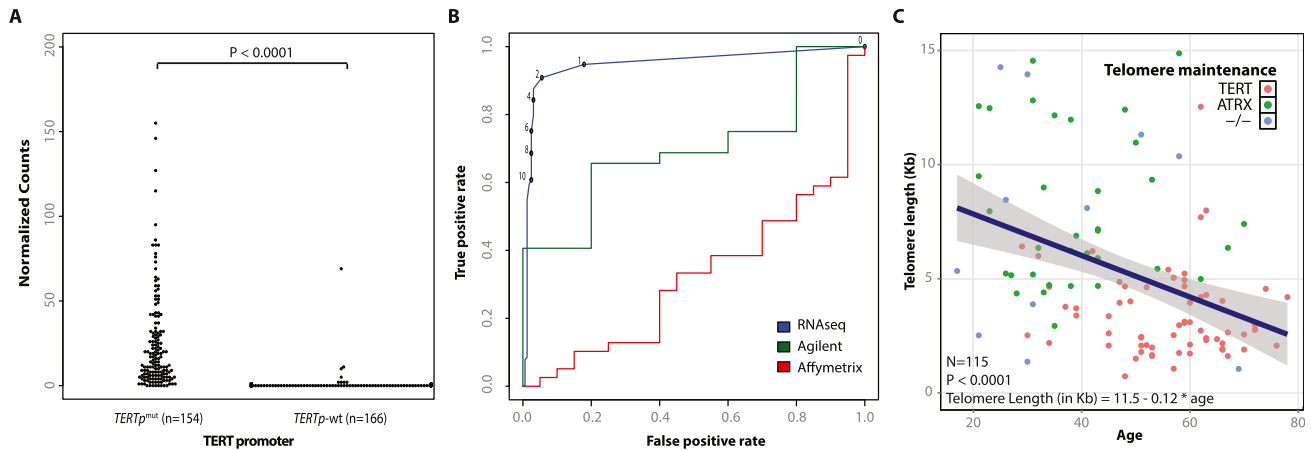
## REFERENCES

- Bailey, M.L., O'Neil, N.J., van Pel, D.M., Solomon, D.A., Waldman, T., and Hieter, P. (2014). Glioblastoma cells containing mutations in the cohesin component STAG2 are sensitive to PARP inhibition. *Mol. Cancer Ther.* *13*, 724–732.
- Berman, B.P., Weisenberger, D.J., Aman, J.F., Hinoue, T., Ramjan, Z., Liu, Y., Noushmehr, H., Lange, C.P., van Dijk, C.M., Tollenaar, R.A., et al. (2012). Regions of focal DNA hypermethylation and long-range hypomethylation in colorectal cancer coincide with nuclear lamina-associated domains. *Nat. Genet.* *44*, 40–46.
- Borah, S., Xi, L., Zaugg, A.J., Powell, N.M., Dancik, G.M., Cohen, S.B., Costello, J.C., Theodorescu, D., and Cech, T.R. (2015). Cancer. TERT promoter mutations and telomerase reactivation in urothelial cancer. *Science* *347*, 1006–1010.
- Brennan, C.W., Verhaak, R.G., McKenna, A., Campos, B., Noushmehr, H., Salama, S.R., Zheng, S., Chakravarty, D., Sanborn, J.Z., Berman, S.H., et al.; TCGA Research Network (2013). The somatic genomic landscape of glioblastoma. *Cell* *155*, 462–477.
- Cancer Genome Atlas Research Network, Brat, D.J., Verhaak, R.G., Aldape, K.D., Yung, W.K., Salama, S.R., Cooper, L.A., Rheinbay, E., Miller, C.R., Vitucci, M., et al. (2015). Comprehensive, integrative genomic analysis of diffuse lower-grade gliomas. *N. Engl. J. Med.* *372*, 2481–2498.
- Chi, A.S., Batchelor, T.T., Yang, D., Dias-Santagata, D., Borger, D.R., Ellisen, L.W., Iafrate, A.J., and Louis, D.N. (2013). BRAF V600E mutation identifies a subset of low-grade diffusely infiltrating gliomas in adults. *J. Clin. Oncol.* *31*, e233–e236.
- Eckel-Passow, J.E., Lachance, D.H., Molinaro, A.M., Walsh, K.M., Decker, P.A., Sicotte, H., Pekmezci, M., Rice, T., Kosel, M.L., Smirnov, I.V., et al. (2015). Glioma Groups Based on 1p/19q, IDH, and TERT Promoter Mutations in Tumors. *N. Engl. J. Med.* *372*, 2499–2508.
- Flynn, R.L., Cox, K.E., Jeitany, M., Wakimoto, H., Bryll, A.R., Ganem, N.J., Bersani, F., Pineda, J.R., Suvà, M.L., Benes, C.H., et al. (2015). Alternative lengthening of telomeres renders cancer cells hypersensitive to ATR inhibitors. *Science* *347*, 273–277.
- Frattoni, V., Trifonov, V., Chan, J.M., Castano, A., Lia, M., Abate, F., Keir, S.T., Ji, A.X., Zoppi, P., Niola, F., et al. (2013). The integrated landscape of driver genomic alterations in glioblastoma. *Nat. Genet.* *45*, 1141–1149.
- Graham, V., Khudyakov, J., Ellis, P., and Pevny, L. (2003). SOX2 functions to maintain neural progenitor identity. *Neuron* *39*, 749–765.

## Figure 5. Overview of Major Subtypes of Adult Diffuse Glioma

Integrative analysis of 1,122 adult gliomas resulted in 7 different subtypes with distinct biological and clinical characteristics. The groups extend across six DNA methylation subtypes of which the LGM6 cluster was further separated by tumor grade into PA-like and LGM6-GBM. The size of the circles is proportional to the percentages of samples within each group. DNA methylation plot is a cartoon representation of overall genome-wide epigenetic pattern within glioma subtypes. Survival information is represented as a set of Kaplan-Meier curves, counts of grade, histology and LGG/GBM subtypes within the groups are represented as bar-plots, whereas age is represented as density. Labeling of telomere length and maintenance status is based on the enrichment of samples within each column, similarly for the biomarkers and the validation datasets.

- Guintivano, J., Aryee, M.J., and Kaminsky, Z.A. (2013). A cell epigenotype specific model for the correction of brain cellular heterogeneity bias and its application to age, brain region and major depression. *Epigenetics* 8, 290–302.
- Holland, E.C., Celestino, J., Dai, C., Schaefer, L., Sawaya, R.E., and Fuller, G.N. (2000). Combined activation of Ras and Akt in neural progenitors induces glioblastoma formation in mice. *Nat. Genet.* 25, 55–57.
- Horn, S., Figl, A., Rachakonda, P.S., Fischer, C., Sucker, A., Gast, A., Kadel, S., Moll, I., Nagore, E., Hemminki, K., et al. (2013). TERT promoter mutations in familial and sporadic melanoma. *Science* 339, 959–961.
- Jones, D.T., Hutter, B., Jäger, N., Korshunov, A., Kool, M., Warnatz, H.J., Zichner, T., Lambert, S.R., Ryzhova, M., Quang, D.A., et al.; International Cancer Genome Consortium PedBrain Tumor Project (2013). Recurrent somatic alterations of FGFR1 and NTRK2 in pilocytic astrocytoma. *Nat. Genet.* 45, 927–932.
- Killela, P.J., Reitman, Z.J., Jiao, Y., Bettegowda, C., Agrawal, N., Diaz, L.A., Jr., Friedman, A.H., Friedman, H., Gallia, G.L., Giovanella, B.C., et al. (2013). TERT promoter mutations occur frequently in gliomas and a subset of tumors derived from cells with low rates of self-renewal. *Proc. Natl. Acad. Sci. USA* 110, 6021–6026.
- Kim, H., Zheng, S., Amini, S.S., Virk, S.M., Mikkelsen, T., Brat, D.J., Grimsby, J., Sougnez, C., Muller, F., Hu, J., et al. (2015). Whole-genome and multisector exome sequencing of primary and post-treatment glioblastoma reveals patterns of tumor evolution. *Genome Res.* 25, 316–327.
- Kon, A., Shih, L.Y., Minamino, M., Sanada, M., Shiraishi, Y., Nagata, Y., Yoshida, K., Okuno, Y., Bando, M., Nakato, R., et al. (2013). Recurrent mutations in multiple components of the cohesin complex in myeloid neoplasms. *Nat. Genet.* 45, 1232–1237.
- Kurscheid, S., Bady, P., Sciuscio, D., Samarzija, I., Shay, T., Vassallo, I., Crieckinge, W.V., Daniel, R.T., van den Bent, M.J., Marosi, C., et al. (2015). Chromosome 7 gain and DNA hypermethylation at the HOXA10 locus are associated with expression of a stem cell related HOX-signature in glioblastoma. *Genome Biol.* 16, 16.
- Lambert, S.R., Witt, H., Hovestadt, V., Zucknick, M., Kool, M., Pearson, D.M., Korshunov, A., Ryzhova, M., Ichimura, K., Jabado, N., et al. (2013). Differential expression and methylation of brain developmental genes define location-specific subsets of pilocytic astrocytoma. *Acta Neuropathol.* 126, 291–301.
- Lodato, M.A., Ng, C.W., Wamstad, J.A., Cheng, A.W., Thai, K.K., Fraenkel, E., Jaenisch, R., and Boyer, L.A. (2013). SOX2 co-occupies distal enhancer elements with distinct POU factors in ESCs and NPCs to specify cell state. *PLoS Genet.* 9, e1003288.
- Louis, D.N., Ohgaki, H., Wiestler, O.D., Cavenee, W.K., Burger, P.C., Jouvet, A., Scheithauer, B.W., and Kleihues, P. (2007). The 2007 WHO classification of tumours of the central nervous system. *Acta Neuropathol.* 114, 97–109.
- Mazor, T., Pankov, A., Johnson, B.E., Hong, C., Hamilton, E.G., Bell, R.J., Smirnov, I.V., Reis, G.F., Phillips, J.J., Barnes, M.J., et al. (2015). DNA Methylation and Somatic Mutations Converge on the Cell Cycle and Define Similar Evolutionary Histories in Brain Tumors. *Cancer Cell* 28, 307–317.
- Mur, P., Mollejo, M., Ruano, Y., de Lope, A.R., Fiaño, C., García, J.F., Castresana, J.S., Hernández-Lain, A., Rey, J.A., and Meléndez, B. (2013). Codeletion of 1p and 19q determines distinct gene methylation and expression profiles in IDH-mutated oligodendroglial tumors. *Acta Neuropathol.* 126, 277–289.
- Noushmehr, H., Weisenberger, D.J., Diefes, K., Phillips, H.S., Pujara, K., Berman, B.P., Pan, F., Pelloski, C.E., Sulman, E.P., Bhat, K.P., et al.; Cancer Genome Atlas Research Network (2010). Identification of a CpG island methylator phenotype that defines a distinct subgroup of glioma. *Cancer Cell* 17, 510–522.
- Olar, A., Wani, K.M., Alfaro-Munoz, K.D., Heathcock, L.E., van Thuijl, H.F., Gilbert, M.R., Armstrong, T.S., Sulman, E.P., Cahill, D.P., Vera-Bolanos, E., et al. (2015). IDH mutation status and role of WHO grade and mitotic index in overall survival in grade II-III diffuse gliomas. *Acta Neuropathol.* 129, 585–596.
- Ozawa, T., Riester, M., Cheng, Y.K., Huse, J.T., Squatrito, M., Helmy, K., Charles, N., Michor, F., and Holland, E.C. (2014). Most human non-GCIMP glioblastoma subtypes evolve from a common proneural-like precursor glioma. *Cancer Cell* 26, 288–300.
- Peters, J.M., and Nishiyama, T. (2012). Sister chromatid cohesion. *Cold Spring Harb. Perspect. Biol.* 4, a011130.
- Pineda, C.T., Ramanathan, S., Fon Tacer, K., Weon, J.L., Potts, M.B., Ou, Y.H., White, M.A., and Potts, P.R. (2015). Degradation of AMPK by a cancer-specific ubiquitin ligase. *Cell* 160, 715–728.
- Reuss, D.E., Mamatjan, Y., Schrimpf, D., Capper, D., Hovestadt, V., Kratz, A., Sahm, F., Koelsche, C., Korshunov, A., Olar, A., et al. (2015). IDH mutant diffuse and anaplastic astrocytomas have similar age at presentation and little difference in survival: a grading problem for WHO. *Acta Neuropathol.* 129, 867–873.
- Roggendorf, W., Strupp, S., and Paulus, W. (1996). Distribution and characterization of microglia/macrophages in human brain tumors. *Acta Neuropathol.* 92, 288–293.
- Schwartzbaum, J.A., Fisher, J.L., Aldape, K.D., and Wrensch, M. (2006). Epidemiology and molecular pathology of glioma. *Nat. Clin. Pract. Neurol.* 2, 494–503.
- Solomon, D.A., Kim, T., Diaz-Martinez, L.A., Fair, J., Elkahoul, A.G., Harris, B.T., Toretsky, J.A., Rosenberg, S.A., Shukla, N., Ladanyi, M., et al. (2011). Mutational inactivation of STAG2 causes aneuploidy in human cancer. *Science* 333, 1039–1043.
- Sturm, D., Witt, H., Hovestadt, V., Khuong-Quang, D.A., Jones, D.T., Konermann, C., Pfaff, E., Tönjes, M., Sill, M., Bender, S., et al. (2012). Hotspot mutations in H3F3A and IDH1 define distinct epigenetic and biological subgroups of glioblastoma. *Cancer Cell* 22, 425–437.
- Sturm, D., Bender, S., Jones, D.T., Lichter, P., Grill, J., Becher, O., Hawkins, C., Majewski, J., Jones, C., Costello, J.F., et al. (2014). Paediatric and adult glioblastoma: multifactorial (epi)genomic culprits emerge. *Nat. Rev. Cancer* 14, 92–107.
- Suvà, M.L., Rheinbay, E., Gillespie, S.M., Patel, A.P., Wakimoto, H., Rabkin, S.D., Riggi, N., Chi, A.S., Cahill, D.P., Nahed, B.V., et al. (2014). Reconstructing and reprogramming the tumor-propagating potential of glioblastoma stem-like cells. *Cell* 157, 580–594.
- Suzuki, H., Aoki, K., Chiba, K., Sato, Y., Shiozawa, Y., Shiraishi, Y., Shimamura, T., Niida, A., Motomura, K., Ohka, F., et al. (2015). Mutational landscape and clonal architecture in grade II and III gliomas. *Nat. Genet.* 47, 458–468.
- Turcan, S., Rohle, D., Goenka, A., Walsh, L.A., Fang, F., Yilmaz, E., Campos, C., Fabius, A.W., Lu, C., Ward, P.S., et al. (2012). IDH1 mutation is sufficient to establish the glioma hypermethylator phenotype. *Nature* 483, 479–483.
- van den Bent, M.J. (2010). Interobserver variation of the histopathological diagnosis in clinical trials on glioma: a clinician's perspective. *Acta Neuropathol.* 120, 297–304.
- Verhaak, R.G., Hoadley, K.A., Purdom, E., Wang, V., Qi, Y., Wilkerson, M.D., Miller, C.R., Ding, L., Golub, T., Mesirov, J.P., et al.; Cancer Genome Atlas Research Network (2010). Integrated genomic analysis identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in PDGFRA, IDH1, EGFR, and NF1. *Cancer Cell* 17, 98–110.
- Yan, H., Parsons, D.W., Jin, G., McLendon, R., Rasheed, B.A., Yuan, W., Kos, I., Batonic-Haberle, I., Jones, S., Riggins, G.J., et al. (2009). IDH1 and IDH2 mutations in gliomas. *N. Engl. J. Med.* 360, 765–773.
- Yoshihara, K., Shahmoradgoli, M., Martínez, E., Vegesna, R., Kim, H., Torres-García, W., Treviño, V., Shen, H., Laird, P.W., Levine, D.A., et al. (2013). Inferring tumour purity and stromal and immune cell admixture from expression data. *Nat. Commun.* 4, 2612.



