

1 **Chinese Glioma Genome Atlas (CGGA): A Comprehensive Resource with**  
2 **Functional Genomic Data for Chinese Glioma Patients**

3  
4 Zheng Zhao<sup>1,#,a</sup>, Ke'nan Zhang<sup>1,#,b</sup>, Qiangwei Wang<sup>1,#,c</sup>, Guanzhang Li<sup>1,d</sup>, Fan Zeng<sup>1,e</sup>, Ying Zhang<sup>1,f</sup>,  
5 Fan Wu<sup>1,g</sup>, Ruichao Chai<sup>1,h</sup>, Zheng Wang<sup>2,i</sup>, Chuanbao Zhang<sup>2,j</sup>, Wei Zhang<sup>2,k</sup>, Zhaoshi Bao<sup>2,\*1</sup>, Tao  
6 Jiang<sup>1,2,3,4,\*m</sup>

7  
8 <sup>1</sup>Beijing Neurosurgical Institute, Capital Medical University, Beijing 100070, China

9 <sup>2</sup>Department of Neurosurgery, Beijing Tiantan Hospital, Capital Medical University, Beijing  
10 100070, China

11 <sup>3</sup>Center of Brain Tumour, Beijing Institute for Brain Disorders, Beijing 100069, China

12 <sup>4</sup>China National Clinical Research Center for Neurological Diseases, Beijing 100070, China

13 #Equal contribution

14 \*Corresponding authors:

15 E-mail: [taojiang1964@163.com](mailto:taojiang1964@163.com) (Jiang T), [bzsjoel985@163.com](mailto:bzsjoel985@163.com) (Bao ZS)

16

17 **Running title:** *Zhao et al / Chinese Glioma Genome Atlas*

18

19 <sup>a</sup>ORCID: 0000-0001-8945-9632

20 <sup>b</sup>ORCID: 0000-0001-7270-569X

21 <sup>c</sup>ORCID: 0000-0002-7308-049X

22 <sup>d</sup>ORCID: 0000-0002-0353-5751

23 <sup>e</sup>ORCID: 0000-0001-5351-2155

24 <sup>f</sup>ORCID: 0000-0002-7613-6188

25 <sup>g</sup>ORCID: 0000-0001-9256-0176

26 <sup>h</sup>ORCID: 0000-0003-3451-8871

27 <sup>i</sup>ORCID: 0000-0003-1687-6990

28 <sup>j</sup>ORCID: 0000-0003-2615-4190

29 <sup>k</sup>ORCID: 0000-0001-7800-3189

- 1 <sup>l</sup>ORCID: 0000-0003-4922-4470
- 2 <sup>m</sup>ORCID: 0000-0002-7008-6351
- 3
- 4 Total word counts: 3,030
- 5 Total references: 31
- 6 Total figures: 1
- 7 Total tables: 1
- 8 Total supplementary figures: 0
- 9 Total supplementary tables: 0
- 10 Total supplementary files: 0
- 11

1 **Abstract:**

2 Gliomas are the most common and malignant intracranial tumours in adults. Recent studies have  
3 shown that functional genomics greatly aids in the understanding of the pathophysiology and  
4 therapy of glioma. However, comprehensive genomic data and analysis platforms are relatively  
5 limited. In this study, we developed the Chinese Glioma Genome Atlas (CGGA,  
6 <http://www.cgga.org.cn>), a user-friendly data portal for storage and interactive exploration of multi-  
7 dimensional functional genomic data that includes nearly 2,000 primary and recurrent glioma  
8 samples from Chinese cohorts. CGGA currently provides access to whole-exome sequencing (286  
9 samples), messenger RNA sequencing (1,018 samples) and microarray (301 samples), DNA  
10 methylation microarray (159 samples), and microRNA microarray (198 samples) data, as well as  
11 detailed clinical data (e.g., WHO grade, histological type, critical molecular genetic information,  
12 age, sex, chemoradiotherapy status and survival data). In addition, we developed an analysis tool to  
13 allow users to browse mutational, mRNA/microRNA expression, and DNA methylation profiles and  
14 perform survival and correlation analyses of specific glioma subtypes. CGGA greatly reduces the  
15 barriers between complex functional genomic data and glioma researchers who seek rapid, intuitive,  
16 and high-quality access to data resources and enables researchers to use these immeasurable data  
17 sources for biological research and clinical application. Importantly, the free provision of data will  
18 allow researchers to quickly generate and provide data to the research community.

19

20 **KEYWORDS:** Glioma; Functional genomics; Chinese Glioma Genome Atlas; Database

21

## 1 **Introduction**

2 Gliomas are the most frequent malignant tumours of the adult brain. According to a multi-centre  
3 cross-sectional study on brain tumours in China, the prevalence of primary brain tumours in all  
4 populations is approximately 22.52 per 100,000 persons, with gliomas accounting for 31.1% of the  
5 population aged 20–59 years [1-3]. According to the histopathological classification of the 2016  
6 World Health Organization (WHO) grading system, glioma is classified from grade II to grade IV  
7 by both histological characteristics and several new molecular pathological features, such as *IDH*  
8 mutation status and chromosome 1p/19q co-deletion status [4]. Despite advances in current  
9 treatment standards, the survival rate of patients with glioma has not changed in decades, especially  
10 for aggressive gliomas (with a poor median survival time of only 12 to 14 months) [5, 6]. In addition,  
11 most lower-grade gliomas (grade II and III, LGG) will progress to glioblastoma (grade IV, GBM)  
12 in less than 10 years [4, 7, 8]. At present, the reasons for glioma recurrence or malignant progression  
13 may be as follows: 1) infiltrative tumour cells cannot be completely removed by neurosurgical  
14 resection [9, 10]; 2) retained tumour cells cannot be effectively suppressed by limited postoperative  
15 treatment options [3, 11, 12]; 3) multiple lesions may develop [13, 14]; 4) cell cloning is rapid under  
16 chemotherapy and/or radiotherapy [7, 15]; 5) the adaptive tumour microenvironment permits  
17 tumour cells [16, 17]; and 6) limited data resources lead to limited research. Therefore, it is essential  
18 to collect clinical specimens and generate genomic data for the glioma research community.

