

MYC amplification at diagnosis drives therapy-induced hypermutation of recurrent glioma

Jiguang Wang (✉ jgwang@ust.hk)

The Hong Kong University of Science and Technology <https://orcid.org/0000-0002-6923-4097>

Quanhua Mu

The Hong Kong University of Science and Technology

Ruichao Chai

Beijing Neurosurgical Institute

Hanjie Liu

Beijing Neurosurgical Institute

Yingxi Yang

Division of Life Science and Department of Chemical and Biological Engineering, Hong Kong University of Science and Technology, Hong Kong

Zheng Zhao

Beijing Neurosurgical Institute

Ming Hong Lui

Division of Life Science and Department of Chemical and Biological Engineering, Hong Kong University of Science and Technology, Hong Kong

Zhaoshi Bao

Beijing Neurosurgical Institute

Dong Song

hong kong university of science and technology <https://orcid.org/0000-0002-5213-5121>

Biaobin Jiang

Division of Life Science and Department of Chemical and Biological Engineering, Hong Kong University of Science and Technology, Hong Kong

Jason Sa

Korea University <https://orcid.org/0000-0002-3251-5004>

Hee Jin Cho

Samsung Medical Center

Yuzhou Chang

Beijing Neurosurgical Institute

Kaitlin Hao Yi Chan

HKUST <https://orcid.org/0000-0001-7363-2082>

Danson Shek Chun Loi

HKUST

Sindy Sing Ting Tam

HKUST

Aden Ka Yin Chan

CUHK

Angela Wu

Division of Life Science and Department of Chemical and Biological Engineering, Hong Kong University of Science and Technology, Hong Kong <https://orcid.org/0000-0002-3531-4830>

Wai San Poon

The Chinese University of Hong Kong

H.K. Ng

The Chinese University of Hong Kong

Danny Chan

Prince of Wales Hospital

Antonio Iavarone

Columbia University <https://orcid.org/0000-0002-0683-4634>

Do-Hyun Nam

Samsung Medical Center

Tao Jiang

Beijing Neurosurgical Institute <https://orcid.org/0000-0002-7008-6351>

Biological Sciences - Article

Keywords: Clonal Evolution, Adult Diffuse Gliomas, Matched Initial-recurrent Tumor Pairs, Grade Progression, Loss-of-function Mutations, Cancer Evolution

Posted Date: January 5th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-138020/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **MYC amplification at diagnosis drives therapy-induced hypermutation**
2 **of recurrent glioma**

3

4 Quanhua Mu^{1,#}, Ruichao Chai^{2,#}, Hanjie Liu², Yingxi Yang¹, Zheng Zhao², Ming Hong
5 Lui¹, Zhaoshi Bao^{1,2}, Dong Song¹, Biaobin Jiang¹, Jason K. Sa^{4,5}, Hee Jin Cho⁴, Yuzhou
6 Chang², Kaitlin Hao Yi Chan³, Danson Shek Chun Loi³, Sindy Sing Ting Tam³, Aden Ka
7 Yin Chan⁷, Angela Ruohao Wu³, Wai Sang Poon⁶, Ho Keung Ng⁷, Danny Tat Ming Chan⁶,
8 Antonio Iavarone^{8,*}, Do-Hyun Nam^{4,9,10,11,*}, Tao Jiang^{2,11,*}, Jiguang Wang^{1,11,*}

9

10 ¹ Division of Life Science, Department of Chemical and Biological Engineering, Center of
11 Systems Biology and Human Health, State Key Laboratory of Molecular Neuroscience,
12 The Hong Kong University of Science and Technology, Hong Kong SAR, China.

13 ²Beijing Neurosurgical Institute, Capital Medical University, Beijing, China.

14 ³Division of Life Science, Department of Chemical and Biological Engineering, The Hong
15 Kong University of Science and Technology, Hong Kong SAR, China.

16 ⁴Institute for Refractory Cancer Research, Samsung Medical Center, Seoul, South Korea.

17 ⁵Department of Biomedical Sciences, Korea University College of Medicine, Seoul, Korea.

18 ⁶CUHK Otto Wong Brain Tumour Centre, Department of Surgery, Prince of Wales
19 Hospital, The Chinese University of Hong Kong, Hong Kong SAR, China.

20 ⁷Department of Anatomical and Cellular Pathology, Prince of Wales Hospital, The Chinese
21 University of Hong Kong, Hong Kong SAR, China.

22 ⁸Institute for Cancer Genetics, Columbia University, New York, New York, USA.

23 ⁹Department of Neurosurgery, Samsung Medical Center, Sungkyunkwan University
24 School of Medicine, Seoul, Korea.

25 ¹⁰Department of Health Science & Technology, Samsung Advanced Institute for Health
26 Sciences & Technology, Sungkyunkwan University School of Medicine, Seoul, Korea.

27 ¹¹Chinese Glioma Genome Atlas (CGGA) and Asian Glioma Genome Atlas (AGGA)
28 Research Networks.

29

30 # These authors contributed equally

31 * Emails:

32 jgwang@ust.hk; ai2102@cumc.columbia.edu; nsnam@skku.edu; taojiang1964@163.com.

1 **Abstract**

2 Clonal evolution drives cancer progression and therapeutic resistance^{1,2}. Recent
3 longitudinal analyses revealed divergent clonal dynamics in adult diffuse gliomas³⁻¹¹.
4 However, the early genomic and epigenomic factors that steer post-treatment molecular
5 trajectories remain unknown. To track evolutionary predictors, we analyzed sequencing
6 and clinical data of matched initial-recurrent tumor pairs from 511 adult diffuse glioma
7 patients. Using machine learning we developed methods capable of predicting grade
8 progression and hypermutation from tumor characteristics at diagnosis. Strikingly, *MYC*
9 copy number gain in initial tumors emerged as a key factor predicting development of
10 hypermutation under temozolomide (TMZ) treatment. The driving role of *MYC* in TMZ-
11 associated hypermutagenesis has been experimentally validated in a model of TMZ-
12 induced hypermutator using both patient-derived gliomaspheres and established glioma
13 cell lines. Subsequent studies showed that c-Myc binding to open chromatin and
14 transcriptionally active regions increases the vulnerability of genomic regions to TMZ-
15 induced mutagenesis. Consequently, *MYC* target genes, including the key mismatch repair
16 genes, develop loss-of-function mutations, thus triggering the hypermutation process. This
17 study reveals *MYC* as an early predictor of cancer evolution and provides a machine
18 learning platform for predicting cancer dynamics to improve patient management.

1 **Main**

2 Malignant diffuse gliomas (WHO grade II to IV) constitute the most common primary
3 brain tumors in adults¹². Despite surgery combined with radiotherapy plus alkylator
4 chemotherapy, aggressive gliomas inevitably recur. At relapse, a proportion of lower-grade
5 gliomas (LGG) remain grade II-III but a remarkable number of cases progress into grade
6 IV glioblastoma (GBM), leading to worse prognosis³. Molecular mechanisms that drive
7 this process remain elusive.

8

9 To characterize the evolutionary landscape of glioma under therapy, enormous efforts were
10 devoted to sequence cancer genomes at multiple time points³⁻¹¹. The hypermutation (HM)
11 phenomenon was detected in temozolomide-treated samples in several independent cohorts,
12 highlighting the role of stress-induced mutagenesis in shaping adaptive glioma
13 evolution^{13,14}. Our recent study revealed a mutually exclusive pattern between HM and an
14 alternative TMZ-resistant mechanism driven by *MGMT*¹⁵. Moreover, group efforts of
15 Glioma Longitudinal AnalySiS (GLASS) consortium reconstructed molecular trajectories
16 of 222 individual cases, demonstrating a heterogeneous and largely stochastic pattern of
17 glioma evolution after treatment⁸. Yet, it is still unknown how the evolutionary routes were
18 determined and whether the early predictors of cancer evolution exist.

19

20 Here we analyzed longitudinal sequencing data from a large number of diffuse glioma
21 cases that cover all three molecular subtypes—IDH wildtype (IDHwt), IDH-mutant
22 without chromosome 1p/19q co-deletion (IDHmut-noncode1), and IDH-mutant with
23 chromosome 1p/19q co-deletion (IDHmut-codel). Leveraging information theory and

24 machine learning, we aim to quantify clonal stochasticity and track early molecular
25 predictors that drive cancer evolution.

26

27 **Longitudinal sequencing reveals mutational dynamics in glioma**

28 We assembled 511 matched initial and recurrent diffuse glioma pairs, among which 95 out
29 of 219 LGG cases progressed to GBM at recurrence. Out of these samples, 276 were
30 profiled by whole-genome sequencing (WGS), 989 by whole exome sequencing (WES),
31 394 by RNA sequencing, and 201 by clinical panel sequencing (**Extended Data Fig. 1a-**
32 **b, Extended Data Table 1-2**). Using calibrated computational pipelines, we portrayed the
33 dynamic landscape of somatic mutations in key driver genes (**Fig. 1a, Extended Data Fig.**
34 **1c-d**). Consistent with previous reports^{3,8,9,16}, we found that 17.0% (62 out of 365) of TMZ-
35 treated cases in our cohort gained hypermutation at recurrence. Excluding hypermutated
36 samples, elevated mutational burden at recurrence was observed in each of the three
37 molecular subtypes of glioma with *P* values of 8.0×10^{-6} , 2.9×10^{-6} and 2.0×10^{-3} for IDHwt,
38 IDHmut-noncodel, and IDHmut-codel (paired *t*-test), respectively (**Extended Data Fig.**
39 **1e**). Ternary plot showed that unlike IDHwt gliomas, IDH-mutant tumors had a remarkable
40 number of key alterations enriched in recurrent-specific area (**Extended Data Fig. 2a-c**).
41 In particular, 82 out of 128 IDHmut-noncodel and 13 out of 59 IDHmut-codel gliomas
42 developed recurrence-specific alterations in canonical GBM drivers such as *CDKN2A*
43 deletion and/or *CDK4* amplification and deletions or mutations in *PTEN* (**Extended Data**
44 **Fig. 2b-c**). Interestingly, six recurrent IDHwt gliomas lost their initial alteration in
45 *EGFRvIII* and nine recurrent IDHwt gliomas lost the copy number amplification in *MDM4*.
46 Moreover, development of HM at recurrence was observed in all three subtypes, while one

47 IDHwt and two IDHmut-noncodel cases developed MGMT translocation. These results
48 collectively implied that diffuse glioma recurred and developed chemoresistance in
49 response to current therapeutic intervention via divergent molecular mechanisms.

50

51 We next systematically compared genetic alterations of initial and recurrent tumors in all
52 patients and found that 70.2% of cases (269 out of 383 patients sequenced by WES/WGS)
53 lost >5 coding mutations (single nucleotide variants and small indels) while 81.9% (313
54 out of 383) cases gained >5 new mutations (**Fig. 1b**). For each molecular subtype, we
55 observed significantly larger number of recurrence-specific compared to initial-specific
56 mutations ($P = 8.4 \times 10^{-7}$, 6.3×10^{-12} and 1.5×10^{-4} for IDHwt, IDHmut-noncodel, and
57 IDHmut-codel, respectively) (**Fig. 1b**). To quantify the changeability of each glioma
58 feature during evolution, we defined its longitudinal divergence by binarizing the change
59 of its mutational status over time and calculating the information entropy (**Fig. 1c**). As
60 expected, no longitudinal change was observed for mutational status of *IDH1* or co-
61 deletion of 1p/19q and hence these features contained zero entropy, implying that these
62 alterations either achieved fixation due to their overwhelming survival advantage or
63 occurred extremely early during glioma evolution (**Fig. 1d, Extended Data Fig. 2d-f**). In
64 contrast, the pathologic grade, somatic hypermutation (HM), together with genomic
65 deletions in *CDKN2A*, *RBI* and *PTEN* contained the highest entropy, implying these
66 features changed frequently over time (**Fig. 1d**). We reason that the change of mutational
67 status of a glioma feature was not only determined by stochasticity, but also related to
68 potential predisposition factors observed early. Therefore, we calculated mutual
69 information between each feature's longitudinal divergence and data acquired from initial

70 tumor to quantify predictability of the corresponding feature (Methods). Interestingly,
71 pathological grading, and HM are two features with the highest predictability, implying
72 these two events might be deployed earlier (**Fig. 1d, Extended Data Fig. 2g-h**).

73

74 **Genomic characteristics of initial glioma informs grade progression**

75 We then investigated the longitudinal change of pathologic grade in our cohort. Not
76 surprisingly, we found that patients who progressed to GBM showed significantly worse
77 survival compared to those who recurred as LGG, regardless of their molecular subtypes
78 (**Fig. 2a**). The above information theory analysis suggested that the risk of glioma grade
79 progression was related to early genetic and clinical factors. We therefore developed a
80 machine learning model to predict grade progression (LGG to GBM) from genomic and
81 clinical characteristics collected at diagnosis.

