

Clinical Feature-Related Single-Base Substitution Sequence Signatures Identified with an Unsupervised Machine Learning Approach

Ji Hongchen

Xijing Hospital <https://orcid.org/0000-0003-3838-7354>

Li Junjie

Fourth Military Medical University: Air Force Medical University

Zhang Qiong

Fourth Military Medical University: Air Force Medical University

Yang Jingyue

Fourth Military Medical University: Air Force Medical University

Duan Juanli

Fourth Military Medical University: Air Force Medical University

Wang Xiaowen

Fourth Military Medical University: Air Force Medical University

Ma Ben

Chinese PLA General Hospital

Zhang Zhuochao

Chinese PLA General Hospital

Pan Wei

Fourth Military Medical University: Air Force Medical University

Zhang Hongmei (✉ zhm_fmму@163.com)

Fourth Military Medical University <https://orcid.org/0000-0002-3991-3750>

Research article

Keywords: Mutation sequence, unsupervised learning, cancer, clinical feature, prognosis

Posted Date: March 22nd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-344127/v1>

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

[Read Full License](#)

1 Clinical feature-related single-base substitution sequence signatures identified with an
2 unsupervised machine learning approach

3

4 Ji Hongchen^{1,2,*}, Li Junjie^{3,*}, Zhang Qiong^{1,*}, Yang Jingyue^{1,*}, Duan Juanli⁴, Wang
5 Xiaowen¹, Ma Ben², Zhang Zhuochao², Pan Wei¹, Zhang Hongmei¹

6

7 ¹ Department of Oncology, Xijing Hospital, Fourth Military Medical University. No. 127
8 West Changle Road, Xi'an, China

9 ² Faculty of Hepatopancreatobiliary Surgery, Chinese PLA General Hospital. No. 28 Fuxing
10 Road, Beijing, China

11 ³ Department of Emergency, Xijing Hospital, Fourth Military Medical University. No. 127
12 West Changle Road, Xi'an, China

13 ⁴ Department of Hepatobiliary Surgery, Xijing Hospital, Fourth Military Medical
14 University. No. 127 West Changle Road, Xi'an, China

15 * These authors contributed equally to this work.

16

17 **Corresponding author**

18 Dr. Zhang Hongmei, Department of Oncology, Xijing Hospital, Fourth Military Medical
19 University, No. 127 West Changle Road, Xi'an, 710032, China. Tel: (+86) 13991293309, E-
20 mail: zhm@fmmu.edu.cn

21

22

23

24

25

26 **Abstract**

27 **Background:** Mutation processes leave different signatures in genes. For single-base
28 substitutions, previous studies have suggested that mutation signatures are not only reflected
29 in mutation bases but also in neighboring bases. However, because of the lack of a method to
30 identify features of long sequences next to mutation bases, the understanding of how flanking
31 sequences influence mutation signatures is limited.

32 **Methods:** We constructed a long short-term memory – self organizing map (LSTM-SOM)
33 unsupervised neural network. By extracting mutated sequence features via LSTM and
34 clustering similar features with the SOM, single-base substitutions in The Cancer Genome
35 Atlas database were clustered according to both their mutation site and flanking sequences.
36 The relationship between mutation sequence signatures and clinical features was then
37 analyzed. Finally, we clustered patients into different classes according to the composition of
38 the mutation sequence signatures by the K-means method and then studied the differences in
39 clinical features and survival between classes.

40 **Results:** Ten classes of mutant sequence signatures (mutation blots, MBs) were obtained
41 from 2,141,527 single-base substitutions via LSTM-SOM machine learning approach.
42 Different features in mutation bases and flanking sequences were revealed among MBs. MBs
43 reflect both the site and pathological features of cancers. MBs were related to clinical
44 features, including age, gender, and cancer stage. The class of an MB in a given gene was
45 associated with survival. Finally, patients were clustered into 7 classes according to the MB
46 composition. Significant differences in survival and clinical features were observed among
47 different patient classes.

48 **Conclusions:** We provided a method for analyzing the characteristics of mutant sequences.
49 Result of this study showed that flanking sequences, together with mutation bases, shape the
50 signatures of SBSs. MBs were shown related to clinical features and survival of cancer

51 patients. Composition of MBs is a feasible predictive factor of clinical prognosis. Further
52 study of the mechanism of MBs related to cancer characteristics is suggested.

53

54

55 **Keywords**

56 Mutation sequence; unsupervised learning; cancer; clinical feature; prognosis

57

58 **Background**

59 The stability of the cell genome is continually threatened by endogenous and
60 exogenous factors that may lead to DNA damage [1, 2]. If not repaired properly, DNA
61 damage may result in genetic mutations [3, 4]. The development of cancers involves a series
62 of genetic mutations [5]. A number of internal and external factors underlying genetic
63 mutations have been identified, such as smoking, alcohol consumption and mismatch repair
64 deficiency [5, 6]. In some kinds of cancers, such as colon cancer and breast cancer, there has
65 been a great deal of research elucidating the relationship between genetic mutations and
66 cancer-related processes [7]. However, in most cases, the patten of genetic mutations and its
67 role in tumor progression are still poorly understood.