19         Recent high-throughput technologies have enabled extensive characterization of genomic  
20 status, including but not limited to DNA methylation modification, genetic alteration, and gene  
21 expression regulation. In the cancer research community, major large-scale projects, such as The  
22 Cancer Genome Atlas (TCGA, including 516 LGGs and 617 GBMs before Oct. 18, 2019) [18] and  
23 the International Cancer Genome Consortium (ICGC, excluding TCGA samples, including 80 adult  
24 GBMs and 50 paediatric GBMs before April. 3, 2019) [19, 20], have generated an unparalleled  
25 amount of functional genomic data. These projects have begun to transform our understanding of  
26 cancer and even lead to improvements in our ability to diagnose, treat, and prevent human cancers.  
27 Importantly, they have provided an opportunity to make and validate important discoveries for  
28 cancer genomic researchers around the globe. However, the data resources generated by these  
29 projects are often not easy to access directly, analyse or visualize, especially for researchers with no

1 bioinformatics skills, thus preventing the translation of functional genomics results into novel  
2 findings of biological significance for drug development and clinical treatment. Although several  
3 webservers, such as cBioportal [21, 22] and GlioVis [23], have been built to integrate analysed  
4 multi-dimensional glioma data, they have ignored the presence of cancer heterogeneity in gliomas,  
5 which cannot be examined in specific subtypes and is rarely found in recurrent glioma samples.

6 Here, we introduce the CGGA (Chinese Glioma Genome Atlas, <http://www.cgga.org.cn>)  
7 database, which is an open-access and easy-to-use platform for interactive exploration of multi-  
8 dimensional functional genomic datasets for nearly 2,000 primary and recurrent glioma samples  
9 from Chinese cohorts. CGGA currently contains whole-exome sequencing (286 samples),  
10 messenger RNA (mRNA) sequencing (1,018 samples), microarray (301 samples), DNA methylation  
11 microarray (159 samples), microRNA microarray (198 samples) and comprehensive clinical data.  
12 We also developed an analysis module to allow users to browse the mutational landscape profile,  
13 mRNA/microRNA expression profile and DNA methylation profile as well as to perform survival  
14 and correlation analyses for specific glioma subtypes. We believe that this website will greatly  
15 reduce the barriers between complex functional genomic data and glioma researchers who seek  
16 rapid, intuitive, and high-quality access to data resources.

## 17 **Results**

### 18 **Database content and usage**

19 The CGGA database was designed to store functional genomic data and to allow interactive  
20 exploration of multi-dimensional datasets from primary and recurrent gliomas in Chinese cohorts;  
21 it is available at <http://www.cgga.org.cn/>. Currently, CGGA contains whole-exome sequencing data  
22 (286 samples), messenger RNA sequencing data (total: 1,018 samples, batch 1 with 693 samples  
23 and batch 2 with 325 samples), microarray data (301 samples), DNA methylation microarray data  
24 (159 samples), and microRNA microarray data (198 samples) for glioma. The database also contains  
25 detailed clinical data (including WHO grade and histological type, critical molecular genetic  
26 information, age, sex, chemoradiotherapy status and survival data). Detailed statistical information  
27 for each dataset is provided in **Table 1**. We organized the web interface of CGGA according to the  
28 three main functional features: (i) Home, (ii) Analyse, and (iii) Download. In the following context,  
29 we provide an example for using CGGA.

## 1 **The home page**

2 On the ‘Home’ page, CGGA provides a statistical table for a glioma dataset, including the dataset  
3 name, data type, number of samples in each subgroup, clinical data and analysis purposes. For  
4 instance, we performed messenger RNA sequencing on 1,018 glioma samples included in two  
5 datasets (693 samples in batch 1 and 325 samples for batch 2, including 282 primary LGGs, 161  
6 recurrent LGGs, 140 primary GBMs and 109 recurrent GBMs in batch 1 and 144 primary LGGs,  
7 38 recurrent LGGs, 85 primary GBMs, 24 recurrent GBMs and 30 secondary GBMs in batch 2). To  
8 the best of our knowledge, CGGA is the first database to store the functional genomic data for both  
9 LGG and GBM recurrent gliomas. In addition, users can obtain a visualized result for the analysis  
10 of each dataset for a specific glioma subtype by clicking on a hyperlink on the ‘Home’ page. The  
11 ‘Download’ and ‘Help’ pages can also be accessed directly from the ‘Home’ page.

## 12 **Overall analyses and results**

13 To facilitate analysis of the CGGA data by researchers, we developed four online modules in the  
14 ‘Analyse’ tab, including ‘WESeq data’, ‘mRNA data’, ‘methylation data’, and ‘microRNA data’, to  
15 analyse whole-exome, mRNA expression, DNA methylation and microRNA expression data,  
16 respectively (**Figure 1A**). A key feature of CGGA is that it is easy to use. In the context below, we  
17 demonstrate the use of the ‘Analyse’ tab in CGGA.

18 On the ‘WESeq data’ page, users are allowed to visualize the mutational profile of a gene set  
19 of interest and survival analysis of a specific gene of interest in a specific glioma subtype. In the  
20 ‘Oncoprint’ section, users are guided to a) input a gene set of interest (*IDHI TP53 ATRX* for  
21 example), and b) select a dataset of interest (‘All’ for example). Based on user input, this tool  
22 automatically generates visualized results. In this result, each case or patient is represented as  
23 columns, each gene is displayed as rows, and a colour map on the bottom is used to depict specific  
24 clinical information (**Figure 1B**). This ‘Oncoprint’ can be very useful for visualizing the mutational  
25 profile for a gene set of interest in a specific glioma subtype and for intuitively validating trends  
26 such as mutational frequency and mutual exclusivity or co-occurrence for a gene pair. In the above  
27 example, mutations in the *IDHI* (47%), *TP53* (46%) and *ATRX* (30%) genes were the most common  
28 mutations in all gliomas. In the ‘Survival’ section, users are allowed to a) input a specific gene of  
29 interest (*IDHI* for example), and b) select a dataset of interest (‘Primary LGG’ for example) to  
30 investigate the association of the mutation with severe functional consequences. Consistent with

1 previous studies [24], primary LGG cases with *IDHI* mutations have a better overall survival than  
2 do cases with *IDHI* wild-type tumours ( $p < 0.0001$ , **Figure 1C**, left). These analysis results from  
3 the ‘WEseq data’ section can be exported as a PDF file. For the sake of reproducibility, we provide  
4 the analysis data (**Figure 1C**, middle) and R code (**Figure 1C**, right), which allow users to reproduce  
5 the figure to be able to modify or adapt each figure according to each researcher’s demands.