82

83 To develop the machine learning model, we compiled data of 116 LGG cases as the training
84 cohort (**Fig. 2b, Extended Data Table 1**). Particularly, the training cohort included 53
85 LGG-GBM longitudinal pairs together with 63 LGG-LGG pairs. Overall, three
86 classification models were trained using different sets of features: Model 1 used the initial
87 grade as the single predictor, Model 2 used the initial grade, age, gender, and treatment as
88 predictors, while Model 3 included the initial grade, age, gender, treatment, and genomic
89 alterations as predictors. The performance of the models was measured by the area under
90 the receiver operating characteristic curve (AUC). Based on five-fold cross validation we
91 found that the model using integrated features including patient clinical information,
92 applied treatment plus genomic alterations achieved the highest AUC (0.87, 95% CI 0.81-

93 0.94, DeLong test, **Fig. 2c**). The SHAP (Shapley additive explanations) score, a metric to
94 quantify the importance of each feature^{17,18}, was then calculated and used to prioritize the
95 relevant features. On the top of the list is IDH status, where IDH mutation contributed
96 negatively to grade progression, and IDH wildtype are more prone to progress. Following
97 IDH mutations are 17p copy neutral loss of heterozygosity (CNLOH), 1p/19q codeletion,
98 alkylator treatment and *CDKN2A* deletion (**Fig. 2d**). Other relevant genetic factors include
99 amplification/gain of *MET*, chr7p and *MYC*, deletions of *TP53* and *RBI*, as well as *TERT*
100 promoter mutations. Notably, *MYC*, a master regulator of stemness and proliferation, has
101 been recently revealed to regulate glioma progression¹⁹.

102

103 To test the model performance, 50 independent LGG cases from GLASS were used as the
104 hold-out testing cohort. In this cohort, 21 LGGs progressed to GBM at recurrence and the
105 other 29 did not. Our model achieved an AUC of 0.79 (95% confidence interval 0.67-0.92,
106 DeLong test; **Fig. 2e**). We then applied the model on 496 TCGA primary LGGs to assess
107 whether the predicted progression risk to GBM at recurrence might be associated with
108 patient survival outcome. As a result, remarkable worse overall survival (OS, $P = 1.34 \times$
109 10^{-8} , log-rank test) and progression-free survival (PFS, $P = 1.15 \times 10^{-6}$, log-rank test)
110 outcomes were observed for cases predicted to progress to GBM, compared to those
111 predicted to recur as LGG (**Fig. 2f**). More importantly, when glioma patients were
112 evaluated in each of the three molecular groups, poorer PFS and OS were consistently
113 associated with predicted progression risk (**Extended Data Fig. 3a-b**). These findings
114 demonstrated that our machine learning model was able to integrate clinical information

115 and molecular features in order to evaluate the risk of glioma progression and early predict
116 patient prognosis.

117

118 **Features of initial tumors predispose hypermutation at recurrence**

119 We adopted CELLO toolkit²⁰ to infer HM from multi-platform sequencing data and
120 identified 68 HM samples with average mutation load 157.7 mutations per Mb and average
121 HM score 1.39 (**Extended Data Fig. 4a**). In our dataset, we demonstrated that HM was
122 associated with glioma progression and impacts post-progression survival (**Extended Data**
123 **Fig. 4b-c**). Consistent with previous studies^{16,21,22}, most (16 out of 18 samples with
124 available MGMT methylation status, 88.9%) of HM samples were MGMT methylated, and
125 59 of 65 (90.8%) HM cases gained loss-of-function mutations in the mismatch repair
126 (MMR) pathway genes such as *MSH6* and *MSH2* (**Extended Data Fig. 4d**). We then
127 compared the transcriptional profile of TMZ-treated HM versus NHM recurrent gliomas.
128 As expected, *MGMT* expression is significantly lower in HM samples compared to the
129 alkylator-treated non-hypermethylated glioma, while genes such as cyclin dependent kinase 2
130 (*CDK2*) and polo-like kinase 1 (*PLK1*) had significantly higher expression in HM
131 (**Extended Data Fig. 4e**). The elevated *PLK1* level in HM glioma is compatible with
132 results of a previous study that applied PLK1 inhibitors in MMR-deficient glioma cell
133 lines²³. Subsequent gene set enrichment analysis revealed MYC targets V2 as the most
134 significantly up-regulated gene set (normalized enrichment score (NES) = 3.32, $P = 5.2 \times$
135 10^{-4} , **Extended Data Fig. 4f**) in hypermutated recurrent glioma. Other most upregulated
136 pathways included DNA repair, E2F targets, G2M checkpoint, and MYC targets V1,
137 implying HM samples had more active transcription and proliferation²⁴.

138

139 To identify the molecular features at diagnosis that drive the occurrence of HM at
140 recurrence, we compared the transcriptome of initial gliomas that recurred as hypermutated
141 (pre-HM, n=21) versus the initial gliomas that recurred as non-HM tumors (pre-NHM, n =
142 143). Mimicking the recurrent tumors, we found significant upregulation of MYC targets,
143 G2M checkpoint, E2F targets and mitotic spindle pathways in pre-HM samples (Figure
144 S5a-b). MYC pathway activation is known to be associated with cell proliferation. Indeed,
145 the commonly used cell proliferation marker, MKI67, is significantly up regulated genes
146 in pre-HM gliomas ($P = 1.54 \times 10^{-3}$, Wilcoxon's rank-sum test, **Extended Data Fig. 5c**).
147 Immunohistochemistry staining of 67 gliomas showed that gliomas with high Ki-67
148 staining have significantly higher probability to develop HM after TMZ treatment ($P =$
149 0.04, Fisher's exact test, **Extended Data Fig. 5d**). *MGMT* was down-regulated in pre-HM
150 samples, which is consistent with the observation that MGMT-methylated gliomas, when
151 treated with TMZ, were prone to develop hypermutation ($P = 0.021$, Fisher's exact test,
152 **Extended Data Fig. 5e**). Strikingly, comparing genomic alterations of pre-HM and pre-
153 non-HM gliomas we found that genomic gain of *MYC* was the most enriched feature in the
154 pre-HM tumors (fold change = 5.31, $P = 3.23 \times 10^{-5}$, Fisher's exact test), well in line with
155 the upregulation of MYC pathways (**Fig. 3a-b**). Specifically, about 46.9% of *MYC*-gained
156 gliomas developed HM after TMZ treatment, while the proportion in gliomas without *MYC*
157 gain was only 12.7% (**Fig. 3c**). Following *MYC* gain, other genomic alterations enriched
158 in initial tumors of hypermutators include *ATRX* mutation and *RB1* deletion (**Fig. 3a**). Loss-
159 of-function mutations in *ATRX* have been shown to impair DNA repair, leading to
160 genetically unstable tumors that rapidly accumulate oncogenic mutations²⁵. *RB1* encodes

161 Rb protein and plays an essential role as a cell-cycle regulator acting at the G1/S cell cycle
162 checkpoint²⁶. To make use of the comprehensive information for hypermutation prediction,
163 machine learning models were trained using genomic features of initial glioma plus
164 alkylator treatment information and achieved an Area Under Precision-Recall Curve
165 (AUPRC) of 0.88. When the expression of *MGMT* and *MKI67* were included, the AUPRC
166 further improved to 0.95 (**Fig. 3d**). SHAP analysis showed that alkylator treatment was the
167 most predictive feature in the model, followed by *MYC* gain, *ATRX*, *FUBP1* and *IDH*
168 mutations, and *EGFR* amplification (**Fig. 3e**).

169

170 To interrogate the association between *MYC* and HM in glioma, we grew a *MYC*-
171 amplified glioma cell line U251 and a *MYC*-wild-type glioma cell line U87 in medium
172 containing escalating concentrations of TMZ (**Extended Data Fig. 5f-g**), until they
173 developed elevated TMZ resistance (TR; **Extended Data Fig. 5h**). During the whole
174 process, U251 maintained its high copy number and high expression of *MYC* gene, while
175 U87 had relatively low *MYC* expression (**Extended Data Fig. 5i**). Upon resistance,
176 U87TR showed elevated *MGMT* expression (**Extended Data Fig. 5j**). On the contrary,
177 U251TR maintained low level of *MGMT* expression but gained thousands of mutations.
178 Over 90% of U251TR mutations were C>T (G>A) single nucleotide variants, including an
179 *MSH6* T1219I hotspot mutation with allele frequency 38%, demonstrating the emergence
180 of the TMZ-induced hypermutation *in vitro* (**Extended Data Fig. 5i-j**). Considering U87
181 and U251 might not perfectly represent glioma patients, we further carried out the TMZ
182 inducing experiment in the newly derived gliomasphere (T2-4) sample from a patient with
183 high *MYC* copy number and expression (**Fig. 3f**). This sample developed dramatic TMZ

184 resistance after the induction experiment (**Fig. 3g**). Interestingly, T2-4TR developed
185 hypermutation with a stop codon mutation *MLH1* W597* (AF = 23%) and a splicing-donor
186 mutation (*MLH1* c.453+1G>A, AF = 32%, **Fig. 3h**), while maintaining a low level of
187 *MGMT* expression and a high level of MYC copy number and expression. Notably, the
188 splicing-donor mutation in *MLH1* caused an aberrant splicing which included an additional
189 nine basepairs in the spliced RNA (**Fig. 3i**). We next tracked the development of
190 hypermutator phenotype during T2-4 cell line passaging, and found the percent spliced in
191 (PSI) of the aberrant isoform emerged at passage six and increased over time (**Fig. 3i**).
192 Whole-exome sequencing of different passages demonstrated that hypermutation occurred
193 almost simultaneously with loss-of-function mutations in *MLH1* at the sixth-seventh
194 passage of cell culture (**Fig. 3i-j**), confirming the underlying association between a burst
195 of somatic mutation and the MMR deficiency in glioma cells.

196

197 **Hypermutation occurs frequently in actively transcribed genomic regions with MYC** 198 **binding**

199 To explore the mechanism of how MYC activation promotes somatic hypermutation, we
200 first characterized the genomic distribution of hypermutation sites in relation to epigenetic
201 regulatory elements. The human genome was annotated using chromosomal binding peaks
202 from chromatin immunoprecipitation sequencing (ChIP-seq) data (**Extended Data Fig.**
203 **Table 3, Methods**), and then the mutational density of the TMZ-associated hypermutations
204 (HM) and non-TMZ-associated conventional somatic mutations (NHM) were calculated
205 and normalized for comparison. Consistent with previous observations that conventional
206 somatic mutations in the cancer genomes were not randomly distributed but enriched in

207 transcriptionally inactive regions²⁷⁻³¹, we found remarkably higher density of NHM
208 mutations around the H3K9me3 modification sites which are often associated with
209 heterochromatin (**Fig. 4a**, upper panel). In contrast, the density of hypermutations in these
210 regions were close to random. Conversely, while NHM mutations were significantly lower
211 than expected in the open chromatin regions marked by DNase hypersensitivity, HM had
212 remarkably higher mutation (**Fig. 4a**, lower panel). More importantly, the HM mutation
213 was significantly enriched in the periphery of H3K27ac and H3K4me3 sites that
214 respectively mark active enhancers and promoters, while NHM was of lower density (**Fig.**
215 **4b**, upper panel). In regions surrounding binding sites of RNA polymerase II (RPII), a key
216 enzyme in the transcriptional machinery, the density of HM mutations was almost two-fold
217 of NHM mutations (**Fig. 4b**, lower left panel). It has been shown that in normal cells, the
218 actively transcribed genes recruit mismatch repair complex to protect the gene body from
219 mutation by H3K36me3^{32,33}. We investigated the mutation densities around H3K36me3
220 modification sites and noted that the mutation density of both HM and NHM mutations
221 were lower than random control, but the density of HM was much higher than NHM (**Fig.**
222 **4b**, lower right panel). Altogether, these results demonstrated that in TMZ-associated
223 hypermutators, mutations were enriched to the proximity of active transcription sites.

224

225 To explore the role of MYC gain in hypermutagenesis, we compared MYC gene expression
226 between samples with and without MYC gain and found significantly higher MYC
227 expression in gliomas with MYC gain ($P = 1.4 \times 10^{-7}$, Wilcoxon rank-sum test; **Extended**
228 **Data Fig. 6a**). In tumor cells with elevated MYC expression, c-Myc forms a heterodimer
229 with MAX and binds to the promoters and enhancers of active genes, causing global

230 transcription amplification³⁴. To quantify the global transcriptional activity, we used
231 fourteen housekeeping genes as reference and counted the number of high-expression
232 genes (**Methods**). As expected, more high expression genes were observed in MYC-gain-
233 positive gliomas ($P = 1.0 \times 10^{-3}$, Wilcoxon rank-sum test, **Fig. 4c**), indicating globally
234 amplified transcription in these samples. Comparison of the expression level of MYC direct
235 targets also showed consistent results ($P = 1.3 \times 10^{-4}$, Wilcoxon rank-sum test, **Extended**
236 **Data Fig. 6b**).