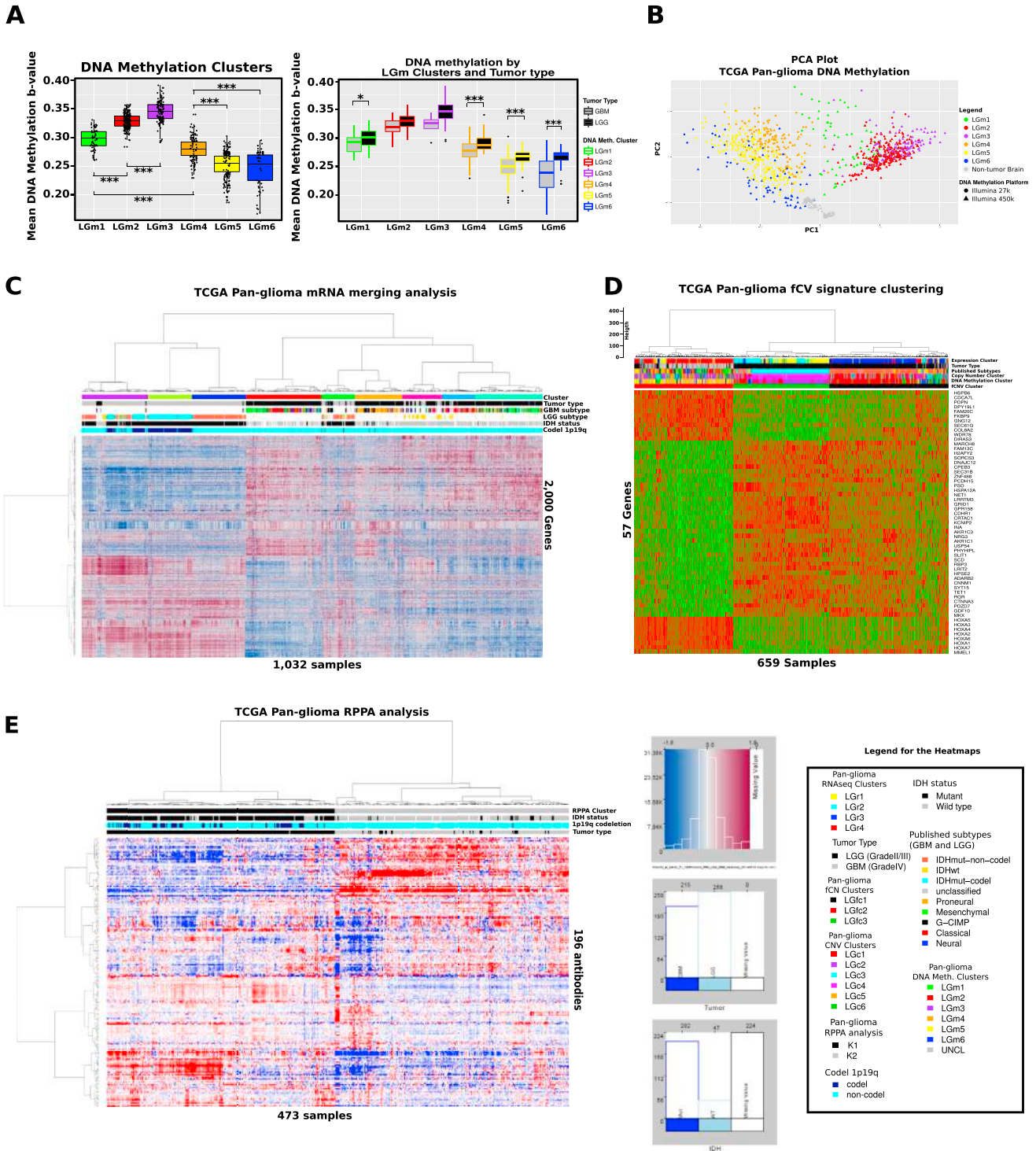
**Figure S1. Telomere Length Quantification of 120 Gliomas with Known *TERT* Promoter and *ATRX* Mutational Status, Related to Figure 1**

(A) RNaseq *TERT* expression is upregulated in *TERTp* mutant cases, but not in *ATRX* and double negative cases ( $p < 0.0001$ ).

(B) *TERT* expression as quantified by RNA sequencing is a highly sensitive and specific marker for the absence or presence of the *TERTp* mutation (AUC 0.95). Using a cutoff value of 2, sensitivity and specificity are 91% and 95%, respectively. Interestingly, microarray data is poor substitute for the *TERT* promoter mutation with an AUC of 0.70 and 0.32 for the Agilent and Affymetrix microarray respectively.

(C) Telomeres gradually shorten with increasing age in tumor samples ( $p < 0.0001$ ). Note the steeper decline relative to Figure 2B and that *ATRX* mutant patients are in the younger age range whereas *TERTp* mutant patients are in the younger age range. This suggests an independent contribution of telomere maintenance to telomere length.





**Figure S2. Pan-glioma DNA Methylation and Transcriptome Subtypes, Related to Figure 2**

(A) Boxplot of the mean DNA methylation beta-values genome-wide (20,036 probes) for each sample distributed by the six Pan-glioma DNA methylation clusters (left) and divided by tumor type (right). Significant differences are highlighted with \* (p-value < 0.01) and \*\*\* (p-value < 1e-04).

(B) Principal component analysis of 932 TCGA glioma samples and 77 non-tumor brain samples (Guintivano et al., 2013) performed on 19,520 CpG probes (genome-wide).

(C) LGG-GBM mRNA merging analysis. Clustered heatmap of merged data with 569 GBM and 463 LGG non-duplicate samples, and 2000 most variable genes. Consensus clustering revealed 9 clusters. The 3 left-most clusters show predominantly LGG samples, 3 clusters show predominantly GBM samples, whereas 3

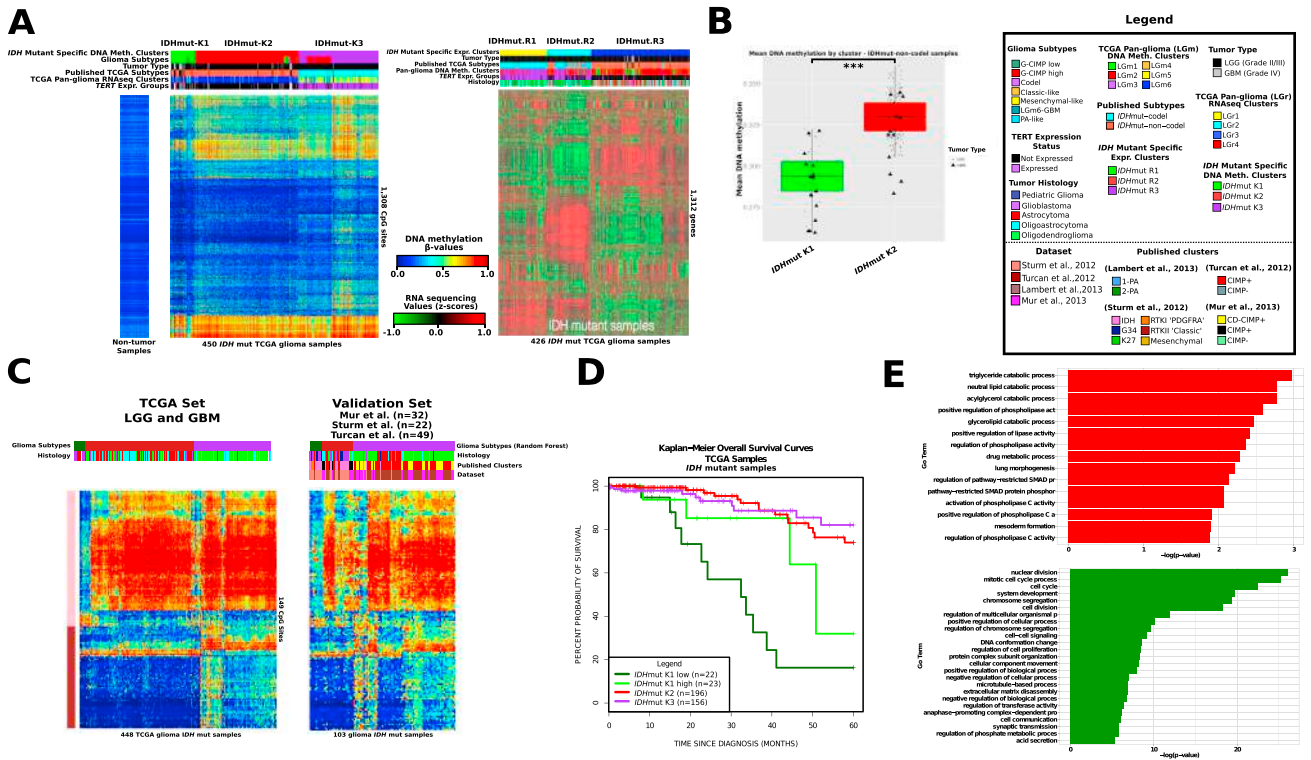
(legend continued on next page)

---

clusters show a mixture of GBM and LGG samples. The LGG IDH-wild-type samples clustered mostly with the GBM classical samples, whereas many of the LGG IDH mutant-non-codel samples cluster with the GBM G-CIMP samples.

(D) Functional Copy Number (fCN) gene signature Heatmap. Genes with Spearman's correlation between CN and Expression above 0.5,  $\text{abs}(FC > 1.5)$  and  $\text{abs}(\Delta CN > 0.5)$  define the fCN signature. The Heatmap illustrate the samples unsupervised clustering given the fCN signature. RNA expression levels range from green (low) to red (high). Each row reports the annotation of a different analysis performed in the paper. Last row reports the fCN annotation.

(E) Clustered heatmap of unsupervised hierarchical clustering of 473 samples (columns) and 196 antibodies (rows). The annotation bars (shown on top) were not used for clustering. The legend for the annotation bars is shown on the left. Two clusters can be found that largely correspond to tumor type. The left cluster has largely LGG samples and one GBM sample. However, the right cluster has mostly GBM samples but 26 LGG samples, 17 of which have no mutations in *IDH1/2*. In the heatmap, low, medium, and high expression is represented by blue, white, and red colors, respectively.



**Figure S3. Identification of a Distinct *IDH* mutant Subtype Defined by Epigenomics, Related to Figure 3**

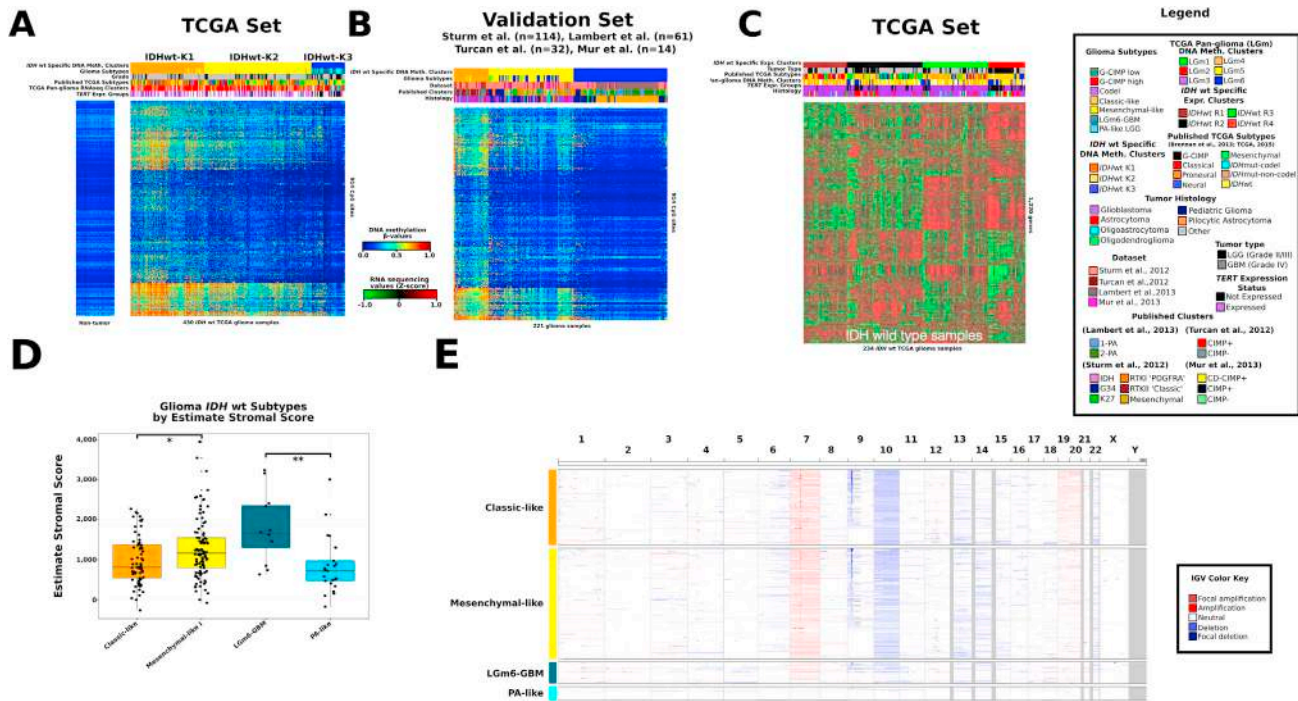
(A) Left, Heatmap of DNA methylation data. Unsupervised consensus clustering analysis using 1,308 CpG probe specific CpG probes defined among the TCGA *IDH* mutant gliomas. Column-wise represents 450 *IDH* mutant glioma samples, row-wise represents probes. Samples are ordered according to the consensus cluster output, and rows are ordered by hierarchical clustering. DNA methylation beta-values ranges from 0 (low) to 1 (high). Three clusters were defined, each cluster separated and labeled. Non-tumor brain samples are represented on the left of the heatmap (Guintivano et al., 2013). Additional tracks are included at the top of the heatmaps to identify each sample membership within separate cluster analysis (Glioma subtypes, tumor type, previous published subtypes (Brennan et al. Cell, 2013, TCGA Research Network, NEJM, 2015), RNA sequencing and *TERT* expression). Legend is provided for the heatmap. Right, Clustering of *IDH* mutant samples transcriptional profiles. Unsupervised clustering of gene expression separated by *IDH* status 426 samples confirming the presence of three main groups resembling the clusters reported in (TCGA Network, New Eng J Med 2015) where all GBM G-CIMP cluster together with the LGG *IDH* mutant-non-codel.