68 Genetic mutations include single-base substitutions (SBSs), small insertions and
69 deletions (indels), genome rearrangement and chromosome copy-number changes [8]. SBSs
70 contribute the largest proportion of genetic mutations. Mathematical methods have been used
71 to decipher mutation signatures from somatic mutation catalogs [2, 8-15]. At present, large
72 amounts of mutation data from cancer patients have been obtained and made available in
73 relevant databases, such as The Cancer Genome Atlas (TCGA) database. In the context of
74 increasing sample sizes, a number of mutation signatures that are correlated with certain
75 mutation processes have been identified [16, 17]. The clustering methods for SBSs applied in
76 some studies have included 1-2 bases next to mutated bases, and the results have suggested
77 that bases next to the mutation site influence mutation signatures [2, 8]. However, the
78 inclusion of adjacent genes in such analyses leads to an exponential increase in the number of
79 possible classifications. Because of the lack of a highly efficient method to identify features
80 of long sequences next to mutation bases, the understanding of how flanking sequences
81 influence somatic mutation characteristics is limited.

82 The application of machine learning, especially neural networks, makes it possible to
83 effectively mine information from large amounts of data. A long short-term memory (LSTM)
84 network is a special kind of recurrent neural network (RNN). Compared with a naive RNN,
85 LSTM performs better in extracting features from long sequences, such as sentences [18, 19].
86 LSTM has been used to analyze DNA or RNA sequence information in some studies [20-22].
87 A self-organizing map (SOM) algorithm is an unsupervised clustering algorithm. The method
88 of "competitive learning" can identify interconnections between samples and present their
89 categories in a lower-dimensional form [23, 24]. The use of LSTM to extract the features of
90 mutated sequences and the identification of similar features with the SOM algorithm
91 provided an approach for analyzing the characteristics of mutated sequences and their
92 relationship with cancer development. In this study, we established an LSTM-SOM
93 unsupervised learning network to include long flanking sequences into the analysis of mutant
94 sequence signatures. Via the LSTM-SOM method, we clustered the mutation sequences in
95 the TCGA database into different classes (for a clear understanding, mutant sequence
96 signatures clustered by the LSTM-SOM are referred to as mutation blots, MBs) and then
97 analyzed the relationships among MBs, clinical features, and cancer patient survival.

98

99 **Methods**

100

101 *Data sources*

102 SBS data and clinical data of patients enrolled in this study were obtained from the
103 TCGA database. First, the SBS information includes the sample barcode, chromosomal
104 location, mutant allele, reference allele, Hugo gene symbol, etc. Clinical data, including age,
105 gender, weight, cancer stage, and survival time or time to the last follow-up, were extracted
106 according to the sample barcode. In the LSTM-SOM model, 100 flanking bases were

107 included in the analysis, and the flanking sequence was obtained from the Genome Reference
108 Consortium human genome build 38 (GRCh38) based on the mutation sites of SBSs in
109 TCGA data.

110

111 *LSTM-SOM model building*

112 In brief, our LSTM-SOM model works via a cycle of 3 steps: 1. extraction of the
113 feature vector of the mutant sequence by LSTM; 2. clustering of feature vectors by the SOM,
114 and feature vectors are updated at the same time to bring vectors with similar features closer
115 together; and 3. use of the updated feature vectors for the labeling and training of the LSTM
116 model.

117 **Step 1. Obtaining feature vectors with LSTM.** Mutant sequences are represented in
118 the form of a matrix. A 1×2 vector is used to represent different bases (A: [0, 0]; T: [0, 1]; C:
119 [1, 0]; G: [1, 1]; N:[-1, -1]). When placing the reference sequence in the corresponding
120 position, mutated bases can be recorded as a 1×4 vector. When the flanking bases are
121 included, a mutated sequence can be represented by an n×4 matrix. For example, CATTG >
122 CACTG can be expressed as follows:

$$123 \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

124 RNNs have long been used in the analysis of sequence data. A naive RNN effectively
125 analyzes short sequences. An LSTM network is based on the network structure of RNNs [21].
126 The LSTM approach introduces the mechanisms of "forgetting" and "memory". Thus, the
127 capacity of the LSTM network to analyze long sequences is improved by controlling the
128 long-term state [22]. As the "forgetting" mechanism of LSTM, the unit closer to the end of the
129 sequences has a greater influence on the output of LSTM. In our study, LSTM was designed

130 to read from both ends of the mutated sequence.''' In this way, the mutation site is placed at
131 the ends of both sequences to reinforce its influence on the LSTM output.

