

1     **Title page**

2     **Title:**

3     Integrating gene mutation spectra from tumors and the general population with gene  
4     expression topological networks to identify novel cancer driver genes

5

6     **Author affiliations**

7     Dan He<sup>1,2#</sup>, Ling Li<sup>1,2#</sup>, Zhiya Lu<sup>1,2</sup>, Shaoying Li<sup>1,2</sup>, Tianjun Lan<sup>3</sup>, Feiyi Liu<sup>1,2</sup>,  
8     Huasong Zhang<sup>1,2</sup>, Bingxi Lei<sup>4</sup>, David N. Cooper<sup>5</sup>, Huiying Zhao<sup>1,2,\*</sup>

9     <sup>1</sup>Department of Medical Research Center, Sun Yat-Sen Memorial Hospital, Sun  
10    Yat-Sen University, Guangzhou 510006, China

11    <sup>2</sup>Guangdong Provincial Key Laboratory of Malignant Tumor Epigenetics and Gene  
12    Regulation, Guangzhou 510006, China

13    <sup>3</sup>Department of Oral and Maxillofacial Surgery, Sun Yat-Sen Memorial Hospital, Sun  
14    Yat-Sen University, Guangzhou 510010, China

15    <sup>4</sup>Department of Neurosurgery, Sun Yat-Sen Memorial Hospital, Sun Yat-Sen  
16    University, Guangzhou 510006, China

17    <sup>5</sup>Institute of Medical Genetics, Cardiff University, Heath Park, Cardiff CF14 4XN,  
18    UK

19    # Contributed equally

## 20 Correspondence email

21 Huiying Zhao, PhD

22 Department of Medical Research Center, Sun Yat-Sen Memorial Hospital, 107 Yan

23 Jiang West Road, Guangzhou, P.R. China, 500001

24 Zhaohy8@mail.sysu.edu.cn

25

## 26 Abstract

27 **Background:** Understanding the genetics underlying cancer development and  
 28 progression is the most important goal of biomedical research to improve patient  
 29 survival rates. Recently, researchers have proposed computationally combining the  
 30 mutational burden with biological networks as a novel means to identify cancer driver  
 31 genes. However, these approaches treated all mutations as having the same functional  
 32 impact on genes and incorporated gene-gene interaction networks without considering  
 33 tissue specificity, which may have hampered our ability to identify novel cancer  
 34 drivers. **Methods:** We have developed a framework, DGAT-cancer that integrates the  
 35 predicted pathogenicity of somatic mutation in cancers and germline variants in the  
 36 healthy population, with topological networks of gene expression in tumor tissues,  
 37 and the gene expression levels in tumor and paracancerous tissues in predicting cancer  
 38 drivers. These features were filtered by an unsupervised approach, Laplacian selection,  
 39 and those selected were combined by Hotelling and Box-Cox transformations to score  
 40 genes. Finally, the scored genes were subjected to Gibbs sampling to determine the

41 probability that a given gene is a cancer driver.

42 **Results:** This method was applied to nine types of cancer, and achieved the best area  
43 under the precision-recall curve compared to three commonly used methods, leading  
44 to the identification of 571 novel cancer drivers. One of the top genes, *EEF1A1* was  
45 experimentally confirmed as a cancer driver of glioma. Knockdown of *EEF1A1* led to  
46 a ~ 41-50% decrease in glioma size and improved the temozolomide sensitivity of  
47 glioma cells.

48 **Conclusion:** By combining the pathogenic status of mutational spectra in tumors  
49 alongside the spectrum of variation in the healthy population, with gene expression in  
50 both tumors and paracancerous tissues, DGAT-cancer has significantly improved our  
51 ability to detect novel cancer driver genes.

52

53 **Keywords:** Cancer drivers, Pathogenic status of mutations, Laplacian selection,  
54 Hotelling and Box-Cox transformations, Gibbs sampling.

55

## 56 **Background**

57 The identification of cancer driver genes is important for the early diagnosis of cancer,  
58 for identifying efficacious anti-cancer therapeutics and for investigating the  
59 underlying mechanisms of tumorigenesis. Traditionally, cancer driver genes have  
60 been recognized on the basis of their being recurrently altered in tumors[1, 2].  
61 However, the ability of many somatic mutations to alter gene function is often

uncertain, making it hard to identify cancer driver genes unambiguously. Clearly, we require information other than somatic mutation data in order to reliably detect cancer driver genes.

Traditional approaches to identifying cancer driver genes have relied upon the statistical testing of the mutational burden of individual genes[1, 3], a strategy that assumes that driver mutations occur more frequently than expected by chance alone[4]. This type of approach, exemplified by MutSigCV[4] and MuSiC[5], is in common use and has been successful in identifying many genes that harbor recurrent mutations in cancer(s). Such approaches are useful for identifying those genes which are mutated across a large number of samples but can easily miss genes mutated in only a small number of samples. However, most cancer driver genes are only mutated in a small proportion of patients[6]. Indeed, somatic mutation frequencies are influenced by specific characteristics of the gene, such as length, replication time and mutation rate in the healthy population[7-10]. A high gene mutation rate in the healthy population suggests that many somatic mutations within that gene are likely to be neutral passengers rather than drivers[11, 12]. Various other methods have been developed that improve our ability to recognize genes characterized by a higher-than-usual rate of somatic mutation[4, 5, 11-13]. OncodriveFML[14] is one such method designed to use the patterns of mutations across tumors in coding and non-coding regions to identify cancer driver genes. Similarly, OncodriveCLUSTL[15] is a method that detects cancer driver genes by clustering the mutations in cancer cohorts according to the number of the mutation distribution.

84 Mutations residing within transcribed regions of the genome are known to be more  
85 likely to influence the gene expression profiles of tumors[4]. The mutation rates of  
86 specific genes can be cross-compared with tumor expression signatures[16, 17].  
87 Cancer driver genes could in principle be identified by integrating gene mutation data  
88 with gene expression data. As an example, a previous study performed enrichment  
89 analysis to integrate genomic and transcriptomic alterations from whole-exomes and  
90 functional data from protein function predictions with gene interaction networks to  
91 reveal breast cancer driver genes[18]. Recent approaches have been developed by  
92 combining mutation scores with biological network protein-protein interaction (PPI)  
93 data to predict cancer driver genes[19, 20]. These approaches did not consider  
94 functional data for mutations and cancer tissue-specific PPI networks, which may  
95 reduce the ability in predicting novel cancer drivers. Thus, no approach to identifying  
96 cancer driver genes has yet been devised that fully considers the pathogenic status of  
97 mutational spectra in tumors alongside the spectrum of variation in the healthy  
98 population, in combination with gene expression in both tumors and normal tissues.

99 To address these shortcomings, we devised a new model, DGAT-cancer  
100 (Distinguish cancer drivers using Genomics and Transcriptome data), which integrates  
101 the predicted pathogenicity scores of somatic mutations in cancers and germline  
102 mutations in the healthy population, with gene expression in tumors and  
103 paracancerous tissues in order to detect cancer driver genes. The work scheme of  
104 DGAT-cancer is shown in Fig. 1. First, for each gene, DGAT-cancer calculated a  
105 unidirectional Earth Mover's Difference score (uEMD) to evaluate the difference

106 between the predicted pathogenic scores (obtained from 19 predictors in dbNSFP[21])  
107 of somatic mutational spectra in cancers and that of germline variants in healthy  
108 populations. The influence of mutation on gene expression has been evaluated by  
109 integrating gene expression data in tumors with somatic mutations through  
110 topological data analysis (TDA)[22]. Briefly, using TDA, we have constructed a gene  
111 expression topological network by clustering samples with similar gene expression  
112 profiles. The gene expression topological network was then used to evaluate the  
113 frequencies of mutational spectra occurring in different sample clusters. For those  
114 mutations occurring in adjacent sample clusters, we calculated the divergence of their  
115 frequencies across clusters, which was used as one feature of the gene in which the  
116 mutations resided. The Jensen–Shannon divergence between the gene expression  
117 profile and mutation profile of each gene on the topological network was further  
118 computed and used as another feature of the gene. In total, 24 features for each gene  
119 were included in DGAT-cancer. The most effective features were identified by  
120 Laplacian selection in an unsupervised way. The selected features were integrated by  
121 means of the Hotelling and Box-Cox transformations to score the genes. Finally, by  
122 using gene scores as weights, we performed Gibbs sampling to identify cancer drivers.  
123 This method was then applied to 6,643 samples containing mutation and gene  
124 expression data from tumors and paracancerous tissues derived from 9 cancer cohorts  
125 and succeeded in identifying 734 genes as being significant cancer drivers. Of these,  
126 571 were previously unreported as cancer-associated genes. Further, these genes were  
127 found to be highly enriched in pathways related to cancer, as well as significantly

128 enriched in drug-response genes, demonstrating that this new approach facilitates the  
129 identification of clinically relevant genes. One of the top genes, *EEF1A1*, was  
130 predicted to be a driver of glioma. This is the first time that *EEF1A1* has been  
131 considered to be a driver of glioma. Its relationship with glioma was confirmed by  
132 analysis of surgical specimens, a cell model and an animal model.

133

## 134 **Methods**

### 135 **Data collection**

136 ***Somatic mutation data.*** We collected simple somatic mutation (SSM) data from 12  
137 cancer cohorts from the Broad Institute GDAC Firehose Portal, the International  
138 Cancer Genome Consortium (ICGC) Data Portal and The Cancer Genome Atlas  
139 (TCGA) (Additional file 1: Table S1). The coordinates of the data are by reference to  
140 the genome assembly version hg19. The germline mutation data of 2,557 individual  
141 samples were collected from Phase 3 of the 1000 Genomes Project (GRCh38).  
142 Although gnomAD[23] contains germline variants from a larger cohort, it does not  
143 provide genetic data from each sample and was therefore inappropriate for use in  
144 generating mutation score profiles for this study. Somatic mutation entries that were  
145 duplicated between multiple databases were removed such that only one  
146 non-redundant entry was retained.

147

148 ***The pathogenicity of mutations.*** There are many methods available with which to

149 predict the pathogenicity of mutations. We employed dbNSFP (version: dbnsfp30a,  
150 also called LJB\*)[21] of ANNOVAR[24] to annotate mutations from the 1000  
151 Genomes Project (1000GP) and the somatic mutations (Additional file 1: Table S1).  
152 The dbNSFP includes 19 predictors which provide scores representing the probability  
153 of the non-synonymous variants being pathogenic (Additional file 2: Table S2). The  
154 scores given by each predictor were normalized into the range of 0 to 1 in this study.

155

156 ***Known cancer driver genes.*** The accuracy of the predictions made by DGAT-cancer  
157 was examined using a set of known cancer driver genes. This set of 1,168 unique  
158 cancer driver genes comprised a non-redundant set of 723 genes from the COSMIC  
159 Cancer Gene Census[25] (CGC) and 1,064 genes from OncoKB[26]. OncoKB  
160 contains more cancer driver genes than CGC because it collects cancer genes from  
161 various panels, including class Tier 1 genes in CGC (576), the MSK-IMPACT™  
162 panel (505), the MSK-IMPACT™ Heme and HemePACT panels (575), the  
163 FoundationOne CDx panel (324), the FoundationOne Heme panel (593), and data  
164 from Vogelstein et al.[125][6]. We also collated gene sets containing cancer  
165 type-specific driver genes, which were collected from IntOGen[27]. The number of  
166 known driver genes for each cancer type are shown in Additional file 1: Table S3.  
167 These genes were used as gold standard cancer-associated genes to evaluate the  
168 accuracy of the predictions made in this study.