(B) Boxplot of the average DNA methylation beta-value genome-wide (20,000 probes) for each sample grouped by *IDHmut* K1 and *IDHmut* K2. Dots represent LGG tumors and triangles represent GBM tumors. Significant difference is highlighted with \*\*\* ( $p$ -value  $< 2.2 \times 10^{-16}$ )

(C) Left, Heatmap of DNA methylation data. Supervised statistical analysis using 149 CpG tumor specific CpG probes that define each TCGA *IDH* mutant glioma subtype. Column-wise represents 448 *IDH* mutant (codels and non codels) TCGA glioma samples, row-wise represents probes. DNA methylation beta-values ranges from 0 (low) to 1 (high). Right, Heatmap of DNA methylation data for the validation dataset (Sturm et al., 2012; Turcan et al., 2012; Mur et al., 2013), using the 149 CpG tumor specific probes that define each TCGA *IDH* mutant glioma subtype. Non-TCGA glioma samples were classified into one of the three *IDH* mutant type specific clusters using the random forest machine learning method. DNA methylation beta-values ranges from 0 (low) to 1 (high). Additional tracks are included at the top of the heatmap to identify tumor histology, published clusters (Published Clusters) and each sample membership according to its dataset (Study). Legend is provided for the heatmap.

(D) Kaplan-Meier survival curves showing samples separated by *IDHmut* K1 low, *IDHmut* K1 high, *IDHmut* K2 and *IDHmut* K3. Tick represent censorship.

(E) Pathway analysis of differentially expressed genes between *IDHmut* K1, *IDHmut* K2, ranked by  $p$ -value. The top red panel shows categories enriched in *IDHmut*K2; the bottom green panel displays categories enriched in *IDHmut*K1.



**Figure S4. Identification of a Distinct Subgroup of IDH-Wild-Type Gliomas with Pilocytic Astrocytoma Features, Related to Figure 4**

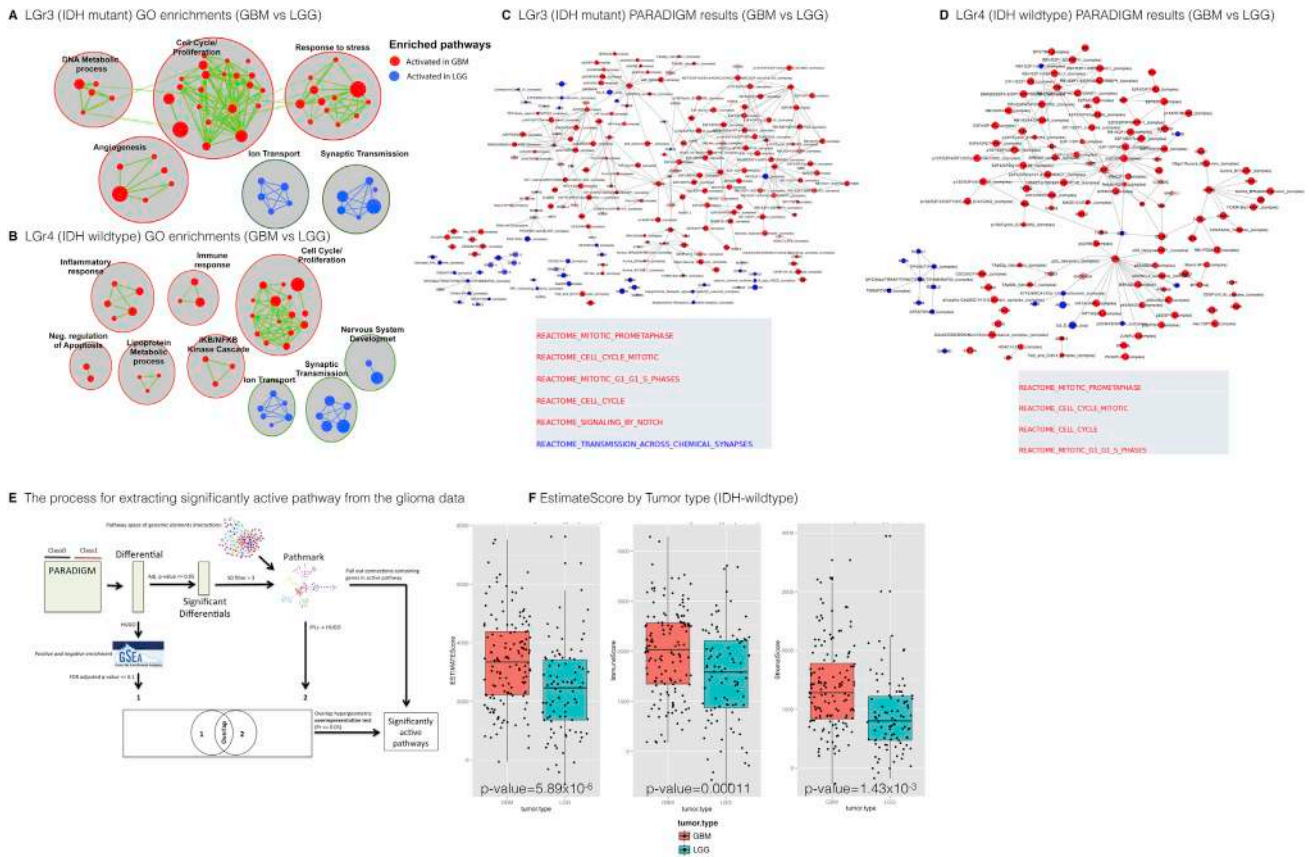
(A) Heatmap of DNA methylation data. Unsupervised consensus clustering analysis using 914 CpG tumor specific probes defined among the TCGA IDH-wild-type gliomas. Column-wise represents 430 IDH-wild-type TCGA glioma samples, row-wise represents probes. Samples are ordered according to the consensus cluster output, and rows are ordered by hierarchical clustering. DNA methylation beta-values ranges from 0 (low) to 1 (high). Three clusters were defined, each cluster separated and labeled. Non-tumor brain samples are represented on the left of the heatmap (Guintivano et al., 2013). Additional tracks are included at the top of the heatmaps to identify each sample membership within separate cluster analysis (Glioma subtypes, tumor type, previous published subtypes Brennan et al. Cell, 2013, TCGA Research Network, NEJM, 2015), RNA sequencing and TERT expression). Legend is provided for the heatmap.

(B) Heatmap of DNA methylation data for the validation dataset (Sturm et al., 2012; Turcan et al., 2012; Lambert et al., 2013; Mur et al., 2013), using the 914 CpG tumor specific probes defined in panel S4A. Non-TCGA glioma samples were classified into one of the three IDH-wild-type specific clusters using the random forest machine learning method. The second track from top to bottom shows the classification of non-TCGA glioma samples into one of the seven glioma subtypes also using the random forest machine learning method. DNA methylation beta-values ranges from 0 (low) to 1 (high). Additional tracks are included at the top of the heatmap to identify each sample membership according to its dataset (Dataset), to previous published clusters (Published Clusters) and to tumor histology. Legend is provided for the heatmap.

(C) Clustering of IDH-wild-type samples transcriptional profiles. Unsupervised clustering of gene expression separated by IDH status showed that the LGr4 cluster identified in the pan-glioma unsupervised analysis splits into four mixed LGG/GBM clusters (234 samples), where the first two, although separated by a relatively small number of genes, are respectively enriched with Classical subtype (59%) and LGm4 samples and the second with Mesenchymal (75%) subtype and LGm5 samples, the third enriched with Proneural subtype (85%) and a fourth mostly containing LGG IDH-wild-type samples.

(D) Boxplot of the estimate stromal score for each sample distributed by the four glioma IDH wild-type subtypes. Significant differences are highlighted with \* (p-value < 0.05) and \*\* (p-value < 0.005).

(E) IGV screenshot demonstrating differences in copy number landscape across glioma subtypes.



**Figure S5. Progression of LGG to GBM Is Marked by Cell Cycle/Proliferation or Invasion/Microenvironmental Changes, Related to the Transcriptome Clusters Shown in Figure 2**

The pathways involved with progression from LGG to GBM were identified through supervised analysis of co-clustered LGG and GBM using Gene Set Enrichment Analysis. Gene sets were compiled from the Gene Ontology pathway database. Significantly enriched gene sets (FDR < 0.1, p-value < 0.005) were depicted as an annotation module network using Cytoscape and EnrichmentMap. Nodes represent enriched gene sets, which are grouped and annotated by their similarity. Node size is proportional to the total number of genes within each gene set. Proportion of shared genes between gene sets is represented as the thickness of the line between nodes.

(A) Progression of LGG IDH mutant-non-codel to GBM G-CIMP in LGr3 was strongly marked by a hyper-proliferation signature and revealed four major gene sets groups related to cell cycle and hyperproliferation, DNA metabolic processes, response to stress and angiogenesis.

(B) Similar analysis of the gene sets activated in the GBM compared to the LGG component of LGr4 (IDH-wild-type) identified an inflammation and immunologic response signature characterized by the activation of several chemokines and interleukins enriching sets involved in inflammatory and immuno response, negative regulation of apoptosis, cell cycle and proliferation, IKB/NFKB kinase cascade.

(C) Differential regulatory networks describing differential molecular activities between GBM and LGG in LGr3. Dichotomies were selected by only choosing those where samples form tight linearly separable clusters in the high dimensional genomic space. The size of the node is inversely proportional to the magnitude of the p-value computed by LIMMA for each differential. Curated canonical MSigDB pathways significantly represented in these networks are listed below each network, following the same color scheme as described above.

(D) Same as in C. for LGr4

(E) Overview of the adopted pipeline for extracting significant pathways.

(F) Distribution of Estimate, Immuno and Stromal score by tumor type in the IDH-wild-type samples.

Cell

Supplemental Information

## **Molecular Profiling Reveals Biologically Discrete Subsets and Pathways of Progression in Diffuse Glioma**

**Michele Ceccarelli, Floris P. Barthel, Tathiane M. Malta, Thais S. Sabedot, Sofie R. Salama, Bradley A. Murray, Olena Morozova, Yulia Newton, Amie Radenbaugh, Stefano M. Pagnotta, Samreen Anjum, Jiguang Wang, Ganiraju Manyam, Pietro Zoppoli, Shiyung Ling, Arjun A. Rao, Mia Grifford, Andrew D. Cherniack, Hailei Zhang, Laila Poisson, Carlos Gilberto Carlotti, Jr., Daniela Pretti da Cunha Tirapelli, Arvind Rao, Tom Mikkelsen, Ching C. Lau, W.K. Alfred Yung, Raul Rabadan, Jason Huse, Daniel J. Brat, Norman L. Lehman, Jill S. Barnholtz-Sloan, Siyuan Zheng, Kenneth Hess, Ganesh Rao, Matthew Meyerson, Rameen Beroukhim, Lee Cooper, Rehan Akbani, Margaret Wrensch, David Haussler, Kenneth D. Aldape, Peter W. Laird, David H. Gutmann, TCGA Research Network, Houtan Noushmehr, Antonio Iavarone, and Roel G.W. Verhaak**

# Molecular profiling refines the classification of adult diffuse lower- and high-grade glioma

## Supplemental Information

Supplemental Information content:

### Supplemental Experimental Procedures

<b>1. Biospecimens</b> .....	<b>3</b>
<b>2. DNA sequencing</b> .....	<b>4</b>
2.1 DNA sequencing data production .....	4
2.2 Identification of somatic mutations .....	4
2.4 Identification of TERT promoter mutations.....	6
2.5 Mutation significant analysis.....	6
2.6 Telomere quantification.....	6
2.7 Whole genome mutation calling.....	7
<b>3. DNA copy number analysis</b> .....	<b>7</b>
3.1 Preprocessing and peak calling.....	7
3.2 Functional Copy Number (CN) analysis .....	8
3.3 Mutations with Common Focal Alterations (MutComFocal) .....	8
<b>4. mRNA Expression</b> .....	<b>9</b>
4.1 Data preparation and gene selection .....	9
4.2 Classification of Affymetrix samples .....	10
4.3 Tumor Map and Pathway Activity Analysis .....	10
4.3.1 Combining multi-platform multi-tumor datasets .....	10
4.3.2 Tumor Map method (manuscript in preparation) .....	10
4.3.3 Multi-platform maps using Bivariate Standardization similarity space Transformation (BST).....	11
4.3.4 Extracting significantly active pathways .....	12
4.4 Combining GBM Agilent G4502A mRNA data with LGG Illumina Hi-Seq RNA-seq data.....	13
4.5 RNA Fusion analysis .....	13
4.5.1 Fusion transcript detection using PRADA.....	13
4.5.2 Fusion transcript detection using deFuse.....	14
4.6 Identification of Transcriptional Regulatory Factors underlying IDH wild type and IDH mutant phenotypes in Glioma .....	14
<b>5. DNA methylation profiling</b> .....	<b>15</b>
5.1 Preprocessing and clustering.....	15
5.2 Unsupervised clustering analysis of DNA methylation data .....	16
5.3 Supervised analysis of DNA methylation.....	17

5.4 Identification of Epigenetically Regulated Genes .....	18
5.5 Classification of new glioma samples based on DNA methylation glioma subtypes .....	20
5.6 Patient centric table (DNA methylation).....	20





5.7 Homer de novo motif searches .....	22
<b>6. Reverse phase protein array (RPPA).....</b>	<b>22</b>
6.1 Data Processing .....	22
6.2 Data normalization.....	23
6.3 Clustering .....	24
<b>7. Regulome Explorer .....</b>	<b>25</b>
7.1. Feature Matrix .....	25
7.2. All-by-all Pairwise Associations .....	26
<b>8. Supplemental References .....</b>	<b>27</b>



# Supplemental Experimental Procedures

## 1. Biospecimens

**Authors:** Jay Bowen, Kristen M. Leraas, Tara M. Lichtenberg

**Correspondence and questions should be directed to:** Jay Bowen (jay.Bowen@nationwidechildrens.org)

Biospecimens were collected from patients diagnosed with low grade gliomas (LGG) and glioblastoma multiforme (GBM) undergoing surgical resection.