132 We used the torch.nn package in PyTorch to construct a neural network. The LSTM
133 procedure that we used consists of two hidden layers, each with 64 nodes. The data
134 subsequently entered a full connection layer, and a 1×8 vector was finally output as the
135 feature vector of a single mutated sequence.

136 **Step 2. Clustering with the SOM.** The SOM consists of two kinds of layers: an input
137 layer and a competition layer. The randomized units in the competition layer were trained to
138 describe the distribution of units in the input layer via the mechanism of "competitive
139 learning" [25]. In the SOM process of the LSTM-SOM model, the feature vector obtained
140 from the LSTM process is used as the input. Units in the competition layer are adjusted
141 continuously according to their distance to the input unit. For one input unit, the unit in the
142 competition layer nearest to it is regarded as the "winning unit", which will move the
143 maximal distance to the input unit (target), and for the other units, their travel distance to the
144 target decreases with the increase in the distance to the winning unit. To avoid an excessive
145 concentration of the results, we set a threshold value in the model. When the distance
146 between the competition layer unit and the target is over the threshold value, the unit will
147 move in the opposite direction to the target. In particular, not only will units in the
148 competition layer be updated in our SOM model, but the input unit will also be updated in the
149 opposite direction of the vector sum of the competition layer unit movement. Then, the
150 updated input unit will be used as a label to train the LSTM model.

151 First, we obtained feature vectors of 100 samples from LSTM in one batch, and they
152 were used as the input units of the SOM. The settings included 200 units in the SOM
153 competition layer. For each input vector, the Euclidean distance between it (x) and each unit
154 in the competition layer (w_j) was calculated as follows:

155
$$d_j(x) = \sqrt{\sum_{i=1}^D (x_i - w_{ji})^2}$$

156 The unit closest to x is recorded as w_{min} , and the distance between w_{min} and each
 157 other competition layer unit is calculated as follows:

158
$$d_j(w_{min}) = \sqrt{\sum_{i=1}^D (w_{ji} - w_{min_i})^2}$$

159 A threshold of S was set in the process of training. If $d_j(w_{min}) \leq S$, w_j will move in
 160 the direction of x ; otherwise, w_j will move in the opposite direction. The transportation
 161 distance decays with an increase in $d_j(w_{min})$. The neighborhood function refers to the
 162 Gaussian function [25]:

163
$$D(w_j) = e^{-\frac{d_j(w_{min})^2}{2\pi\sigma^2}}$$

164 In the decay function, σ is a constant that affects the amplitude of transportation
 165 distance decay. The update vector is as follows (where L is the learning rate of the SOM):

166
$$\Delta(w_j) = \begin{cases} L \times D(w_j) \times (w_j - x) & d_j(w_{min}) \leq S \\ -L \times D(w_j) \times (w_j - x) & d_j(w_{min}) > S \end{cases}$$

167 When the distance between w_j and the target x is less than S , they will approach each
 168 other. Otherwise, they will pull away from each other. Due to the existence of the decay
 169 function, the influence of distant units on each other is very small, and no excessive
 170 dispersion of units was observed in training. To avoid overfitting, the units in the SOM
 171 competition layer are updated after each training batch of 100 samples. The samples in each
 172 batch are selected randomly from different cancers. To change the discrete status of the input
 173 vectors and cause similar input vectors to aggregate, the input units are updated in the
 174 opposite direction (x is the input vector):

175
$$x(new) = x + \sum_{j=1}^{200} \Delta(w_j)$$

176 **Step 3. Training the LSTM model.** The updated $x(new)$ is used as the label to train
177 the LSTM network. In this way, the output feature vectors of LSTM with similar features can
178 be gradually closed.

179 The above three steps are repeated until a clear, stable classification is obtained.

180

181 *Obtain the classification*

182 During training, the units in the competition layer of the SOM were sorted according
183 to the distance to w_{min} . S was set as the distance of unit rank 40 (5% of entire competition
184 layer units) to w_{min} . After each iteration of SOM analysis, the updated input data were used
185 as labels to train the LSTM model for 2 iterations. The LSTM learning rate was set as 0.001.
186 The SOM learning rate was set as 0.005.

187 Two classes were obtained after one round of training. After 3 rounds of training, a
188 total of 8 clustered classes were obtained. It was observed that there were 2 classes showing
189 significantly larger sample sizes than the other classes. Therefore, an additional round of
190 clustering was carried out in the 2 classes. Finally, we obtained 10 classes of mutated
191 sequences.