169

170 ***Constrained genes.*** A previous study[28] provided a set of 1,003 genes that are



171 significantly lacking in missense variations (NHLBI's Exome Sequencing Project),  
172 suggesting that they have a high intolerance to germline mutation.

173

174 ***Genes affecting cell proliferation and/or viability.*** These genes were identified by a  
175 previous study that performed shRNA screens in 216 cancer cell lines[29]. We  
176 downloaded the final shRNA quality file from the study and selected 687 genes as  
177 significantly affecting cell survival if the genes had shown significant gene  
178 suppression in cells according to the evaluation of ATARiS[30] ( $q < 0.05$ ) in at least  
179 one half of the shRNA screens.

180

181 ***Drug response genes.*** From OncoKB[26], we collected 43 genes that have been  
182 reported to respond to FDA-approved drugs (Level 1) and 17 genes that respond to  
183 drugs used in standard care (Level 2). From IntOGen[27], we identified 51 genes  
184 whose protein products interact with FDA-approved drugs[31]. From these datasets,  
185 we removed off-target genes, gene therapy targets in the IntOGen list, and drugs  
186 targeting fusion driver genes, as well as 15 genes associated with drug resistance in  
187 COSMIC[32]. Finally, 46 actionable genes were obtained which could be used for  
188 drug target enrichment analysis.

189

190 ***Pathology data.*** Genes whose mRNA expression significantly correlated with the  
191 survival of cancer patients were downloaded from The Human Protein Atlas  
192 (HPA)[33]. The HPA performed Kaplan-Meier analysis to estimate the correlation

193 between mRNA expression and patient survival for each gene (a total of 20,090 genes)  
 194 for 20 cancer types. We selected genes with expression levels significantly (log-rank  
 195  $p < 0.05$ ) correlating with patient survival in at least one cancer type investigated by  
 196 this study. The proportion of genes predicted to be cancer drivers by DGAT-cancer in  
 197 these cancer survival-related genes was compared to those genes predicted to be  
 198 non-cancer drivers by a one-sided Fisher's Exact test.

199

## 200 **Pathogenicity score profiles of mutations in cancers and in the healthy** 201 **population cohort**

202 Each tumor sample or specific individual from the 1000GP may harbor more than one  
 203 mutation in each gene. The pathogenic status of these mutations was scored by 19  
 204 distinct approaches each representing different detrimental effects of mutations on  
 205 gene function. The mutation score of a given gene for a specific sample determined by  
 206 a particular predictive approach was defined as the average mutation score across all  
 207 mutations in that gene. Mutation scores of a given gene from all tumor samples or all  
 208 samples from 1000GP (Additional file 1: Table S1) were considered to represent the  
 209 mutation score profiles of that gene calculated by one particular predictive approach.  
 210 Then, for each gene, we constructed a density distribution for mutation scores in each  
 211 cancer cohort and in the healthy population cohort, respectively. The density  
 212 distribution was divided into 100 evenly spaced bins. A gene was filtered out from the  
 213 construction of mutation score profiles if all the mutations in the gene were derived

214 from fewer than five samples.

215 In order to compare the mutation score profiles in tumors with that of the healthy  
216 population cohort, we calculated the difference between the two profiles by means of  
217 a unidirectional Earth Mover's Difference score (uEMD, *Equation (1)*). For each gene,  
218 we repeated this calculation for all 19 types of functional score and obtained 19  
219 uEMD scores. These uEMD scores were termed uEMD-Mut scores.

$$220 \quad uEMD_g = \sum_{B=100}^1 \max \{ \sum_{b=100}^B (M_{b,g} - N_{b,g}), 0 \} \quad (1),$$

221 where  $M_{b,g}$  is the fraction of normalized scores in the  $b$ -th bin for gene  $g$  in the  
222 density distribution of a cancer cohort,  $N_{b,g}$  is the fraction of normalized scores in  
223 the  $b$ -th bin for gene  $g$  in the density distribution of the general population, and  $B$   
224 is the index of bins in the density distribution. Thus, genes with a significantly  
225 different distribution in cancer tissues from the general cohort would be given higher  
226 uEMD scores.

227

## 228 Comparing gene expression profiles of tumors and paracancerous tissues

229 RNA-seq data (RSEM normalized count, log2 transformed) of tumors from 12 cancer  
230 types in TCGA were downloaded from the UCSC Xena platform[34]. The  
231 corresponding RNA-seq data of paracancerous tissues from eight types of cancer  
232 including Bladder urothelial carcinoma (BLCA), Breast invasive carcinoma (BRCA),  
233 Cervical and endocervical cancers (CESC), Colon adenocarcinoma (COAD),  
234 Glioblastoma multiforme (GBM), Head and neck squamous cell carcinoma (HNSC),

235 Lung adenocarcinoma (LUAD) and Stomach adenocarcinoma (STAD) were  
236 downloaded from TCGA. RNA-seq data from paracancerous tissues or tumors with a  
237 sample size fewer than five were not included in the study. The sample sizes of  
238 RNA-seq data for each type of cancer are shown in Additional file 1: Table S4. For  
239 each gene, we calculated the median expression level across samples. Then, we  
240 computed the uEMD score for each gene to measure the difference in the density  
241 distribution of gene expression in tumor and paracancerous tissues. This uEMD score  
242 of the gene was termed uEMD-Ex.

243

## 244 **Topological network constructed from gene expression data and somatic** 245 **mutations in tumors**

246 By using gene expression data derived from tumor samples, we constructed a  
247 topological network for each type of cancer (Additional file 1: Table S4 and Table S5)  
248 using Mapper algorithm, an R package in TDAmapper (the detailed input arguments  
249 are listed in Additional file 1: Table S6)[35]. Briefly, tumor samples with similar  
250 expression profiles are clustered into one node. Where two nodes have at least one  
251 tumor sample in common, they are connected by an edge. For each gene in the  
252 topological network, we evaluated the divergence of the mutation frequency of the  
253 gene in samples across similar gene-expression clusters (*Equation (2)*). The strategy  
254 used here is similar to that described in previous studies[22, 36]. Prior to the  
255 calculation, we used a threshold of  $MAF > 0.001$  to filter out mutations occurring at

low frequencies in tumors.

$$C_g = \frac{N}{N-1} \frac{\sum_{i,j \in \Gamma} e_{i,g} A_{ij} e_{j,g}}{(\sum_{k \in \Gamma} e_{k,g})^2} \quad (2),$$

where  $\Gamma$  denotes the set of nodes in the topological network,  $A$  is the adjacency matrix of the topological network,  $N$  is the number of nodes in  $\Gamma$ , and  $e_{i,g}$  is the average frequency of non-synonymous mutations of gene  $g$  for samples in the node  $i$ . A lower  $C_g$  represents a higher divergence of the mutation frequency between similar clusters.

To evaluate the similarity between the profiles of mutation frequency and mRNA expression, we computed the Jensen–Shannon divergence between the expression and mutation profiles of each gene based on the topological network (Equation (3)).  $C_g$  and  $JSD_g$  were termed MutExTDA scores.

$$JSD_g = \frac{1}{2} \sum_{i \in \Gamma} [-(\widetilde{e}_{i,g} + \widetilde{r}_{i,g}) \log \left( \frac{\widetilde{e}_{i,g} + \widetilde{r}_{i,g}}{2} \right) + \widetilde{e}_{i,g} \log (\widetilde{e}_{i,g}) + \widetilde{r}_{i,g} \log (\widetilde{r}_{i,g})] \quad (3),$$

where  $\widetilde{e}_{i,g}$  denotes the fraction of tumors with gene  $g$  somatically mutated in node  $i$ ,  $\widetilde{r}_{i,g}$  denotes the average expression of gene  $g$  in the tumors associated with node  $i$ , and  $\Gamma$  denotes the set of nodes in the topological network. Prior to calculation, they were normalized to meet  $\sum_{i \in \Gamma} \widetilde{e}_{i,g} = \sum_{i \in \Gamma} \widetilde{r}_{i,g} = 1$ . A lower  $JSD_g$  denotes less difference between the two distributions.

## Data preprocessing

We obtained 24 features to describe each gene, comprising 19 uEMD-Mut scores, one uEMD-Ex score, two MutExTDA scores and two features representing the median

expression levels of the gene in tumor and paracancerous tissues, respectively (Additional file 2: Table S2). We filtered out genes lacking more than half the number of features (Additional file 3: Table S7 and Additional file 1: Table S8). For the remaining genes, we used k-nearest neighbors (KNN) imputation to fill in the missing features. More than 99% of genes for cancer types Pheochromocytoma and paraganglioma (PCPG), Testicular germ cell tumors (TGCT) and Thyroid carcinoma (THCA) were lacking uEMD-Mut scores. For these cancer types, we constructed the predictor only using gene expression-based features (uEMD-Ex or MutExTDA) to perform the prediction.

## Using Laplacian Score to select features

We used an unsupervised method, the Laplacian Score, to select features that have high power to preserve the local geometric structure of the feature space[37]. The detailed steps for the application of the Laplacian Score in feature selection were as follows:

- (1) A network  $N$  was constructed to connect all  $m$  candidate genes ( $N \in \mathbb{R}^{m \times m}$ ). For each pair of genes, we calculated the Euclidean distance,  $|x_i - x_j|$  between their feature vectors ( $x_i$  is the feature vector of gene  $i$ , and  $x_j$  is the feature vector of gene  $j$ ). Based on the Euclidean distance, if gene  $i$  is among the top  $k$  (here we set  $k = [0.01 \times m]$ ) of the nearest genes to gene  $j$ , or the gene  $j$  is among the top  $k$  of the nearest gene to the gene  $i$  ( $i \neq j$ ), we set the connection

298 between  $i$  and  $j$  as  $N_{ij} = 1$ . Otherwise,  $N_{ij} = 0$ .

299 (2) In the network  $N$ , if  $N_{ij} = 1$ , we weigh the edge by  $W_{ij} = e^{-||x_i - x_j||^2}$   
 300 ( $W \in \mathbb{R}^{m \times m}$ ). Otherwise,  $W_{ij} = 0$ . The weighted network reflects the local  
 301 structure of the  $m$  genes in the feature space.

302 (3) Computing the Laplacian Score of each feature. Let  $y_l = [y_{l1}, y_{l2}, \dots, y_{lm}]^T$   
 303 denote the  $l$ -th feature values for all  $m$  genes. We redefine  $y_l$  by removing the  
 304 mean from the samples as in Equation (4). The Laplacian Score ( $LS$ ) is computed  
 305 by Equation (5). The lower the  $LS$  value, the more important the feature is.  
 306 According to the scores of each feature, we finally selected the top 20 features  
 307 with the smallest  $LS$  to be included in the prediction model for each type of  
 308 cancer.

$$309 \quad \tilde{y}_l = y_l - \frac{y_l^T D I}{I^T D I} I, \quad (4)$$

310 where  $D = \text{diag}(\sum_{j=1}^m W_{1j}, \sum_{j=1}^m W_{2j}, \dots, \sum_{j=1}^m W_{mj})$ ,  $I = [1, 1, \dots, 1]^T$ .

$$311 \quad LS_l = \frac{\tilde{y}_l^T L \tilde{y}_l}{\tilde{y}_l^T D \tilde{y}_l}, \quad (5)$$

312 where  $L = D - W$ .