The case list freeze included 1122 cases comprising 516 LGG and 606 GBM. Samples were from the following 32 tissue source sites: Asterand (n=2); Case Western (n=188); Cedars Sinai (n=34); CHI-Penrose Colorado (n=2); Christiana Healthcare (n=12); Cureline (n=26); Dept of Neurosurgery at University of Heidelberg (n=48); Duke University (n=90); Emory University (n=44); Fondazione-Besta (QH) (n=38); Greenville Health System (n=1); Hartford (n=2); Henry Ford Hospital (n=243); Huntsman Cancer Institute (n=8); International Genomics Consortium (n=2); John Wayne Cancer Center (n=2); Johns Hopkins (n=7); Mayo Clinic (n=39); MD Anderson Cancer Center (n=101); Memorial Sloan Kettering Cancer Center (n=15); Northwestern University (n=2); St. Joseph AZ (n=30); Swedish Neurosciences (n=6); The University of New South Wales (n=19); Thomas Jefferson University (n=44); Toronto Western Hospital (n=14); University of California San Francisco (n=50); University of Florida (n=30); University of Kansas (n=1); University of Miami (n=3); University of North Carolina (n=2); University of Sao Paulo (n=17).

Samples were acquired and processed according to previous descriptions (Brennan et al., 2013; TCGA\_Network, 2015).

A detailed list of clinical and molecular data elements is included in Table S1 and reflects the clinical data package frozen on 05/01/2015. Clinical data elements comprise histology, grade, gender, age at diagnosis/surgery, treatments, vital status, overall and progression-free survival. Clinical data available at the BCR was manually curated. Where possible, additional de-identified follow-up data were requested from TSSs through BCR and manually added into the clinical data freeze package.

Overall survival was defined as the time from surgical diagnosis until death. Cases that were still alive at the time of this study have overall survival time censored at the time of last follow-up. Survival curves were estimated and plotted using the Kaplan-Meier method. Log-rank tests were used to compare curves between groups. Single-predictor and multiple-predictor models were fit using Cox regression under the proportional hazards assumption. Hazard ratios and 95% confidence intervals are reported. Nested models were compared using the likelihood ratio test (LRT). Harrell's concordance index (C-index) was used to assess and report model performance



(Harrell et al., 1982). These analyses were conducted in R (v 3.1.2) using the survival package (Therneau, 2014; Therneau and Grambsch, 2000).

## 2. DNA sequencing

**Authors:** Floris Barthel, Bradley Murray, Siyuan Zheng, Roel Verhaak

**Correspondence and questions should be directed to:** Roel Verhaak (rverhaak@mdanderson.org)

### 2.1 DNA sequencing data production

Whole exome, whole genome and targeted validation and TERT promoter sequencing (including low-pass sequencing) was performed as previously described (Brennan et al., 2013; Cancer Genome Atlas Research, 2015; Verhaak et al., 2010).

Platform	Center	Disease	Exome capture kit	Read length	Paired samples
Illumina HiSeq	BI	GBM	Agilent Sure-Select Hun All Exon v2.0, 44Mb kit	2 x 76 bp	307
Illumina HiSeq	BI	LGG	Agilent Sure-Select Hun All Exon v2.0, 44Mb kit	2 x 76 bp	513
<b>Union</b>					<b>820</b>

#### Whole exome sequencing

Platform	Center	Disease	Libraries	Read length	Paired samples
Illumina HiSeq	BI	GBM	2-59	2 x 101 bp	38
Illumina HiSeq	BI	LGG	3-11	2 x 101 bp	20
Illumina HiSeq	WUGSC	GBM	16-167	100 bp	13
Illumina HiSeq	HMS-RK	LGG	1	2 x 51 bp	52
<b>Union</b>					<b>123</b>

#### Whole genome sequencing (including low-pass)

### 2.2 Identification of somatic mutations

The Broad Institute's Firehose cancer genome analysis pipeline used BAM files for tumor and matched normal samples to perform quality control, local realignment coverage calculations and others on whole exome sequencing (Table 1) as described (Imielinski et al., 2012). For the identification of somatic single nucleotide variations we used a multicenter approach integrating the output of three different somatic mutation algorithms: MuTect (Cibulskis et al., 2013), RADIA (Radenbaugh et al., 2014) and Varscan (Koboldt et al., 2012). MAF files from each mutation calling algorithm were integrated in a unique MAF file considering those mutations that were called at least by two of the three considered methods. The integrated MAF contains 28637 somatic mutation called by all the methods, 5559 called by MuTect and VarScan, 7971 called by MuTect and RADIA

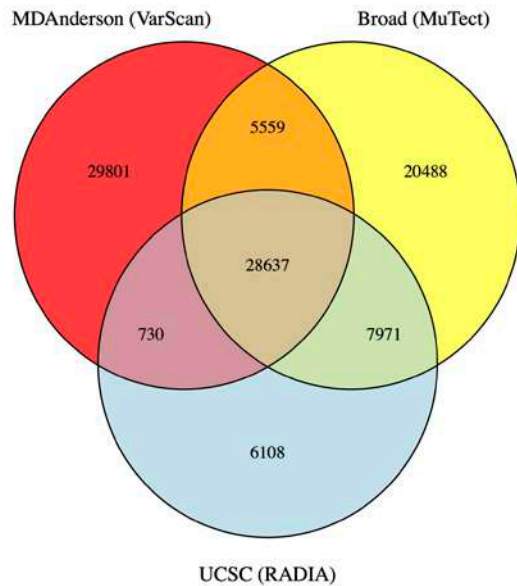


and 730 called by VarScan and RADIA. Similarly, for the detection of somatic insertions and deletions we intersected the calls produced by Indelocator and Varscan algorithms obtaining 1956 high confidence indels.

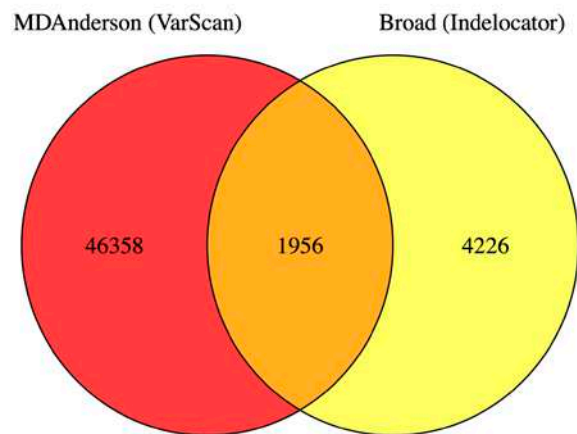
### 2.3 Identification of IDH mutations

In order to expand the annotation of IDH status in our cohort, previously reported (Cancer Genome

**A.** Multicenter Somatic Single Nucleotide Variations



**B.** Multicenter Somatic Small Insertions and Deletions



Atlas Research, 2008) mutation calls on Sanger sequenced DNA and exome sequencing of whole genome amplified DNA were used. Sanger sequencing and whole exome sequencing of whole genome amplified DNA was performed as previously described (Brennan et al., 2013; Cancer Genome Atlas Research, 2008; Verhaak et al., 2010). Except for bona fide IDH1/2 mutations, no other mutations were called on these platforms.

Platform	Center	Aliquot	Disease	Paired samples
ABI	WUGSC	DNA	GBM	158
Illumina	BI	WGA	GBM	163
<b>Union*</b>				<b>174</b>

**Additional data used to determine IDH mutation status.**



## 2.4 Identification of TERT promoter mutations

Targeted sequencing at the TERT promoter region (Chr5:1295150-1295300) was performed on a subset of 287 cases as previously described (Cancer Genome Atlas Research, 2015). Additionally, we evaluated whole genome sequencing (including low-pass) for the presence of somatic variants using GATK pileup. We required a minimum coverage of at least 6 bp and a minimum variant allele fraction of 15% for detection of TERT promoter mutations. A total of 328 cases had sufficient coverage to detect a mutation and 162 cases showed a somatic mutation at one of three sites.

Nucleotide change	Site	Paired samples
A161C	Chr5:1295161	2
C228T	Chr5:1295228	121*
C250T	Chr5:1295250	39*

*\*One case showed mutations in both C250T and C228T*

## 2.5 Mutation significant analysis

Significantly mutated genes were identified using the MutSigCV algorithm. Analyses were conducted on the entire sample set (n=820) except a single hypermutator phenotype (TCGA-DU-6392). Intronic mutations were excluded. A mutation blacklist was applied for remove potential technical artifacts (Lawrence et al., 2013b). Genes with a q-value less than 0.1 were considered significant.

## 2.6 Telomere quantification

Quantification of telomere length was performed using the TelSeq tool (Ding et al., 2014). This tool counts the number of reads containing any (range 0 to  $k$ ) amount of telomeric repeats ( $n_k$ ), or TTAGGG, and then computes the estimated telomere length in bp  $l$  further based on the average chromosome length in bp  $c$  and the total coverage  $s$ .

$$1) l = c \times \frac{n_k}{s}$$

The authors recommend a  $k$  of 7 based on their experimentally validated results. Given that TelSeq was not designed for cancer, it does not take into account tumor ploidy and purity. We have therefore modified the TelSeq computation to consider tumor purity  $p$  and ploidy  $\tau$ :

$$2) \frac{n_k}{s} = \frac{l_t \times \tau \times p + l_n \times (1-p)}{\tau \times c \times p + c \times (1-p)}$$

Because  $p$  and  $\tau$  are given by the ABSOLUTE analysis (Carter et al., 2012), solving  $l_t$  is straightforward, whereas  $l_n$  can be calculated using 1) above.

The average chromosome length  $c$  is calculated as follows:

$$3) c = 46/G$$



Here  $G$  is the total genome length and 46 is the expected number of chromosomes. Because GC content is a potential confounding factor,  $G$  was set to the genome length in bp with GC content between 48% and 52%. The average coverage  $s$  is adjusted in a similar fashion.

## 2.7 Whole genome mutation calling

MuTect (Cibulskis et al., 2013) was used to call somatic mutations on 89 matched primary tumor-normal pairs. We required a minimum coverage of 14 in the tumor sample and 8 in the normal sample. Variants known to dbSNP v132 and unknown to COSMIC v54 were filtered resulting in 714,305 variants. Using these samples we used overlapping RNA-seq expression data to form an integrated dataset of 67 pairs (29 GBM, 38 LGG). In order to identify potential promoter sites we used the GENCODE v19 transcript annotation ( $n=196,520$  transcripts) and used a subset of 24,001 transcripts that have an exact UniProt database match and has been curated according to known clinically relevant protein changes (Ramos et al., 2015). We then reduced the transcripts down to one transcript per gene ( $n=17,722$  transcripts). For each remaining transcript we then took a region spanning from 2,000 bp upstream of the transcription start site and 200 bp into the coding region. We then determined overlapping mutations for each region using the Bioconductor package "GenomicRanges" (Lawrence et al., 2013a). We removed regions with hits from less than 7 unique samples, removed regions which were upstream of genes lacking RNA-seq counts or counts that were lacking any variability, removed regions in which the variants had a median of read count of 1 or more alternate reads in the matching normal. This filtering resulted in 141 mutations across 12 putative promoter regions (Table S2E). For each of the remaining gene promoter regions we then performed a t-test and a mann-whitney-U test comparing the log<sub>2</sub> normalized gene expression counts in mutant cases to wild type cases. When we subsequently filtered out promoter regions with a Benjamini-Hochberg adjusted gene expression correlation Q-value  $< 0.05$  only three promoter regions remained including TERT, TRIM28 and CACNG6.

## 3. DNA copy number analysis

**Authors:** Bradley Murray, Floris Barthel, Roel Verhaak

**Correspondence and questions should be directed to:** Roel Verhaak  
(rverhaak@mdanderson.org)

### 3.1 Preprocessing and peak calling

Tumor and normal samples were profiled on Affymetrix SNP6.0 GeneChip arrays and subsequently processed into genome segmentation files (McCarroll et al., 2008). The tool GISTIC 2.0 was then used to identify significantly reoccurring focal and broad copy number changes (Mermel et al.,



2011). Events with a Q-value  $< 0.10$  were considered significant. In order to identify low-frequent subtype specific events, we ran GISTIC both across the entire cohort ( $n=1084$ ) and smaller subsets within DNA methylation clusters ( $n=6$  groups), RNA expression clusters ( $n=4$  groups) and IDH-codel subtypes ( $n=3$  groups). For each statistically significant peak, GISTIC 2.0 indicates a narrow focal peak and a wider surrounding peak. We intersected all overlapping focal peaks across all GISTIC run and identified 57 disjoint amplified regions and 105 deleted regions. Using this method, while drastically limiting the number of genes compared to using the wide peak boundaries, we were still about to find 80% of genes that were considered as potential tumor drivers in previous studies. Genes previously suggested as tumor drivers not found using this method include IRS2 gain, LSAMP loss and KDR/KIT gain (the neighboring oncogene PDGFRA however was still found). In order to further narrow down the list of genes per peak and to identify potential tumor drivers, we sought to correlate copy number change to gene expression and prioritized genes in which we found significant mutations. Using this method, we were able find evidence for several new tumor drivers including GIGYF2 loss, ERRF11 loss, ARID2 loss and FGFR2 gain. For the complete list of peaks, genes and their mutation and expression correlates see Table S2B.

### **3.2 Functional Copy Number (CN) analysis**

**Authors:** Pietro Zoppoli, Antonio Iavarone

**Correspondence and questions should be directed to:** Pietro Zoppoli (zoppoli@icg.cpmc.columbia.edu)

In order to define the functional copy number (fCN) genes we calculated the Spearman's correlation between the copy number and the expression of each gene in the dataset. We selected all the genes with  $p$ -value  $< 0.05$  and  $cor > 0.5$ .

In order to highlight the different behavior between the four expression groups, we selected only the differentially expressed ( $abs(FC) > 1.5$ ) and aberrated ( $abs(\Delta CN) > 0.5$ ) fCN genes obtaining a list of 57 genes (the fCN signature).

### **3.3 Mutations with Common Focal Alterations (MutComFocal)**

**Authors:** Raul Rabadan, Jiguang Wang, Antonio Iavarone

**Correspondence and questions should be directed to:** Antonio Iavarone (ai2102@cumc.columbia)

By considering both copy number and somatic mutation data of LGG/GBM samples, we applied the algorithm of MutComFocal (Trifonov et al., 2013). Particularly, focality score and recurrence score were calculated based on samples with at least 10 and at most 1,000 copy number segments. The



focality score assigns equal weight to all genes participating in a genomic alteration inversely proportional to the size of that alteration, while recurrence score assigns equal weight to all genes altered in a sample inversely proportional to the total number of gene altered in the sample (Frattini et al., 2013; Trifonov et al., 2013).