237

238 We next divided the open chromatin regions and active transcription sites based on whether
239 c-Myc also binds to the site, and then compared the density of HM mutations in these two
240 types of regions (**Fig. 4d**). The density of HM mutations was higher in open chromatin
241 with c-Myc binding compared to those without c-Myc binding. Similarly, when checking
242 H3K27ac modification sites with and without c-Myc binding, we found the mutation
243 density of HM to be much higher in sites with c-Myc binding than that in sites without c-
244 Myc binding (**Fig. 4d**). The same phenomenon was observed for H3K4me3 modification
245 sites, RPII binding sites and H3K36me3 modification sites with and without c-Myc binding
246 (**Fig. 4d**).

247

248 **c-Myc binds MMR genes and increases probability of hypermutation**

249 The eight MMR genes, namely *PMS1*, *PMS2*, *MSH2*, *MSH3*, *MSH5*, *MSH6*, *MLH1* and
250 *MLH3*, all had H3K4me3 modification, c-Myc/MAX and RPII binding, especially in the
251 periphery of the promoter regions (**Fig. 4e**). In U87MG cell line, we experimentally
252 induced MYC overexpression and confirmed that MYC overexpression upregulated the

253 expression of MSH6 and MSH2, the two most frequently altered MMR genes in HM (**Fig.**
254 **4f, Extended Data Fig. 6c**). Chi-square test with Yates' correction revealed significant
255 enrichment of MMR genes in c-Myc binding targets ($P = 0.0002$, **Fig. 4g**), a group which
256 was highly enriched for temozolomide-driven mutations. Therefore, we propose that MYC
257 over activation increased the density of TMZ-induced mutations in the genomic regions of
258 MMR genes. With an elevated risk of developing loss-of-function mutations in MMR
259 genes, MYC-amplified gliomas hereby tend to develop HM in response to TMZ treatment.
260 Mechanistically, the double strand DNA is opened during transcription, and MAX and c-
261 Myc binding at the active promoters prolonged the time window of opening status. Under
262 TMZ treatment, the drug molecules attack the single stranded DNA, causing more damages
263 in the transcription active regions. When such damages occurred in MMR genes and failed
264 for repair, mutations will accumulate in other open sites and finally lead to hypermutation
265 (**Fig. 4h**).

266 **Discussion**

267 The high rate of recurrence represents the major obstacle in improving survival of adult
268 diffuse glioma patients. Here, we have assembled and analyzed a mostly original cohort of
269 511 cases with longitudinal matched samples of gliomas. This work highlighted the
270 longitudinal trajectory of evolution for each of the three main subtypes of gliomas (IDH
271 mutant codel, IDH mutant non-codel, IDH wild type). Through the development of novel
272 computational approaches, we found that grade progression and gain of hypermutation are
273 largely predictable from the molecular features of the initial glioma. Following training on
274 clinical and genomic features, our machine learning model robustly distinguished LGGs
275 that tend to progress to GBM from those that recur as LGG. Among the features with
276 highest importance, *CDKN2A* deletion, a key cell-cycle related alteration has been
277 correlated with aggressive disease and tumor enhancing.^{8,35} Another important feature that
278 emerged from the model was *TP53* copy-neutral loss of heterozygosity. Despite the
279 importance and prevalence of TP53 inactivation in glioma, little is known about how TP53
280 loss of function alterations contribute to glioma malignancy. Other important contributing
281 factors in the model include alkylator treatment and MYC gain, both of which are
282 associated with hypermutation. The relevance of these alterations as drivers of glioma
283 aggressiveness is underscored by the significant association between hypermutation and
284 grade progression.

285

286 The most important finding that emerged from our work is the discovery of the gain-of-
287 function genomic alterations of MYC as drivers of the HM state at recurrence. Several
288 reports including our own suggested that therapy-associated hypermutations in recurrent

289 glioma are driven by MMR deficiency^{9,13,16}. In our cohort, MMR deficiency was found in
290 86% of hypermutated gliomas at recurrence. Indeed, whereas MMR deficiency is the direct
291 causal factor leading to hypermutation, the absence of mutations in MMR genes in
292 untreated tumors suggested that other, upstream alterations might be necessary to set in
293 motion the machinery leading to the hypermutated state. The key early event we uncovered,
294 was gain of MYC, an event significantly detected in untreated glioma primed to recur as
295 hypermutated tumors. The role of MYC is especially significant in the context of the IDH
296 mutant non-codel subtype, which is also the glioma group with the highest likelihood to
297 undergo transformation to the hypermutated state at recurrence. Through the integration of
298 exome and transcriptomic analyses, we found that gain of MYC results in marked
299 activation in downstream pathways associated with proliferation and deregulated cell cycle
300 progression, which are the primary biological functions of MYC in mammalian cells. The
301 role of MYC as unexpected driver of the hypermutation state has been experimentally
302 validated in human glioma cells. Furthermore, by integrating genomic and epigenetic data,
303 we confirmed our previous hypothesis that mutagenic mechanisms related to TMZ
304 treatment and subsequent MMR alteration act more efficiently in highly expressed regions
305 of open chromatin¹⁶, and further showed that c-Myc binding augmented the chance of
306 MMR mutations. Interestingly, highly expressed regions are protected by transcription-
307 coupled repair (TCR)³⁶, and hypermutated tumors exhibited an enrichment of somatic
308 mutations in TCR genes in the hypermutated glioma samples. However, whether the TCR
309 pathway is functional in these samples will have to be established by future studies.

310

311 Collectively, this study painted the evolution routes of three molecular types of glioma,
312 trained machine learning models to predict grade progression based on clinical and
313 genomic features of initial glioma and demonstrated MYC as a potential predictor of
314 hypermutation after treatment. In addition, we have developed an interactive, publicly
315 available web resource (**Extended Data Fig. 7**; <https://wanglab.shinyapps.io/cello/>) to
316 explore the longitudinal glioma dataset and make predictions of developing treatment-
317 induced hypermutation and grade progression.
318

319 **Acknowledgements**

320 This work was supported by NSFC Excellent Young Scientists Fund (Hong Kong and
321 Macau) (No. 31922088), RGC grants (N_HKUST606/17, 26102719), ITC grant
322 (ITCPD/17-9), National Natural Science Foundation of China (NSFC)/Research Grants
323 Council (RGC) Joint Research Scheme (No. 81761168038), National Natural Science
324 Foundation of China (No. 81903078, No. 81972816), National Key Research and
325 Development Project of China (No. 2019YFE0109400), Hong Kong Epigenomics Project
326 (LKCCFL18SC01-E), and a grant from the Korea Health Technology R&D through the
327 Korea Health Industry Development Institute funded by the Ministry of Health & Welfare,
328 Republic of Korea (HI14C3418). The authors would like to thank all contributors to the
329 Chinese Glioma Genome Atlas. In addition, the authors acknowledge data generators of
330 the published datasets, especially TCGA research network, the GLASS Consortium,
331 Gregory M. Kiez and Mehmet Kutman Foundation, and the Yale University Department
332 of Neurosurgery for providing access to the raw and/or processed sequencing data.

333

334 **Author Contributions**

335 J.W., T.J., D.H.N., and A.I. conceptualized and supervised the project. Q.M. performed the
336 computational and statistical analyses; R.C., H.L. and Y.C. carried out the cell line and
337 PDC experiments; Y.Y. and Q.M. developed the CELLO webserver; Z.Z., Z.B., and T.J.
338 contributed to sample preparation, genomic sequencing, and clinical data collection of the
339 CGGA cohort; M.H.L., D.S. and B.J. contributed to some of the bioinformatics analyses;
340 J.K.S., H.J.C., and D.H.N. contributed to sample preparation, genomic sequencing and
341 clinical data collection of the SMC cohort; A.C., W.S.P., H.K.N., and D.M.C. collected the

342 CUHK samples; K.C., D.L, S.T, and A.W. prepared libraries for samples of the CUHK
343 cohort; J.W. and Q.M. wrote the manuscript, which was then revised and proof-read by all
344 authors.