313

## 314 Data transformation

315 In order to integrate multiple features of each gene into a risk score (Fig. 1), we  
 316 combined the features by Hotelling and Box-Cox transformations, which converted  
 317 the feature values into  $p$ -values. For a given scaled matrix of  $m$  genes with  $n$   
 318 features ( $P \in \mathbb{R}^{m \times n}$ , here  $n = 20$ ), the Hotelling transformation is performed as:

$$P' = U(P - IM)^T,$$

where  $U = [v_1; v_2; \dots; v_{n'}]^T$  with  $n' \leq n$  as the number of chosen principal components of  $P$  and  $V = [v_1; v_2; \dots; v_n]$  are eigenvectors for the covariance matrix of  $P$  corresponding to decreasing eigenvalues with  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ .  
 $M = [\frac{1}{m} \sum_{i=1}^m P_{i1}, \frac{1}{m} \sum_{i=1}^m P_{i2}, \dots, \frac{1}{m} \sum_{i=1}^m P_{in}]$ . Thus, transformed  $P' \in \mathbb{R}^{n' \times m}$ . Then, the Box-Cox transformation is performed as follows,

$$pi' = \begin{cases} \frac{pi^{\beta_i-1}}{\beta_i}, \beta_i \neq 0 \\ \log(pi), \beta_i = 0 \end{cases},$$

where  $pi$  is the  $i$ -th row vector of  $P'$  and all elements of  $pi$  vector are forced to be positive before being transformed, and  $\beta_i$  is the parameter for transforming  $pi$  to  $pi'$ .

Finally, we standardized each  $pi'$  and calculated P-values for elements of  $pi'$  that is in a standard Gaussian distribution. The P-values of the elements were combined as  $S(j)$  by Fisher's method (Fisher's combined probability test).  $S(j)$  is termed the score of gene  $j$ .

$$S(j) = -\ln \left( \chi_{n'}^2{}^{-1} \left( -2 \sum_{i=1}^{n'} \log(pi'_j) \right) \right)$$

### Gibbs sampling

The scores ( $S$ ) of genes were used as conditional probabilities in Gibbs sampling to obtain a convergent probability distribution of candidate genes. In the first round of sampling, Gibbs sampling was initiated by randomly selecting  $m'$  ( $m' \leq m$ ) genes



337 from the candidate genes. The  $m'$  genes were assumed to have equal probabilities of  
 338 being selected. Then, in the second round of sampling, another set of  $m'$  genes was  
 339 sampled from the remaining  $m-m'$  genes weighted by their scores ( $S$ ). The  
 340 following rounds were the same as the second round. In each round, the selected  
 341 frequency ( $\frac{\# of times the gene is selected}{\# of sampling}$ ) of each gene was updated. All the selected  
 342 frequencies of candidate genes in the  $i$ -th round were denoted as a vector ( $m \times$   
 343  $1$ ),  $Freq_i$ . When the Euclidean norm of  $Freq_i - Freq_{i-1}$  was smaller than  $E_{Gibbs}$   
 344 ( $E_{Gibbs}$  was set as 0.01), the iteration was stopped.  $Freq_{last}$  was assigned as the  
 345 posterior probabilities (PP) of candidate genes. Then, we constructed a null  
 346 distribution of PP in order to obtain the likelihood of a given gene being a cancer  
 347 driver gene. The null distribution was generated by giving genes randomly weighted  
 348 scores that were derived from the uniform distribution with the same range as the true  
 349 scores. We generated 1,000,000 sets of null distributions of PPs for the  $m$  genes by  
 350 running Gibbs sampling. For each gene, we obtained 1,000,000 random PPs. By  
 351 counting the number of times that a random PP of a gene was larger than the real PP  
 352 of the gene, an experience  $p$ -value was estimated. The  $p$ -values were adjusted by  
 353 Bonferroni correction and we selected those genes with  $p_{adj} < 0.01$  as being  
 354 significant.

355

## 356 **Gene set enrichment analysis**

357 The significance of the enrichment was evaluated by the null hypothesis that the

proportion of predicted cancer driver genes in a given gene set is equal to the expected proportion of the protein-coding genes in the same gene set. These human protein-coding genes were downloaded from GENCODE[38] if genes were annotated as with biotype “protein\_coding”. In total, 19,350 unique human protein-coding genes were obtained. The gene set enrichment analysis was performed using a one-sided Fisher’s Exact test, and the enrichment  $p$ -values were corrected for multiple testing using the Bonferroni correction. We defined  $p_{adj} < 0.05$  as a significant result.

## Comparison with other cancer driver prediction methods

The performance of DGAT-cancer was evaluated by the area under the precision–recall curve (AUPRC) as calculated by the perfMeas package in R[39]. The AUPRC was calculated by using known cancer driver genes collected from CGC, OncoKB and IntOGen (1,199 genes, Methods) as a positive gene set. The negative gene set contained the human genes after removing the positive genes and the genes directly interacting with the positive genes according to protein-protein interaction data in BIOGRID[40] (version 4.4.214, organism: human sapiens). In total, 8,153 genes were allocated to the negative gene set. The numbers of positive and negative cancer driver genes for each specific cancer are shown in Additional file 1: Table S9.

DGAT-cancer was compared to three classical methods, MutSigCV[4] (<https://software.broadinstitute.org/cancer/cga/mutsig>), OncodriveFML[14] (<http://bbgglab.irbbarcelona.org/oncodrivefml/home>) and OncodriveCLUSTL[15]

(<http://bbglab.irbbarcelona.org/oncodriveclustl/home>). MutSigCV (version 1.3.5) is a method that considers the heterogeneity of mutations as a means to identify genes that are significantly more highly mutated in cancer than expected by chance alone given background mutation processes. We ran MutSigCV using default parameters and the mutation data from Broad GDAC Firehose with the gene covariates provided by the original article[4]. OncodriveFML is a method designed to use the somatic mutation pattern across cancers to identify cancer driver genes; it was run using the online version by inputting files in “maf” format obtained from ICGC and Broad GDAC Firehose (Additional file 1: Table S1). The detailed parameters for running OncodriveFML are shown in Additional file 1: Fig. S1. OncodriveCLUSTL is a sequence-based clustering algorithm designed to detect cancer drivers. We used the online version of OncodriveCLUSTL to predict cancer driver genes based on default parameters.

392

### 393 **Tissue specimens and patient information**

394 91 surgical specimens were collected from the Department of Neurosurgery, Sun  
395 Yat-sen Memorial Hospital, Sun Yat-sen University. These samples include 55  
396 glioblastoma multiforme (GBM) specimens, 31 brain lower grade glioma (LGG)  
397 specimens and 5 non-tumor brain tissues, which had been diagnosed between 2012  
398 and 2022. The study was approved by the Ethics Committee of Sun Yat-sen University,  
399 and informed consent was obtained from all subjects. The non-tumor brain tissues

were obtained from patients with non-tumor diseases and required partial brain excision from patients with traumatic brain injury, or other diseases such as cerebral angiomas or vascular malformations.

## **Experimental validation of the role of novel driver genes in cancer**

In order to examine the performance of DGAT-cancer, we experimentally validated the roles of the predicted cancer drivers by using surgical specimens, a cell model and an animal model. All experimental methods are provided in the Additional file 1: Supplementary Methods.

## **Results**

### **Application of DGAT-cancer to multiple types of cancer**

DGAT-cancer was applied to the prediction of cancer driver genes in nine cancer types, Bladder urothelial carcinoma (BLCA), Breast invasive carcinoma (BRCA), Cervical and endocervical cancers (CESC), Colon adenocarcinoma (COAD), Glioblastoma multiforme (GBM), Head and neck squamous cell carcinoma (HNSC), Brain lower grade glioma (LGG), Lung adenocarcinoma (LUAD) and Stomach adenocarcinoma (STAD) whose mutation (uEMD-Mut) and gene expression (uEMD-Ex, gene expression level and MutExTDA) features were all available in TCGA database (Methods). The numbers of genes predicted to be cancer drivers by

420 DGAT-cancer are shown in Additional file 1: Table S8.

421 We found significant ( $p_{adj} \leq 0.032$ ) overlaps between the predicted cancer drivers  
422 and the cancer gene sets, CGC, OncoKB and IntOGen, respectively (Fig. 2a)  
423 compared to the background genes (19,350 protein-coding genes in ENSEMBL[41]  
424 biotype). Since cancer genes are likely to have experienced a slower evolutionary rate  
425 and stronger purifying selection than those of non-cancer, Mendelian disease, and  
426 orphan disease genes[42], we tested the enrichment of the predicted genes in the  
427 genes under selective constraint[28]. We found that the predicted cancer drivers of the  
428 nine cancer types were significantly ( $p_{adj} \leq 4.18 \times 10^{-4}$ ) enriched in genes that  
429 have been under selective constraint (Fig. 2a). We further evaluated the enrichment of  
430 the predicted cancer drivers in the shRNA gene set by performing shRNA screens in  
431 216 cancer cell lines[29] (Methods). As shown in Fig. 2a, the predicted cancer driver  
432 genes in the nine cancer types were significantly ( $p_{adj} \leq 0.034$ ) enriched in the  
433 shRNA gene set, illustrating the potential roles of the predicted cancer drivers in  
434 cancer cell survival. The predicted cancer drivers were also enriched in cancer-related  
435 pathways defined in the Kyoto Encyclopedia of Genes and Genomes (KEGG)[43]  
436 database, e.g. Melanoma ( $2.90 \times 10^{-5} \leq FDR \leq 0.046$ ), Proteoglycans in cancer  
437 ( $5.13 \times 10^{-6} \leq FDR \leq 0.033$ ), cancer immunotherapy ( $FDR = 6.16 \times 10^{-3}$ ), cell  
438 cycle ( $0.017 \leq FDR \leq 0.024$ ), ErbB signalling pathway ( $5.63 \times 10^{-5} \leq FDR \leq$   
439  $0.024$ ), MAPK signalling pathway ( $3.35 \times 10^{-3} \leq FDR \leq 0.010$ ), Wnt signalling  
440 ( $1.78 \times 10^{-3} \leq FDR \leq 0.044$ ), p53 signalling pathway ( $6.40 \times 10^{-3} \leq FDR \leq$   
441  $0.046$ ), and TGF-beta signalling pathway ( $1.61 \times 10^{-6} \leq FDR \leq 0.044$ ) (Additional

442 file 1: Fig. S2).

443 There were 20 genes predicted as cancer drivers by DGAT-cancer in multiple  
 444 cancer types (Fig. 2b). Among them, *TP53* was predicted to be a cancer driver in nine  
 445 types of cancer, with predicted scores ranking between the top 1 to the top 24. The  
 446 *COL1A2* gene was predicted as a cancer driver in BLCA and HNSC, with predicted  
 447 scores ranking in the top 10 and the top 4, respectively. The *COL1A2* gene encodes  
 448 the pro-alpha2 chain of type I collagen which is involved in the TGF-beta signalling  
 449 pathway and has been reported to be associated with the migration of chondrosarcoma  
 450 and fibrosarcoma cells[44]. Another gene, *PTEN*, was predicted to be a cancer driver  
 451 in BRCA, CESC, COAD, GBM, LGG and STAD. *PTEN* is a tumour suppressor[45,  
 452 46] that is involved in the p53 and PI3K-Akt signalling pathways and has been found  
 453 to be frequently mutated in a large number of different cancers[47, 48].