## 4. mRNA Expression

**Authors:** Michele Ceccarelli, Stefano M. Pagnotta, Antonio Iavarone

**Correspondence and questions should be directed to:** Michele Ceccarelli (ceccarelli@unisannio.it)

### 4.1 Data preparation and gene selection

RNA-seq raw counts of 667 cases (513 LGG and 154 GBM) were downloaded, normalized and filtered using the Bioconductor package TCGAbiolinks (Colaprico et al., 2015) using TCGAquery(), TCGAdownload() and TCGAprepare() for both tumor types ("LGG" and "GBM", level 3, and platform "IlluminaHiSeq\_RNASeqV2"). The union of the two matrices was then normalized using within-lane normalization to adjust for GC-content effect on read counts and upper-quantile between-lane normalization for distributional differences between lanes applying the TCGAanalyze\_Normalization() function encompassing EDASeq protocol. Gene selected for clustering were chosen by applying two filters, the first was aimed at reducing the batch effect between the two tumor cohorts. We computed differentially expressed genes with TCGAanalyze\_DEA() (implementing the EdgeR protocol (Robinson et al., 2010)), and filtered out genes differentially expressed between the two sets ( $\alpha = 10^{-10}$ ), obtaining 10,389 genes. We then applied variability filters that select genes having a sufficiently high variation (100%) between the mean of top 5% and the mean of the bottom 5% values and having these means respectively above and below the overall median value of the data matrix. The filtering steps resulted in 2,275 genes that were used for the consensus clustering. ConsensusClusterPlus Bioconductor package was used to perform the clustering with hierarchical clustering as inner method and 1000 resampling steps (epsilon=0.8). Number of cluster ( $n = 4$ ) was used as local maxima of the Calinsky-Harabasz curve. Within cluster analysis was done generating differentially expressed genes between GBM and LGG cohorts (log fold change greater and 1.0 and FDR less than 0.05), lists were then analyzed using DAVID functional annotation tool (Huang et al., 2009) and ClueGO (Bindea et al., 2009).

### 4.2 Classification of Affymetrix samples

Once the four RNA-seq clusters were obtained, we reclassified 378 GBM samples for which no RNA-seq data were available using their Affymetrix profiles. We used the 151 GBM samples (20 in LGr1,





4 in LGr2, 10 in LGr3 and 117 in LGr4) having both the Affymetrix and RNA-seq profiles as training set of a kNN classifier ( $k = 3$ ) to assign LGr cluster memberships to the remaining 378 Affymetrix samples. The feature set of the classifier was based on a signature of 327 probesets obtained by selecting up-regulated and down-regulated genes for the training samples in each cluster.

### **4.3 Tumor Map and Pathway Activity Analysis**

**Authors:** Yulia Newton, Olena Morozova, Sofie Salama

**Correspondence and questions should be directed to:** Sofie Salama (ssalama@soe.ucsc.edu)

#### **4.3.1 Combining multi-platform multi-tumor datasets**

We utilized the ComBat batch effect removal method (Johnson et al., 2007) in order to combine mRNA expression data from the GBM RNA-seq (n=154), GBM Agilent (n=525), LGG RNA-seq (n=513), and LGG Agilent (n=27) datasets. We chose to use data generated using Agilent microarray platform over those generated using Affymetrix because such data were available for both tumor types, while Affymetrix data were only available for GBM samples. We combined the 4 datasets and ran ComBat. We flagged 4 batches, one for each dataset, to be removed by the ComBat method. One hundred and forty nine GBM samples were analyzed using both Agilent and RNA-seq platforms. Twenty seven LGG samples were analyzed using both Agilent and RNA-seq platforms. We utilized these matched samples as biological covariates in the ComBat method. Upon completion of the data transformation, we removed all redundant samples analyzed using the Agilent platform whenever the sample was also analyzed using RNA-seq. This combined mRNA expression dataset (n=1043) was used for Tumor Map analysis.

#### **4.3.2. Tumor Map method (manuscript in preparation)**

Tumor Map is a dimensionality reduction and visualization method for high dimensional genomic data. It allows viewing and browsing relationships between high dimensional heterogeneous genomic samples in a two-dimensional map, in a manner much like exploring geo maps in Google Maps web application.

Prior to the analysis, technical and batch effects in the gene expression data were mitigated as a preprocessing step and as described above. We computed sample-by-sample pair-wise similarities. From RNA expression data, we selected 6002 genes whose expression was the most variable based on the variance distribution curve. The 1301 most important methylation probes were selected by manual curation of the probe list as described in the DNA methylation analysis section. We used Spearman rank correlation (Spearman, 1904) on these continuous variable data (mRNA and methylation). To build maps based on a single data type, for each sample the closest



neighborhood of 10 samples is selected. The Tumor Map method represents these local neighborhoods as a graph. The edge weight in this graph is proportional to the magnitude of the similarity metric. Then spring-embedded graph layout (Golbeck and Mutton, 2005) algorithm is applied to the constructed graph. The spring-embedded layout algorithm treats edges as springs and allows the springs to oscillate for a fixed amount of time with the energy inversely proportional to the edge weights. Under these conditions, springs with large weights do not oscillate much, causing those vertices to stay together. However, springs with small weights oscillate more and end up farther away from each other. The method then projects the positions of all the vertices in the resulting graph layout onto a two-dimensional grid. Each cell in the grid allows only one vertex to be placed into it. If multiple vertices contest for the same grid cell, a random vertex selection is made and placed into the cell; and the other competing vertices are placed into the nearest empty cell, snapping around the original cell in a spiral-like manner. Thus, dense clumps of samples are separated so that they can be viewed at approximately the same scale as the distances that separate them. After computing pairwise sample similarities in the gene expression and DNA methylation space separately, the two similarity spaces are combined after standardizing each space was standardized.

#### **4.3.3 Multi-platform maps using Bivariate Standardization similarity space Transformation (BST)**

We computed sample pairwise similarities for each data type separately, producing a square samples-by-samples similarity matrix. For each of the similarity matrices, we perform bivariate standardization by transforming each value to be an arithmetic mean of the z-scores of this value within both the row and the column empirical distributions. This method is an adaptation of the approach by Faith et. al (2007). Once each of the similarity matrices is transformed into a z-score space, we combine each available z-scores (from N platforms) for each pair of samples by taking a weighted average of the z-scores, where the weights indicate the importance of each of the N platforms being combined. When genomic data for a given platform is not available for at least one of the samples from the pair, a pairwise similarity for this pair will not be available for this platform. Our method allows such omissions, as it will only combine similarity z-scores from those platforms that are available for any pair of samples. The resulting BST matrix is a square samples-by-samples matrix that contains a union of samples in all the platforms.



#### 4.3.4 Extracting significantly active pathways

We used mRNA expression for samples available through RNA-seq platform only and the CNV data to transform the data into inferred pathway activity levels using PARADIGM (Vaske et al., 2010). We then considered a number of dichotomies, such as LGm1 GBM vs. LGG (see Table S5). Some of the dichotomies we considered have significantly different numbers of samples in each class (see Table S5). In order to make statistically strong inferences about pathway activities we only considered those dichotomies in which both classes are well represented by their members and the variance within the classes is much smaller than the variance between the classes. In other words, we selected those dichotomies where sample scatter is small within the classes and classes are separable in the pathway space. Based on the PARADIGM IPLs (Inferred Pathway Levels) we computed pair-wise Spearman rank correlation for each pair of samples. We then computed within-class and between-class variance of the correlations, first for the first class and then for the second class. We then computed the F-statistic for each of the classes in the dichotomy and the p-value based on the F-distribution. We aggregated the p-value for the dichotomy by computing the mean p-value. We selected those dichotomies that had an aggregated p-value of  $\leq 0.05$ . Table S5 shows final dichotomies analyzed for the differential pathway activities. For each dichotomy selected, we computed differential activity levels using the linear models for microarrays and RNA-seq data (LIMMA) method (Smyth, 2005). We then applied Gene Set Enrichment Analysis (GSEA) (Subramanian et al., 2005) to the HUGO members of the full differential vector. We extracted only those pathways that had FDR-adjusted q-value of  $\leq 0.1$ . At the same time, we extracted statistically significant differentials (multiple hypothesis adjusted p-value  $\leq 0.05$ ). We ran PATHMARK (Cancer Genome Atlas Research, 2013) on the statistically significant differential activities obtained from LIMMA to extract connected components of the global PARADIGM regulatory network. An additional filter of 3 standard deviations was applied to the PATHMARK method. This means only those activities that fall outside 3 standard deviations of the empirical distribution of the statistically significant differentials pass through the filter. A network connection is extracted if both vertices connected by that connection pass the filter. For each pathway gene set that passed the GSEA q-value of 0.1 we computed the overlap of the pathway genes and those that survived the PATHMARK filter as well as the over-representation hypergeometric p-value. We then extracted those pathways that passed with the p-value of  $\leq 0.05$ . Figure S5E shows an overview of the process for extracting significantly active pathway from the glioma data. Figures S5C-D show pathway views of the significant IPLs from Table S5 in which IPLs representing families, complexes, phospho-events and redundant complexes were removed for better visualization.



## **4.4 Combining GBM Agilent G4502A mRNA data with LGG Illumina Hi-Seq RNA-seq data**

**Authors:** Shiyun Ling, Rehan Akbani

**Correspondence and questions should be directed to:** Rehan Akbani (rakbani@mdanderson.org)

Approximately 15,700 genes were common between the two platforms and a total of 185 pairs of GBM and LGG sample replicates were run on both platforms. Initial tests by combining the GBM and LGG replicates and clustering them showed two clusters based entirely on platform differences and the replicates didn't merge with each other. To remove the platform effect, we developed a novel algorithm that randomly divided the 185 replicate pairs into training, testing and validation sets. The training set was used to train an Empirical Bayes (Johnson et al., 2007) based model, which was then applied to the testing set. The testing set was used to figure out which genes didn't merge well by using a *t*-test to find the genes with the most differences between the platforms. The process was repeated 1000 times by using a bootstrapping approach for the training set. The top 3000 genes that were consistently found to be the most variable in the testing set were removed from the data set. The resulting model was then applied to the validation set, after removing those 3000 genes, to evaluate the algorithm. The evaluations showed that all 43 of the replicate pairs in the validation set clustered in matched pairs. The median of Pearson's correlations between the matched pairs was 0.23 before adjustment and 0.93 after adjustment, indicating very successful merging. We then applied the model to the full GBM and LGG dataset to perform overall merging, and then removed duplicates by randomly keeping one sample from the pairs. The final dataset had 1032 samples and 12,717 genes.

## **4.5 RNA Fusion analysis**

**Authors:** Olena Morozova, Floris Barthel, Sofie Salama, Roel Verhaak

**Correspondence and questions should be directed to:** Roel Verhaak (rverhaak@mdanderson.org)

### **4.5.1 Fusion transcript detection using PRADA**

Transcript fusions were detected in 665 samples using the Pipeline for RNA-seq Data Analysis (PRADA) fusion detection tool (Torres-Garcia et al., 2014). We classified fusions to one of four tiers based on the number of junction spanning reads and discordant read pairs, gene partner uniqueness, gene homology, whether the fusion preserves the open reading frame, transcript allele fraction and DNA breakpoints in SNP6 array data, as previously described (Yoshihara et al., 2014). Briefly, tier one fusions are the highest confidence fusions and tier four fusions are the lowest



confidence ones. For the purpose of this analysis we chose to include tiers one and two. A summary of included fusions can be found in Table S2C.

#### **4.5.2 Fusion transcript detection using deFuse**

RNA-seq reads were analyzed using deFuse package version 0.6.0 (McPherson et al., 2011). Fusions involving receptor tyrosine kinase genes were manually reviewed using blat analysis (Kent, 2002) of the breakpoint sequence in the UCSC Genome Browser (Kent et al., 2002). Candidate fusions were filtered based on the following deFuse parameters:

- Splitr\_count > = 5 (5 or more split reads supporting the fusion)
- Span\_count > = 10 (10 or more spanning reads supporting the fusion)
- Read\_through ~ "N" (fusion is not a readthrough)
- Adjacent ~ "N" (fusion does not involve adjacent genes)
- Altsplice ~ "N" (fusion cannot be explained by alternative splicing)
- Min\_map\_count = 1 (at least one spanning read supporting the fusion is uniquely mapped)
- ORF ~ "Y" (fusion preserves the open reading frame)

deFuse and PRADA fusion predictions were combined to generate a list of 204 events identified by both methods (Table S2C).

## **4.6 Identification of Transcriptional Regulatory Factors underlying IDH wild type and IDH mutant phenotypes in Glioma**

**Authors:** Ganiraji Manyam, Arvind Rao, Ganesh Rao

**Correspondence and questions should be directed to:** Ganesh Rao (grao@mdanderson.org)

Batch-corrected expression data from Agilent Microarray and Illumina HiSeq RNA-seq platforms using MBatch was used for differential expression and transcription factor analysis. Linear regression was used to find the genes that are differentially expressed between IDH wild type and IDH-mutant groups after adjusting for the effect of expression platform (RNA-seq or microarray) in the model. The p-values are adjusted for multiple testing using the Bonferroni method. Genes with adjusted p-value less than 0.01 are considered significant.

Transcription Factor (TFs) Analysis was performed using the Match Algorithm of Biobase (TRANSFAC) system to identify TFs enriched in promoters of genes differentially expressed between IDH wild type and mutant groups. This algorithm compares the number of TF binding sites found in a query sequence set against a background set and identifies factors whose frequencies are enriched in the query compared to the background. Genes significantly upregulated in the IDH



mutant group are considered as the background for TF analysis of genes upregulated in IDH wild type group and vice-versa. The TFs enriched with p-value less than 0.05 are considered significant. Differential expression analysis was used to assess the expression differences of the enriched TFs themselves between the two groups (IDHmut vs wt). The transcription factors with Bonferroni-adjusted p-value less than 0.05 are defined as significant candidates (Excel file).

Ingenuity Pathway Analysis (IPA) was used to generate downstream networks for the top ranking transcription factors. Rank of the transcription factor is defined based on fold change between the two groups and the number of transcription factor binding sites in the promoter region of the target genes. Twelve transcription factor families were found to have log fold change of >1 between the IDH mut and IDHwt groups. The ones with the highest number of target genes are NKX2-5, PAX8, ETV7, CEBPD, ETV4, ELF4, and NFE2L3. Several of these TFs have been shown to be important for carcinogenesis. For example, Pax8 has been shown to be minimally expressed in LGG and normal brain but highly expressed in glioblastoma (Hung et al., 2014) and plays a role in telomerase regulation (Chen et al., 2008). Similarly, enrichment of the pro-proliferative TF ETV4 in 1p/19q codeleted gliomas has been demonstrated (Gleize et al., 2015).