6 On the ‘mRNA data’ page, users are allowed to perform gene expression distribution,  
7 correlation and survival analyses for a specific gene of interest in a specific glioma subtype. Three  
8 mRNA datasets are available for users, including two batch RNA-seq datasets (batch 1: 693 samples;  
9 batch 2: 325 samples) and one microarray dataset (301 samples). In the ‘Distribution’ section, users  
10 can display one gene distribution pattern for each glioma subtype by selecting a dataset  
11 (‘mRNAseq\_325’ for example) and inputting a gene name of interest (‘*ADAMTSL4*’ for example).  
12 The results show the gene expression pattern in each glioma subtype classified by clinical  
13 information. Similar to our previous studies [25], the *ADAMTSL4* gene was shown to be  
14 differentially expressed according to the WHO 2016 classification based on the *IDH* mutation  
15 and/or 1p/19q co-deletion status (**Figure 1D**, left). Moreover, a critical feature of the CGGA dataset  
16 is the inclusion of recurrent gliomas. This module allows users to infer whether a gene may be a  
17 candidate factor that drives malignant progression if it is differentially expressed in primary and  
18 recurrent gliomas. In the ‘Correlation’ section, the user is allowed to validate the co-expression  
19 pattern by selecting a dataset (‘mRNAseq\_325’ for example) and inputting a gene pair  
20 (‘*ADAMTSL4*’ and ‘*CD274*’ for example). As a result, the co-expression patterns in each glioma  
21 subtype will be displayed with the results of Pearson’s test and the p value (**Figure 1D**, middle). In  
22 the ‘Survival’ section, users can perform survival analysis based on gene expression by selecting a  
23 dataset (‘mRNAseq\_325’ for example) and inputting a gene of interest (‘*ADAMTSL4*’ for example).  
24 All primary glioma patients with low *ADAMTSL4* expression showed better overall survival than  
25 did those with high *ADAMTSL4* expression ( $p < 0.0001$ , **Figure 1D**, right). The above results from  
26 the ‘mRNA data’ section are consistent with our previous study [25]. Similar to the ‘mRNA data’  
27 page, users can also display the methylation/microRNA distribution and perform correlation and  
28 survival analyses on the ‘methylation data’ page and the ‘microRNA data’ page, respectively.

## 29 **Data acquisition**

1 All the data sets in CGGA can be downloaded on the ‘Download’ page by both the community and  
2 researchers. Each data type is saved at the gene and/or probe level and is then combined with the  
3 available clinical data, including basic clinical information, survival and therapy information.

#### 4 **Perspectives and concluding remarks**

5 The current version of the CGGA is the first release of our database, and it incorporates multi-  
6 dimensional functional genomic glioma data, including whole-exome sequencing, mRNA and  
7 microRNA expression, and DNA methylation data for nearly 2,000 samples from Chinese cohorts.  
8 Considering the importance of these data for glioma research, CGGA is publicly available. To the  
9 best of our knowledge, CGGA is the first database to store the functional genomic data for both  
10 recurrent LGGs and GBMs. In addition, CGGA provides several tools that allow users to analyse  
11 these datasets, including mutational profile, distribution pattern, correlation and survival analysis  
12 tools. These tools will be useful for users to generate or validate findings of novel biological  
13 significance.

14 We anticipate several future directions for our CGGA database. First, through the Beijing  
15 Neurosurgical Institute, Beijing Tiantan Hospital and Chinese Glioma Cooperative Group (CGCG)  
16 Research Network, we will continue to collect glioma samples and perform multiple ‘Omics’  
17 sequencing/microarray analyses, and we will continue to update this database regularly in the future.  
18 Second, we also plan to add image-genomic data that match the ‘Omics’ data in CGGA. Third, we  
19 will develop more advanced features, including data for other ‘Omics’ analyses, search functions  
20 for clinical information on a patient of interest, and further extensions for the data analysis tools. In  
21 summary, CGGA facilitates access to functional genomic data for Chinese cohorts for the entire  
22 glioma community. It provides an easy-to-use, user-friendly interface for obtaining integrated data  
23 sets, performing intuitive visualized analysis, and downloading these datasets. CGGA greatly  
24 reduces the barriers between complex functional genomic data and glioma researchers, which  
25 empowers researchers to use functional genomic data into important biological insights and  
26 potential clinical applications.

#### 27 **Materials and methods**

##### 28 **Clinical specimen collection**



1 Glioma tissues, corresponding genomic data and patient follow-up information were obtained from  
2 Beijing Tiantan Hospital at Capital Medical University, Tianjin Medical University General  
3 Hospital, Sanbo Brain Hospital at Capital Medical University, the Second Affiliated Hospital of  
4 Harbin Medical University, the First Affiliated Hospital of Nanjing Medical University, and the  
5 First Hospital of China Medical University. All research performed was approved by the Beijing  
6 Tiantan Hospital Capital Medical University Institutional Review Board (IRB) and was conducted  
7 according to the principles of the Helsinki Declaration. According to the central pathology reviews  
8 of independent committee certified neuropathologists, all the subjects were consistently diagnosed  
9 with glioma and further classified according to the 2007/2016 WHO classification system. All  
10 patients provided written informed consent. The specimens were collected under IRB KY2013-017-  
11 01 and frozen in liquid nitrogen within 5 min of resection.