345

346 **Competing Interests**

347 The authors declare no potential conflict of interests.

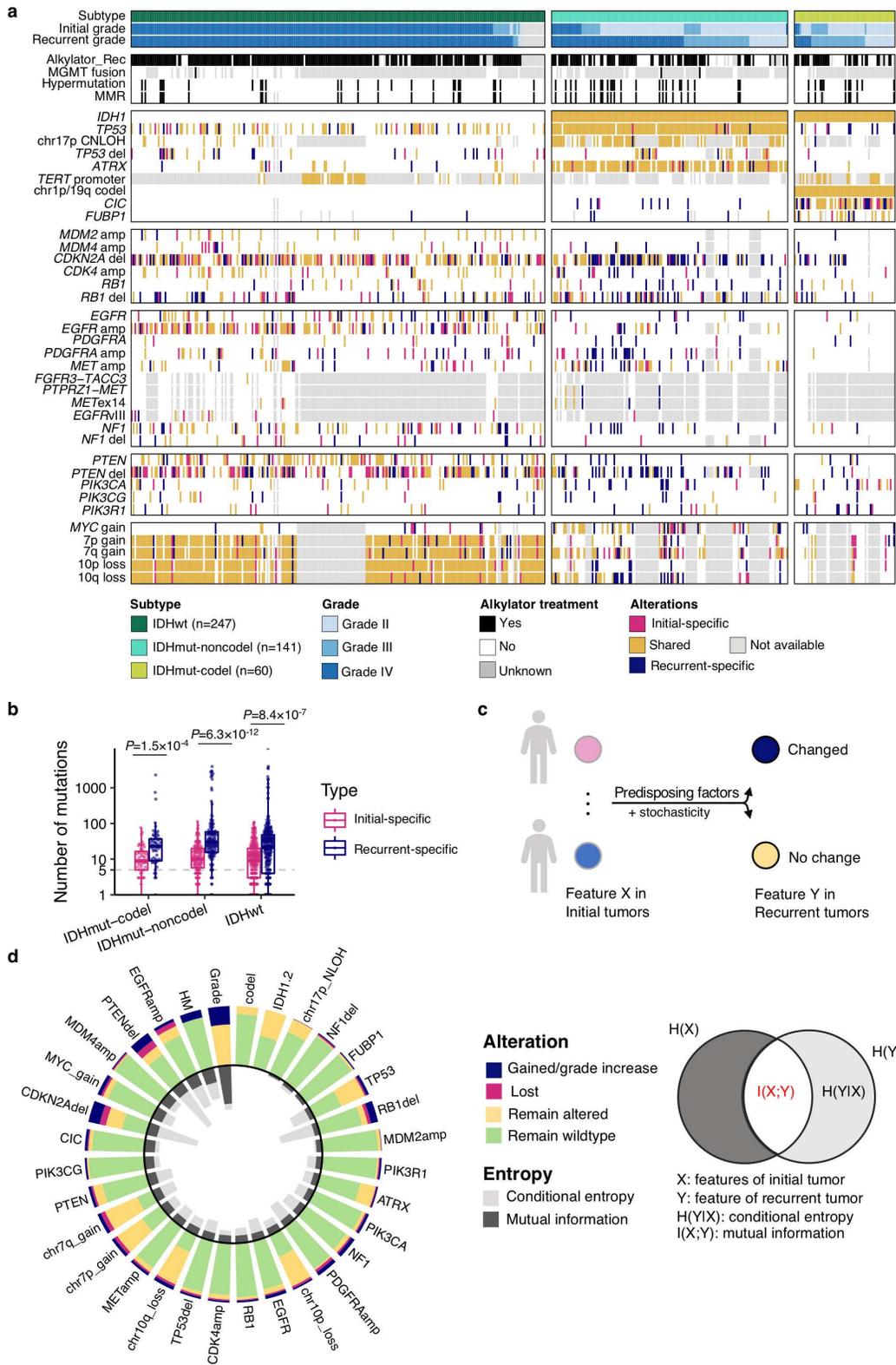
348 **References**

- 349 1. Greaves, M. & Maley, C. C. Clonal evolution in cancer. *Nature* **481**, 306–313
350 (2012).
- 351 2. McGranahan, N. & Swanton, C. Clonal Heterogeneity and Tumor Evolution: Past,
352 Present, and the Future. *Cell* (2017). doi:10.1016/j.cell.2017.01.018
- 353 3. Hu, H. *et al.* Mutational Landscape of Secondary Glioblastoma Guides MET-
354 Targeted Trial in Brain Tumor. *Cell* **175**, 1665-1678.e18 (2018).
- 355 4. Lee, J.-K. *et al.* Spatiotemporal genomic architecture informs precision oncology
356 in glioblastoma. *Nat. Genet.* **49**, 594–599 (2017).
- 357 5. Kim, J. *et al.* Spatiotemporal Evolution of the Primary Glioblastoma Genome.
358 *Cancer Cell* **28**, 318–328 (2015).
- 359 6. Mazor, T. *et al.* DNA Methylation and Somatic Mutations Converge on the Cell
360 Cycle and Define Similar Evolutionary Histories in Brain Tumors. *Cancer Cell* **28**,
361 307–317 (2015).
- 362 7. Suzuki, H. *et al.* Mutational landscape and clonal architecture in grade II and III
363 gliomas. *Nat. Genet.* **47**, 458–468 (2015).
- 364 8. Barthel, F. P. *et al.* Longitudinal molecular trajectories of diffuse glioma in adults.
365 *Nature* **576**, 112–120 (2019).
- 366 9. Johnson, B. E. *et al.* Mutational Analysis Reveals the Origin and Therapy-Driven
367 Evolution of Recurrent Glioma. *Science* (80-.). **343**, 189–193 (2014).
- 368 10. Bai, H. *et al.* Integrated genomic characterization of IDH1-mutant glioma
369 malignant progression. *Nat. Genet.* **48**, 59–66 (2016).
- 370 11. Wang, Q. *et al.* Tumor Evolution of Glioma-Intrinsic Gene Expression Subtypes
371 Associates with Immunological Changes in the Microenvironment. *Cancer Cell*
372 **32**, 42-56.e6 (2017).
- 373 12. Schwartzbaum, J. A., Fisher, J. L., Aldape, K. D. & Wrensch, M. Epidemiology
374 and molecular pathology of glioma. *Nature Clinical Practice Neurology* **2**, 494–
375 503 (2006).
- 376 13. Touat, M. *et al.* Mechanisms and therapeutic implications of hypermutation in
377 gliomas. *Nature* **580**, 517–523 (2020).
- 378 14. Cipponi, A. *et al.* MTOR signaling orchestrates stress-induced mutagenesis,
379 facilitating adaptive evolution in cancer. *Science* **368**, 1127–1131 (2020).
- 380 15. Oldrini, B. *et al.* MGMT genomic rearrangements contribute to chemotherapy
381 resistance in gliomas. *Nat. Commun.* **11**, 1–10 (2020).
- 382 16. Wang, J. *et al.* Clonal evolution of glioblastoma under therapy. *Nat. Genet.* **48**,
383 768–776 (2016).
- 384 17. Lundberg, S. M. *et al.* From local explanations to global understanding with
385 explainable AI for trees. *Nat. Mach. Intell.* **2**, 56–67 (2020).
- 386 18. Lundberg, S. M. & Lee, S.-I. A Unified Approach to Interpreting Model
387 Predictions. in *Advances in Neural Information Processing Systems* (eds. Guyon, I.
388 *et al.*) **30**, 4765–4774 (Curran Associates, Inc., 2017).
- 389 19. Bai, H. *et al.* Integrated genomic characterization of IDH1-mutant glioma
390 malignant progression. *Nat. Genet.* **48**, 59–66 (2016).
- 391 20. Jiang, B., Song, D., Mu, Q. & Wang, J. CELLO: a longitudinal data analysis
392 toolbox untangling cancer evolution. *Quant. Biol.* **8**, 256–266 (2020).

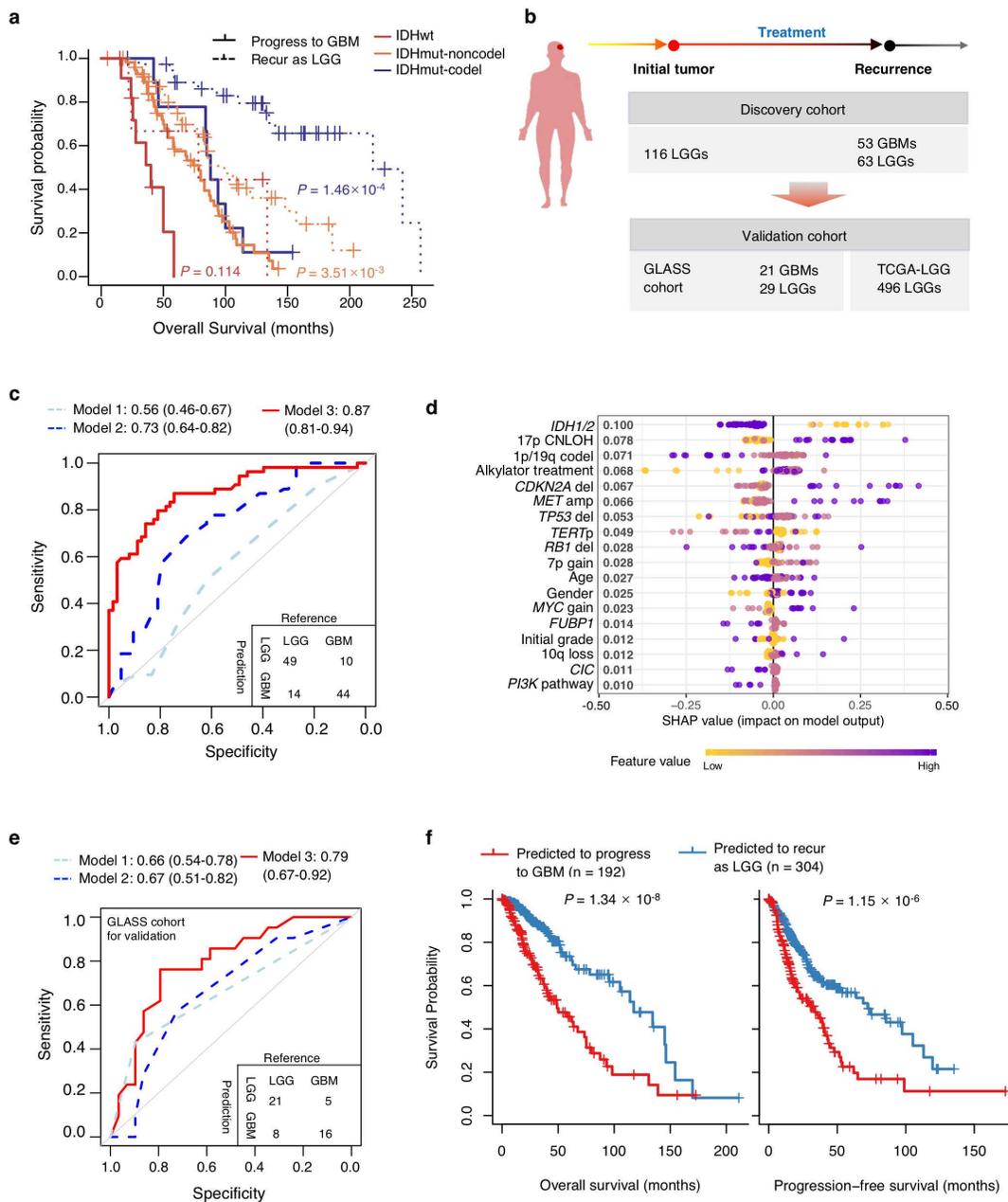
- 393 21. van Thuijl, H. F. *et al.* Evolution of DNA repair defects during malignant
394 progression of low-grade gliomas after temozolomide treatment. *Acta*
395 *Neuropathol.* **129**, 597–607 (2015).
- 396 22. Mathur, R. *et al.* MGMT promoter methylation level in newly diagnosed low-
397 grade glioma is a predictor of hypermutation at recurrence. *Neuro. Oncol.* (2020).
398 doi:10.1093/neuonc/noaa059
- 399 23. Higuchi, F. *et al.* PLK1 inhibition targets MYC-activated malignant glioma cells
400 irrespective of mismatch repair deficiency–mediated acquired resistance to
401 temozolomide. *Mol. Cancer Ther.* **17**, 2551–2563 (2018).
- 402 24. Liberzon, A. *et al.* The Molecular Signatures Database Hallmark Gene Set
403 Collection. *Cell Syst.* **1**, 417–425 (2015).
- 404 25. Koschmann, C. *et al.* ATRX loss promotes tumor growth and impairs
405 nonhomologous end joining DNA repair in glioma. *Sci. Transl. Med.* **8**, 328ra28-
406 328ra28 (2016).
- 407 26. Bertoli, C., Skotheim, J. M. & De Bruin, R. A. M. Control of cell cycle
408 transcription during G1 and S phases. *Nature Reviews Molecular Cell Biology* **14**,
409 518–528 (2013).
- 410 27. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new
411 cancer-associated genes. *Nature* **499**, 214–218 (2013).
- 412 28. Polak, P. *et al.* Cell-of-origin chromatin organization shapes the mutational
413 landscape of cancer. *Nature* **518**, 360–364 (2015).
- 414 29. Gonzalez-Perez, A., Sabarinathan, R. & Lopez-Bigas, N. Local Determinants of
415 the Mutational Landscape of the Human Genome. *Cell* **177**, 101–114 (2019).
- 416 30. Makova, K. D. & Hardison, R. C. The effects of chromatin organization on
417 variation in mutation rates in the genome. *Nature Reviews Genetics* **16**, 213–223
418 (2015).
- 419 31. Akdemir, K. C. *et al.* Somatic mutation distributions in cancer genomes vary with
420 three-dimensional chromatin structure. *Nat. Genet.* **52**, 1178–1188 (2020).
- 421 32. Huang, Y., Gu, L. & Li, G. M. H3K36me3-mediated mismatch repair
422 preferentially protects actively transcribed genes from mutation. *J. Biol. Chem.*
423 **293**, 7811–7823 (2018).
- 424 33. Li, F. *et al.* The histone mark H3K36me3 regulates human DNA mismatch repair
425 through its interaction with MutS α . *Cell* **153**, 590–600 (2013).
- 426 34. Lin, C. Y. *et al.* Transcriptional amplification in tumor cells with elevated c-Myc.
427 *Cell* **151**, 56–67 (2012).
- 428 35. Jonsson, P. *et al.* Genomic correlates of disease progression and treatment
429 response in prospectively characterized gliomas. *Clin. Cancer Res.* **25**, 5537–5547
430 (2019).
- 431 36. Hanawalt, P. C. & Spivak, G. Transcription-coupled DNA repair: Two decades of
432 progress and surprises. *Nature Reviews Molecular Cell Biology* **9**, 958–970 (2008).
- 433 37. Kim, H. *et al.* Whole-genome and multisector exome sequencing of primary and
434 post-treatment glioblastoma reveals patterns of tumor evolution. *Genome Res.* **25**,
435 316–327 (2015).
- 436 38. Ceccarelli, M. *et al.* Molecular Profiling Reveals Biologically Discrete Subsets and
437 Pathways of Progression in Diffuse Glioma. *Cell* **164**, 550–563 (2016).
- 438 39. Zhao, J. *et al.* Immune and genomic correlates of response to anti-PD-1

- 439 immunotherapy in glioblastoma. *Nat. Med.* **25**, 462–469 (2019).
- 440 40. Hu, H. *et al.* Mutational Landscape of Secondary Glioblastoma Guides MET-
441 Targeted Trial in Brain Tumor. *Cell* **175**, 1665-1678.e18 (2018).
- 442 41. Lee, J.-K. *et al.* Pharmacogenomic landscape of patient-derived tumor cells
443 informs precision oncology therapy. *Nat. Genet.* **50**, 1399–1411 (2018).
- 444 42. Trifonov, V., Pasqualucci, L., Tiacci, E., Falini, B. & Rabadan, R. SAVI: A
445 statistical algorithm for variant frequency identification. *BMC Syst. Biol.* **7**, S2
446 (2013).
- 447 43. Cingolani, P. *et al.* A program for annotating and predicting the effects of single
448 nucleotide polymorphisms, SnpEff. *Fly (Austin)*. **6**, 80–92 (2012).
- 449 44. Sherry, S. T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids*
450 *Res.* **29**, 308–11 (2001).
- 451 45. Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer.
452 *Nucleic Acids Res.* **47**, D941–D947 (2019).
- 453 46. Talevich, E., Shain, A. H., Botton, T. & Bastian, B. C. CNVkit: Genome-Wide
454 Copy Number Detection and Visualization from Targeted DNA Sequencing. *PLOS*
455 *Comput. Biol.* **12**, e1004873 (2016).
- 456 47. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics*
457 **25**, 2078–2079 (2009).
- 458 48. Haas, B. J. *et al.* Accuracy assessment of fusion transcript detection via read-
459 mapping and de novo fusion transcript assembly-based methods. *Genome Biol.* **20**,
460 213 (2019).
- 461 49. Wang, J. *et al.* Tumor evolutionary directed graphs and the history of chronic
462 lymphocytic leukemia. *Elife* **3**, (2014).
- 463 50. Chen, T. & Guestrin, C. XGBoost: A Scalable Tree Boosting System. in *Proc.*
464 *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data*
465 *Mining* 785–794 (2016). doi:10.1145/2939672.2939785
466
467
468

469 **Figure and legend**



471 **Figure 1. Identification of predictable events from pan-glioma evolutionary landscape.**
472 **(a)** Longitudinal molecular landscape of glioma patients with paired DNA sequencing data.
473 Each column represents a patient which were stratified by the molecular subtype, and each
474 row represents a clinical feature or molecular variant. The color of each variant indicates
475 whether it is specific to initial, recurrence or shared. MMR: mismatch repair; CNLOH:
476 copy-neutral loss of heterozygosity. **(b)** Comparison of the number of initial- and recurrent-
477 specific somatic mutations during evolution in each glioma subtype. The horizontal dashed
478 line represents five mutations. The *P* values were calculated by Wilcoxon's rank-sum test.
479 **(c)** Proposed model explaining the changes of the clinical and molecular features over time.
480 Predisposing factors in the initial tumor, together with stochasticity, determine whether the
481 status of the feature will change at recurrence. **(d)** Alteration status, entropy and mutual
482 information of clinical and genomic features. Stacked bars outside the circle summarizes
483 the frequency of a clinical feature or molecular variant to be lost, gain, remaining wildtype
484 or remaining altered. Inside the circle shows the overlapped bar plots of mutual information
485 (black) and conditional entropy (grey), where high proportion of overlap indicates better
486 predictability. The features are ranked by mutual information.