## 5. DNA methylation profiling

**Authors:** Thais S. Sabedot, Tathiane M. Malta, Simon G. Coetzee, Peter W. Laird, Houtan Noushmehr

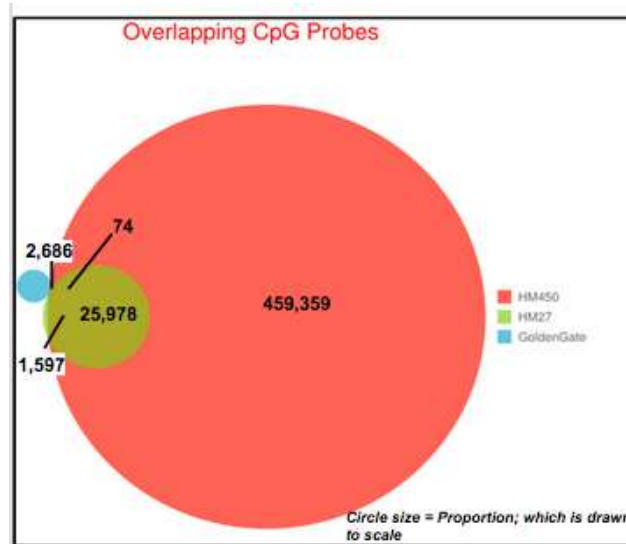
**Correspondence and questions should be directed to:** Houtan Noushmehr (houtan@usp.br)

### 5.1 Preprocessing and clustering

For data acquisition, we used the the Bioconductor package TCGAAbiolinks (Colaprico et al., 2015). First, TGCAquery() was used to search the samples of “GBM” and “LGG” tumors in TCGA repository using the following parameters: data level = 3, platform type = “HumanMethylation450” and “HumanMethylation27”, version 12 for LGG and version 5 GBM samples. Second, TCGAdownload() was used to download the data; and, finally, TCGAprepare() was used to read the data into a dataframe. A total of 932 TCGA glioma samples assessed for DNA methylation, including 516 LGG and 416 GBM samples, profiled using 2 different Illumina platforms, were included. During the initial phase of the TCGA project, 287 GBM samples (batches 1 to 9) were profiled using the Illumina HumanMethylation 27 platform (HM27), which interrogates 27,578 CpG probes. As a new platform became available, the TCGA LGG (batches 1 to 16) and 129 GBM (batches 1 to 12) samples were transitioning into the larger more comprehensive Illumina platform known as the HumanMethylation450 (HM450), which interrogates 485,421 CpG sites. The DNA methylation score for each locus is presented as a beta ( $\beta$ ) value ( $\beta = (M/(M+U))$ ) in which M and U indicate the mean



methyated and unmethyated signal intensities for each locus, respectively.  $\beta$ -values range from zero to one, with scores of zero indicating no DNA methylation and scores of one indicating complete DNA methylation. A detection p-value also accompanies each data point and compares the signal intensity difference between the analytical probes and a set of negative control probes on the array. Any data point with a corresponding p-value greater than 0.01 is deemed not to be statistically significantly different from background and is thus masked as “NA” in TCGA level 3 data packages. The data levels and the files contained in each data level package are present on the TCGA Data Portal website (<http://tcga-data.nci.nih.gov/tcga/>). Please note that as continuing updates of genomic databases and data archive revisions frequently become available, the data packages on TCGA Data Portal are updated accordingly. Data of the two platforms (HM450 and HM27) were merged as previously described (Brennan et al., 2013) and we ended with 25,978 probes that match both 27k and 450k platforms, as illustrated in the following Venn diagram. Duplicated samples and secondary tumors were excluded. The 932 sample IDs used for DNA methylation analysis are listed in Table S1.



## 5.2 Unsupervised clustering analysis of DNA methylation data

Methods to capture tumor-specific DNA methylation probes were used as recently described (Cancer Genome Atlas Research, 2014b) and is provided here as reference, with slight modifications to the total numbers. We used the Level 3 DNA methylation data contained in the packages listed above for analyses. We first removed probes which had any “NA”-masked data points and probes that were designed for sequences on X and Y chromosomes. We selected CpG sites that are located in high CpG density regions (top 25% of the sites with the highest observed/expected CpG ratio around their 3kb regions spanning from 1,500 bp upstream to 1,500



bp downstream of the transcription start sites) and CpGs associated with CpG islands extracted from the UCSC Genome Browser (<http://genome.ucsc.edu>). To capture cancer-specific DNA hypermethylation events, we further eliminated sites that were methylated (mean  $\beta$ -value  $\geq 0.3$ ) in histologically non-tumor brain tissues (Guintivano et al., 2013). This selection method reduced the initial 25,978 probes to 1,300 glioma-specific CpG probes, which corresponded to 6.5% of the full available data. However, a clustering analysis can be strongly confounded by the purity of tumor samples. To alleviate the potential influence of variable levels of tumor purity in our sample set on our clustering result, we dichotomized the data using a  $\beta$ -value of  $>0.3$  as a threshold for positive DNA methylation. We then performed unsupervised hierarchical clustering on 1,300 CpG sites with this threshold that are methylated in at least 10% of the tumors using a binary distance metric for clustering and Ward's method for linkage. The cluster assignments were generated by cutting the resulting dendrogram. The probes are arranged based on the order of unsupervised hierarchical clustering of the dichotomous data using a binary distance metric and Ward's linkage method. We identified six groups (LGm1-LGm6) shown in Figure 2A generated based on the original  $\beta$ -values to visualize 1,300 CpG sites used in the clustering.

The approach described above to capture tumor-specific DNA methylation probes was used to select glioma-specific CpG probes and perform unsupervised clustering separated by IDH status. We identified 1,308 tumor specific CpG probes for IDH-mutant analysis and identified three IDH-mutant-specific clusters (Figure S3A). Likewise, we identified 914 tumor specific CpG probes for IDH-wild type samples and identified three IDH-wildtype-specific clusters (Figure S4A).

In order to classify the newly acquired TCGA samples (not included in the previous studies; LGG = 227; GBM = 20) into the context of previously published DNA methylation clusters (Brennan et al., 2013; TCGA\_Network, 2015), we randomly selected a set of 80% of TCGA samples to train a random forest machine-learning. We then evaluated the performance on the remaining 20% of samples and got an accuracy of more than 88% on average. We then tested the new TCGA samples and classified them into the previously DNA methylation clusters.

### **5.3 Supervised analysis of DNA methylation**

We used Wilcoxon test followed by multiple testing using the Benjamini and Hochberg (BH) method for false discovery rate estimation (Benjamini and Hochberg, 1995) to identify differentially DNA methylated probes between two groups of interest.

The 131 probes presented in Figure 3A were defined comparing samples from IDHmut-K1 (n=53) to IDHmut-K2 (n=221), using the following criteria: FDR  $< 10e-15$ , absolute difference in mean methylation beta-value  $> 0.27$ .





The 90 probes presented in Figure 3H were identified comparing samples from G-CIMP-low (n=25) to G-CIMP-high (n=249), in order to identify probes defining the G-CIMP-low group, using the following criteria: FDR < 10e-13, difference in mean methylation beta-value > 0.3 and < -0.4.

The 149 probes presented in Figure 3H were a combination of the 90 probes described above with 73 probes identified from the comparison between non-codons (from LGm2, n=210) and codons (from LGm3, n=120), using the following criteria: FDR < 10e-30, absolute difference in mean methylation beta-value > 0.25, removing probes with NA values. All probeset lists are provided on the publication portal accompanying this publication ([https://tcga-data.nci.nih.gov/docs/publications/lgggbm\\_2015/](https://tcga-data.nci.nih.gov/docs/publications/lgggbm_2015/)).

#### **5.4 Identification of Epigenetically Regulated Genes**

To increase our statistical power, we decided to evaluate epigenetically regulated genes using the Pan-glioma subtypes, which allowed us to use the entire TCGA glioma cohort. We selected tumor samples that have both DNA methylation and RNA-sequencing based gene expression data to do this analysis, resulting in 636 samples (513 LGG and 123 GBM). We also randomly selected 110 non-tumor TCGA samples from 11 different tissues ([https://tcga-data.nci.nih.gov/docs/publications/lgggbm\\_2015/](https://tcga-data.nci.nih.gov/docs/publications/lgggbm_2015/)), profiled using the same platforms. Each DNA methylation probe was mapped to the nearest UCSC gene, and after merging the DNA methylation and gene expression data, we retained a total of 19,530 pairs of DNA methylation and gene expression probes. We organized the tumor samples as either methylated ( $\beta \geq 0.3$ ) or unmethylated ( $\beta < 0.3$ ) for each probe. We selected the pair of DNA methylation and gene expression probes for which the mean expression in the methylated group was lower than 1.28 standard deviation (bottom 10%) of the mean expression in the unmethylated group, and in which >80% of the samples in the methylated group have expression levels lower than the mean expression in the unmethylated group. We labeled each tumor sample as epigenetically silenced for a specific probe/gene pair if: it belonged to the methylated group and the gene expression level was lower than the mean of the unmethylated group silenced (Cancer Genome Atlas Research, 2014a), resulting in 3,806 probes/genes identified as epigenetically regulated. A Fisher test was used to detect if these 3,806 pairs were enriched in a DNA methylation cluster. For each probe, tumor samples labeled as methylated and downregulated by cluster, while non-tumor samples labeled as unmethylated and upregulated, were counted and arranged into a contingency table for a Fisher test, using 50% as a cutoff. p-value was calculated for each probe/gene pair and then was adjusted for multiple testing using the BH method for false discovery rate estimation (Benjamini and Hochberg, 1995). This analysis identified 3 Epigenetically Regulated groups (EReg): EReg2 with 233 genes enriched in LGm2 (resembling G-CIMP high), EReg3 with 15 genes enriched in LGm3



(resembling Codels) and EReg4 with 14 genes enriched in LGm4 (resembling Classic-like) and 1 gene enriched in LGm5 (resembling Mesenchymal-like). Since LGm1 (enriched for G-CIMP-low) and LGm6 (comprising LGm6-GBM and PA-like) are heterogeneous clusters, we applied a different approach in order to identify epigenetically regulated genes for these groups. For EReg1, we compared the DNA methylation and gene expression levels for G-CIMP-low samples (n=25) with G-CIMP-high samples (randomly selected 50 samples out of 249) and those probes/genes with Wilcoxon BH adjusted p-value less than  $1e-10$ , methylation difference greater than 0,25 and RNA expression log Fold Change greater than 0,85 were selected, resulting in 15 epigenetically regulated genes enriched in G-CIMP-low. For EReg5, we compared the DNA methylation levels for LGm6 samples (n=77) with a subset of randomly selected samples from the 855 remaining TCGA glioma samples (n=140) and those probes with Wilcoxon BH adjusted p-value less than  $1e-21$  and methylation difference greater than 0,33 were selected, resulting in 12 epigenetically regulated genes enriched in LGm6.

To validate the EReg genes in order to confirm the existence of these signatures in an independent, non-TCGA data, we downloaded 4 different and publicly available datasets (Lambert et al., 2013; Mur et al., 2013; Sturm et al., 2012; Turcan et al., 2012), comprising 324 samples with distinct histology and clinical attributes. These samples included adult, pediatric gliomas of both low and high grade, reported with codel, IDH status and G-CIMP status. Our independent data set included a pool of 61 pilocytic astrocytomas defined as grade I gliomas (Lambert et al., 2013). In order to classify the additional non-TCGA gliomas into our LGm clusters, we selected a random set of 80% TCGA samples to train a random forest machine-learning model and evaluated the performance on the remaining 20%. Given the high specificity and sensitivity of our model (accuracy > 88% on average), we, then, tested the LGm cluster prediction model on the additional non-TCGA samples using the random forest method. Data were visualized using the same 45 pairs of CpG probes/genes that define the epigenetically regulated genes for IDH mutant samples (Figure 3F) and the same 27 pairs of CpG probes/genes that define the epigenetically regulated genes for IDH wild type samples (Figure 4D). Applying a similar ordering in the validation set and accounting for differences in sample size, we recapitulated the five EReg groups both for IDH mutant samples (Figure 3G) and IDH wild type samples (Figure 4E) in molecular level. The list of epigenetically regulated genes can be found at [https://tcga-data.nci.nih.gov/docs/publications/lgggbm\\_2015/](https://tcga-data.nci.nih.gov/docs/publications/lgggbm_2015/).

The same random forest machine learning model approach was used for the IDH-mutant samples (using the 1,308 IDH-mutant tumor specific CpG probes) and for the IDH-wildtype samples (using the 914 IDH-wildtype tumor specific CpG probes), separately. We then tested the models in the IDH-



mutant and IDH-wildtype samples from the validation set (Lambert et al., 2013; Mur et al., 2013; Sturm et al., 2012; Turcan et al., 2012) (Figure S4B).

## **5.5 Classification of new glioma samples based on DNA methylation glioma subtypes**

New glioma samples can be classified into one of our glioma subtypes using our CpG probe methylation signatures provided on the publication portal accompanying this publication ([https://tcga-data.nci.nih.gov/docs/publications/lggqbm\\_2015/](https://tcga-data.nci.nih.gov/docs/publications/lggqbm_2015/)).

First, all glioma samples should be divided by their known IDH status, separated into either IDH-mutant and IDH-wildtype. IDH-mutant is defined as those samples harboring any type of known IDH1 or IDH2 mutation as described recently (TCGA\_Network, 2015). IDH-wildtype refers to those samples with an intact IDH1 or IDH2. Samples as either IDH-mutant or IDH-wildtype are then further classified accordingly:

### **IDH-mutant:**

In order to define newly diagnosed glioma samples into one of the 3 glioma subtypes within IDH-mutants, we recommend applying Random Forest in a two-step process. 1) using the 1,308 tumor specific CpG probes which defines the IDHmut specific clusters (Fig S3A) and 2) using the 163 CpG probes which defines each TCGA IDH-mutant glioma subtype (Fig S3C).

1. If the sample was classified as IDHmut-K1 or IDHmut-K2 using the 1,308 tumor specific CpG probes for IDH-mutant and as G-CIMP-low using the 163 CpG probes defined by a supervised analysis across IDH-mutant subgroups, we classify the sample as G-CIMP-low;
2. If the sample was classified as IDHmut-K1 or IDHmut-K2 using the 1,308 tumor specific CpG probes for IDH-mutant and as G-CIMP-high using the 163 CpG probes defined by a supervised analysis across IDH-mutant subgroups, we classify the sample as G-CIMP-high;
3. If the sample was classified as IDHmut-K3 using the 1,308 tumor specific CpG probes for IDH-mutant, we classify the sample as Codel.

### **IDH-wildtype:**

Likewise, IDH-wildtype can be classified using a single random forest machine-learning model applied with a signature defined by the 914 tumor specific CpG probes for IDH-wildtype (Figures S4A-B). Samples following into IDHwt-K3 (aka LGm6), we recommend subdividing this group based on grade, resulting in either LGm6-GBM and PA-like (LGG).