#### 12 **Data processing for whole-exome sequencing data**

13 Genomic DNA from tumours and the matched blood samples was extracted, and high integrity was  
14 confirmed by 1% agarose gel electrophoresis. The DNA was subsequently fragmented and quality-  
15 controlled, and paired-end libraries were prepared. Agilent SureSelect kit v5.4 was used for target  
16 capture. Sequencing was performed using the Illumina HiSeq 4000 platform with a paired-end  
17 sequencing strategy. Valid DNA sequencing data were mapped to the reference human genome  
18 (UCSC hg19) using Burrows-Wheeler Aligner (v0.7.12-r1039, bwa mem) [26] with default  
19 parameters. SAMtools (version 1.2) [27] and Picard (version 2.0.1, Broad Institute) were then used  
20 to sort the reads by coordinates and mark duplicates. Statistics such as sequencing depth and  
21 coverage were calculated based on the resulting BAM files. SAVI2 was used to identify somatic  
22 mutations (including single-nucleotide variations and short insertions/deletions) as previously  
23 described [7, 8]. Briefly, in this pipeline, SAMtools mpileup and bcftools (version 0.1.19) [28] were  
24 employed to perform variant calling, and the preliminary variant list was filtered to remove positions  
25 with no sufficient sequencing depth, positions with only low-quality reads, and positions biased  
26 toward either strand. Somatic mutations were identified and evaluated by an empirical Bayesian  
27 method. In particular, mutations with a significantly higher mutation allele frequency in tumours  
28 than in normal controls were selected.

#### 29 **Data processing for mRNA sequencing data**

1 Prior to library preparation, total RNA was isolated using RNeasy Mini Kit (Qiagen) according to  
2 the manufacturer's instructions. A pestle and QIAshredder (Qiagen) were used to disrupt and  
3 homogenize frozen tissue. The RNA integrity was checked using a 2,100 Bioanalyzer (Agilent  
4 Technologies), and only high-quality samples with an RNA integrity number (RIN) value greater  
5 than or equal to 6.8 were used to construct the sequencing library. Typically, 1 µg of total RNA was  
6 used with the TruSeq RNA library preparation kit (Illumina) in accordance with the low-throughput  
7 protocol, except that SuperScript III reverse transcriptase (Invitrogen) was used to synthesize first-  
8 strand cDNA. After PCR enrichment and purification of adapter-ligated fragments, the  
9 concentration of DNA with adapters was determined by quantitative PCR (Applied Biosystems  
10 7,500) using primers QP1 5'-AATGATACGGCGACCACCGA-3' and QP2 5'-  
11 CAAGCAGAAGACGGCATAACGAGA-3'. The length of the DNA fragment was measured using a  
12 2,100 Bioanalyzer, with median insert sizes of 200 nucleotides. The RNA-seq libraries were  
13 sequenced using the Illumina HiSeq 2,000, 2,500 or 4,000 Sequencing System. The libraries were  
14 prepared using the paired-end strategy with read lengths of 101 bp, 125 bp or 150 bp. Base calling  
15 was performed by the Illumina CASAVA v1.8.2 pipeline. RNA-seq mapping and quantification were  
16 processed by using STAR (version v2.5.2b) [29] and RSEM (version 1.2.31) software [30]. Briefly,  
17 reads were aligned to the human genome reference (GENCODE v19, hg19) with STAR, and then  
18 sequencing read counts for each GENCODE gene were calculated using RSEM. The expression  
19 levels of different samples were merged into an FPKM (fragments per kilobase transcriptome per  
20 million fragments) matrix. We defined a gene as expressed only if its expression level was greater  
21 than 0 in half of the samples. Finally, we retained only expressed genes in the mRNA expression  
22 profile.

### 23 **Data processing for mRNA microarray data**

24 A rapid haematoxylin & eosin stain for frozen sections was performed on each sample to assess the  
25 tumour cell proportion before RNA extraction. RNA was extracted from only samples with >80%  
26 tumour cells. Total RNA was extracted from frozen tumour tissue with the mirVana miRNA Isolation  
27 Kit (Ambion), as described previously [31]. A NanoDrop ND-1000 spectrophotometer (NanoDrop  
28 Technologies) was used to evaluate the quality and concentration of extracted total RNA and an  
29 Agilent 2100 Bioanalyzer (Agilent) to assess the integrity. The qualified RNA was collected for  
30 further processing. cDNA and biotinylated cRNA were synthesized and hybridized to Agilent Whole

1 Human Genome Array according to the manufacturer's instructions. Finally, the array-generated  
2 data were analyzed by the Agilent G2565BA Microarray Scanner System and Agilent Feature  
3 Extraction Software (Version 9.1). GeneSpring GX11.0 was applied to calculate the probe intensity.

#### 4 **Data processing for methylation microarray data**

5 A haematoxylin and eosin-stained frozen section was prepared for assessment of the percentage of  
6 tumour cells before RNA extraction. Only samples with greater than 80% tumour cells were selected.  
7 Genomic DNA was isolated from frozen tumour tissues using the QIAamp DNA Mini Kit (Qiagen)  
8 according to the manufacturer's protocol. The DNA concentration and quality were assessed using  
9 a NanoDrop ND-1000 spectrophotometer (NanoDrop Technologies, Houston, TX). The microarray  
10 analysis was performed using Illumina Infinium HumanMethylation27 Bead-Chips (Illumina Inc.),  
11 which contains 27,578 highly informative CpG sites covering more than 14,000 human RefSeq  
12 genes. This allows researchers to investigate all sites per sample at a single-nucleotide resolution.  
13 Bisulfite modification of DNA, chip processing and data analysis were performed following the  
14 manufacturer's manual at Wellcome Trust Centre for Human Genetics Genomics Lab, Oxford, UK.  
15 The array results were examined with the BeadStudio software (Illumina).

#### 16 **Data processing for microRNA microarray data**

17 Total RNA (tRNA) was extracted from frozen tissues by using the mirVana miRNA Isolation Kit  
18 (Ambion, Inc., Austin, Tex), and the concentration and quality were determined with a NanoDrop  
19 ND-1000 spectrophotometer (NanoDrop Technologies, Wilmington, Del). microRNA expression  
20 profiling was performed using the human v2.0 microRNA Expression BeadChip (Illumina, Inc.,  
21 San Diego, Calif) with 1146 microRNAs covering 97% of the miRBase 12.0 database according to  
22 the manufacturer's instructions.