192

193 *Analysis of clinical features*

194 In the analysis of clinical features, measurement data were expressed as the mean \pm
195 standard deviation. In the analysis of differences between groups, an independent-samples T
196 test (number of groups = 2) or analysis of variance (ANOVA) (number of groups > 2) was
197 used. Enumeration data were expressed as count data, and chi-square analysis was used for
198 difference testing. A sample was removed if the data of an item required for statistics were
199 missing. $P < 0.05$ was considered to indicate a statistically significant difference. In the

200 survival analysis, the log-rank test was used to analyze the difference in survival between
201 different groups.

202

203 ***Clustering of patients according to the MB composition***

204 Patients were clustered according to their MB composition. In the clustering method
205 according to the MB composition, each kind of MB was reflected as the percentage of the
206 entire MB in one patient. The K-means method was used for clustering performed by the K-
207 means method in the scikit-learn package. An "elbow method" was used to evaluate the K
208 value (number of clustered groups) [26]. The K value evaluated in different cancers, and the
209 entire sample was generally between 5-8. After comparing the clustering results, K=7 was
210 selected as the class number for K-means clustering.

211

212 ***Code available***

213 All mathematical methods were performed with Python. The code for the
214 pretreatment of TCGA data and the construction, training and testing of the model is stored at
215 https://github.com/FruedDolce/SBS_CLUSTER/. For clinical data analysis, patient clustering,
216 survival analysis and drawing, the code is stored at <https://github.com/FruedDolce/SATA/>.
217 All the code is open source and freely available.

218

219 **Results**

220

221 ***SBS clustering via the LSTM-SOM unsupervised machine learning approach***

222 A total of 2,141,527 somatic SBS data points from 9596 patients were collected from
223 the TCGA database. For each SBS sample, 100 flanking bases (50 bases at the 5' end and 50
224 at the 3' end) were included in the LSTM training data.

225 In brief, our LSTM-SOM model functions by extracting the features of mutant
226 sequences via the LSTM network and then taking the generated feature vector as the input
227 data for the SOM. Units in competitive layers of the SOM are then refreshed to edges closer
228 to the distribution of the input data. After each iteration of the SOM in our LSTM-SOM
229 model, not only will the units in the competitive layer of the SOM be refreshed, but the input
230 data generated by LSTM will also be adjusted in the opposite direction (Fig. 1a). Then, the
231 refreshed input data are used as the labels to train the LSTM model. The above steps were
232 repeated until the LSTM outputs formed clear classifications.

233 One hundred samples from patients with different cancers were selected randomly in
234 each training iteration. In the LSTM process, the influence of unit data on the LSTM output
235 results decreased with increasing distance to the ending unit. The LSTM process was carried
236 out on both sides of the mutation site in opposite directions. Thus, the mutation site was
237 placed at the end of both sequences to expand its influence on LSTM output and to reflect the
238 difference between the reference allele and mutant allele. Mutated sequences were clustered
239 into 2 types after one stage of training. Thus, we obtained 8 classes of MBs after 3 stages of
240 training. Then, an additional stage of training was performed for 2 classes of MB with a
241 significantly larger number of samples and ultimately revealed 10 classes of MBs, recorded
242 as MB 1-MB 10 (Fig. 1b).

243

244 *Features of mutation bases and flanking sequences in different MBs*

245 Following the principle of complementary base pairing, 4 kinds of bases form 6
246 classes of base substitutions: C>A, C>G, C>T, T>A, T>C, and T>G, where base substitutions
247 are represented by the pyrimidine residue of the base pair. Among the 10 classes of MBs
248 clustered by the LSTM-SOM, 4 contained a single kind of mutation (MB 7: C>A; MB 8:
249 C>T; MB 9: T>A; MB 10: T>C). The other 5 classes contained multiple types of mutations

250 (MB 1: C>G, T>C, and T>G; MB 2: C>A, C>T, T>A and T>G; MB 3: T>A and T>C; MB 4:
251 C>G and C>T; MB 5: C>A, C>G, T>C and T>G; MB 6: C>A, C>T, T>A and T>G) (Fig. 2,
252 Additional file 2: Table S1).

253 The clustering results were strongly influenced by the flanking bases of the mutation
254 site. For example, both MB 5 and MB 7 exhibited C>A mutations, and the flanking bases of
255 MB 5 were dominated by T bases, but MB 7 was dominated by A bases. Differences in
256 flanking bases could also be observed in other classes of MBs with similar mutation features,
257 such as MB 2 and MB 6, MB 4 and MB 8 (Fig. 2). With an increase in the distance from the
258 mutation site, the proportions of the four bases tended to become balanced. In the analysis of
259 cancers with high incidence (lung, breast, prostate, colon, stomach, bladder, ovary, cervix
260 uteri, liver, thyroid, skin and kidney cancers), the composition of the bases in the mutation
261 site and the flanking sites of each MB basically followed that in the entire sample (Additional
262 file 1: Figure S1).