454

# 455 **Expression of cancer drivers predicted by DGAT-cancer correlates significantly** 456 **with the level of drug activity**

457 As shown in Fig. 2c, the predicted cancer drivers in BRCA, CESC, COAD, GBM,  
 458 LGG and STAD were significantly ( $7.65 \times 10^{-8} \leq p_{adj} \leq 0.046$ ) enriched in drug  
 459 response genes. We further explored the correlation between the expression patterns  
 460 of the predicted cancer driver genes and drug activities, expressed as 50% growth  
 461 inhibitory levels (GI50) in the NCI-60 cell line, which were derived from  
 462 CellMinerCDB[49]. The analysis was performed by calculating Pearson correlation

coefficients between gene expression levels and the z-scores of negative log 10 (GI50) in all NCI-60 cell lines. We found that the expression levels of many predicted cancer drivers were significantly ( $p_{adj} \leq 0.05$ ) correlated with drug activities (Fig. 2d, thereby illustrating the potential of these genes for clinical treatment. Briefly, the genes predicted to be cancer drivers in four out of nine cancer types (BLCA, COAD, HNSC and STAD) were significantly enriched (one-sided Fisher's Exact test  $p_{adj} < 0.036$ ) in genes whose expression was correlated with drug activity comparing to random genes selected from the human protein gene sets.

## **DGAT-cancer outperformed other methods**

DGAT-cancer was compared with three other cancer driver identification methods, MutSigCV, OncodriveFML and OncodriveCLUSTL with respect to the area under the precision-recall curve (AUPRC). First, DGAT-cancer was compared with other methods of predicting cancer drivers from a set of genes (a total of 21,664 genes, Additional file 1: Table S8) having prediction scores given by at least one of the four methods. Additional file 1: Table S9 shows the numbers of positive and negative genes used in evaluating the methods. As shown in Fig. 3a, the AUPRC (in the range of 0.336 to 0.469) of DGAT-cancer in predicting cancer drivers for nine types of cancer (BLCA, BRCA, CESC, COAD, HNSC, GBM, LGG, LUAD and STAD) were the highest by comparison with the other three methods, MutSigCV (AUPRC ranged in 0.225 to 0.282), OncodriveFML (AUPRC ranged in 0.283 to 0.365) and

484 OncodriveCLUSTL (AUPRC ranged in 0.289 to 0.356) (Fig. 3a). When assessing the  
485 AUPRC of these methods for predicting cancer drivers from genes with prediction  
486 scores provided by all four methods (Additional file 1: Table S8), DGAT-cancer also  
487 performed the best in the eight cancer types (BLCA, BRCA, CESC, COAD, HNSC,  
488 GBM, LUAD and STAD) (Additional file 1: Fig. S3a).

489 Each of the four methods yielded a  $p$ -value to represent the probability of the  
490 predicted cancer driver being a false positive. The  $p$ -value distributions of genes  
491 generated by these four methods were then compared with those expected  $p$ -values  
492 from a uniform distribution using quantile-quantile plots. As shown in Additional file  
493 1: Fig. S3b, the  $p$ -values of genes not predicted to be cancer drivers by DGAT-cancer  
494 ( $p > 0.05$ ) in BLCA, LUAD and STAD exhibited better agreement with the expected  
495  $p$ -values than those predicted by other methods. Additionally, the genes predicted ( $p <$   
496  $0.05$ ) by DGAT-cancer to be cancer drivers in all nine cancer types showed higher  
497 inflation from the expected  $p$ -values than the genes predicted by other methods  
498 (Additional file 1: Fig. S3c). This suggested that DGAT-cancer has enhanced potential  
499 to distinguish novel cancer drivers from random genes.

500 We next explored the consistency of DGAT-cancer with respect to the other  
501 methods. The posterior probabilities (PPs) given by DGAT-cancer for genes predicted  
502 to be cancer drivers ( $q_{risk} < 0.05$ ) by the other methods were compared to the PP  
503 scores of the genes not predicted to be cancer drivers ( $q_{risk} \geq 0.05$ ) by the other  
504 methods, MutSigCV, OncodriveFML and OncodriveCLUSTL, respectively. The  
505 differences were evaluated by the Wilcoxon rank-sum test. The total number of cancer



506 drivers with  $q_{risk} < 0.05$  identified by MutSigCV, OncodriveFML and  
 507 OncodriveCLUSTL are shown in Additional file 1: Table S10. As depicted in Fig. 3b,  
 508 those genes predicted to be cancer drivers in BLCA, BRCA, CESC, COAD, LGG and  
 509 LUAD by the three methods were given significantly higher PPs ( $p < 0.033$ ) by  
 510 DGAT-cancer than the genes not predicted to be cancer drivers by the other methods.  
 511 These results suggested a high degree of consistency between DGAT-cancer and the  
 512 other methods.

513 There were 67, 58, 50, 166, 41, 87, 9, 39 and 112 genes predicted to be cancer  
 514 drivers of nine types of cancer (BLCA, BRCA, CESC, COAD, GBM, HNSC, LGG,  
 515 LUAD and STAD) by DGAT-cancer, respectively but not predicted to be cancer  
 516 drivers by the other three methods (OncodriveCLUSTL, OncodriveFML, MutSigCV).  
 517 Moreover, 91, 152, 32, 135, 13, 177, 11, 67 and 78 genes were predicted as cancer  
 518 drivers in nine types of cancer (BLCA, BRCA, CESC, COAD, GBM, HNSC, LGG,  
 519 LUAD and STAD) by at least of one of the three methods, were not predicted to be  
 520 cancer drivers by DGAT-cancer. We compared these genes in relation to their scores  
 521 of features used in DGAT-cancer, and found that the gene expression-related features  
 522 such as uEMD-Ex scores, expression values and MutExTDA scores showed a  
 523 significant difference in the two groups of genes across multiple cancers (Fig. 3c and  
 524 Additional file 1: Fig. S4). Specifically, the genes missed by the other three methods  
 525 were expressed more highly in both tumors and paracancerous tissues than the genes  
 526 that were missed by DGAT-cancer in BRCA, COAD, HNSC and LUAD. This  
 527 illustrated that DGAT-cancer was more likely to detect active genes in cancer and

paracancerous tissues than other methods without considering gene expression. DGAT-cancer also used features based on mutation frequency and expression profiles in tumor cohorts, JSD and C score (Methods), and fin predicting cancer drivers. For these two scores, genes missed by other methods had relatively lower JSD scores in three cancers (BRCA, COAD and STAD) and lower C scores in four cancers (BRCA, COAD, LUAD and STAD) than genes missed by DGAT-cancer (Fig. 3c), suggesting a preference for DGAT-cancer to identify genes that have highly correlated gene expression and mutation profiles.

# **Gene expression-based features in tumour and paracancerous tissue served to improve DGAT-cancer**

DGAT-cancer integrated both mutation-based features (uEMD-Mut) and gene expression-based features (uEMD-Ex, MutExTDA and expression values) in order to identify cancer drivers. The importance of these two types of features in the prediction was then evaluated. First, we used only uEMD-Mut to predict cancer drivers for each of the nine cancer types (Methods); this predictor was termed DGAT-Mut. The comparison between DGAT-Mut with other methods was based on a total of 23,433 genes that have the prediction score from at least one of the methods. Additional file 1: Table S11 shows the numbers of positive and negative genes used for evaluating the methods. As shown in Fig. 4, DGAT-Mut yielded higher AUPRC values (ranging from 0.272 to 0.370) for predicting cancer drivers in BLCA, BRCA, CESC, COAD,

549 GBM, LGG and STAD than OncodriveCLUSTL (ranging from 0.288 to 0.442),  
550 OncodriveFML (ranging from 0.283 to 0.406) and MutSigCV (ranging from 0.221 to  
551 0.282). However, the AUPRC values achieved by DGAT-Mut in predicting cancer  
552 drivers for nine cancer types were lower than with DGAT-cancer. These results  
553 illustrated the importance of integrating uEMD-Mut with the gene expression-related  
554 features for improving the prediction of cancer drivers.

555 When we only used gene expression-based features (uEMD-Ex, MutExTDA and  
556 gene expression values) to construct the predictive model, DGAT-Exp (Methods), it  
557 provided lower AUPRC values than DGAT-cancer in predicting cancer drivers for five  
558 out of nine types of cancer (BLCA, BRCA, COAD, HNSC and LUAD) (Fig. 4).  
559 DGAT-cancer and DGAT-Mut were not applied to the prediction of cancer drivers for  
560 three cancer types (PCPG, TGCT and THCA) because limited mutational information  
561 was available from the TCGA or ICGC databases for these three types of cancer  
562 (Methods). When DGAT-Exp was used to predict cancer drivers in these types of  
563 cancer (BLCA, CESC, COAD, GBM, LGG, STAD and TGCT), it yielded the highest  
564 AUPRC values compared to the other three methods (MutSigCV, OncodriveCLUSTL  
565 and OncodriveFML). In short, DGAT-Exp performed less well than DGAT-cancer  
566 with respect to the prediction of cancer drivers in five types of cancer (BLCA, BRCA,  
567 COAD, HNSC, and LUAD) although it yielded the best AUPRC in predicting the  
568 cancer drivers than other four methods (OncodriveCLUSTL, OncodriveFML,  
569 MutSigCV, and DGAT-Mut). This result confirms the key role of gene  
570 expression-based features in the model. Thus, combining both uEMD-Mut and

571 expression-based features improved the performance of DGAT-cancer.

572

### 573 **DGAT-cancer identifies novel cancer drivers**

574 DGAT-cancer identified many novel cancer driver genes that have not hitherto been  
575 reported to be related to cancer in CGC, OncoKB or IntOGen gene sets. Briefly, it  
576 predicted 71, 131, 60, 99, 148, 117, 20, 46 and 105 novel cancer drivers in GBM,  
577 BLCA, BRCA, CESC, COAD, HNSC, LGG, LUAD and STAD, respectively  
578 (Additional file 4: Table S12).

579 We next wondered if these predicted novel cancer drivers might correlate with the  
580 prognosis of the cancer patients (Methods). Relationships between the survival time  
581 of patients and gene expression were explored using a Kaplan-Meier model to analyze  
582 data from BRCA, BLCA, CESC, COAD, GBM, HNSC, LGG, LUAD and STAD in  
583 The Human Protein Atlas[33]. By comparing the number of genes whose mRNA  
584 expression level significantly (log-rank  $p < 0.05$ ) correlated with patient survival in  
585 at least one cancer type, we observed that the predicted novel cancer drivers contain a  
586 significantly higher proportion of genes (total number 571 and proportion 96.67%)  
587 (one-sided Fisher's Exact test,  $p = 1.03 \times 10^{-5}$ ,  $OR = 2.48$ ) that correlated with  
588 patient survival than the genes predicted to be of low probability to be cancer drivers  
589 (removing genes contained in CGC, OncoKB, IntOGen and predicted cancer drivers,  
590 total number 7,723 and proportion 92.14%). For example, *CD44*, a gene that was  
591 predicted to be a novel cancer driver in LUAD (with scores ranked in the top 48 and

592  $p_{adj} < 10^{-6}$ ) was found to be significantly correlated with patient survival for COAD  
593 ( $p = 5.00 \times 10^{-2}$ ), GBM ( $p = 3.50 \times 10^{-3}$ ), LGG ( $p = 3.50 \times 10^{-3}$ ) and HNSC  
594 ( $p = 1.46 \times 10^{-2}$ ). It has been reported that the expression level of *CD44* exhibits a  
595 positive correlation with PD-L1 protein in LUAD patients [52]. *TTN* was another  
596 novel gene predicted to be a cancer driver in six types of cancer, BRCA, CESC,  
597 COAD, HNSC, LUAD and STAD (with scores ranked in the top 2 to top 39, and  
598  $p_{adj} < 10^{-6}$ ) and the expression level was found to correlate with survival in BRCA  
599 ( $p = 1.00 \times 10^{-2}$ ), BLCA ( $p = 1.34 \times 10^{-2}$ ), GBM ( $p = 3.71 \times 10^{-2}$ ), LGG  
600 ( $p = 3.71 \times 10^{-2}$ ) and HNSC ( $p = 1.49 \times 10^{-2}$ ). Finally, *EEF1A1* was predicted to  
601 be a cancer driver in BLCA, BRCA, GBM, HNSC and LUAD (ranked 21 to ~68,  
602  $p_{adj} < 10^{-6}$ ). Further analysis indicated that the expression of *EEF1A1* was  
603 significantly correlated with the survival of patients in CESC ( $p = 7.86 \times 10^{-3}$ ),  
604 COAD ( $p = 2.53 \times 10^{-2}$ ), GBM ( $p = 1.39 \times 10^{-2}$ ) and LGG ( $p = 1.39 \times 10^{-2}$ ).  
605 A previous study has indicated that *EEF1A1* is expressed at a significantly higher  
606 level in glioblastomas than in non-neoplastic white matter[53]. However, no studies  
607 have so far been performed to establish a relationship between *EEF1A1* and GBM.  
608 We, therefore, performed further experiments to validate a possible role for *EEF1A1*  
609 in glioma.