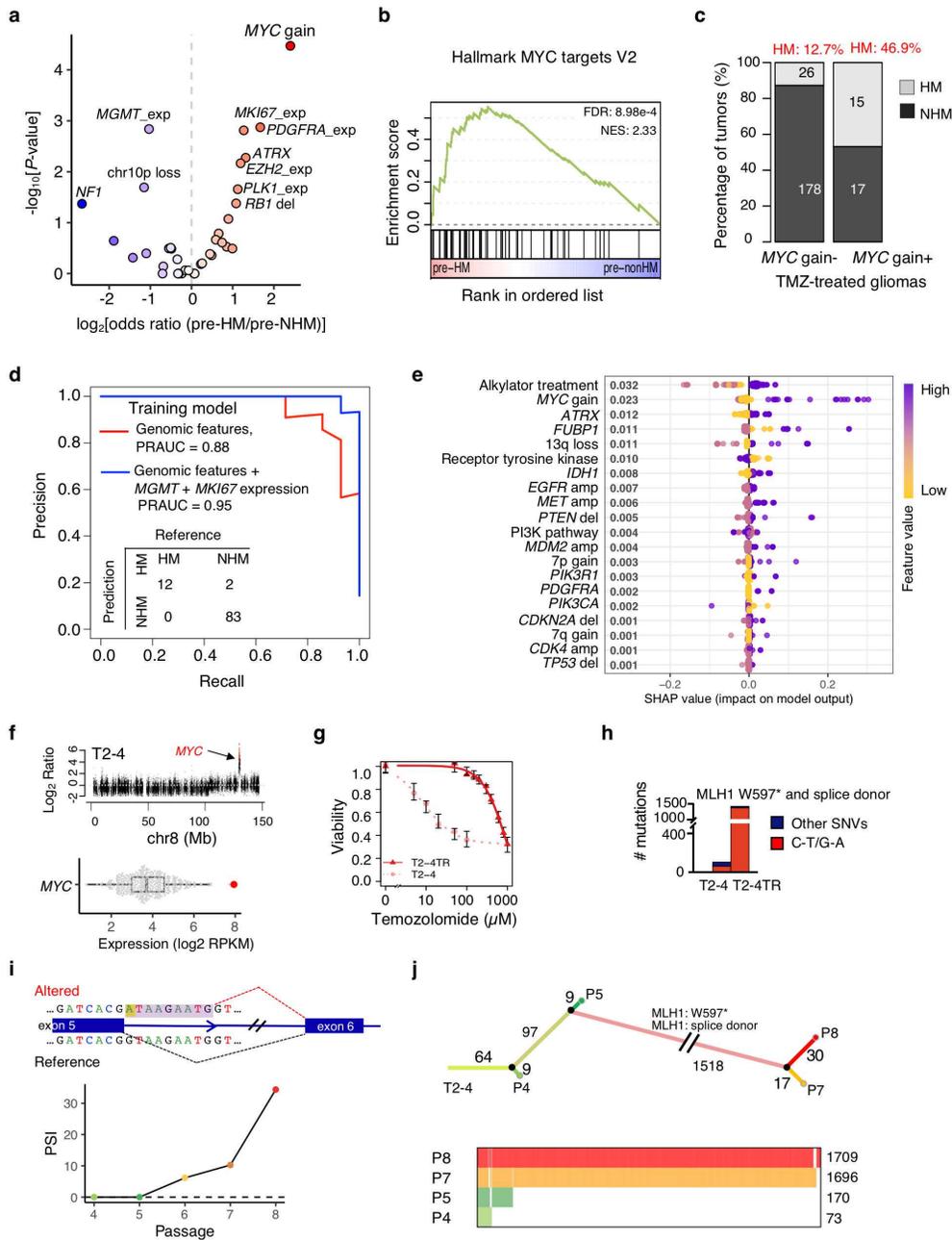


487

488 **Figure 2. Genomic characteristics of initial glioma informs grade progression after**
 489 **treatment.**

490 **(a)** Overall survival of low-grade gliomas that progressed to glioblastoma (GBM) and those
 491 that recurred as low-grade glioma (LGG). The colors represent glioma subtypes, while the
 492 line type (solid or dashed) show whether the patient progressed to GBM at recurrence. The

493 *P* values were calculated by log-rank test. **(b)** Design of the machine learning tasks. The
494 discovery cohort, including 117 LGGs, was used to train machine learning models to
495 predict whether the initial LGG patients would progress to GBM at recurrence. The models
496 were then applied to two independent cohorts for validation. **(c)** Receiving operation
497 characteristic (ROC) curves showing the cross-validation performance of the models
498 trained using different combinations of features. Model 1: the initial grade as the single
499 predictor; Model 2: the initial grade, age, gender, and treatment as predictors; Model 3: the
500 initial grade, age, gender, treatment, and genomic alterations as predictors. **(d)** Shapley
501 additive explanation (SHAP) scores of features in the prediction model of grade
502 progression. The color of the points represents value of each feature. **(e)** ROCs showing
503 the three models' prediction performance in the GLASS cohort. A confusion matrix is
504 shown at the right bottom. **(f)** Overall survival (left panel) and progression-free survival
505 (right panel) of LGGs from TCGA-LGG cohort, stratified by their predicted risk of
506 progressing to GBM. Patients that were in this longitudinal cohort were excluded. The *P*
507 values were calculated by log-rank test.

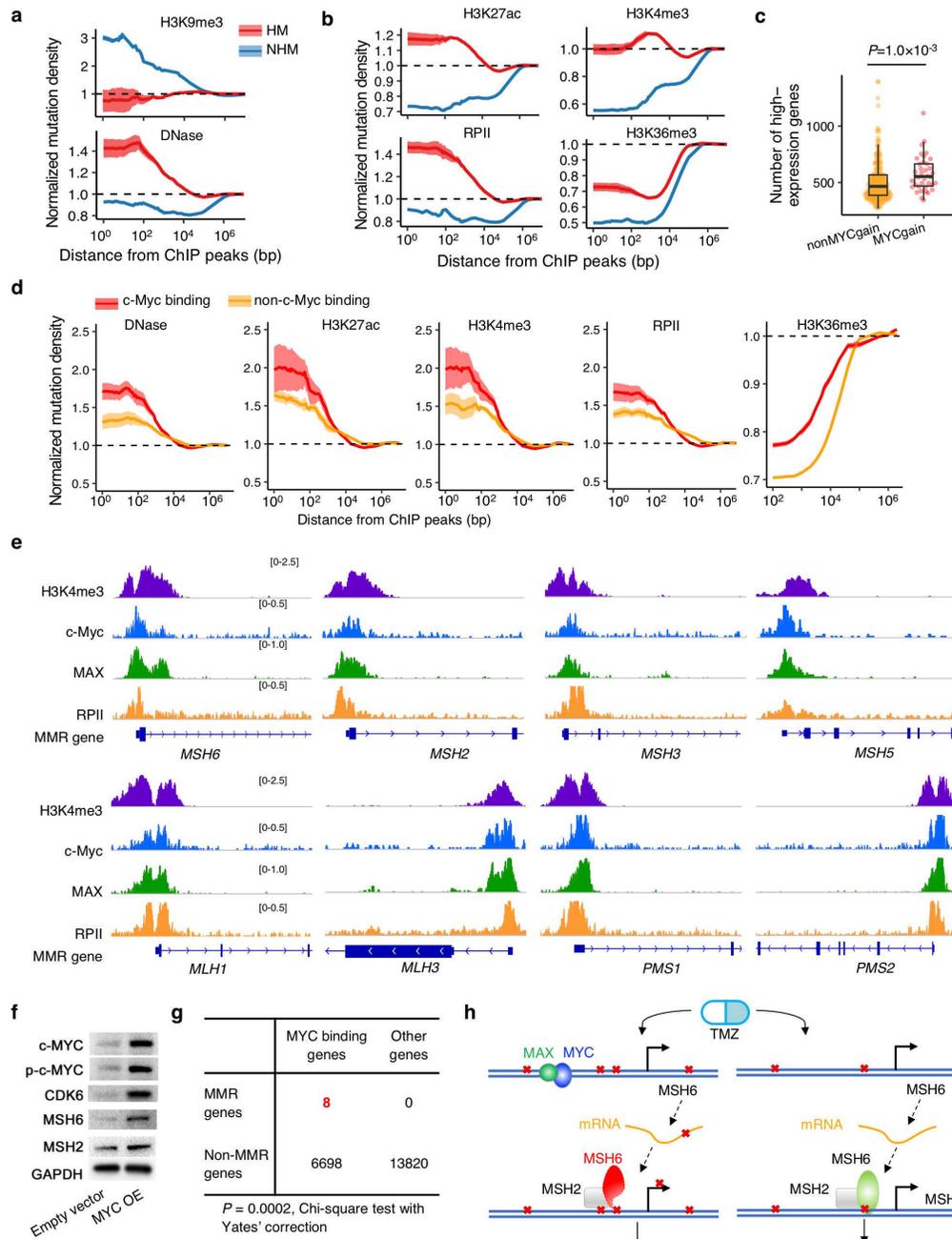


508

509 **Figure 3. Features of initial tumors predispose TMZ-associated hypermutation at**
 510 **recurrence.**

511 **(a)** Comparison of genomic and transcriptomic features of the TMZ-treated initial gliomas
 512 preceding hypermutators (pre-HM) and initial tumors preceding non-hypermutators (pre-
 513 NHM). **(b)** Gene set enrichment analysis. Gene expression profile of pre-HM gliomas were

514 compared to pre-NHM gliomas. FDR: false-discovery rate; NES: normalized enrichment
515 score. **(c)** The frequency of developing hypermutation in TMZ treated gliomas with or
516 without *MYC* copy number gain. **(d)** Precision-recall curve of HM prediction the model
517 trained using genomic features (red curve) and genomic features plus expression of *MGMT*
518 and *MKI67* (blue curve). **(e)** SHAP value of each feature in the model trained with genomic
519 features plus expression of *MGMT* and *MKI67* in (c). **(f)** The DNA copy number of *MYC*
520 in T2-4 (upper panel), and the expression level (lower panel) of T2-4 (the red dot) and other
521 RNA sequenced gliomas in this study (grey dots). **(g)** Response curves of T2-4 and the
522 induced T2-4TR cell line to TMZ. **(h)** Number of somatic mutations in T2-4 and T2-4 TR.
523 Dark red represents C to T or G to A mutations, while dark blue represents other somatic
524 mutations. **(i)** Aberrant splicing related to the MLH1 splicing donor mutation in T2-4. The
525 upper panel shows the aberrant splicing which included nine additional base pairs in the
526 spliced mRNA. The lower panel shows the percentage of spliced-in (PSI) reads of the
527 aberrant isoform in different passages of T2-4. **(j)** Mutations in T2-4 at different time points
528 during the TMZ-inducing experiment. The upper panel shows the phylogenetic tree
529 constructed from the mutations in passages 4, 5, 7 and 8 as compared to the untreated T2-
530 4 PDC, and the lower panel shows a heatmap of the mutations.



531

532 **Figure 4. TMZ-induced mutations are enriched in open active regions, while c-Myc**
 533 **binding further increases mutation density**

534 **(a)** Normalized mutation density of TMZ-associated hypermutations (HM) and non-TMZ
 535 associated recurrence-specific somatic mutations (NHM) around H3K9me3 modification
 536 sites (upper panel) and DNase hypersensitive sites (lower panel). The dashed horizontal

537 line represents the random level. **(b)** Normalized mutation density of HM and NHM around
538 H3K27ac modification sites (upper left), H3K4me3 modification sites (upper right), RNA
539 polymerase II (RPII) binding sites (lower left) and H3K36me3 modification sites (lower
540 right). **(c)** Number of high-expression genes in gliomas with and without MYC gain. The
541 *P* value was calculated by Wilcoxon's rank-sum test. **(d)** Normalized mutation density of
542 HM around DNase hypersensitive sites, H3K27ac modification sites, H3K4me3
543 modification sites, RPII binding sites and H3K36me3 modification sites with (red color)
544 and without (orange color) c-Myc binding. **(e)** Chromatin immunoprecipitation (ChIP)
545 intensity of H3K4me3 modification, c-Myc, MAX and RPII binding at eight mismatch
546 repair (MMR) genes. **(f)** Contingent table of c-Myc binding and MMR genes. **(g)** Western
547 blot showing c-MYC, p-c-MYC, CDK6, MSH6, MSH2, and MLH1 in U87 cell lines with
548 empty vector or MYC overexpression. OE: overexpression. **(h)** Proposed working model
549 of how MYC activation increases chance of hypermutation.