263

264 ***MBs reflect the difference in cancers according to both location and pathological type***

265 Significant differences in the composition of MBs existed among cancers with
266 different pathologies. Overall, MB 4, MB 5, MB 7 and MB 8 accounted for much greater
267 percentages of the MBs than the other classes of MBs, especially MB 4 and MB 8 (Fig. 3a).
268 Malignant mesenchymal tumors seemed to present a higher percentage of MB 2 and MB 6
269 than epithelial malignant tumors. Transitional cell carcinoma of the urinary tract showed a
270 distinctly higher MB 1 incidence than other cancers. Cancers of germ cells and the glomus
271 (paragangliomas) exhibited a high proportion of MB 10. An obvious feature of melanomas
272 was the dominance of MB 4 and MB 8. This finding suggested that these classes of MBs may
273 be correlated with ultraviolet light exposure.

274 The components of MBs varied in different cancers, and some cancers presented
275 distinct features. The proportions of MBs in different cancers were influenced by the
276 pathological type to some extent (Fig. 3b). For example, cancers of the skin and lymph nodes
277 showed extraordinarily high proportions of MB 4 and MB 8 but small proportions of other
278 MBs. In both types of cancers, melanoma is the major pathologic type. Lung cancer
279 presented high proportions of MB 5 and MB 7. Among the 2 major pathological types of lung
280 cancer, adenocarcinoma (AC) exhibited much higher proportions of MB 5 and MB 7 than did
281 squamous cell carcinoma (SCC). This was consistent with the MB composition in the two
282 pathological types. However, for the same pathological type, differences in the MB
283 composition could be observed in different cancers. For example, AC of the colon presented
284 higher proportions of MB 4 and MB 8 than did AC of the lung. SCC of the lung exhibited
285 more MB 5 and MB 7 than SCC in the head and neck (Additional file 1: Figure S2).

286 In most classes of MBs, the frequency of genes that were commonly mutated in
287 malignant tumors, such as TTN, TP53, and MUC16, was relatively high. Distinct features
288 existed in some classes of MB. The proportion of TP53 mutations was generally high, but it
289 was relatively low in MB 7 and MB 10. Remarkably, BRAF was the most common mutated
290 gene in MB 9. A higher proportion of PIK3CA mutations was observed in MB 8 and MB 10
291 than in the other classes of MBs (Additional file 2: Table S2). More distinct features could be
292 observed when considering specific cancers. For example, in pancreatic cancer, MB 4 and
293 MB 5 contained a higher frequency of KRAS mutations than did the other classes of MBs. In
294 kidney cancer, the frequency of VHL mutations ranked high in MB 3, MB 5, MB 7 and MB
295 9. In skin and thyroid cancers, BRAF mutations were common in MB 9 but not in the other
296 classes of MBs (Additional file 1: Figure S3).

297

298 *Survival analysis of patients with different MBs in the same mutation gene.*

299 In the mutated genes with a high frequency, the composition of MBs varied between
300 different kinds of cancers. Such differences reflected the overall MB composition of each
301 cancer (Fig. 4). To further study the influence of certain genes with different MBs on
302 survival, we analyzed the survival of patients who exhibited mutations in genes with high
303 mutation frequencies (TTN, MUC16, TP53, DNAH5, USH2A, PIK3CA, SYNE1, etc.).
304 Patients were grouped according to the MB classification of specific genes. Patients carrying
305 genes with MB 4 and MB 8 mutations usually showed better survival. In contrast, MB 1, 6,
306 and 9 in a gene could predict worse survival (Fig. 4 and Additional file 1: Figure S4 and S5).

307

308 *Relationship between MBs and clinical features of cancer patients*

309 Analysis was performed to determine the relationship between MBs and the clinical
310 features of tumor patients, including their age, gender, weight, AJCC stage and TNM stage.
311 The change in MBs showed a nonmonotonic trend with patient age. The percentages of MB
312 2, MB 5, and MB 7 in single patients increased with age within the first interval (<70 for MB
313 2; <75 for MB 5 and MB 7) but decreased when age exceeded the threshold. This trend was
314 reversed for MB 4 and MB 8. An exception was observed for MB 9, whose proportion in
315 single patients decreased monotonically with age. The proportion of MBs in a single patient
316 generally varied between the genders. Female patients were likely to show higher percentages
317 of MB 2, MB 3, MB 5, MB 6 and MB 10, while male patients exhibited higher percentages
318 of MB 4, MB 8 and MB 9. The difference was not significant in MB 1 and MB 7. No
319 apparent rule regarding the relationship between the weight and MB composition of a patient
320 was observed (Fig. 5).