#### 23 **Implementation**

24 In CGGA, all data are organized with MySQL 14.14 based on relational schema, which will be  
25 supported by future CGGA updates. The website code was written based on Java Server Pages using  
26 the Java Servlet framework. The website is deployed on the Tomcat 6.0.44 web server and runs on  
27 a CentOS 5.5 Linux system. JQuery was used to generate, render and manipulate data visualization.  
28 The 'Analyse' module was realized with Perl and R scripts. The CGGA website has been fully tested  
29 in Google Chrome and Safari browsers.

## 1 **Data Availability**

2 All data for this article can be found online at <http://www.cgga.org.cn>.

## 3 **Authors' contributions**

4 TJ, ZB, WZ, and ZZ conceived and supervised this study. ZZ, KZ and QW designed the research.  
5 ZZ, GL, FZ, YZ, FW, RC, ZW, CZ performed data analysis. ZZ developed CGGA web server. TJ,  
6 ZB, and ZZ wrote the manuscript. All authors read and approved the final manuscript.

## 7 **Competing interests**

8 The authors have declared no competing interests.

## 9 **Acknowledgments**

10 This work was supported by the National Natural Science Foundation of China (NSFC) fund (Nos.  
11 81702460 and 81802994).

## 12 **References**

- 13 [1] Jiang T, Tang GF, Lin Y, Peng XX, Zhang X, Zhai XW, et al. Prevalence estimates for primary  
14 brain tumors in China: a multi-center cross-sectional study. *Chin Med J (Engl)* 2011;124:2578-83.
- 15 [2] Zhao Z, Meng F, Wang W, Wang Z, Zhang C, Jiang T. Comprehensive RNA-seq transcriptomic  
16 profiling in the malignant progression of gliomas. *Sci Data* 2017;4:170024.
- 17 [3] Jiang T, Mao Y, Ma W, Mao Q, You Y, Yang X, et al. CGCG clinical practice guidelines for the  
18 management of adult diffuse gliomas. *Cancer Lett* 2016;375:263-73.
- 19 [4] Louis DN, Perry A, Reifenberger G, von Deimling A, Figarella-Branger D, Cavenee WK, et al.  
20 The 2016 World Health Organization Classification of Tumors of the Central Nervous System: a  
21 summary. *Acta Neuropathol* 2016;131:803-20.
- 22 [5] Stupp R, Mason WP, van den Bent MJ, Weller M, Fisher B, Taphoorn MJ, et al. Radiotherapy  
23 plus concomitant and adjuvant temozolomide for glioblastoma. *N Engl J Med* 2005;352:987-96.
- 24 [6] Van Meir EG, Hadjipanayis CG, Norden AD, Shu HK, Wen PY, Olson JJ. Exciting new advances  
25 in neuro-oncology: the avenue to a cure for malignant glioma. *CA Cancer J Clin* 2010;60:166-93.
- 26 [7] Wang J, Cazzato E, Ladewig E, Frattini V, Rosenbloom DI, Zairis S, et al. Clonal evolution of  
27 glioblastoma under therapy. *Nat Genet* 2016;48:768-76.
- 28 [8] Hu H, Mu Q, Bao Z, Chen Y, Liu Y, Chen J, et al. Mutational Landscape of Secondary  
29 Glioblastoma Guides MET-Targeted Trial in Brain Tumor. *Cell* 2018;175:1665-78 e18.
- 30 [9] Chaichana KL, Jusue-Torres I, Navarro-Ramirez R, Raza SM, Pascual-Gallego M, Ibrahim A,  
31 et al. Establishing percent resection and residual volume thresholds affecting survival and  
32 recurrence for patients with newly diagnosed intracranial glioblastoma. *Neuro Oncol* 2014;16:113-  
33 22.
- 34 [10] Aldape K, Brindle KM, Chesler L, Chopra R, Gajjar A, Gilbert MR, et al. Challenges to curing  
35 primary brain tumours. *Nat Rev Clin Oncol* 2019;16:509-20.
- 36 [11] Yi GZ, Huang G, Guo M, Zhang X, Wang H, Deng S, et al. Acquired temozolomide resistance