## 5.6 Patient centric table (DNA methylation)

To generate DNA methylation calls for each sample per gene per overlapping platforms (HM27, HM450), we began by first collapsing multiple CpGs to one representative gene. Using the associated gene expression data (organized as one gene - one expression value per sample), we merged the samples and CpG probes with gene expression data for each platform. We next calculated the spearman correlation ( $\rho$ ) across all samples for all CpG probes for each gene to one gene expression value. For multiple CpGs for each annotated gene promoter, we selected one CpG probe with the lowest correlation rho value to the associated gene expression profile to capture the most biologically representative event (epigenetic silencing). This effectively reduced the number of CpG probes from N:1 to 1:1. Our data set was then reduced down to 636 samples x 19,486 CpG:Gene.

Next, we assigned discrete categories based on the spearman correlation rho value according to the following criteria:

1. Strongly negatively correlated (SNC) when  $\rho$  value is less than 0.5;
2. Weakly negatively correlated (WNC) when  $\rho$  value is between 0.5 and 0.25;
3. No negative correlation (NNC) when  $\rho$  value is greater than 0.25.

Next, we assigned samples to either the 10th (T10 or N10) or 90th (T90 or N90) percentile based on the observed beta-value across tumor samples (T) and normal samples (N). For the normal samples, we used 110 non-tumor TCGA samples from 11 different tissues previously described. We assigned labels for each gene per platform per tissue type (tumor and normal) according to the following rules:

1. If percentile 90 < 0.25, we assign it as CUN or CUT (constitutively unmethylated in normal or tumor);
2. If percentile 10 > 0.75, we assign it as CMN or CMT (constitutively methylated in normal or tumor);
3. If percentile 10 > 0.25 and percentile 90 < 0.75, we assign it as IMN or IMT (intermediate methylated in normal or tumor);
4. If it doesn't fall in any of the above categories, it is assign VMN or VMT (variably methylated in normal or tumor).

Next we assigned a 'call' and a confidence 'score' for each possible combinations (48) [3 (SNC, WNC, NNC) x 4 (CUN, CMN, VMN, IMN) x 4 (CUT, CMT, VMT, IMT)]. We created the following relationship for each call and score based on our interpretation of the most informative epigenetic event (e.g. promoter DNA hypermethylation and low expression). Users should understand that the selection and criteria performed were done to the best of our knowledge at the time. We felt most



confident with calling epigenetically silenced events and this is reflected in the confidence score.

The methylation calls are as follows:

MG: Methylation gain compared to normal

ML: Methylation loss compared to normal

MT: Methylated in tumor

UT: Unmethylated in tumor

ES: Epigenetically silenced

UC: Unable to make call

Methylation class confidence scores vary from 0 (no call) to 4 (high confidence). Patient centric table can be accessed at [https://tcga-data.nci.nih.gov/docs/publications/lgggbm\\_2015/](https://tcga-data.nci.nih.gov/docs/publications/lgggbm_2015/).

## 5.7 Homer de novo motif searches

De novo Motif discovery was performed using HOMER (script v4.4 (8-25-2014)), an algorithm previously described (Heinz et al., 2010). Briefly, differentially methylated probes were classified according to genomic location into CpG island, CpG shores, and open seas as follow: CpG islands were defined based on UCSC annotation and as per the criteria previously described (Gardiner-Garden and Frommer, 1987; Takai and Jones, 2002). Coverage of CpG island regions was further enhanced by including the 2 kb regions flanking CpG island, referred to here as CpG shores. CpGs isolated in the genome were defined as open seas. Probes mapped to each region were used to performed de novo motif analysis using HOMER (HOMER perl script 'findMotifsGenome.pl'). To increase sensitivity of the method, up to two mismatches were allowed in each oligonucleotide sequence and distributions of CpG content in 'target' and 'background' sequences were selectively weighted to equalize the distributions of CpG content in both sets. Raw outputs from HOMER can be found at [https://tcga-data.nci.nih.gov/docs/publications/lgggbm\\_2015/](https://tcga-data.nci.nih.gov/docs/publications/lgggbm_2015/).

## 6. Reverse phase protein array (RPPA)

**Authors:** Rehan Akbani, Zhenlin Ju, Yiling Lu, Gordon Mills

**Correspondence and questions should be directed to:** (rakbani@mdanderson.org)

### 6.1 Data Processing

Protein was extracted using RPPA lysis buffer (1% Triton X-100, 50 mmol/L Hepes (pH 7.4), 150 mmol/L NaCl, 1.5 mmol/L MgCl<sub>2</sub>, 1 mmol/L EGTA, 100 mmol/L NaF, 10 mmol/L NaPPi, 10% glycerol, 1 mmol/L phenylmethylsulfonyl fluoride, 1 mmol/L Na<sub>3</sub>VO<sub>4</sub>, and aprotinin 10 ug/mL) from human tumors and RPPA was performed as described previously (Coombes, 2011; Hennessy et al., 2007; Hu et al., 2007; Liang et al., 2007; Tibes et al., 2006). Lysis buffer was used to lyse frozen



tumors by Precellys homogenization. Tumor lysates were adjusted to 1  $\mu\text{g}/\mu\text{L}$  concentration as assessed by bicinchoninic acid assay (BCA) and boiled with 1% SDS. Tumor lysates were manually serially diluted in two-fold of 5 dilutions with lysis buffer. An Aushon Biosystems 2470 arrayer (Burlington, MA) printed 1,056 samples on nitrocellulose-coated slides (Grace Bio-Labs). Slides were probed with 196 validated primary antibodies (Cancer Genome Atlas Research, 2015) followed by corresponding secondary antibodies (Goat anti-Rabbit IgG, Goat anti-Mouse IgG or Rabbit anti-Goat IgG). Signal was captured using a DakoCytomation-catalyzed system and DAB colorimetric reaction. Slides were scanned in a CanoScan 9000F. Spot intensities were analyzed and quantified using Array-Pro Analyzer (Media Cybernetics Washington DC) to generate spot signal intensities (Level 1 data). The software SuperCurveGUI (Coombes, 2011; Hu et al., 2007), available at <http://bioinformatics.mdanderson.org/Software/supercurve/>, was used to estimate the EC50 values of the proteins in each dilution series (in log2 scale). Briefly, a fitted curve ("supercurve") was plotted with the signal intensities on the Y-axis and the relative log2 concentration of each protein on the X-axis using the non-parametric, monotone increasing B-spline model (Tibes et al., 2006). During the process, the raw spot intensity data were adjusted to correct spatial bias before model fitting. A QC metric (Coombes, 2011) was returned for each slide to help determine the quality of the slide: if the score is less than 0.8 on a 0-1 scale, the slide was dropped. In most cases, the staining was repeated to obtain a high quality score. If more than one slide was stained for an antibody, the slide with the highest QC score was used for analysis (Level 2 data). Protein measurements were corrected for loading as described (Coombes, 2011; Gonzalez-Angulo et al., 2011; Hu et al., 2007) using median centering across antibodies (level 3 data). In total, 196 antibodies and 473 samples were used. Final selection of antibodies was also driven by the availability of high quality antibodies that consistently pass a strict validation process as previously described (Hennessy et al., 2010). These antibodies are assessed for specificity, quantification and sensitivity (dynamic range) in their application for protein extracts from cultured cells or tumor tissue. Antibodies are labeled as validated and use with caution based on degree of validation by criteria previously described (Hennessy et al., 2010).

Two RPPA arrays were quantitated and processed (including normalization and load controlling) as described previously, using MicroVigene (VigeneTech, Inc., Carlisle, MA) and the R package SuperCurve (version-1.3), available at <http://bioinformatics.mdanderson.org/OOMPA> (Hu et al., 2007; Tibes et al., 2006). Raw data (level 1), SuperCurve nonparametric model fitting on a single array (level 2), and loading corrected data (level 3) were deposited at the DCC.



## 6.2 Data normalization

We performed median centering across all the antibodies for each sample to correct for sample loading differences. Those differences arise because protein concentrations are not uniformly distributed per unit volume. That may be due to several factors, such as differences in protein concentrations of large and small cells, differences in the amount of proteins per cell, or heterogeneity of the cells comprising the samples. By observing the expression levels across many different proteins in a sample, we can estimate differences in the total amount of protein in that sample vs. other samples. Subtracting the median protein expression level forces the median value to become zero, allowing us to compare protein expressions across samples.

Surprisingly, processing similar sets of samples on different slides of the same antibody may result in datasets that have very different means and variances. Neely et al. (2009) processed clinically similar ALL samples in two batches and observed differences in their protein data distributions. There were additive and multiplicative effects in the data that could not be accounted by biological or sample loading differences. We observed similar effects when we compared the two batches of GBM and LGG tumor protein expression data. A new algorithm, replicates-based normalization (RBN), was therefore developed using replicate samples run across multiple batches to adjust the data for batch effects. The underlying hypothesis is that any observed variation between replicates in different batches is primarily due to linear batch effects plus a component due to random noise. Given a sufficiently large number of replicates, the random noise is expected to cancel out (mean=zero by definition). Remaining differences are treated as systematic batch effects. We can compute those effects for each antibody and subtract them out. Many samples were run in both batches. One batch was arbitrarily designated the “anchor” batch and was to remain unchanged. We then computed the means and standard deviations of the common samples in the anchor batch, as well as the other batch. The difference between the means of each antibody in the two batches and the ratio of the standard deviations provided an estimate of the systematic effects between the batches for that antibody (both location-wise and scale-wise). Each data point in the non-anchor batch was adjusted by subtracting the difference in means and multiplying by the inverse ratio of the standard deviations to cancel out those systematic differences. Our normalization procedure significantly reduced technical effects, thereby allowing us to merge the datasets from different batches.

## 6.3 Clustering

We used consensus clustering to cluster the samples in an unsupervised way, with Pearson correlation as the distance metric and Ward as the linkage algorithm. A total of 473 samples and



196 antibodies were used in the analysis. Two clusters were observed that largely corresponded with tumor type (Figure S3E), however, there were a few notable exceptions. Whereas only one GBM sample clustered with the LGG samples, twenty-six LGG samples were found to cluster with the GBM samples. Seventeen of those twenty-six samples had no mutations in IDH1/2, similar to the GBM samples. Furthermore, compared to the LGG-like cluster, the GBM-like cluster had elevated expression of IGFBP2, fibronectin, PAI1, HSP70, EGFR, phosphoEGFR, phosphoAKT, Cyclin B1, Caveolin, Collagen VI, Annexin1 and ASNS, whereas it had low expression of PKC (alpha, beta and delta), PTEN, BRAF, and phosphoP70S6K.

## 7. Regulome Explorer

**Authors:** Geetika Sethi, Brady Bernard, Vesteynn Thorsson, Sheila Reynolds, Lisa Lype, Ilya Shmulevich

**Correspondence and questions should be directed to:** [ilya.shmulevich@systemsbiology.org](mailto:ilya.shmulevich@systemsbiology.org)

### 7.1. Feature Matrix

Associations among the diverse clinical and molecular data are identified through construction of a “feature matrix” (FM) by integrating values from all data types. Each column in the FM represents one of the 1123 tumor samples. Each row in the FM represents a single clinical, sample or molecular data element (mRNA expression levels, microRNA expression levels, protein levels (RPPA), copy number alterations, DNA methylation levels and somatic mutations), and the individual data values may be numerical (continuous or discrete) or categorical, as appropriate. Missing values are indicated within the FM by “NA”, and the number of non-NA data values varies significantly across the different data types (rows). Data were retrieved from the DCC on November 18, 2015 and further processed as follows. Clinical and sample data (633 features): DCC clinical and sample data were processed into a matrix. Cluster assignments: Cluster memberships resulting from unsupervised clustering for each of the individual molecular data types: SCNA (Supplement 3), RNAseq (Supplement 4), DNA methylation (Supplement 5), and RPPA (Supplement 6) were incorporated into the FM. Mutation





rates and categories (Supplement 2) were included in the FM as well. Molecular datasets include Gene expression (15,401 features): Gene level RSEM values from RNA-seq (Supplement SA) were log2 transformed, and filtered to remove low-variability genes (bottom 25% removed, based on interdecile range). MicroRNA expression (692 features): The summed and normalized microRNA quantification files were log2 transformed, and filtered to remove low-variability microRNAs (bottom 25% based on zero-count). Somatic copy number alterations: Copy number and focal copy number changes were obtained for peaks identified by GISTIC as described above (Supplement 3, 6318 features). DNA methylation (19,727 features): Probe-specific level-3  $\beta$ -values were obtained as described above (Supplement). We started with the probes common between the two methylation platforms, and then removed the bottom 25% based on interdecile range. Somatic mutations (2842): The Mutations Annotation Format file (Supplement 1), was used to generate a binary indicator vector indicating whether a particular non-silent mutation is present in a specific sample. Mutation features found in fewer than two tumor samples were removed. Overall, the `gbm_lgg` feature matrix has 45839 features (inclusive of the above mentioned analysis platforms) for all the 1122 patients (data freeze list) resulting in 51477197 matrix elements (48501 x 38), with approximately 89% non-NA elements (197478 out of 1843038).

The Synapse platform by Sage Bionetworks ([www.sagebase.org](http://www.sagebase.org), [1]) was used during the development of this project for distributing versioned data to project researchers and as a staging area for assembling files into the Feature Matrix.

## **7.2. All-by-all Pairwise Associations**

Statistical association among the diverse data elements in this study was evaluated by comparing pairs of columns in the feature matrix. Hypothesis testing was performed by testing



against null models for absence of association, yielding a p-value. P-values for the association between and among clinical and molecular data elements were computed according to the nature of the data levels for each pair: discrete vs. discrete (Fisher's exact test); discrete vs. continuous (ANOVA F- test, equivalently t-test for binary vs. continuous) or continuous vs. continuous (F-test). Ranked data values were used in each case. To account for multiple-testing bias, the p-value was adjusted using the Bonferroni correction. Exploring potentially interesting genomic relationships have been of interest to researchers previously [2]. In order to allow researchers to further explore genomic associations in TCGA gbm\_lgg dataset, including primary data, the statistically significant pairs of associations were loaded into the Regulome Explorer web application, which is designed to enable researchers to explore associations among multiple data types in cancer genomics. Prior to loading, a p-value threshold was chosen specific to each pair of data types in such a way as to strike a balance between making potentially interesting associations available to queries by the tool, while still allowing the tool to be responsive, since the number of loaded graph edges (each corresponding to a statistically significant relationship) is in the millions. All identified pairwise relationships, including those described in this manuscript can be found at <http://explorer.cancerregulome.org>.