321 Although the detailed methods of AJCC staging in different cancers are not the same,
322 they generally follow similar principles [27]. Therefore, we merged the subdivisions of the
323 stages in some cancers to analyze cancer stage. The proportions of MB 3, MB 7 and MB 9

324 showed a decreasing trend with increasing T and N stages. In contrast, MB 4 and MB 8 had a
325 positive relationship with T and N stages. For some MBs, their relationship with cancer
326 staging was complicated. MB 5 decreased with the progression of T and N stages, but M1
327 patients presented more MB 5 than M0 patients. MB 2 and MB 6 exhibited a remarkably high
328 prevalence in N3 patients. (Fig. 5).

329 In most cancers, the MB composition at different ages basically followed the pattern
330 shown in the total samples. The proportion of MB 2 in most cancers was significantly higher
331 in males than in females. Regarding cancer staging, T and M stages showed obvious
332 tendencies in most kinds of cancers, and their trends were basically consistent with those for
333 the total sample. Stomach cancer and colon cancer, in particular, showed opposite MB
334 tendencies in T and N stages compared with the entire sample and with other cancers with
335 high incidence (Additional file 1: Figure S6-S9).

336

337 *Composition of MBs in cancer patients is related to clinical prognosis*

338 To further analyze the influence of the MB composition on the clinical features of
339 patients, a K-means clustering method was used to classify patients according to MB
340 composition. Different kinds of MBs were recorded according to their proportion rather than
341 their number in a single patient. K=7 was selected as the number of classes to be
342 distinguished. Clustered patients were designated as Classes 1-7. The compositions of MBs in
343 different cancers are shown in Fig. 6a.

344 In the survival analysis, significant differences in survival curves were observed in
345 different classes of patients (Fig. 6b). In the pairwise survival analysis, patients in Classes 2,
346 4, and 5 showed better survival, and patients in Classes 1, 3, 6, and 7 showed worse survival
347 (Fig. 6c). In the analysis of specific cancers, survival in different classes of patients generally
348 followed the results obtained for the total sample but with some discrepancies that were not

349 significant. Class 3 patients, in particular, seemed to show poor survival for most of the
350 analyzed cancers (Additional file 1: Figure S10).

351 Patients of different classes showed distinct clinical features (Fig. 6d and e).
352 According to AJCC staging, a significantly lower proportion of stage IV patients and a higher
353 proportion of stage I patients were observed in Classes 4 and 5, which may be related to the
354 better survival of these 2 classes of patients. Interestingly, Class 4 included significantly more
355 T4 patients but hardly any M1 patients. This suggests that the MB composition of Class 4
356 may be associated with the local progression of cancers. Class 6 patients showed the highest
357 percentage of AJCC stage 4 and lowest percentage of AJCC stage I, which may be the reason
358 for the poor survival of these patients. In the analysis of age, patients of Class 3 were found to
359 present significantly greater ages. At the same time, the weight of Class 3 patients was also
360 high. Class 1 patients exhibited a high percentage of AJCC stage 1 and a low percentage of
361 stage IV. Moreover, the proportion of N0 patients in Class 1 was significantly higher than
362 that in other classes.

363

364 **Discussion**

365 Several studies on mutation signatures have been published. Most of the studies were
366 based on the TCGA or other databases. Several mathematical methods are now used to
367 cluster the mutation signature [2, 4, 8, 10, 11]. Some of the studies have suggested that
368 adjacent bases may affect the characteristics of the mutation signature. However, because of
369 the lack of a method to analysis the long sequences near the mutation base, the studies have
370 been limited to 1-2 bases next to the mutation base, and the understanding of the influence of
371 flanking sequences on mutation characteristics is still limited.

372 An increase in the number of included flanking bases leads to an exponential increase
373 in the number of possible classifications. In our study, together with the 50 flanking bases on

374 both sides, there were theoretically 6×4^{100} possible classes, making it nearly impossible to
375 analyze such classes with classical statistical methods. LSTM is a machine learning approach
376 that is good at extracting the features of long sequences [28]. This approach provided us with
377 a method for extracting the features of mutated sequences across a wider spatial scope. A
378 follow-up SOM method can then be used to discover internal relationships between the
379 extracted features and ultimately obtain different categories of mutant sequences. To avoid
380 overfitting of the model, the weight of the vectors in the competitive layer was updated after
381 all input data were trained in one batch. Each iteration of training included 2 LSTM iterations
382 and 2 SOM iterations. In this way, we identified 10 classes of mutation sequences. No one
383 kind of mutation was contained in a single class of MBs. The composition of the bases
384 flanking the mutation sites differed considerably. Generally, units located far from the
385 endpoint had less influence on the LSTM output than those located close to the endpoint [22].
386 This characteristic was reflected in the flanking bases of the mutation site. In all kinds of
387 MBs, the proportions of A, T, C, and G were quite different among the bases near the
388 mutation site. With an increase in the distance from the mutation site, the proportions of the
389 four bases tended to become balanced.