610

## 611 **Experimental validation of a role for *EEF1A1* in glioma**

### 612 ***EEF1A1* expression is increased in glioma and correlates with a poor prognosis**

613 We compared the mRNA expression levels of *EEF1A1* in 698 glioma samples from  
614 the TCGA GBM/LGG dataset and 1,157 normal brain samples from GTEx[50]. The  
615 results showed that *EEF1A1* expression was upregulated in tumor tissue relative to  
616 normal brain tissue, as well as significantly increased in GBM relative to LGG (Both  
617  $P < 0.001$ , Fig. 5a). Moreover, the expression level of *EEF1A1* in the glioma samples  
618 from TCGA was found to negatively correlate with the overall survival of patients  
619 ( $P = 0.026$ , Fig. 5b).

620 The *EEF1A1* protein level was examined in 91 glioma specimens by  
621 immunohistochemical analysis; these included 55 GBM tissues, 31 LGG tissues and 5  
622 non-tumor brain tissues. Relative to non-tumor brain tissues, LGG exhibited a  
623 significantly higher *EEF1A1* level in LGG ( $p = 6.59 \times 10^{-3}$ , Fig. 5c) and GBM  
624 ( $p = 2.12 \times 10^{-3}$ , Fig. 5c). Notably, the expression of *EEF1A1* protein was  
625 significantly increased in GBM as compared with LGG ( $p = 4.02 \times 10^{-3}$ , Fig. 5c).  
626 When we cultured GBM cell lines (U251 and U87 cells) and LGG cell lines (Hs683),  
627 we found that *EEF1A1* protein was more highly expressed in GBM cells (U251 and  
628 U87) than in LGG cells (Hs683) ( $p = 1.01 \times 10^{-3}$  for U87 vs. Hs683,  $p = 6.82 \times$   
629  $10^{-4}$  for U251 vs. Hs683, Additional file 1: Fig. S5). Thus, we surmised that  
630 *EEF1A1* may be involved in glioma tumorigenesis.

631

## 632 **Knockdown of *EEF1A1* inhibited the proliferation and migration of glioma cells**

633 To examine the role of *EEF1A1* in glioma tumorigenesis, we used shRNA to knock  
634 down *EEF1A1* expression in U251 and U87 cells. To assess knockdown efficiency,  
635 we performed real-time quantitative PCR (RT-qPCR) and Western blot assays. The  
636 results showed that compared with the control group, *EEF1A1* mRNA in the  
637 knockdown group decreased by about 77% and 86% in U87 and U251 glioma cells  
638 ( $p = 2.82 \times 10^{-17}$  for U87 and  $p = 8.35 \times 10^{-11}$  for U251, Additional file 1: Fig.  
639 S6a), respectively. Western blot experiments showed that *EEF1A1* protein expression  
640 in the knockdown groups decreased by 65% and 45% in U87 and U251 glioma cells  
641 ( $p = 1.67 \times 10^{-2}$  for U87 and  $p = 2.10 \times 10^{-3}$  for U251), respectively  
642 (Additional file 1: Fig. S6b).

643 To assess the effect of *EEF1A1* knockdown on the proliferation of U251 and U87  
644 cells, we performed a Cell Counting Kit-8 (CCK-8) assay, an EdU flow cytometry  
645 assay and a colony formation assay. The results showed that U251 and U87 cells with  
646 reduced *EEF1A1* expression exhibited a significant decrease in cell viability within  
647 72h according to the CCK-8 assay (-25% and  $p = 9.19 \times 10^{-4}$  for U251, -11% and  
648  $p = 1.49 \times 10^{-2}$  for U87, Additional file 1: Fig. S6c), whilst the proportions of  
649 EdU-positive cells were decreased by 45% for U251 ( $p = 1.32 \times 10^{-5}$ ) and by 38%  
650 for U87 ( $p = 7.74 \times 10^{-9}$ ) (Fig. 5d), with the numbers of colonies being decreased  
651 by 54% for U251 ( $p = 5.80 \times 10^{-4}$ ) and by 69% for U87 ( $p = 1.61 \times 10^{-4}$ ) (Fig. 5e)  
652 compared with the control group. These results demonstrated that *EEF1A1*  
653 knockdown significantly inhibited the proliferation of U251 and U87 cells.

654 The migration capacity of glioma cells also decreased significantly after *EEF1A1*  
655 knockdown. We found that after *EEF1A1* knockdown, the percentage of cells  
656 migrating through the transwell plate significantly decreased (-69% and  $p = 4.31 \times$   
657  $10^{-11}$  for U251, -71% and  $p = 3.96 \times 10^{-10}$  for U87, Fig. 5f and Additional file 1:  
658 Fig. S6d). A scratch-wound healing assay yielded a similar result. A significantly  
659 shorter migration distance was observed in U251 cells with *EEF1A1* knockdown  
660 (-44%,  $p = 2.62 \times 10^{-2}$ , Additional file 1: Fig. S6e) and U87 cells (-30%,  
661  $p = 2.97 \times 10^{-3}$ , Additional file 1: Fig. S6f). To assess the proliferation of glioma  
662 cells *in vivo*, U251 and U87 cells transfected with *EEF1A1* shRNA or control shRNA  
663 were injected into zebrafish and the areas of GFP fluorescent foci were measured. As  
664 shown in Fig. 5g, U251 and U87 cells with *EEF1A1* knockdown exhibited a  
665 significant decrease in the areas of fluorescent foci (-50% and  $p = 1.96 \times 10^{-3}$  for  
666 U251, -41% and  $p = 3.20 \times 10^{-2}$  for U87, Additional file 1: Fig. S6g). Taken  
667 together, these results suggest that *EEF1A1* plays a role in regulating the proliferation  
668 and migration of glioma cells.

669

# 670 **Knockdown of *EEF1A1* increased temozolomide (TMZ) sensitivity in glioma** 671 **cells**

672 In order to explore the role of *EEF1A1* in the sensitivity of glioma cells to  
673 temozolomide (TMZ), U251 and U87 cells with *EEF1A1* knockdown, and control  
674 glioma cells were cultured in different concentrations of TMZ. The results showed



that knockdown of *EEF1A1* significantly decreased cell viability at different concentrations of TMZ ( $p = 2.66 \times 10^{-2}$  for U251 cells and  $p = 9.10 \times 10^{-3}$  for U87 cells, Fig. 6a), as well as at the half maximal inhibitory concentration (IC50) of TMZ in glioma cells (-41% and  $p = 2.94 \times 10^{-3}$  for U251, -44% and  $p = 6.44 \times 10^{-3}$  for U87, Fig. 6b). Next, a colony formation assay was performed in cells with treatment of TMZ (50 $\mu$ M). In the presence of TMZ, knockdown of *EEF1A1* significantly reduced colony formation (-64% and  $p = 2.93 \times 10^{-5}$  for U251 cells, -58% and  $p = 1.60 \times 10^{-4}$  for U87 cells, Fig. 6c and d).

Culturing U251 and U87 cells in each group with 200 $\mu$ M TMZ, we observed that the cell viability of the *EEF1A1*-knockdown group was significantly lower ( $p < 0.01$ ) than that of the control group according to CCK-8 assays (Additional file 1: Fig. S7a). These results suggest that the knockdown of *EEF1A1* is able to inhibit the proliferation of glioma cell lines in the presence of TMZ.

To investigate the effect of *EEF1A1* on the apoptosis of glioma cells treated by TMZ, we cultured U251 and U87 cells with 200 $\mu$ M TMZ for 24h, and assessed cellular apoptosis by means of a flow cytometer. The results showed that the proportion of apoptotic cells in the *EEF1A1*-knockdown group was significantly higher than that in the control group (+127% and  $p = 1.60 \times 10^{-5}$  for U251, +129% and  $p = 6.50 \times 10^{-3}$  for U87, Fig. 6e and f). Moreover, apoptosis-related proteins also showed significant changes after the knockdown of *EEF1A1*. Thus, cleaved caspase-3, the effector of apoptotic activity, was significantly increased, whilst bcl-2, an inhibitor of apoptosis, was significantly decreased in *EEF1A1*-knockdown glioma

697 cells (Both  $p = 2.86 \times 10^{-2}$ ) (Fig. 6g and h). These findings indicated that the  
698 decreased expression of *EEF1A1* can promote apoptosis in glioma cells, thereby  
699 increasing the TMZ sensitivity of glioma cells.

## 700 Discussion

701 The computational identification of cancer driver genes is key to improving our  
702 understanding of the underlying mechanisms of tumorigenesis. Most of the current  
703 methods for the identification of cancer drivers rely on the use of a single type of  
704 genomic data, such as mutations or gene expression[51]. Recently, approaches have  
705 been developed that integrate somatic mutation rates with biological networks for the  
706 prediction of cancer driver genes[52]. However, none of these approaches have  
707 considered the functional impact of mutations and tissue type-specific gene  
708 expression networks, potentially resulting in the introduction of false positives in the  
709 prediction. Here, we have developed a method, DGAT-cancer, which integrates the  
710 predicted pathogenicity of mutation profiles from cancer cohorts and a healthy  
711 population with the gene expression profiles of tumors and paracancerous tissues into  
712 a risk score that is capable of predicting cancer drivers. The method takes advantage  
713 of the huge amount of mutation data from individual samples generated by the 1000  
714 Genomes Project, the TCGA project and the ICGC project to obtain the distribution  
715 differences of the predicted pathogenic scores of mutations in the healthy population  
716 and in tumor tissues. These differences were then integrated with the topological  
717 network of gene expression in tumor tissues and gene expression data from

718 paracancerous tissues. To filter out the non-redundant features, we selected those  
719 features that allowed the preservation of the local geometric structure of the feature  
720 space by an unsupervised method, the Laplacian Score. DGAT-cancer is capable of  
721 detecting cancer drivers for any type of cancer by inputting mutation and/or gene  
722 expression data.