- 1 in MGMT-deficient glioblastoma cells is associated with regulation of DNA repair by DHC2. *Brain*  
2 2019;142:2352-66.
- 3 [12] Frosina G. DNA Repair and Resistance of Gliomas to Chemotherapy and Radiotherapy.  
4 *Molecular Cancer Research* 2009;7:989-99.
- 5 [13] Lee JK, Wang J, Sa JK, Ladewig E, Lee HO, Lee IH, et al. Spatiotemporal genomic architecture  
6 informs precision oncology in glioblastoma. *Nat Genet* 2017;49:594-9.
- 7 [14] Liu Q, Liu Y, Li W, Wang X, Sawaya R, Lang FF, et al. Genetic, epigenetic, and molecular  
8 landscapes of multifocal and multicentric glioblastoma. *Acta Neuropathol* 2015;130:587-97.
- 9 [15] Barthel FP, Johnson KC, Varn FS, Moskalik AD, Tanner G, Kocakavuk E, et al. Longitudinal  
10 molecular trajectories of diffuse glioma in adults. *Nature* 2019.
- 11 [16] Wang Q, Hu B, Hu X, Kim H, Squatrito M, Scarpace L, et al. Tumor Evolution of Glioma-  
12 Intrinsic Gene Expression Subtypes Associates with Immunological Changes in the  
13 Microenvironment. *Cancer Cell* 2017;32:42-56 e6.
- 14 [17] Quail DF, Joyce JA. The Microenvironmental Landscape of Brain Tumors. *Cancer Cell*  
15 2017;31:326-41.
- 16 [18] Tomczak K, Czerwinska P, Wiznerowicz M. The Cancer Genome Atlas (TCGA): an  
17 immeasurable source of knowledge. *Contemp Oncol (Pozn)* 2015;19:A68-77.
- 18 [19] Zhang J, Bajari R, Andric D, Gerthoffert F, Lepsa A, Nahal-Bose H, et al. The International  
19 Cancer Genome Consortium Data Portal. *Nat Biotechnol* 2019;37:367-9.
- 20 [20] Zhang J, Baran J, Cros A, Guberman JM, Haider S, Hsu J, et al. International Cancer Genome  
21 Consortium Data Portal--a one-stop shop for cancer genomics data. *Database (Oxford)*  
22 2011;2011:bar026.
- 23 [21] Cerami E, Gao J, Dogrusoz U, Gross BE, Sumer SO, Aksoy BA, et al. The cBio cancer  
24 genomics portal: an open platform for exploring multidimensional cancer genomics data. *Cancer*  
25 *Discov* 2012;2:401-4.
- 26 [22] Gao J, Aksoy BA, Dogrusoz U, Dresdner G, Gross B, Sumer SO, et al. Integrative analysis of  
27 complex cancer genomics and clinical profiles using the cBioPortal. *Sci Signal* 2013;6:pl1.
- 28 [23] Bowman RL, Wang Q, Carro A, Verhaak RG, Squatrito M. GlioVis data portal for visualization  
29 and analysis of brain tumor expression datasets. *Neuro Oncol* 2017;19:139-41.
- 30 [24] Cancer Genome Atlas Research N, Brat DJ, Verhaak RG, Aldape KD, Yung WK, Salama SR,  
31 et al. Comprehensive, Integrative Genomic Analysis of Diffuse Lower-Grade Gliomas. *N Engl J*  
32 *Med* 2015;372:2481-98.
- 33 [25] Zhao Z, Zhang KN, Chai RC, Wang KY, Huang RY, Li GZ, et al. ADAMTSL4, a Secreted  
34 Glycoprotein, Is a Novel Immune-Related Biomarker for Primary Glioblastoma Multiforme. *Dis*  
35 *Markers* 2019;2019:1802620.
- 36 [26] Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
37 *Bioinformatics* 2009;25:1754-60.
- 38 [27] Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence  
39 Alignment/Map format and SAMtools. *Bioinformatics* 2009;25:2078-9.
- 40 [28] Narasimhan V, Danecek P, Scally A, Xue Y, Tyler-Smith C, Durbin R. BCFtools/RoH: a hidden  
41 Markov model approach for detecting autozygosity from next-generation sequencing data.  
42 *Bioinformatics* 2016;32:1749-51.
- 43 [29] Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal  
44 RNA-seq aligner. *Bioinformatics* 2013;29:15-21.

1 [30] Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without  
2 a reference genome. *BMC Bioinformatics* 2011;12:323.

3 [31] Yan W, Zhang W, You G, Zhang J, Han L, Bao Z, et al. Molecular classification of gliomas  
4 based on whole genome gene expression: a systematic report of 225 samples from the Chinese  
5 Glioma Cooperative Group. *Neuro Oncol* 2012;14:1432-40.

6

7

1 **Figures**

2 Figure 1. An overview of the CGGA database.

3 A. The CGGA contains whole-exome sequencing, mRNA and microRNA expression, and DNA  
4 methylation data, clinical data, and several analysis modules; B. The mutation profile in all gliomas  
5 (in the 'WSseq\_286' dataset); C. left: the overall survival of glioma patients with IDH1 mutation  
6 and the wild-type gene from primary LGGs (in the 'WSseq\_286' dataset); middle: the data was used  
7 to generate the plot; right: the R code was used to generate the plot; D. left: the ADAMTSL4 gene  
8 expression distribution in primary gliomas based on 2016 WHO grading system (in the  
9 'mRNAseq\_325' dataset); middle: the gene expression correlation between ADAMTSL4 and  
10 CD274 genes (using 'mRNAseq\_325' dataset); right: the overall survival of glioma patients with  
11 low and high ADAMTSL4 gene expression (in the 'mRNAseq\_325' dataset).

12

13

1 **Table 1.** Clinical and Phenotypical Characteristics of Data Set in CGGA database

	All	Primary LGG	Recurrent LGG	Primary GBM	Recurrent GBM	Secondary GBM
<b>Wtseq_286 dataset</b>						
No. of samples –no. (%)	286	126 (44%)	58 (20%)	54 (19%)	48 (17%)	0 (0%)
Age at diagnosis – yr.						
Mean	42.0±12.3	39.6±10.3	37.3±8.7	50.2±14.7	44.5±13.3	-
Range	10-76	10-69	15-61	19-76	19-69	-
Male sex – no. (%)	168	78 (46%)	35 (21%)	29 (17%)	26 (15%)	-
Therapy						
Radiotherapy only	62	52 (84%)	4 (6%)	4 (6%)	2 (3%)	-
Chemotherapy	13	8 (62%)	2 (15%)	0 (0%)	3 (23%)	-
Cheomoradiotherapy	144	49 (34%)	27 (19%)	42 (29%)	26 (18%)	-
No therapy	23	9 (39%)	8 (35%)	4 (17%)	2 (9%)	-
Unknown	44	8 (18%)	17 (39%)	4 (9%)	15 (34%)	-
Survival – month						
Median (95% CI)	51.0 (37.2-98.1)	117.2 (99.4-)	28.5 (20.9-76.0)	16.5 (10.2-28.7)	14.7 (8.9-)	-
IDH_mut_status						
Mutant	161	88 (55%)	45 (28%)	12 (7%)	16 (10%)	-
Wildtype	125	38 (30%)	13 (10%)	42 (34%)	32 (26%)	-
1p19q_codeletion_status						
Codel	51	28 (55%)	17 (33%)	1 (2%)	5 (10%)	-
Non-codel	139	48 (35%)	33 (24%)	23 (17%)	35 (25%)	-
Unknown	96	50 (52%)	8 (8%)	30 (31%)	8 (8%)	-
<b>RNAseq_1018 dataset</b>						
No. of samples –no. (%)	1,018	426 (42%)	199 (20%)	225 (22%)	133 (13%)	30 (3%)
Age at diagnosis – yr.						
Mean	43.2±12.3	40.2±10.8	40.2±9.6	51.0±12.9	45.0±13.2	38.8±11.4
Range	8-79	10-74	15-64	11-79	14-71	8-57