1 **Methods**

2 **Human Patients and Samples**

3 In this study we collected longitudinally paired glioma from at least two time points,
4 namely the initial tumor and the recurrent tumor. Where possible, blood samples were also
5 collected. The newly collected samples are from Capital Medical University Beijing
6 Tiantan Hospital (China), Samsung Medical Center (SMC, Korea) and CUHK Otto Wong
7 Brain Tumour Centre in Hong Kong.

8 Glioma patients from Beijing Tiantan Hospital were enrolled as part of the Chinese Glioma
9 Genome Atlas project (CGGA, <http://cgga.org.cn>). Ethical approval was obtained from the
10 institutional review board in Capital Medical University (IRB ID: KYSQ2019-200-01),
11 and informed consent for each patient was obtained before surgery. For each specimen, the
12 pathological diagnosis is confirmed independently by two neuropathologists based on the
13 2016 edition of WHO classification of central nervous system tumors. The specimen was
14 flash-frozen within 5 mins after resection and restored in liquid nitrogen until being used
15 for DNA/RNA extraction and other experiments.

16

17 The SMC cohort includes both published cases and newly enrolled samples. The study is
18 approved by the institutional review board of SMC (IRB file No. 2010-04-004 & 2005-04-
19 001), and written consent is obtained from each patient. The glioma specimens were
20 immediately snap-frozen after surgical resection and stored in in liquid nitrogen for further
21 analysis. This study also included samples from six patients with recurrent glioma from
22 CUHK Otto Wong Brain Tumour Centre. Ethical approval was obtained from the Joint

23 Chinese University of Hong Kong-New Territories East Cluster Clinical Research Ethics
24 Committee with Reference No. 2018.389.

25

26 In addition to these newly collected specimens, we also curated sequencing data of
27 longitudinally paired glioma from published cohorts. This include the Istituto Neurologico
28 C. Besta cohort (INCB, n=19)¹⁶, the MD Anderson Cancer Center cohort (MD Anderson,
29 n=10)³⁷, the Kyoto University cohort (KU, n=10)⁷, the University of California San
30 Francisco cohort (UCSF, n=23)⁹, The Cancer Genome Atlas cohort (TCGA, n=27)³⁸, the
31 Yale cohort (n=41)¹⁹, the Columbia University cohort (Columbia, n=14)³⁹, the Memorial
32 Sloan Kettering Cancer Center cohort (MSKCC, n=67)³⁵. The original Glioma
33 Longitudinal AnalySiS (GLASS) cohort⁸ included some samples from the above collection,
34 so we included the rest 139 patients for which sequencing data of both initial and recurrent
35 tumors were available and designated as GLASS cohort.

36

37 **Sample processing and Illumina sequencing**

38 For the samples from Tiantan hospital, DNA and/or RNA from the tumor were extracted
39 using the same protocol as described previously^{15,40}. Where available, DNA from matched
40 blood was also extracted. For patients PS115-PS132, the exome regions of both tumor and
41 blood DNA were captured using the Agilent SureSelect XT Human All ExonV5 kit, and
42 then sequenced using Illumina HiSeq 2500 platform. For patients PS132-PS157, whole
43 genome libraries were prepared and sequenced using Illumina HiSeq 2500 platform,
44 generating about 180Gb and 90Gb 150 bp paired end reads for tumor and normal DNA,
45 respectively. For RNA, the extracted total RNA was firstly depleted for tRNA and rRNA,

46 then reverse-transcribed to cDNA, and sequenced using Illumina HiSeq 2500 platform,
47 generating around 6 Gb 150 bp paired end reads for each sample.

48

49 Samples from SMC were processed and sequenced as previously described⁴¹. For samples
50 from CUHK, both total DNA and RNA were isolated from tumor samples while only total
51 DNA was purified from blood samples. Frozen tumors were homogenized in Buffer RLT
52 (Qiagen) with 40mM DTT and subsequently filtered by QIAshredder (Qiagen). Total
53 DNA/RNA from tissue lysate and from blood were then isolated using AllPrep DNA/RNA
54 Mini kit (Qiagen) according to manufacturer's protocol. The quantity and purity of the
55 nucleic acids were determined using Qubit 3 Fluorometer with either dsDNA HS Assay
56 Kit or RNA HS Assay Kit. For DNA samples with sub-optimal purity determined by
57 A260/280 ratio, they were further cleaned up using DNA Clean & Concentrator-5 kit
58 (ZYMO Research). RNA integrity was assessed by RNA Quality Number (RQN) using
59 Fragment Analyzer (Advanced Analytical) with HS RNA Kit (15NT). DNA samples with
60 a A260/280 of ~1.8, and RNA samples with a A260/280 of ~2.0 and RQN > 5 were
61 subjected to sequencing. Exome DNA were subsequently captured using Agilent
62 SureSelect XT Human All ExonV5 kit. DNA and RNA were sequenced using the same
63 protocol as samples from Tiantan Hospital, but generated 24Gb reads for tumor DNA, 12
64 Gb reads for blood DNA, and 12 Gb reads for RNA.

65

66 For the published samples, a detailed list of the sequencing types and protocols were
67 summarized in Table S1.

68

69 **Sequencing data preprocessing and mapping**

70 The quality of the sequencing reads was first checked by running FastQC v0.11.5. The low-
71 quality reads (containing at least one ambiguous base, or average quality <20) were then
72 removed using fastp v0.20.1. For DNA sequencing, the clean reads were mapped to the
73 human reference genome hg19 using Burrows Wheeler Aligner (BWA, v0.7.15-r1140)
74 mem algorithm, then sorted by coordinates using samtools v1.2. Subsequently, duplicate
75 reads were marked by running Picard MarkDuplicates 2.9.2 tool
76 (<https://broadinstitute.github.io/picard/>), generating the final bam file for mutation and
77 copy number identification. For RNA-seq, the reads were mapped to hg19 using STAR
78 2.6.1d. The gene annotation is ENSEMBL GRCh37.75.

79 For published samples, sequencing reads were either downloaded from NCBI SRA or
80 extracted from the published bam files that we downloaded from EGA, and then realigned
81 to hg19 using the same parameters as the newly sequenced samples.

82

83 **Identification of somatic mutations in longitudinal samples**

84 Somatic mutations in the newly sequenced samples were identified using SAVI⁴². Briefly,
85 raw mutation calls as compared to the normal sample were generated from samtools 1.2
86 mpileup, and then annotated using SNPEff for their impacts⁴³. SNPSIFT was used to add
87 further annotation such as whether the mutations are common SNPs (per dbSNP⁴⁴), cancer
88 related somatic mutations (per COSMIC⁴⁵), or observed in normal patients. A statistical
89 test was performed for each mutation call to determine the significance of being a somatic
90 mutation. SAVI is able to analyze mutations from multiple samples in one single run, which
91 is especially useful for longitudinal studies.

92 For the published data, we identified somatic mutations using the same pipeline as the
93 newly sequenced ones when raw sequencing data were available. For the YALE, MSKCC
94 and GLASS cohort, we retrieved mutation call information (including mutation loci,
95 reference and alternative alleles, reference and alternative depth) from individual sources,
96 annotated with SAVI using the same parameters, and then applied the same filters to
97 identify somatic mutations. Mutations in known glioma driver genes were reported.

98

99 **Copy number alteration detection**

100 Copy number variations in the newly sequenced samples were detected using CNVkit
101 0.9.5⁴⁶, which is suitable for both WES and WGS data. In both cases, only the exon regions
102 were considered. The BED file containing information about the captured regions was
103 downloaded from Agilent website (<https://earray.chem.agilent.com/suredesign/index.htm>).
104 The final results were segmented into SEG files.

105

106 **Detection of genes fusions, *MET*_{ex14} and *EGFR*_{vIII} from RNA-seq data**

107 RNA sequencing data from previous publications were downloaded, and the reads were
108 extracted using samtools 1.2⁴⁷. Starting from FASTQ files, STAR-fusion 1.5.0⁴⁸ was
109 utilized to identify and annotate gene fusion candidates, using the fastq files as input. The
110 fusion candidates were then filtered by removing fusions that were present in normal
111 tissues, fusions involving mitochondria genes and uncharacterized genes, and fusions of
112 two paralog genes.

113

114 *MET*ex14 and *EGFR*vIII were identified by counting the spanning reads over the junction
115 of *MET* exon 13 and exon 15, and *EGFR* exon 1 and exon 8, respectively. Briefly, RNA
116 sequencing reads were aligned to the reference genome (hg19), and then the spanning reads
117 were extracted based on the CIGAR record. The spanning reads were manually checked in
118 Integrative Genomic Viewer to remove false positives such as PCR artifacts and potential
119 mapping errors. Only samples with at least three supporting reads were considered as
120 positive.

121

122 **Integration of mutation profiles**

123 For samples from the Yale, MSKCC and GLASS cohorts, all the mutations were obtained
124 from the original publications, which contains information about the coordinates of the
125 mutations, the reference and alternative alleles, as well as the read count of the reference
126 and alternative alleles. All the mutations were reannotated using SAVI, and the same
127 mutation filter as those used in the new samples were applied. This unified mutation
128 annotation and filtering method ensures the results from multiple cohorts comparable.

129 The copy number alteration calls from all samples were first transformed to SEG format
130 which included the coordinate of each segment and the estimated segment mean value
131 (reflecting the magnitude of copy number change). Subsequently, each segment was
132 annotated for the carried genes using the reference human genome GRCh37.75. The same
133 cutoffs for copy number gain, amplification and deletion were applied.

134 For gene fusions, where available, the raw sequencing reads were downloaded or extracted,
135 and then subjected to gene fusion detection using the same protocol.

136

137 **Construction of Tumor Evolution Directed Graphs**

138 Tumor Evolution Directed Graphs⁴⁹ were constructed using CELLO²⁰. The input to
139 CELLO was the mutation calls from SAVI, and the samples were organized in the normal-
140 initial-recurrence order. CELLO automatically extracted the shared, private to initial and
141 private to recurrence mutations in each patient, which will be used to construct the tree.
142 The evolution trees were visualized in Cytoscape.

143

144 **Calculation of entropy and mutual information**

145 For each patient, an alteration might be gained or lost during evolution, both were counted
146 as one “status change event”. If the alteration was present at both time points, or absent at
147 both time points, it is counted one “status stable event”. For a given alteration X, suppose
148 the number of status change events in the cohort is m , and the number of “status stable
149 event” is n , then the entropy is

$$150 H(X) = -\frac{m}{m+n} \log_2 \left(\frac{m}{m+n} \right) - \frac{n}{m+n} \log_2 \left(\frac{n}{m+n} \right).$$

151 The mutual information between two alterations, X and Y, is calculated as below:

$$152 I(X, Y) = H(X) + H(Y) - H(X, Y).$$

153 In our analysis, the maximum number of features was set to four. All the possible
154 combinations were enumerated to find the four features with the highest mutual
155 information.

156

157 **Machine learning model training, validation, and feature importance**

158 XGBoost was used to train and test the models⁵⁰. XGboost is a boosting machine learning
159 methodology, supporting influential supervised classification tasks. The dataset was

160 divided into training set and test set two times for the prediction of HM and Grade with the
161 ratio of 3:1. For the training dataset, we leveraged five-fold cross validation to select the
162 optimal parameters, which were applied to the final configuration of XGBoost with best
163 validation AUC (Area Under receiver operating characteristic Curve). Then, the
164 independent test dataset was imported into the model to evaluate the performance and
165 generalization ability of the model for unseen data. The model training and testing was
166 implemented in R using the *xgboost* package, and SHAP (SHapley Additive exPlanations)
167 score¹⁷ was extracted from the XGBoost models using *SHAPforxgboost* package.

168

169 **Survival analysis**

170 Survival analysis was performed using the *survival* package in R. The significance levels
171 were calculated by two-sided log-rank test, and the Kaplan-Meier curves were plotted using
172 *survminer* package.

173

174 **Identification of hypermutation**

175 Hypermutation was detected using CELLO, with default parameters. CELLO determines
176 the hypermutation status based on mutation load and a hypermutation score which we
177 defined previously²⁰. For the MSKCC cohort which includes panel sequencing data, the
178 mutation load reported in the publication was directed used, while the hypermutation score
179 was calculated using the same methodology as other samples.