723 Comparing DGAT-cancer with three existing methods, it achieved the highest  
724 AUPRC in predicting cancer drivers from five out of nine types of cancer while  
725 OncodriveFML achieved the highest AUPRC in the prediction of cancer drivers for  
726 four types of cancer. The reliability of DGAT-cancer was further evaluated by  
727 comparing the  $p$ -values of the predicted cancer drivers to the expected  $p$ -values. The  
728 result showed that DGAT-cancer is a powerful tool for discriminating cancer drivers  
729 from random genes. Compared to other methods, DGAT-cancer has the ability to  
730 recognize genes with a higher expression level in tumors and paracancerous tissues  
731 (Fig. 3c), suggesting the importance of introducing tissue-specific gene expression as  
732 a feature in predicting cancer driver genes. Moreover, for the cancer types with larger  
733 sample numbers, such as BLCA, CESC, HNSC and STAD, DGAT-cancer exhibits a  
734 preference for the prediction of genes with higher uEMD-mut scores calculated using  
735 the information on tumor mutations and germline variants in the healthy population  
736 (Additional file 1: Fig. S4). Additional analysis indicated that employing only  
737 mutational information in cancer driver prediction was insufficient to achieve  
738 enhanced performance. As indicated in Fig. 4, DGAT-Mut only performed better than  
739 DGAT-Exp in its prediction of cancer drivers in BRCA, HNSC and LUAD. One likely

740 reason is that the accuracy of DGAT-mut is influenced by the 19 features for the  
741 prediction of the pathogenicity of mutations. Especially, DGAT-cancer performed the  
742 best compared to all other methods in predicting cancer drivers for three types of  
743 cancers, LGG, PCGC and TGCT which are cancer types lacking gene expression data  
744 in paracancerous tissues (Fig. 4). This result may reflect the important roles of gene  
745 expression data in tumor tissue in predicting cancer drivers.

746 DGAT-cancer predicted many novel candidate genes that require further  
747 experimental validation. Thus, *AHNAK* was predicted as a cancer driver in BRCA,  
748 COAD, LUAD and STAD. It encodes a large structural scaffold protein and has been  
749 reported as a tumour suppressor in the proliferation and invasion of triple-negative  
750 BRCA[53] and LUAD[54]. *AHNAK* may function as a tumour suppressor through the  
751 inhibition of p-ERK and ROCK1 in COAD[55]. *DST* was predicted as a cancer driver  
752 in BRCA, COAD, GBM and STAD. Up-regulation of gene *DST* has been observed in  
753 the ductal carcinoma *in situ* component of BRCA[56]. *MUC5B* was predicted as a  
754 cancer driver in BRCA, COAD and STAD. Abnormal expression of *MUC5B* in  
755 COAD has been found to be associated with low expression of p53[57]. *COL1A2* was  
756 predicted to be a cancer driver in BLCA, BRCA, GBM, HNSC and LGG. *COL1A2*  
757 has been found to be associated with tumorigenesis in HNSC[58], and is a hub gene in  
758 the perineural invasion (PNI)-associated co-expression module, where PNI is a key  
759 pathological feature of HNSC[59]. *KRT14* was predicted to be a cancer driver in  
760 HNSC. It has been reported to be upregulated in HNSC[60] and its expression level in  
761 HNSC tumours is associated with patient survival[61]. *LOXL2* was predicted to be a

762 cancer driver of HNSC. A previous study has shown that a novel splice variant in  
763 *LOXL2* upregulated *LOXL2* in HPV-negative HNSC and enhanced proliferation,  
764 migration, and invasion in HPV-negative HNSC cells[62]. The silencing of *LOXL2*  
765 inhibited cell migration and invasion in HNSC cell lines[63].

766 This is the first study to predict *EEF1A1* to be a cancer driver in GBM. *EEF1A1*  
767 protein has been reported to interact with key components of the pathway that  
768 regulates the synthesis of apoptosis-related proteins[64-66]. Additional roles for  
769 *EEF1A1* in apoptosis have been studied in gastric cancer[67], ovarian cancer[68],  
770 renal cell carcinoma[69] and lung cancer[70]. Here, we found that the knockdown of  
771 *EEF1A1* inhibits apoptosis in glioma cells when the cells are treated with TMZ. The  
772 underlying mechanisms require further investigation. However, when we explore the  
773 role of *EEF1A1* in glioma cells without TMZ, we did not observe any change in  
774 apoptosis when *EEF1A1* was knocked down (Additional file 1: Fig. S7b). Further  
775 experiments are required to establish the precise nature of the role that *EEF1A1* plays  
776 in glioma tumorigenesis.

777 There are many ways to further improve the accuracy of DGAT-cancer. First,  
778 DGAT-cancer is dependent on mutation data from the general population and cancer  
779 population. When the sample size is increased, more mutational information will  
780 become available which will help to improve the accuracy of predictions. Moreover,  
781 DGAT-cancer only takes account of mutation and gene expression data in its  
782 predictions. Other features, such as the protein or RNA structures around the  
783 mutations could also contribute important information that might be of use in

784 discriminating cancer drivers. Finally, since the performance of DGAT-cancer is  
785 influenced by pathogenicity predictions for mutations, more accurate predictors for  
786 determining the pathogenicity of mutations could help to improve DGAT-cancer.

## 787 **Conclusions**

788 We demonstrate that DGAT-cancer is powerful in predicting cancer drivers using  
789 mutation and/or gene expression data, and has a superior performance compared to  
790 three commonly used methods. The importance of gene expression data and mutation  
791 information in predicting cancer drivers was evidenced. DGAT-cancer has broadened  
792 our path to detect cancer driver genes and shed a light on cancer therapy.

## 793 **Declarations**

### 794 **Ethics approval and consent to participate**

795 Not applicable.

796

### 797 **Consent for publication**

798 Not applicable.

799

### 800 **Availability of data and materials**

801 The DGAT-cancer method is available as an open-source software package on the

802 GitHub repository (<https://github.com/Dan-He/DGAT-cancer>). Simple somatic  
803 mutation (SSM) data from 12 cancer cohorts were downloaded from the Broad  
804 Institute GDAC Firehose Portal (<http://gdac.broadinstitute.org/>), the International  
805 Cancer Genome Consortium (ICGC) Data Portal  
806 (<https://dcc.icgc.org/releases/current/Projects/>) and The Cancer Genome Atlas (TCGA)  
807 (<https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga>  
808 ). The germline mutation data of a healthy population were collected from Phase 3 of  
809 the 1000 Genomes Project (<https://www.internationalgenome.org/data>, GRCh38).  
810 RNA-seq data of tumors from 12 cancer types in TCGA were downloaded from the  
811 UCSC Xena platform (<http://xena.ucsc.edu/>).  
812 Additional information is available at the website.

813

## 814 **Competing interests**

815 All authors declared that they have no competing interests.

816

## 817 **Funding**

818 This work was supported by the National Key Research and Development Program of  
819 China (2020YFB0204803), the Natural Science Foundation of China (81801132,  
820 81971190, 61772566), Guangdong Key Field Research and Development Plan  
821 (2019B020228001, 2018B010109006, and 2021A1515010256), Introducing

822 Innovative, Guangzhou Science and Technology Research Plan (202007030010).

823

## 824 **Author's contributions**

825 HY.Z designed the study. D.H performed the model construction and analyses, with

826 assistance from Z.L and T.L. L.L performed the experiments, with assistance from

827 HS.Z, F.L, S.L and B.L. D.H, L.L and Z.L wrote the manuscript. HY.Z and D.N.C

828 supervised the study. All authors discussed the results and interpretation and

829 contributed to the final version of the paper.

830

## 831 **Acknowledgments**

832 The authors thank many resources for making data available.

833

## 834 **References**

835 1. Kandoth C, McLellan MD, Vandin F, Ye K, Niu B, Lu C, et al. Mutational landscape and

836 significance across 12 major cancer types. *Nature*. 2013;502(7471):333-9.

837 2. Cheng F, Zhao J, Zhao Z. Advances in computational approaches for prioritizing driver

838 mutations and significantly mutated genes in cancer genomes. *Brief Bioinform*.

839 2016;17(4):642-56.

840 3. Leiserson MD, Vandin F, Wu HT, Dobson JR, Eldridge JV, Thomas JL, et al. Pan-cancer



841 network analysis identifies combinations of rare somatic mutations across pathways and  
842 protein complexes. *Nat Genet.* 2015;47(2):106-14.

843 4. Lawrence MS, Stojanov P, Polak P, Kryukov GV, Cibulskis K, Sivachenko A, et al.  
844 Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature.*  
845 2013;499(7457):214-8.

846 5. Dees ND, Zhang Q, Kandoth C, Wendl MC, Schierding W, Koboldt DC, et al. MuSiC:  
847 identifying mutational significance in cancer genomes. *Genome Res.* 2012;22(8):1589-98.

848 6. Vogelstein B, Papadopoulos N, Velculescu VE, Zhou S, Diaz LA, Jr., Kinzler KW.  
849 Cancer genome landscapes. *Science.* 2013;339(6127):1546-58.

850 7. Garraway LA, Lander ES. Lessons from the cancer genome. *Cell.* 2013;153(1):17-37.

851 8. Ding L, Wendl MC, McMichael JF, Raphael BJ. Expanding the computational toolbox  
852 for mining cancer genomes. *Nat Rev Genet.* 2014;15(8):556-70.

853 9. Martincorena I, Campbell PJ. Somatic mutation in cancer and normal cells. *Science.*  
854 2015;349(6255):1483-9.

855 10. Gonzalez-Perez A, Mustonen V, Reva B, Ritchie GR, Creixell P, Karchin R, et al.  
856 Computational approaches to identify functional genetic variants in cancer genomes. *Nat*  
857 *Methods.* 2013;10(8):723-9.

858 11. Youn A, Simon R. Identifying cancer driver genes in tumor genome sequencing studies.  
859 *Bioinformatics.* 2011;27(2):175-81.

860 12. Przytycki PF, Singh M. Differential analysis between somatic mutation and germline  
861 variation profiles reveals cancer-related genes. *Genome Med.* 2017;9(1):79.

862 13. Korthauer KD, Kendziorski C. MADGiC: a model-based approach for identifying driver

863 genes in cancer. *Bioinformatics*. 2015;31(10):1526-35.

864 14. Mularoni L, Sabarinathan R, Deu-Pons J, Gonzalez-Perez A, Lopez-Bigas N.

865 OncodriveFML: a general framework to identify coding and non-coding regions with cancer

866 driver mutations. *Genome Biol*. 2016;17(1):128.

867 15. Arnedo-Pac C, Mularoni L, Muinos F, Gonzalez-Perez A, Lopez-Bigas N.

868 OncodriveCLUSTL: a sequence-based clustering method to identify cancer drivers.

869 *Bioinformatics*. 2019;35(22):4788-90.

870 16. Giacomini CP, Leung SY, Chen X, Yuen ST, Kim YH, Bair E, et al. A gene expression

871 signature of genetic instability in colon cancer. *Cancer Res*. 2005;65(20):9200-5.

872 17. Banerjee A, Ahmed S, Hands RE, Huang F, Han X, Shaw PM, et al. Colorectal cancers

873 with microsatellite instability display mRNA expression signatures characteristic of increased

874 immunogenicity. *Mol Cancer*. 2004;3:21.

875 18. Suo C, Hrydziusko O, Lee D, Pramana S, Saputra D, Joshi H, et al. Integration of

876 somatic mutation, expression and functional data reveals potential driver genes predictive of

877 breast cancer survival. *Bioinformatics*. 2015;31(16):2607-13.

878 19. Gumpinger AC, Lage K, Horn H, Borgwardt K. Prediction of cancer driver genes

879 through network-based moment propagation of mutation scores. *Bioinformatics*.

880 2020;36(Suppl\_1):i508-i15.