### **Regulome explorer references**

1. Omberg, L., et al., *Enabling transparent and collaborative computational analysis of 12 tumor types within The Cancer Genome Atlas*. Nat Genet, 2013. **45**(10): p. 1121-6.
2. Sethi, G., et al., *An RNA interference lethality screen of the human druggable genome to identify molecular vulnerabilities in epithelial ovarian cancer*. PLoS One, 2012. **7**(10): p. e47086.

## **8. Supplemental References**

Benjamini, Y., and Hochberg, Y. (1995). Controlling the False Discovery Rate - a Practical and Powerful Approach to Multiple Testing. J Roy Stat Soc B Met 57, 289-300.



Bindea, G., Mlecnik, B., Hackl, H., Charoentong, P., Tosolini, M., Kirilovsky, A., Fridman, W.H., Pages, F., Trajanoski, Z., and Galon, J. (2009). ClueGO: a Cytoscape plug-in to decipher functionally grouped gene ontology and pathway annotation networks. *Bioinformatics* 25, 1091-1093.

Brennan, C.W., Verhaak, R.G., McKenna, A., Campos, B., Noushmehr, H., Salama, S.R., Zheng, S., Chakravarty, D., Sanborn, J.Z., Berman, S.H., *et al.* (2013). The somatic genomic landscape of glioblastoma. *Cell* 155, 462-477.

Cancer Genome Atlas Research, N. (2008). Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature* 455, 1061-1068.

Cancer Genome Atlas Research, N. (2013). Comprehensive molecular characterization of clear cell renal cell carcinoma. *Nature* 499, 43-49.

Cancer Genome Atlas Research, N. (2014a). Comprehensive molecular characterization of gastric adenocarcinoma. *Nature* 513, 202-209.

Cancer Genome Atlas Research, N. (2014b). Comprehensive molecular characterization of urothelial bladder carcinoma. *Nature* 507, 315-322.

Cancer Genome Atlas Research, N. (2015). Comprehensive, Integrative Genomic Analysis of Diffuse Lower-Grade Gliomas. *N Engl J Med*.

Carter, S.L., Cibulskis, K., Helman, E., McKenna, A., Shen, H., Zack, T., Laird, P.W., Onofrio, R.C., Winckler, W., Weir, B.A., *et al.* (2012). Absolute quantification of somatic DNA alterations in human cancer. *Nat Biotechnol* 30, 413-421.

Chen, Y.J., Campbell, H.G., Wiles, A.K., Eccles, M.R., Reddel, R.R., Braithwaite, A.W., and Royds, J.A. (2008). PAX8 regulates telomerase reverse transcriptase and telomerase RNA component in glioma. *Cancer research* 68, 5724-5732.

Cibulskis, K., Lawrence, M.S., Carter, S.L., Sivachenko, A., Jaffe, D., Sougnez, C., Gabriel, S., Meyerson, M., Lander, E.S., and Getz, G. (2013). Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nature biotechnology* 31, 213-219.

Colaprico, A., Silva, T.C., Olsen, C., Garofano, L., Cava, C., Garolini, D., Sabedot, T., Malta, T., Pagnotta, S.M., Castiglioni, I., *et al.* (2015). TCGAbiolinks: An R/Bioconductor package for integrative analysis of TCGA data. doi: 10.1093/nar/gkv1507.

Coomes, K.N., S.; Joy, C.; Hu, J.; Baggerly, K.; *et al.* (2011). SuperCurve Package. R package version 1.4.1.

Ding, Z., Mangino, M., Aviv, A., Spector, T., Durbin, R., and Consortium, U.K. (2014). Estimating telomere length from whole genome sequence data. *Nucleic acids research* 42, e75.

Faith, J.J., Hayete, B., Thaden, J.T., Mogno, I., Wierzbowski, J., Cottarel, G., Kasif, S., Collins, J.J., and Gardner, T.S. (2007). Large-scale mapping and validation of Escherichia coli transcriptional regulation from a compendium of expression profiles. *Plos Biol* 5, 54-66.

Frattini, V., Trifonov, V., Chan, J.M., Castano, A., Lia, M., Abate, F., Keir, S.T., Ji, A.X., Zoppoli, P., Niola, F., *et al.* (2013). The integrated landscape of driver genomic alterations in glioblastoma. *Nature genetics* 45, 1141-1149.

Gardiner-Garden, M., and Frommer, M. (1987). CpG islands in vertebrate genomes. *J Mol Biol* 196, 261-282.

Gleize, V., Alentorn, A., Connen de Kerillis, L., Labussiere, M., Nadaradjane, A., Mundwiller, E., Ottolenghi, C., Mangesius, S., Rahimian, A., Ducray, F., *et al.* (2015). CIC inactivating mutations identify aggressive subset of 1p19q codeleted gliomas. *Annals of neurology*.

Golbeck, J., and Mutton, P. (2005). Spring-Embedded graphs for semantic visualization. *Visualizing the Semantic Web*, 172-182.

Gonzalez-Angulo, A.M., Hennessy, B.T., Meric-Bernstam, F., Sahin, A., Liu, W., Ju, Z., Carey, M.S., Myhre, S., Speers, C., Deng, L., *et al.* (2011). Functional proteomics can define prognosis and predict pathologic complete response in patients with breast cancer. *Clin Proteomics* 8, 11.



Guintivano, J., Aryee, M.J., and Kaminsky, Z.A. (2013). A cell epigenotype specific model for the correction of brain cellular heterogeneity bias and its application to age, brain region and major depression. *Epigenetics* 8, 290-302.

Harrell, F.E., Jr., Califf, R.M., Pryor, D.B., Lee, K.L., and Rosati, R.A. (1982). Evaluating the yield of medical tests. *JAMA* 247, 2543-2546.

Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y.C., Laslo, P., Cheng, J.X., Murre, C., Singh, H., and Glass, C.K. (2010). Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol Cell* 38, 576-589.

Hennessy, B.T., Lu, Y., Gonzalez-Angulo, A.M., Carey, M.S., Myhre, S., Ju, Z., Davies, M.A., Liu, W., Coombes, K., Meric-Bernstam, F., *et al.* (2010). A Technical Assessment of the Utility of Reverse Phase Protein Arrays for the Study of the Functional Proteome in Non-microdissected Human Breast Cancers. *Clin Proteomics* 6, 129-151.

Hennessy, B.T., Lu, Y.L., Poradosu, E., Yu, Q.H., Yu, S.X., Hall, H., Carey, M.S., Ravoori, M., Gonzalez-Angulo, A.M., Birch, R., *et al.* (2007). Pharmacodynamic markers of perifosine efficacy. *Clin Cancer Res* 13, 7421-7431.

Hu, J., He, X., Baggerly, K.A., Coombes, K.R., Hennessy, B.T., and Mills, G.B. (2007). Non-parametric quantification of protein lysate arrays. *Bioinformatics* 23, 1986-1994.

Huang, D.W., Sherman, B.T., and Lempicki, R.A. (2009). Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat Protoc* 4, 44-57.

Hung, N., Chen, Y.J., Taha, A., Olivecrona, M., Boet, R., Wiles, A., Warr, T., Shaw, A., Eiholzer, R., Baguley, B.C., *et al.* (2014). Increased paired box transcription factor 8 has a survival function in glioma. *BMC cancer* 14, 159.

Imielinski, M., Berger, A.H., Hammerman, P.S., Hernandez, B., Pugh, T.J., Hodis, E., Cho, J., Suh, J., Capelletti, M., Sivachenko, A., *et al.* (2012). Mapping the hallmarks of lung adenocarcinoma with massively parallel sequencing. *Cell* 150, 1107-1120.

Johnson, W.E., Li, C., and Rabinovic, A. (2007). Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* 8, 118-127.

Kent, W.J. (2002). BLAT--the BLAST-like alignment tool. *Genome Res* 12, 656-664.

Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and Haussler, D. (2002). The human genome browser at UCSC. *Genome Res* 12, 996-1006.

Koboldt, D.C., Zhang, Q., Larson, D.E., Shen, D., McLellan, M.D., Lin, L., Miller, C.A., Mardis, E.R., Ding, L., and Wilson, R.K. (2012). VarScan 2: somatic mutation and copy number alteration discovery in cancer by exome sequencing. *Genome research* 22, 568-576.

Lambert, S.R., Witt, H., Hovestadt, V., Zucknick, M., Kool, M., Pearson, D.M., Korshunov, A., Ryzhova, M., Ichimura, K., Jabado, N., *et al.* (2013). Differential expression and methylation of brain developmental genes define location-specific subsets of pilocytic astrocytoma. *Acta neuropathologica* 126, 291-301.

Lawrence, M., Huber, W., Pages, H., Aboyoun, P., Carlson, M., Gentleman, R., Morgan, M.T., and Carey, V.J. (2013a). Software for computing and annotating genomic ranges. *PLoS Comput Biol* 9, e1003118.

Lawrence, M.S., Stojanov, P., Polak, P., Kryukov, G.V., Cibulskis, K., Sivachenko, A., Carter, S.L., Stewart, C., Mermel, C.H., Roberts, S.A., *et al.* (2013b). Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* 499, 214-218.

Liang, J., Shao, S.H., Xu, Z.X., Hennessy, B., Ding, Z., Larrea, M., Kondo, S., Dumont, D.J., Gutterman, J.U., Walker, C.L., *et al.* (2007). The energy sensing LKB1-AMPK pathway regulates p27(kip1) phosphorylation mediating the decision to enter autophagy or apoptosis. *Nat Cell Biol* 9, 218-224.

McCarroll, S.A., Kuruvilla, F.G., Korn, J.M., Cawley, S., Nemesh, J., Wysoker, A., Shapero, M.H., de Bakker, P.I., Maller, J.B., Kirby, A., *et al.* (2008). Integrated detection and population-genetic analysis of SNPs and copy number variation. *Nature genetics* 40, 1166-1174.



McPherson, A., Hormozdiari, F., Zayed, A., Giuliany, R., Ha, G., Sun, M.G., Griffith, M., Heravi Moussavi, A., Senz, J., Melnyk, N., *et al.* (2011). deFuse: an algorithm for gene fusion discovery in tumor RNA-Seq data. *PLoS computational biology* 7, e1001138.

Mermel, C.H., Schumacher, S.E., Hill, B., Meyerson, M.L., Beroukhi, R., and Getz, G. (2011). GISTIC2.0 facilitates sensitive and confident localization of the targets of focal somatic copy-number alteration in human cancers. *Genome biology* 12, R41.

Mur, P., Mollejo, M., Ruano, Y., de Lope, A.R., Fiano, C., Garcia, J.F., Castresana, J.S., Hernandez-Lain, A., Rey, J.A., and Melendez, B. (2013). Codeletion of 1p and 19q determines distinct gene methylation and expression profiles in IDH-mutated oligodendroglial tumors. *Acta neuropathologica* 126, 277-289.

Neeley, E.S., Kornblau, S.M., Coombes, K.R., and Baggerly, K.A. (2009). Variable slope normalization of reverse phase protein arrays. *Bioinformatics* 25, 1384-1389.

Radenbaugh, A.J., Ma, S., Ewing, A., Stuart, J.M., Collisson, E.A., Zhu, J., and Haussler, D. (2014). RADIA: RNA and DNA integrated analysis for somatic mutation detection. *PloS one* 9, e111516.

Ramos, A.H., Lichtenstein, L., Gupta, M., Lawrence, M.S., Pugh, T.J., Saksena, G., Meyerson, M., and Getz, G. (2015). Oncotator: cancer variant annotation tool. *Hum Mutat* 36, E2423-2429.

Robinson, M.D., McCarthy, D.J., and Smyth, G.K. (2010). edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 26, 139-140.

Smyth, G.K. (2005). Limma: linear models for microarray data. *Bioinformatics and computational biology solutions using R and Bioconductor*, 397-420.

Spearman, C. (1904). The proof and measurement of association between two things. *Am J Psychol* 15, 72-101.

Sturm, D., Witt, H., Hovestadt, V., Khuong-Quang, D.A., Jones, D.T., Konermann, C., Pfaff, E., Tonjes, M., Sill, M., Bender, S., *et al.* (2012). Hotspot mutations in H3F3A and IDH1 define distinct epigenetic and biological subgroups of glioblastoma. *Cancer cell* 22, 425-437.

Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S., *et al.* (2005). Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* 102, 15545-15550.

Takai, D., and Jones, P.A. (2002). Comprehensive analysis of CpG islands in human chromosomes 21 and 22. *Proc Natl Acad Sci U S A* 99, 3740-3745.

TCGA\_Network (2015). Comprehensive, Integrative Genomic Analysis of Diffuse Lower Grade Gliomas. *New England Journal of Medicine*, in press.

Therneau, T.M. (2014). A package for survival analysis in S. version 2.37-7. <http://CRAN.R-project.org/package=survival>.

Therneau, T.M., and Grambsch, P.M. (2000). Modeling survival data : extending the Cox model (New York: Springer).

Tibes, R., Qiu, Y., Lu, Y., Hennessy, B., Andreeff, M., Mills, G.B., and Kornblau, S.M. (2006). Reverse phase protein array: validation of a novel proteomic technology and utility for analysis of primary leukemia specimens and hematopoietic stem cells. *Mol Cancer Ther* 5, 2512-2521.

Torres-Garcia, W., Zheng, S., Sivachenko, A., Vegesna, R., Wang, Q., Yao, R., Berger, M.F., Weinstein, J.N., Getz, G., and Verhaak, R.G. (2014). PRADA: pipeline for RNA sequencing data analysis. *Bioinformatics* 30, 2224-2226.

Trifonov, V., Pasqualucci, L., Dalla Favera, R., and Rabadan, R. (2013). MutComFocal: an integrative approach to identifying recurrent and focal genomic alterations in tumor samples. *BMC Syst Biol* 7, 25.

Turcan, S., Rohle, D., Goenka, A., Walsh, L.A., Fang, F., Yilmaz, E., Campos, C., Fabius, A.W., Lu, C., Ward, P.S., *et al.* (2012). IDH1 mutation is sufficient to establish the glioma hypermethylator phenotype. *Nature* 483, 479-483.



Vaske, C.J., Benz, S.C., Sanborn, J.Z., Earl, D., Szeto, C., Zhu, J.C., Haussler, D., and Stuart, J.M. (2010). Inference of patient-specific pathway activities from multi-dimensional cancer genomics data using PARADIGM. *Bioinformatics* 26, i237-i245.

Verhaak, R.G., Hoadley, K.A., Purdom, E., Wang, V., Qi, Y., Wilkerson, M.D., Miller, C.R., Ding, L., Golub, T., Mesirov, J.P., *et al.* (2010). Integrated genomic analysis identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in PDGFRA, IDH1, EGFR, and NF1. *Cancer cell* 17, 98-110.

Yoshihara, K., Wang, Q., Torres-Garcia, W., Zheng, S., Vegesna, R., Kim, H., and Verhaak, R.G. (2014). The landscape and therapeutic relevance of cancer-associated transcript fusions. *Oncogene*.