180

181 **Gene expression analysis**

182 Differential gene expression analysis was performed by the DEseq package 1.26.0. Gene
183 Set Enrichment Analysis was performed by the *fgsea* package (version 1.12.0) in R 3.6.3.
184 The gene sets regarding human cancer were gene sets which were downloaded from
185 MSigDB database version 6.1 (<http://software.broadinstitute.org/gsea/msigdb>). To
186 quantify global transcription activity, the mean expression level of 14 housekeeping genes:
187 ACTB, GAPDH, PGK1, PPIA, RPL13A, RPLP0, B2M, YWHAZ, SDHA, TFRC, GUSB,
188 HMBS, HPRT1 and TBP in each sample was calculated and the genes with expression
189 levels higher than this reference value was designated as “high-expression” genes.

190

191 **ChIP-seq data acquisition and intersection with hypermutation loci**

192 ChIP-seq data (listed in Table S3), including both binding density and binding peaks, were
193 downloaded from ChIP-Atlas (<https://chip-atlas.org/>). The peaks were extended by 1-
194 5,000,000 bp upstream and downstream, and the overlaps, regions in the centromeres,
195 telomeres or out of chromosome borders were removed.

196

197 Recurrent-specific mutations from WGS data were divided into MYC-gained HM, non-
198 MYC-gained-HM and non-HM mutations, and randomly down-sampled to 150,000 for
199 each type. The sampling process was repeated for five times. The mutation sites were
200 transformed to BED format, and their intersection with the ChIP peaks were calculated
201 using bedtools v2.26.0. The mutation density was calculated by dividing the total number
202 of mutations with the total width of the regions. The mean and standard deviation of the
203 five replicates were calculated.

204

205 **Cell cultures**

206 The human glioma cell line U87 and U251 were obtained from the Cell Resource Center,
207 Peking Union Medical College (Beijing, China), and the cell has been authenticated by the
208 short tandem repeat analysis Chronic temozolomide treatment cells. Both the two cells
209 were routinely cultured in Dulbecco's modified Eagle's medium (DMEM) supplemented
210 with 10% fetal bovine serum (HyClone, Logan, Utah), 100 units/ml penicillin and 100
211 mg/ml streptomycin (Invitrogen, Carlsbad, CA) at 37°C in a humidified atmosphere of 5%
212 CO₂. The PDC T2-4, was obtained from fresh surgical specimens of a human primary GBM
213 and cultured as tumor spheres in DMEM/F12 medium supplemented with B27 supplement
214 (Life Technologies), bFGF and EGF (20 ng/ml each) at 37°C in a humidified atmosphere
215 of 5% CO₂.

216

217 **Chronic temozolomide treatment of glioma cell lines**

218 The human glioma cell line U87MG and U251MG were obtained from the Cell Resource
219 Center, Peking Union Medical College (Beijing, China). U87MG and U251MG cells were
220 treated with TMZ on Initial concentration of 100 μM for 72h. TMZ-sensitive cells
221 generated death, and the survived cells are collected for further resistant induction with
222 increasing TMZ concentration (≤ 200 μM). After more than 20 generations and 6 months
223 later, each cell line's cells survived and growth in DMEM culture medium containing TMZ
224 of final concentration of 200 μM, were referred as U87TR and U251TR cells.

225 The PDCs were treated with TMZ on an initial concentration of 50 μM for 5 days, and the
226 remain cells were incubated in medium with increasing concentration of TMZ until to 200
227 μM. The IC₅₀ values of TMZ were determined in the following situations: 1) after 6

228 generations of growth in medium with 200 μ M TMZ; 2) the PDCs reached the status that
229 they have similar growth rates in normal medium and medium containing 200 μ M TMZ.
230 Finally, the PDCs were considered TMZ-resistant when they have significantly elevated
231 IC50 values compared with their corresponding initial cells.

232

233 **Temozolomide sensitivity test**

234 The temozolomide sensitivity were studied by CCK-8 kit (Dojindo Laboratories,
235 Kumamoto, Japan) according to the manufacturer's protocol. 3000 cells were grown in
236 each well of 96-well plates with normal medium for 24hr. Then the cells were incubated at
237 medium containing TMZ at predesigned concentrations for 72 hrs, and there were six-
238 repeated wells for each predesigned concentrations. Next, 10 μ l CCK-8 reagent was added
239 to the medium of each well, and the absorbances at 450 nm and 630 nm of the medium
240 were measured after incubation for 2 hours. The number of living cells could be reflected
241 by the difference between optical density of 450 (OD₄₅₀) and 630 (OD₆₃₀) nm wavelength.
242 The response curves were fitted using the *drc* package 3.0-1.

243

244 **Western blot**

245 The whole-cell lysates were using RIPA buffer [150 mM NaCl, 0.1% (wt/vol) NP-40, 50
246 mM Tris (pH 8.0), 0.5% (wt/vol) sodium deoxycholate, 1% (wt/vol) sodium dodecyl sulfate,
247 1 mM dithiothreitol, 0.1 mM phenylmethylsulfonyl fluoride]. A bicinchoninic acid (BCA)
248 array was used to measure the total protein content. The equal total proteins (20 μ g) were
249 boiled and then were electrophoresed on a 8% and 12% sodium dodecyl sulfate–
250 polyacrylamide gel electrophoresis gel, respectively. After membrane transfer and

251 blocking, the primary antibodies and used to probe the target proteins in blocking solution
252 overnight at 4°C. The membrane were detected using an ECL Western Blotting Detection
253 System (Bio-Rad) after Goat anti-rabbit IgG-HRP or goat anti-mouse IgG-HRP incubation
254 1 hr at room temperature. Glyceraldehyde 3-phosphate dehydrogenase (GAPDH) was used
255 as the loading control.

256

257 **Generation of lentiviral vectors encoding c-MYC**

258 The sequencing of c-MYC (NM_002467) was cloned from glioma samples. The sequence
259 was inserted into the lentiviral vector plasmid GV358 (Genechem). The GV358 plasmid
260 carrying the c-MYC CDS sequencing was transduced into 293T cells with the lentiviral
261 packaging plasmid mix (Genechem). The culture medium was collected once a day on two
262 consecutive days and centrifuged with an ultracentrifuge at 25,000 rpm for 1.5 h. The
263 precipitate was resuspended, aliquoted and stored at -80°C. This preparation was added to
264 the U87 cell culture or PDCs supplemented with 3 mg/ml polybrene (Sigma) in the medium.
265 The medium was changed after 24 h. After 72 h, cells containing vectors were selected
266 during 3 days of treatment with puromycin. The c-MYC expression in the U87 cells and
267 PDCs were confirmed by Western blot.

268

269 **Validation of MSH6 mutation by Sanger sequencing**

270 DNA extraction and PCR were performed to get the fragment including the MSH6
271 mutation site identified by the whole exon sequencing, the primers sequences were as
272 follows: forward 5'-CTCCCCATGGGCTGCTAAG-3', Reverse 5' -
273 TATGTCCTAGGCGCACAGC-3'. The product bands were extracted from agarose gel
274 after electrophoresis and verified by Sanger sequencing with the forward primer.

275

276 **Code availability**

277 The hypermutation status, and the Tumor Evolution Directed Graphs were determined
278 using CELLO which is available at <https://github.com/WangLabHKUST/CELLO>. All
279 other scripts are available from the authors upon reasonable request.

280

281 **Data availability**

282 The processed genomic and transcriptomic data are available in the CELLO2 website
283 (<https://wanglab.shinyapps.io/cello>). Raw sequencing data of the newly sequenced samples
284 will be deposited to European Genome-Phenome Archive (EGA, <https://ega-archive.org/>).
285 Data from previous SMC samples were available in EGA
286 (<https://www.ebi.ac.uk/ega/datasets/EGAS00001001800>), data from TCGA were
287 downloaded from NCI Genomics Data Commons (GDC) data portal
288 (<https://portal.gdc.cancer.gov>). Previously published CGGA data have been uploaded to
289 the Genome Sequence Archive in BIG Data Center, Beijing Institute of Genomics (BIG),
290 Chinese Academy of Sciences, under accession number BioProject ID: PRJCA001636
291 (<https://bigd.big.ac.cn/bioproject/browse/PRJCA001636>) and PRJCA001747
292 (<https://bigd.big.ac.cn/bioproject/browse/PRJCA001747>). All the other data supporting the
293 findings of this study are available within the article and its information files and from the
294 corresponding author upon reasonable request.

295

296 **Additional information**

297 Supplementary Information is available for this paper. Correspondence and requests for
298 materials should be addressed to Jiguang Wang (jgwang@ust.hk) or Tao Jiang
299 (taojiang1964@163.com).

300 **Supplementary item titles and legends**

301 **Supplementary Table S1:** Clinical characteristics of the longitudinal glioma cohort,
302 **related to Figure 1.**

303 **Supplementary Table S2:** Sequencing platforms used for the samples in the longitudinal
304 glioma cohort, **related to Figure 1.**

305 **Supplementary Table S3:** List of ChIP-seq data used in hypermutation distribution
306 analysis, **related to Figure 4.**

The color of each variant indicates whether it is specific to initial, recurrence or shared. MMR: mismatch repair; CNLOH: copy-neutral loss of heterozygosity. (b) Comparison of the number of initial- and recurrent specific somatic mutations during evolution in each glioma subtype. The horizontal dashed line represents five mutations. The P values were calculated by Wilcoxon's rank-sum test. (c) Proposed model explaining the changes of the clinical and molecular features over time. Predisposing factors in the initial tumor, together with stochasticity, determine whether the status of the feature will change at recurrence. (d) Alteration status, entropy and mutual information of clinical and genomic features. Stacked bars outside the circle summarizes the frequency of a clinical feature or molecular variant to be lost, gain, remaining wildtype or remaining altered. Inside the circle shows the overlapped bar plots of mutual information (black) and conditional entropy (grey), where high proportion of overlap indicates better predictability. The features are ranked by mutual information.

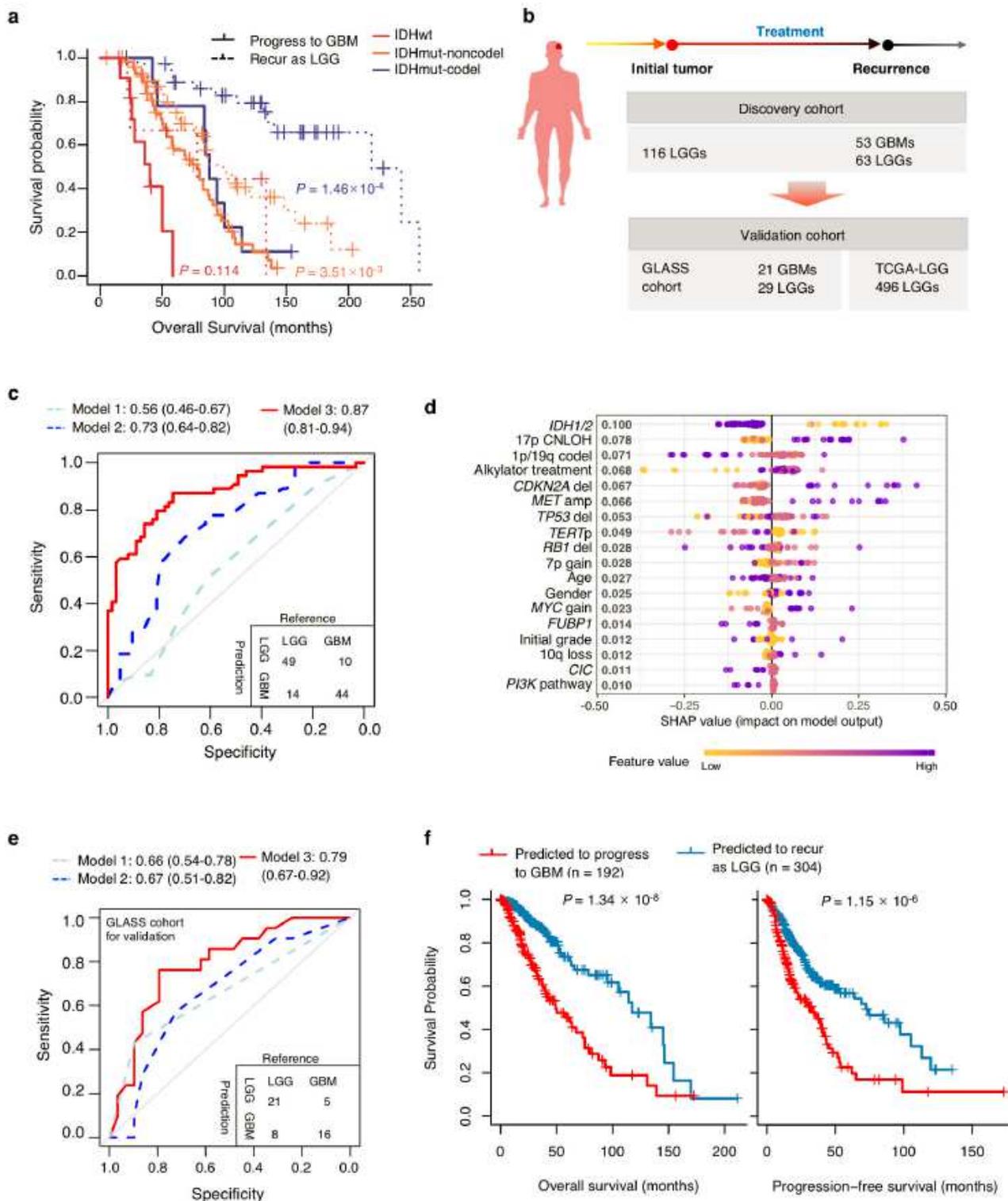


Figure 2

Genomic characteristics of initial glioma informs grade progression after treatment. (a) Overall survival of low-grade gliomas that progressed to glioblastoma (GBM) and those that recurred as low-grade glioma (LGG). The colors represent glioma subtypes, while the line type (solid or dashed) show whether the patient progressed to GBM at recurrence. The P values were calculated by log-rank test. (b) Design of the machine learning tasks. The discovery cohort, including 117 LGGs, was used to train machine learning