881 20. Petrov I, Alexeyenko A. Individualized discovery of rare cancer drivers in global

882 network context. *Elife*. 2022;11.

883 21. Liu X, Wu C, Li C, Boerwinkle E. dbNSFP v3.0: A One-Stop Database of Functional

884 Predictions and Annotations for Human Nonsynonymous and Splice-Site SNVs. *Hum Mutat*.

885 2016;37(3):235-41.

886 22. Rabadan R, Mohamedi Y, Rubin U, Chu T, Alghalith AN, Elliott O, et al. Identification  
887 of relevant genetic alterations in cancer using topological data analysis. *Nat Commun.*  
888 2020;11(1):3808.

889 23. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, et al. The  
890 mutational constraint spectrum quantified from variation in 141,456 humans. *Nature.*  
891 2020;581(7809):434-43.

892 24. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants  
893 from high-throughput sequencing data. *Nucleic Acids Res.* 2010;38(16):e164.

894 25. Sondka Z, Bamford S, Cole CG, Ward SA, Dunham I, Forbes SA. The COSMIC Cancer  
895 Gene Census: describing genetic dysfunction across all human cancers. *Nat Rev Cancer.*  
896 2018;18(11):696-705.

897 26. Chakravarty D, Gao J, Phillips SM, Kundra R, Zhang H, Wang J, et al. OncoKB: A  
898 Precision Oncology Knowledge Base. *JCO Precis Oncol.* 2017;2017.

899 27. Martinez-Jimenez F, Muinos F, Sentis I, Deu-Pons J, Reyes-Salazar I, Arnedo-Pac C, et  
900 al. A compendium of mutational cancer driver genes. *Nat Rev Cancer.* 2020;20(10):555-72.

901 28. Samocha KE, Robinson EB, Sanders SJ, Stevens C, Sabo A, McGrath LM, et al. A  
902 framework for the interpretation of de novo mutation in human disease. *Nat Genet.*  
903 2014;46(9):944-50.

904 29. Cowley GS, Weir BA, Vazquez F, Tamayo P, Scott JA, Rusin S, et al. Parallel  
905 genome-scale loss of function screens in 216 cancer cell lines for the identification of  
906 context-specific genetic dependencies. *Sci Data.* 2014;1:140035.

907 30. Shao DD, Tsherniak A, Gopal S, Weir BA, Tamayo P, Stransky N, et al. ATARiS:  
908 computational quantification of gene suppression phenotypes from multisample RNAi screens.  
909 Genome Res. 2013;23(4):665-78.

910 31. Rubio-Perez C, Tamborero D, Schroeder MP, Antolin AA, Deu-Pons J, Perez-Llamas C,  
911 et al. In silico prescription of anticancer drugs to cohorts of 28 tumor types reveals targeting  
912 opportunities. Cancer Cell. 2015;27(3):382-96.

913 32. Forbes SA, Tang G, Bindal N, Bamford S, Dawson E, Cole C, et al. COSMIC (the  
914 Catalogue of Somatic Mutations in Cancer): a resource to investigate acquired mutations in  
915 human cancer. Nucleic Acids Res. 2010;38(Database issue):D652-7.

916 33. Uhlen M, Zhang C, Lee S, Sjostedt E, Fagerberg L, Bidkhori G, et al. A pathology atlas  
917 of the human cancer transcriptome. Science. 2017;357(6352).

918 34. Goldman MJ, Craft B, Hastie M, Repecka K, McDade F, Kamath A, et al. Visualizing  
919 and interpreting cancer genomics data via the Xena platform. Nat Biotechnol.  
920 2020;38(6):675-8.

921 35. Singh G, Mémoli F, Carlsson GE, editors. Topological Methods for the Analysis of High  
922 Dimensional Data Sets and 3D Object Recognition. PBG@Eurographics; 2007.

923 36. Shuai S, Drivers P, Functional Interpretation Working G, Gallinger S, Stein L,  
924 Consortium P. Combined burden and functional impact tests for cancer driver discovery using  
925 DriverPower. Nat Commun. 2020;11(1):734.

926 37. He X, Cai D, Niyogi P. Laplacian score for feature selection. Proceedings of the 18th  
927 International Conference on Neural Information Processing Systems; Vancouver, British  
928 Columbia, Canada: MIT Press; 2005. p. 507–14.

929 38. Frankish A, Diekhans M, Jungreis I, Lagarde J, Loveland JE, Mudge JM, et al. Gencode  
930 2021. *Nucleic Acids Res.* 2021;49(D1):D916-D23.

931 39. Valentini. G, Re. M. PerfMeas: Performance Measures for ranking and classification  
932 tasks. 2014.

933 40. Oughtred R, Rust J, Chang C, Breitkreutz BJ, Stark C, Willems A, et al. The BioGRID  
934 database: A comprehensive biomedical resource of curated protein, genetic, and chemical  
935 interactions. *Protein Sci.* 2021;30(1):187-200.

936 41. Howe KL, Achuthan P, Allen J, Allen J, Alvarez-Jarreta J, Amode MR, et al. Ensembl  
937 2021. *Nucleic Acids Res.* 2021;49(D1):D884-D91.

938 42. Cheng F, Jia P, Wang Q, Lin CC, Li WH, Zhao Z. Studying tumorigenesis through  
939 network evolution and somatic mutational perturbations in the cancer interactome. *Mol Biol*  
940 *Evol.* 2014;31(8):2156-69.

941 43. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids*  
942 *Res.* 2000;28(1):27-30.

943 44. Omar R, Cooper A, Maranyane HM, Zerbini L, Prince S. COL1A2 is a TBX3 target that  
944 mediates its impact on fibrosarcoma and chondrosarcoma cell migration. *Cancer Lett.*  
945 2019;459:227-39.

946 45. Lee YR, Chen M, Pandolfi PP. The functions and regulation of the PTEN tumour  
947 suppressor: new modes and prospects. *Nat Rev Mol Cell Biol.* 2018;19(9):547-62.

948 46. Chen CY, Chen J, He L, Stiles BL. PTEN: Tumor Suppressor and Metabolic Regulator.  
949 *Front Endocrinol (Lausanne).* 2018;9:338.

950 47. Milella M, Falcone I, Conciatori F, Cesta Incani U, Del Curatolo A, Inzerilli N, et al.

951 PTEN: Multiple Functions in Human Malignant Tumors. *Front Oncol.* 2015;5:24.

952 48. Nassif NT, Lobo GP, Wu X, Henderson CJ, Morrison CD, Eng C, et al. PTEN mutations  
953 are common in sporadic microsatellite stable colorectal cancer. *Oncogene.*  
954 2004;23(2):617-28.

955 49. Luna A, Elloumi F, Varma S, Wang Y, Rajapakse VN, Aladjem MI, et al. CellMiner  
956 Cross-Database (CellMinerCDB) version 1.2: Exploration of patient-derived cancer cell line  
957 pharmacogenomics. *Nucleic Acids Res.* 2021;49(D1):D1083-D93.

958 50. Consortium GT. The Genotype-Tissue Expression (GTEx) project. *Nat Genet.*  
959 2013;45(6):580-5.

960 51. Pham VVH, Liu L, Bracken C, Goodall G, Li J, Le TD. Computational methods for  
961 cancer driver discovery: A survey. *Theranostics.* 2021;11(11):5553-68.

962 52. Wu H, Gao L, Li F, Song F, Yang X, Kasabov N. Identifying overlapping mutated driver  
963 pathways by constructing gene networks in cancer. *BMC Bioinformatics.* 2015;16 Suppl  
964 5:S3.

965 53. Chen B, Wang J, Dai D, Zhou Q, Guo X, Tian Z, et al. AHNAK suppresses tumour  
966 proliferation and invasion by targeting multiple pathways in triple-negative breast cancer. *J*  
967 *Exp Clin Cancer Res.* 2017;36(1):65.

968 54. Park JW, Kim IY, Choi JW, Lim HJ, Shin JH, Kim YN, et al. AHNAK Loss in Mice  
969 Promotes Type II Pneumocyte Hyperplasia and Lung Tumor Development. *Mol Cancer Res.*  
970 2018;16(8):1287-98.

971 55. Cho WC, Jang JE, Kim KH, Yoo BC, Ku JL. SORBS1 serves a metastatic role via  
972 suppression of AHNAK in colorectal cancer cell lines. *Int J Oncol.* 2020;56(5):1140-51.

973 56. Schuetz CS, Bonin M, Clare SE, Nieselt K, Sotlar K, Walter M, et al.  
974 Progression-specific genes identified by expression profiling of matched ductal carcinomas in  
975 situ and invasive breast tumors, combining laser capture microdissection and oligonucleotide  
976 microarray analysis. *Cancer Res.* 2006;66(10):5278-86.

977 57. Walsh MD, Clendenning M, Williamson E, Pearson SA, Walters RJ, Nagler B, et al.  
978 Expression of MUC2, MUC5AC, MUC5B, and MUC6 mucins in colorectal cancers and their  
979 association with the CpG island methylator phenotype. *Mod Pathol.* 2013;26(12):1642-56.

980 58. Misawa K, Kanazawa T, Misawa Y, Imai A, Endo S, Hakamada K, et al.  
981 Hypermethylation of collagen alpha2 (I) gene (COL1A2) is an independent predictor of  
982 survival in head and neck cancer. *Cancer Biomark.* 2011;10(3-4):135-44.

983 59. Zhang Z, Liu R, Jin R, Fan Y, Li T, Shuai Y, et al. Integrating Clinical and Genetic  
984 Analysis of Perineural Invasion in Head and Neck Squamous Cell Carcinoma. 2019;9(434).

985 60. Chung CH, Parker JS, Ely K, Carter J, Yi Y, Murphy BA, et al. Gene expression profiles  
986 identify epithelial-to-mesenchymal transition and activation of nuclear factor-kappaB  
987 signaling as characteristics of a high-risk head and neck squamous cell carcinoma. *Cancer*  
988 *Res.* 2006;66(16):8210-8.

989 61. Zhi X, Lamperska K, Golusinski P, Schork NJ, Luczewski L, Kolenda T, et al. Gene  
990 expression analysis of head and neck squamous cell carcinoma survival and recurrence.  
991 *Oncotarget.* 2015;6(1):547-55.

992 62. Liu C, Guo T, Sakai A, Ren S, Fukusumi T, Ando M, et al. A novel splice variant of  
993 LOXL2 promotes progression of human papillomavirus-negative head and neck squamous  
994 cell carcinoma. *Cancer.* 2020;126(4):737-48.

995 63. Fukumoto I, Kikkawa N, Matsushita R, Kato M, Kurozumi A, Nishikawa R, et al.  
996 Tumor-suppressive microRNAs (miR-26a/b, miR-29a/b/c and miR-218) concertedly  
997 suppressed metastasis-promoting LOXL2 in head and neck squamous cell carcinoma. J Hum  
998 Genet. 2016;61(2):109-18.

999 64. Cancer Genome Atlas Research Network. Electronic address wbe, Cancer Genome Atlas  
1000 Research N. Comprehensive and Integrative Genomic Characterization of Hepatocellular  
1001 Carcinoma. Cell. 2017;169(7):1327-41 e23.

1002 65. Kobayashi Y, Yonehara S. Novel cell death by downregulation of eEF1A1 expression in  
1003 tetraploids. Cell Death Differ. 2009;16(1):139-50.

1004 66. Vera M, Pani B, Griffiths LA, Muchardt C, Abbott CM, Singer RH, et al. The translation  
1005 elongation factor eEF1A1 couples transcription to translation during heat shock response.  
1006 Elife. 2014;3:e03164.