390 The analysis of MBs in different kind of cancers suggested that MBs may
391 comprehensively reflect the difference in cancers according to both location and pathological
392 type. Previous studies have proven that different mutation signatures may be associated with
393 different triggers involved in various mutation processes and result in differing biological
394 behaviors of cancers [2, 8]. A variety of mutation signatures that may be related to the
395 biology and etiology of cancer have been identified [2, 8, 14-20, 29-32]. Our study suggests
396 that a high incidence of MB 4 and MB 8 is associated with pathologic types of cancer that are
397 believed to be caused by external mutagenic exposure, such as SCC, transitional cell
398 carcinoma, malignant mesothelioma and complex epithelial carcinoma. We also found that

399 some kinds of cancer, such as melanoma and transitional cell carcinoma, had distinctive
400 features that are worthy of further study to determine the relationship between each MB and
401 specific cancer processes. In a given gene, SBSs may occur at different bases with different
402 features and present as different kinds of MBs. Each gene with a high mutation frequency
403 contained multiple kinds of MBs. On further study of the MB proportion in genes that are
404 highly frequently mutated, we observed differences in the mutated gene compositions of
405 different MBs. This finding suggests that attention should be paid to the effect of different
406 MBs on the characteristics of cancer when they occur in the same gene.

407 Then, in the subsequent analysis, we focused on the relationship between MBs and
408 clinical features, including survival. First, survival analysis between patients with different
409 MBs in the same gene showed a significant correlation between survival and MBs for
410 specific genes. In the analysis between MBs and clinical features, it was observed that the
411 proportion of MBs generally showed an obvious tendency with a change in clinical features,
412 which suggests that characteristics of MBs reflect the characteristics of cancers. Considering
413 the differences in the clinical significance of staging in different cancers, further analysis was
414 performed on each cancer with high incidence. Generally, in most cancers, the MB
415 composition in patients with different clinical features basically followed the pattern
416 observed in all samples. While there were some exceptions, such as in stomach cancer and
417 colon cancer, MB tendencies in T and N stages were opposite to those in the entire sample
418 and to other cancers with high incidence. This result suggests that local and lymph node
419 progression in gastrointestinal cancers may exhibit distinct mechanisms. In the analysis of
420 age, younger and older patients showed similar MB compositions in the form of a conic
421 structure in the bar graph. This suggests that the similar cancer biologies of young and old
422 patients require further study. Generally, although the results showed a clear relationship

423 between MBs and clinical features, details of the relationship as well as its mechanism still
424 require further study.

425 To further explore the translational relevance of MBs, we then clustered patients into
426 7 classes according to MB composition. Interestingly, patients with a balanced composition
427 of MBs (Classes 1, 3 and 5, especially Class 3) were associated with poor survival for most
428 of the analyzed cancers. These results suggest that a balanced MB composition may predict
429 poor survival in patients and may be related to mixed mutation triggers. Some classes of
430 patients showed typical clinical features. For example, patients in Class 3 were older and
431 weighed more than those in the other classes. These factors may be partly responsible for the
432 poor survival of patients in Class 3. In contrast, although the patients in Class 1 were older,
433 they did not weigh more than those in the other classes. Therefore, further study is still
434 needed to determine the mechanism by which patients in Class 1 experience poor survival.
435 Due to the natural differences in cancer incidence, large differences exist between different
436 cancers. In different cancers, MB may be involved in different kinds of cancer-related
437 processes. Therefore, the analysis of the relationship between MBs and distinctive clinical
438 features in specific kinds of cancer can provide more information about how MBs are related
439 to cancer etiology, processes, prognosis and drug susceptibility.

440 There were still some constraints and limitations to this study. The clustering results
441 obtained from the LSTM-SOM model were largely dependent on the selection of SOM
442 parameters (especially the neighborhood function parameter). There exists the possibility that
443 when training with other parameters, the classification obtained may have been related to
444 clinical features that were not included in this study and thus need further study. Moreover,
445 the mechanism of machine learning models is difficult to explain [33]. It would be
446 meaningful to use a mathematical method to explore the mechanism of the LSTM-SOM
447 functions to improve the interpretability of the LSTM-SOM model and to explain the

448 formation of different classes of MB to determine how sequences of bases affect the
449 characteristics of cancers. Different MBs may also be involved in complex changes in three-
450 dimensional chromosome conformation. Moreover, molecular biology methods are helpful
451 for explaining the different characteristics of MBs.

452

453 **Conclusion**

454 This study provided a method for analyzing the characteristics of mutant sequences. Result of
455 this study showed that flanking sequences, together with mutation bases, shape the signatures
456 of SBSs. The analysis of MBs in different kind of cancers suggested that MBs reflect the
457 difference in cancers according to both location and pathological type. Mutation sequence
458 signatures (MBs) identified via LSTM-SOM method in this study were shown related to
459 clinical features and survival of cancer patients. Composition of MBs is a feasible predictive
460 factor of clinical prognosis. Patients with balanced MB composition seems to have worse
461 survival. Further study on the interpretability of LSTM-SOM network and on the mechanism
462 of MBs related to cancer characteristics is suggested.