models to predict whether the initial LGG patients would progress to GBM at recurrence. The models were then applied to two independent cohorts for validation. (c) Receiving operation characteristic (ROC) curves showing the cross-validation performance of the models trained using different combinations of features. Model 1: the initial grade as the single predictor; Model 2: the initial grade, age, gender, and treatment as predictors; Model 3: the initial grade, age, gender, treatment, and genomic alterations as predictors. (d) Shapley additive explanation (SHAP) scores of features in the prediction model of grade progression. The color of the points represents value of each feature. (e) ROCs showing the three models' prediction performance in the GLASS cohort. A confusion matrix is shown at the right bottom. (f) Overall survival (left panel) and progression-free survival (right panel) of LGGs from TCGA-LGG cohort, stratified by their predicted risk of progressing to GBM. Patients that were in this longitudinal cohort were excluded. The P values were calculated by log-rank test.

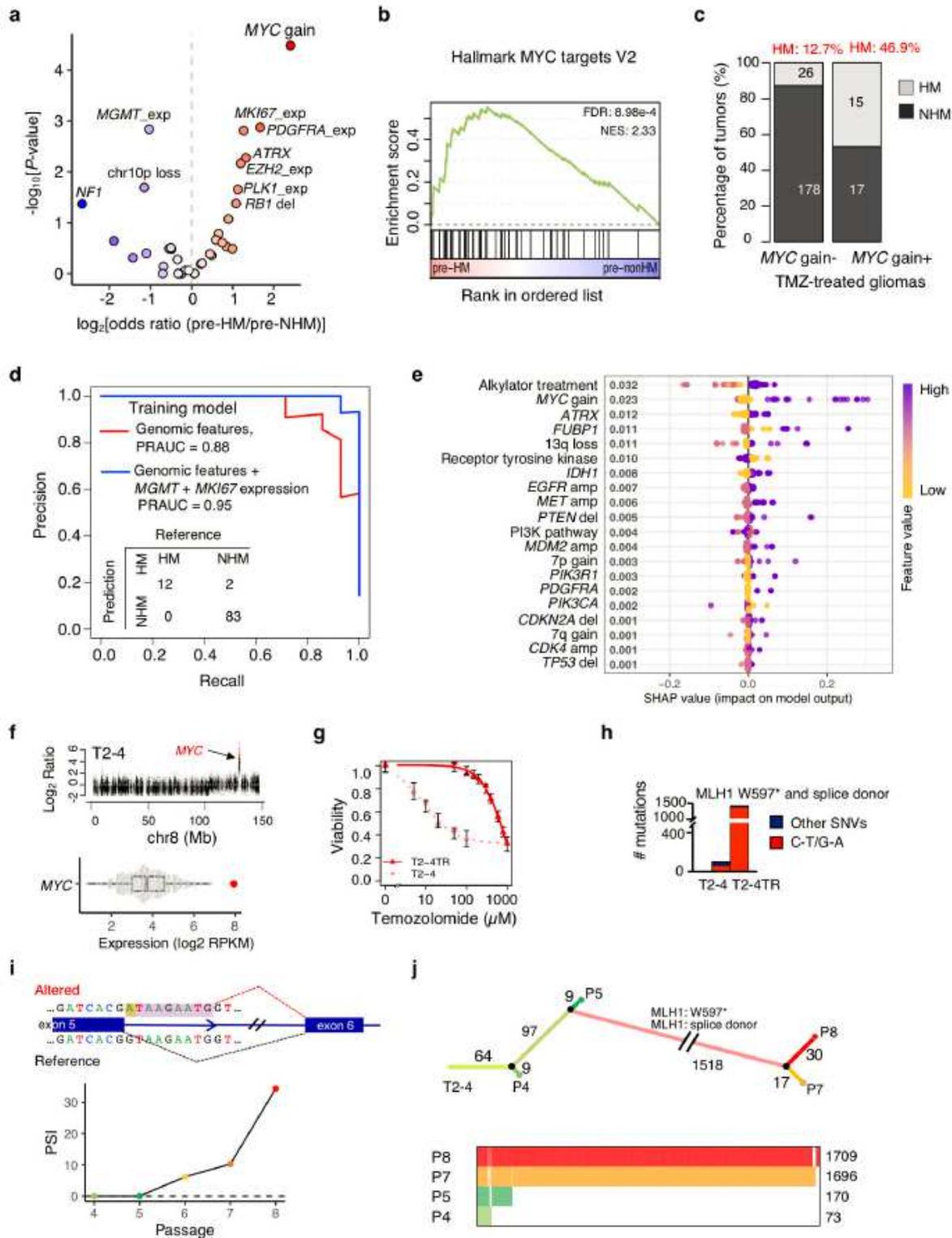


Figure 3

Features of initial tumors predispose TMZ-associated hypermutation at recurrence. (a) Comparison of genomic and transcriptional features of the TMZ-treated initial gliomas preceding hypermutators (pre-HM) and initial tumors preceding non-hypermutators (pre-NHM). (b) Gene set enrichment analysis. Gene expression profile of pre-HM gliomas were compared to pre-NHM gliomas. FDR: false-discovery rate; NES: normalized enrichment score. (c) The frequency of developing hypermutation in TMZ treated gliomas

with or without MYC copy number gain. (d) Precision-recall curve of HM prediction the model trained using genomic features (red curve) and genomic features plus expression of MGMT and MKI67 (blue curve). (e) SHAP value of each feature in the model trained with genomic features plus expression of MGMT and MKI67 in (c). (f) The DNA copy number of MYC in T2-4 (upper panel), and the expression level (lower panel) of T2-4 (the red dot) and other RNA sequenced gliomas in this study (grey dots). (g) Response curves of T2-4 and the induced T2-4TR cell line to TMZ. (h) Number of somatic mutations in T2-4 and T2-4 TR. Dark red represents C to T or G to A mutations, while dark blue represents other somatic mutations. (i) Aberrant splicing related to the MLH1 splicing donor mutation in T2-4. The upper panel shows the aberrant splicing which included nine additional base pairs in the spliced mRNA. The lower panel shows the percentage of spliced-in (PSI) reads of the aberrant isoform in different passages of T2-4. (j) Mutations in T2-4 at different time points during the TMZ-inducing experiment. The upper panel shows the phylogenetic tree constructed from the mutations in passages 4, 5, 7 and 8 as compared to the untreated T2-4 PDC, and the lower panel shows a heatmap of the mutations.

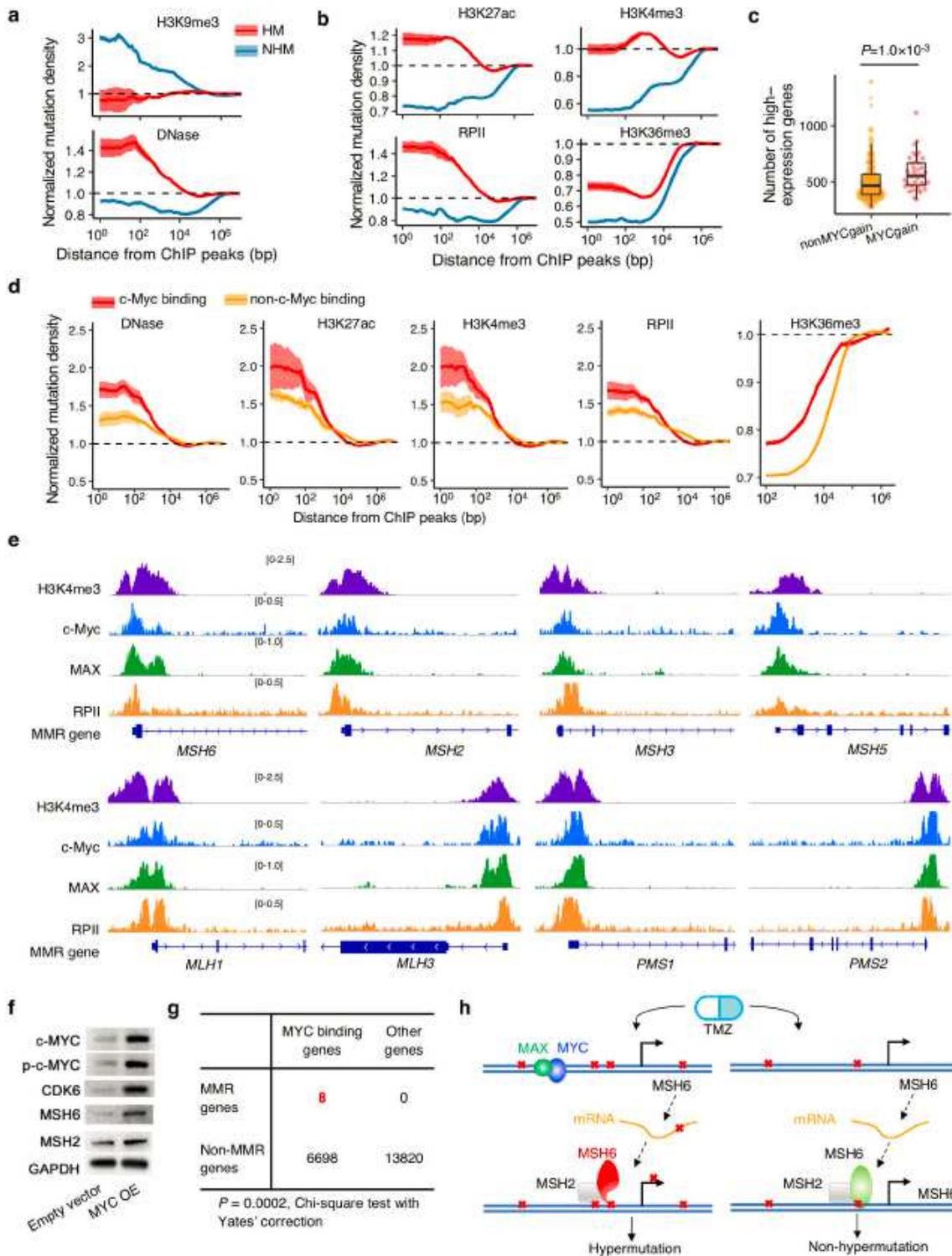


Figure 4

TMZ-induced mutations are enriched in open active regions, while c-Myc binding further increases mutation density (a) Normalized mutation density of TMZ-associated hypermutations (HM) and non-TMZ associated recurrence-specific somatic mutations (NHM) around H3K9me3 modification sites (upper panel) and DNase hypersensitive sites (lower panel). The dashed horizontal line represents the random level. (b) Normalized mutation density of HM and NHM around H3K27ac modification sites

(upper left), H3K4me3 modification sites (upper right), RNA polymerase II (RPII) binding sites (lower left) and H3K36me3 modification sites (lower right). (c) Number of high-expression genes in gliomas with and without MYC gain. The P value was calculated by Wilcoxon's rank-sum test. (d) Normalized mutation density of HM around DNase hypersensitive sites, H3K27ac modification sites, H3K4me3 modification sites, RPII binding sites and H3K36me3 modification sites with (red color) and without (orange color) c-Myc binding. (e) Chromatin immunoprecipitation (ChIP) intensity of H3K4me3 modification, c-Myc, MAX and RPII binding at eight mismatch repair (MMR) genes. (f) Contingent table of c-Myc binding and MMR genes. (g) Western blot showing c-MYC, p-c-MYC, CDK6, MSH6, MSH2, and MLH1 in U87 cell lines with empty vector or MYC overexpression. OE: overexpression. (h) Proposed working model of how MYC activation increases chance of hypermutation.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [TableS1.Patient.ClinicalInformation.xlsx](#)
- [TableS2.Sequencing.Data.xlsx](#)
- [SupplementaryInformationDec30.docx](#)