1007 67. Cui H, Li H, Wu H, Du F, Xie X, Zeng S, et al. A novel 3'tRNA-derived fragment  
1008 tRF-Val promotes proliferation and inhibits apoptosis by targeting EEF1A1 in gastric cancer.  
1009 Cell Death Dis. 2022;13(5):471.

1010 68. Ning X, Shi G, Ren S, Liu S, Ding J, Zhang R, et al. GBAS Regulates the Proliferation  
1011 and Metastasis of Ovarian Cancer Cells by Combining with eEF1A1. Oncologist.  
1012 2022;27(1):e64-e75.

1013 69. Bao Y, Zhao TL, Zhang ZQ, Liang XL, Wang ZX, Xiong Y, et al. High eukaryotic  
1014 translation elongation factor 1 alpha 1 expression promotes proliferation and predicts poor  
1015 prognosis in clear cell renal cell carcinoma. Neoplasma. 2020;67(1):78-84.

1016 70. Wu A, Tang J, Dai Y, Huang H, Nie J, Hu W, et al. Downregulation of Long Noncoding



1017 RNA CRYBG3 Enhances Radiosensitivity in Non-Small Cell Lung Cancer Depending on p53

1018 Status. Radiat Res. 2022;198(3):297-305.

1019

## 1020 **Figure legends**

1021 **Fig. 1. Overview of DGAT-cancer.** First, somatic mutations were collected from two  
 1022 databases, ICGC and TCGA, as well as germline mutations from the 1000 Genomes  
 1023 Project. The predicted pathogenicity scores of these mutations were generated by  
 1024 ANNOVAR. The differences between the pathogenicity scores of mutated genes in  
 1025 cancer (somatic) and those in the general population (germline) were evaluated by  
 1026 DGAT-cancer calculating the uEMD scores. RNA-seq data were collected from  
 1027 tumors and paracancerous tissues. The RNA-seq data from tumors were used to  
 1028 construct a topological network. The gene expression topological network was then  
 1029 used to evaluate the frequencies of mutation spectra occurring in different sample  
 1030 clusters. For the mutations occurring in adjacent sample clusters, we calculated the  
 1031 divergence of their frequencies across clusters, which was used as one feature of the  
 1032 genes harboring the mutations. We also calculated the Jensen–Shannon divergence  
 1033 between the mutation frequency and the mRNA expression of the gene across clusters  
 1034 of the network as one feature. All collected features of genes were filtered using  
 1035 Laplacian scores. We transformed these selected features using Hotelling and  
 1036 Box-Cox transformations in order to integrate them into one composite risk score. By  
 1037 using gene scores as weights, we performed Gibbs sampling to sample genes in order

1038 to identify cancer drivers. Finally, we generated null distributions of PP for each gene  
1039 to compute an experience P-value. The genes with  $p_{adj} < 0.01$  were selected as  
1040 candidate cancer driver genes.

1041

1042 **Fig. 2. Evaluation of the cancer driver genes predicted by DGAT-cancer. a**

1043 Enrichment of predicted cancer drivers in known cancer gene sets (CGC, OncoKB,  
1044 IntOGen) and gene sets related to cancer (constraint: selective constraint genes, and  
1045 genes with expression associated with functions of cancer cells). **b** Top 20 genes most  
1046 frequently predicted as cancer drivers. The color depth denotes the rank order of PP  
1047 (decreasing) for genes generated by DGAT-cancer in that cancer type. **c** Enrichment of  
1048 predicted cancer drivers in drug-targeted genes. **d** The proportions of predicted cancer  
1049 drivers whose expression levels were significantly correlated with drug activity  
1050 (measured in terms of 50% growth inhibitory levels) compared to the background  
1051 genes. The P-values were obtained by using Fisher's Exact test to compare the  
1052 predicted cancer drivers with predicted non-cancer drivers.  $*p_{corrected} < 0.05$ ;  
1053  $**p_{corrected} < 0.001$ ;  $***p_{corrected} < 0.0001$ .

1054

1055 **Fig. 3. Comparison of DGAT-cancer with three methods with respect to their**

1056 **prediction of cancer drivers. a** Comparison of DGAT-cancer with other methods in  
1057 terms of their performance as measured by AUPRC (area under the precision–recall  
1058 curve). **b** Comparing posterior probabilities (PPs) given by DGAT-cancer of genes  
1059 predicted as cancer drivers by MutSigCV, OncodriveCLUSL and OncodriveFML to

the genes predicted as non-cancer drivers by those methods. The P-values were obtained by means of the Wilcoxon rank-sum test. In the boxplot, the center lines represent the median, whilst the boxes represent the first and third quartiles. **c** Comparison between genes missed by DGAT-cancer and genes missed by other methods in terms of their feature scores given by DGAT-cancer. Those scores are uEMD-Ex scores, gene expression level in tumor (tumor-med), gene expression level in paracancerous tissues (normal-med), and MutExTDA scores (JSD and C score). Some boxes are blank due to the missing features in that cancer type. The difference was evaluated by Wilcoxon rank-sum test, \*  $p_{corrected} < 0.05$ ; \*\*  $p_{corrected} < 0.001$ ; \*\*\*  $p_{corrected} < 0.0001$ .

**Fig. 4.** Comparing AUPRC performance of DGAT-cancer, DGAT-Mut, DGAT-Exp, MutSigCV, OncodriveFML and OcodriveCLUST to examine the impact of mutation- and gene expression-based features on improving DGAT-cancer.

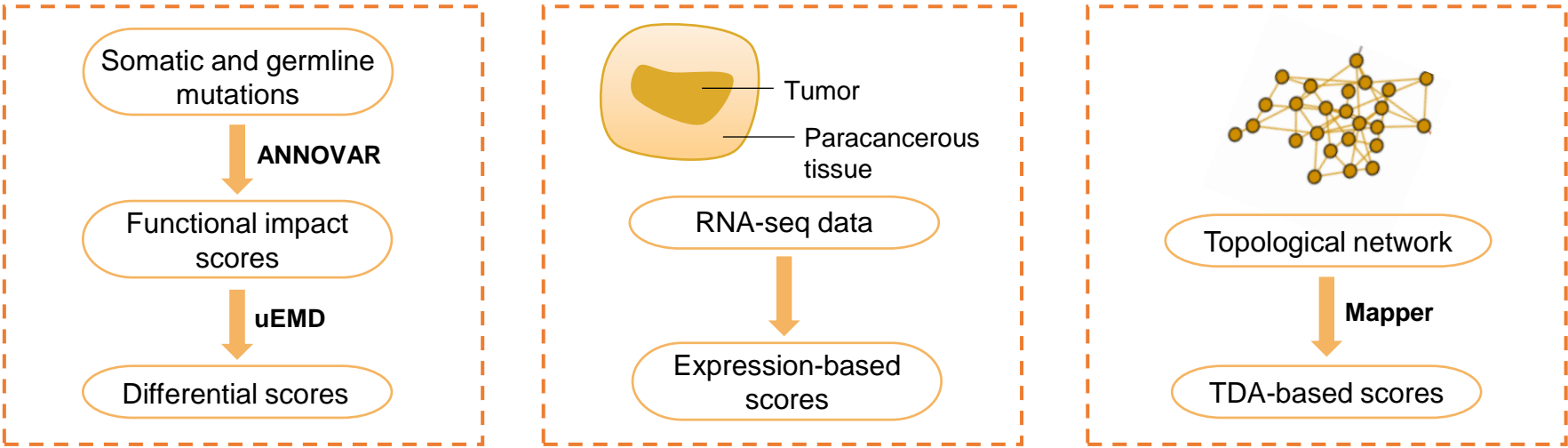
**Fig. 5. *EEF1A1* is highly expressed in glioma and plays a role in regulating the proliferation and migration of glioma cells.** **a** Profile of *EEF1A1* mRNA expression in normal, LGG (WHO grade II–III glioma), or GBM (WHO grade IV glioma) patients in the GTEx and TCGA datasets. *EEF1A1* expression was significantly upregulated in glioma relative to normal brain tissue (left) (Mann-Whitney U test), as well as significantly increased in GBM relative to LGG (right) (Kruskal-Wallis test, adjusted by Bonferroni correction). **b** Kaplan–Meier overall survival plot showing

1082 that the survival rate of glioma patients with high *EEF1A1* expression (red) in the  
 1083 TCGA data set was significantly lower than those with low *EEF1A1* expression (blue)  
 1084 (two-sided log-rank test). **c** *EEF1A1* protein expression in non-tumor brain tissue,  
 1085 LGG and GBM samples using immunohistochemical analysis. Non-tumor brain  
 1086 tissues were obtained from patients with non-tumor brain diseases who had undergone  
 1087 surgical resection. Representative images of immunohistochemical analysis. Scale  
 1088 bar=10μm. The summary of *EEF1A1* protein expression profile in C (Mann-Whitney  
 1089 U test) was on the right. **d** EdU flow cytometry assay showed that the cell  
 1090 proliferation rate was significantly decreased in U251 and U87 cells with the  
 1091 knockdown of *EEF1A1*. **e** *In vitro* colony formation of U251 and U87 cells was  
 1092 decreased after *EEF1A1* knockdown compared to the control. **f** Representative images  
 1093 of the transwell migration assay for migrated cells stained with crystal violet in  
 1094 control cells and sh*EEF1A1* knockdown U251 and U87 cells. Scale bar = 1mm. **g**  
 1095 Representative images of the zebrafish xenograft model used to analyze the  
 1096 proliferation of U251 cells after *EEF1A1* knockdown by measuring GFP fluorescent  
 1097 foci. Scale bar = 0.5mm. Unpaired two-sided t test. \* $p < 0.05$ ; \*\* $p < 0.01$ ;  
 1098 \*\*\* $p < 0.001$ ; \*\*\*\* $p < 0.0001$ .

1099

1100 **Fig. 6. Knockdown of *EEF1A1* improved TMZ sensitivity in glioma cells. a**  
 1101 Knockdown of *EEF1A1* significantly decreased the viability of U251 and U87 cells  
 1102 under different concentrations of TMZ (24h for U251 cells and 48h for U87 cells)  
 1103 ( $n=6$ , paired two-sided t test). **b** Knockdown of *EEF1A1* significantly decreased the

1104 IC<sub>50</sub> of TMZ in U251 and U87 cells ( $n=6$ , unpaired two-sided t test). **c** In the  
 1105 presence of TMZ (50  $\mu$ g/mL), the knockdown of *EEF1A1* significantly reduced  
 1106 colony formation. **d** Quantitation of colony formation in (**c**) ( $n=3$ , unpaired two-sided  
 1107 t test). **e** Flow cytometry results showing the proportions of apoptotic cells (the  
 1108 Annexin V V450-positive cells) in the *EEF1A1*-knockdown and control groups in  
 1109 U251 and U87 cells. **f** The summary of six independent experiments of (**e**) (unpaired  
 1110 two-sided t test). **g** Western-blot showing the cleaved caspase-3 and bcl-2 protein  
 1111 expression in the *EEF1A1*-knockdown and control groups in U251 and U87 cells. **h**  
 1112 Quantitation of four independent experiments of (**g**). (Mann-Whitney U test)  
 1113  $*p < 0.05$ ;  $**p < 0.01$ ;  $***p < 0.001$ ;  $****p < 0.0001$ .



bioRxiv preprint doi: <https://doi.org/10.1101/2023.05.02.539093>; this version posted May 3, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