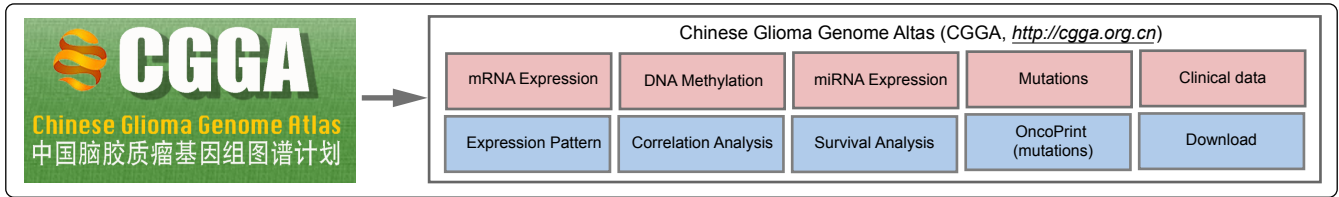


Male sex – no. (%)	601	247 (41%)	115 (19%)	138 (23%)	76 (13%)	21 (3%)
Therapy						
Radiotherapy only	200	128 (64%)	32 (16%)	26 (13%)	10 (5%)	4 (2%)
Chemotherapy	68	30 (44%)	13 (19%)	9 (13%)	11 (16%)	5 (7%)
Cheomoradiotherapy	567	204 (36%)	102 (18%)	159 (28%)	85 (15%)	15 (3%)
No therapy	89	41 (46%)	21 (24%)	18 (20%)	5 (6%)	4 (4%)
Unknown	91	23 (25%)	31 (34%)	13 (14%)	22 (24%)	2 (2%)
Survival – month						
Median (95% CI)	35.0 (30.5-39.9)	108.0 (89.9-)	33.2 (26.1-39.8)	16.1 (13.7-19.7)	9.6 (8.2-11.0)	8.3 (7.1-14.7)
IDH_mut_status						
Mutant	531	289 (54%)	150 (28%)	35 (7%)	34 (6%)	21 (4%)
Wildtype	435	104 (24%)	40 (9%)	183 (42%)	96 (22%)	9 (2%)
Unknown	52	33 (63%)	9 (17%)	7 (13%)	3 (6%)	0
1p19q_codeletion_status						
Codel	212	137 (65%)	54 (25%)	5 (2%)	11 (5%)	4 (2%)
Non-codel	728	254 (35%)	139 (19%)	192 (26%)	118 (16%)	24 (3%)
Unknown	78	35 (45%)	6 (8%)	28 (36%)	4 (5%)	2 (3%)
<b>mRNA-array_301 dataset</b>						
No. of samples –no. (%)	301	156 (52%)	18 (6%)	108 (36%)	5 (2%)	11 (4%)
Age at diagnosis – yr.						
Mean	42.4±11.8	39.6±10.7	38.2±11.2	47.3±12.5	45.6±9.6	38.5±8.6
Range	12-70	17-65	24-62	12-70	36-61	27-51
Male sex – no. (%)	180	93 (52%)	8 (4%)	65 (36%)	2 (1%)	9 (5%)
Therapy						
Radiotherapy only	110	74 (67%)	0	33 (30%)	0	3 (3%)
Chemotherapy	12	1 (8%)	2 (17%)	4 (33%)	3 (25%)	2 (17%)
Cheomoradiotherapy	139	61 (44%)	12 (9%)	60 (43%)	1 (1%)	4 (3%)
No therapy	20	8 (40%)	2 (10%)	6 (30%)	0	2 (10%)
Unknown	20	12 (60%)	2 (10%)	5 (25%)	1 (5%)	0

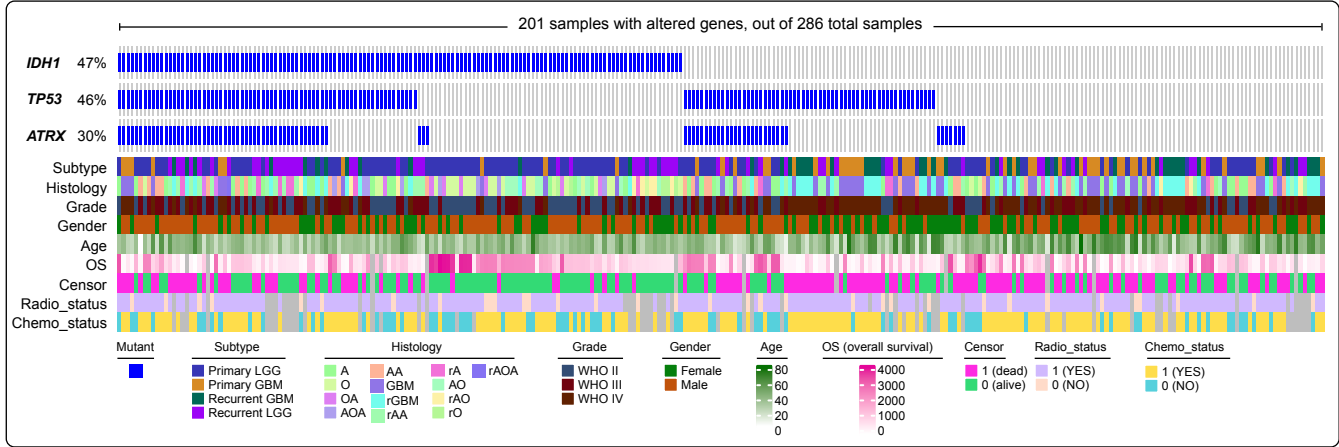
Survival – month						
Median (95% CI)	38.8 (27.2-53.9)	- (99.8- )	39.8 (13.8- )	15.4 (13.3-19.0)	10.5 (7.7- )	7.2 (6.5- )
IDH_mut_status						
Mutant	134	100 (75%)	12 (9%)	14 (10%)	2 (1%)	6 (4%)
Wildtype	165	54 (33%)	6 (4%)	94 (57%)	3 (2%)	5 (3%)
Unknown	2	2 (100%)	0	0	0	0
1p19q_codeletion_status						
Codel	16	14 (88%)	2 (12%)	0	0	0
Non-codel	76	23 (30%)	14 (18%)	27 (36%)	5 (7%)	7 (9%)
Unknown	209	119 (57%)	2 (1%)	81 (39%)	0	4 (2%)
<b>methyl_159_dataset</b>						
No. of samples –no. (%)	159	100 (63%)	8 (5%)	33 (21%)	4 (3%)	6 (4%)
Age at diagnosis – yr.						
Mean	40.2±12.5	39.5±12.2	35.6±12.0	44.2±14.2	41.5±3.7	33.7±7.4
Range	9-70	17-70	24-57	9-70	38-46	27-46
Male sex – no. (%)	89	58 (65%)	4 (4%)	19 (21%)	3 (3%)	5 (6%)
Therapy						
Radiotherapy only	48	39 (81%)	1 (2%)	8 (17%)	0	0
Chemotherapy	10	0	3 (30%)	1 (10%)	3 (30%)	3 (30%)
Cheomoradiotherapy	66	46 (70%)	3 (5%)	16 (24%)	1 (2%)	0
No therapy	12	4 (33%)	1 (8%)	4 (33%)	0	3 (25%)
Unknown	19	11 (58%)	2 (5%)	4 (21%)	3 (16%)	0
Survival – month						
Median (95% CI)	45.8 (36.6-83.9)	107.2 (60.4- )	85.0 (43.8- )	8.5 (6.4-23.1)	16.0 (5.2- )	43.3 (10.6- )
IDH_mut_status						
Mutant	81	65 (80%)	5 (6%)	5 (6%)	2 (2%)	4 (5%)
Wildtype	64	30 (47%)	3 (5%)	27 (42%)	2 (3%)	2 (3%)
Unknown	14	5 (36%)	0	1 (7%)	0	0
1p19q_codeletion_status						

Codel	7	5 (71%)	2 (29%)	0	0	0
Non-codel	18	7 (39%)	3 (17%)	2 (11%)	2 (11%)	4 (22%)
Unknown	134	88 (66%)	3 (2%)	31 (23%)	2 (1%)	2 (1%)
<b>microRNA-array_198 dataset</b>						
No. of samples –no. (%)	198	99 (50%)	8 (4%)	81 (41%)	4 (2%)	6 (3%)
Age at diagnosis – yr.						
Mean	41.9±12.5	39.5±12.3	35.6±12.0	46.1±13.1	41.5±3.7	33.7±7.4
Range	12-70	17-70	24-57	12-70	38-46	27-46
Male sex – no. (%)	123	57 (46%)	4 (3%)	54 (44%)	3 (2%)	5 (4%)
Therapy						
Radiotherapy only	57	38 (67%)	1 (2%)	18 (32%)	0	0
Chemotherapy	12	0	3 (25%)	3 (25%)	3 (25%)	3 (25%)
Cheomoradiotherapy	99	47 (47%)	3 (3%)	48 (48%)	1 (1%)	0
No therapy	15	4 (27%)	1 (7%)	7 (47%)	0	3 (20%)
Unknown	15	10 (67%)	0	5 (33%)	0	0
Survival – month						
Median (95% CI)	28.4(22.1-43.8)	121.6 (60.4- )	85.0 (43.8- )	13.7 (12.7-18.8)	16.0 (5.2- )	43.3 (10.6- )
IDH_mut_status						
Mutant	81	63 (78%)	5 (6%)	7 (9%)	2 (2%)	4 (5%)
Wildtype	106	30 (28%)	3 (3%)	69 (65%)	2 (2%)	2 (2%)
Unknown	11	6 (55%)	0	5 (45%)	0	0
1p19q_codeletion_status						
Codel	7	5 (71%)	2 (29%)	0	0	0
Non-codel	19	7 (37%)	3 (16%)	3 (16%)	2 (11%)	4 (21%)
Unknown	172	87 (51%)	3 (2%)	78 (45%)	2 (1%)	2 (1%)

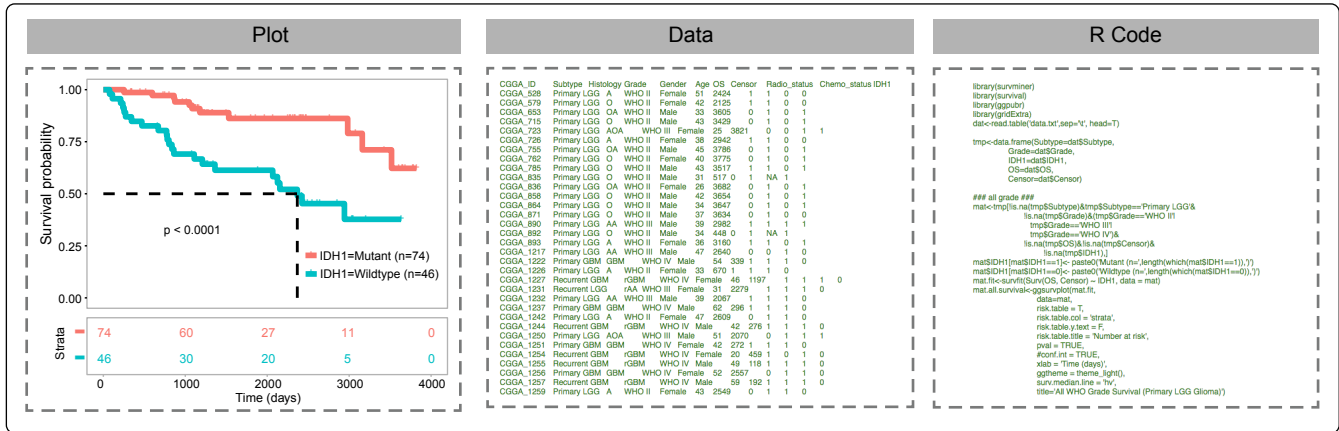
## A CGGA Database Content



## B OncoPrint of Key Gene Alterations in All Grade Gliomas



## C IDH1 Mutation Survival Analysis



## D ADAMTSL4 Expression Analysis