463

464 **List of abbreviations**

LSTM:	long short-term memory
SOM:	self-organizing map
SBS:	single-base substitution
MB:	mutation blot
TCGA:	The Cancer Genome Atlas
GRCh38:	Genome Reference Consortium human genome build 38
AJCC:	American Joint Committee on Cancer
AC:	adenocarcinoma

SCC: squamous cell carcinoma

465

466 **Declarations**

467

468 *Ethics approval and consent to participate*

469 Not applicable

470

471 *Consent for publication*

472 Not applicable

473

474 *Availability of data and materials*

475 SBS data and clinical data of patients enrolled in this study were obtained from the TCGA

476 database. The .maf files including SBS somatic mutation data and the .xml files including

477 clinical data were downloaded from <https://portal.gdc.cancer.gov/repository/>. Reference

478 genome sequences (Genome Reference Consortium human genome build 38, GRCh38) were

479 downloaded from ftp://ftp.ncbi.nih.gov/genomes/H_sapiens/).

480 The code for the pretreatment of TCGA data and the construction, training and testing of the

481 model is stored at https://github.com/FruedDolce/SBS_CLUSTER/ . For clinical data

482 analysis, patient clustering, survival analysis and drawing, the code is stored at

483 <https://github.com/FruedDolce/SATA/> . All the code is open source and freely available.

484

485 *Competing interests*

486 The authors declare that they have no competing interests

487

488 *Funding*

489 This study was supported by the Xijing Hospital Science Foundation (XJZT19ML38).

490

491 ***Author contributions***

492 J. H. C., L. J. J., and Z. H. M. designed this study. J. H. C., Z. Q., M. B., and D. J. L designed
493 and carried out the mathematical method used in this study. J. H. C. and Y. J. Y. wrote the
494 manuscript. W. X. W., M. B., P. W., and Z. Z. C collected and prepared data for analysis. Y.
495 J. Y. created figures and tables. Z. H. M. directed the overall research. All authors read and
496 approved the final manuscript.

497

498 ***Acknowledgements***

499 Not applicable

500

501 **References**

- 502 1. Hanahan D, Weinberg RA. Hallmarks of cancer: the next generation. *Cell*.
503 2011;144:646-74.
- 504 2. Alexandrov LB, Nik-Zainal S, Wedge DC, Campbell PJ, Stratton MR. Deciphering
505 signatures of mutational processes operative in human cancer. *Cell Rep*. 2013;3:246-
506 59.
- 507 3. Cooke MS, Evans MD, Dizdaroglu M, Lunec J. Oxidative DNA damage:
508 mechanisms, mutation, and disease. *FASEB J*. 2003;17:1195-214.
- 509 4. Pfeifer GP. Environmental exposures and mutational patterns of cancer genomes.
510 *Genome Med*. 2010;2:54.
- 511 5. Stratton MR, Campbell PJ, Futreal PA. The cancer genome. *Nature*. 2009;458:719-24.

- 512 6. Peña-Diaz J, Bregenhorn S, Ghodgaonkar M, Follonier C, Artola-Borán M, Castor D,
513 et al. Noncanonical mismatch repair as a source of genomic instability in human cells.
514 Mol Cell. 2012;47:669-80.
- 515 7. Cappell MS. Pathophysiology, clinical presentation, and management of colon cancer.
516 Gastroenterol Clin North Am. 2008;37:1-24, v.
- 517 8. Alexandrov LB, Kim J, Haradhvala NJ, Huang MN, Tian Ng AW, Wu Y, et al. The
518 repertoire of mutational signatures in human cancer. Nature. 2020;578:94-101.
- 519 9. Nik-Zainal S, Alexandrov LB, Wedge DC, Van Loo P, Greenman CD, Raine K, et al.
520 Mutational processes molding the genomes of 21 breast cancers. Cell. 2012;149:979-
521 93.
- 522 10. Poon SL, Pang ST, McPherson JR, Yu W, Huang KK, Guan P, et al. Genome-wide
523 mutational signatures of aristolochic acid and its application as a screening tool. Sci
524 Transl Med. 2013;5:197ra01.
- 525 11. Alexandrov LB, Jones PH, Wedge DC, Sale JE, Campbell PJ, Nik-Zainal S, et al.
526 Clock-like mutational processes in human somatic cells. Nat Genet. 2015;47:1402-7.
- 527 12. Nik-Zainal S, Davies H, Staaf J, Ramakrishna M, Glodzik D, Zou X, et al. Landscape
528 of somatic mutations in 560 breast cancer whole-genome sequences. Nature.
529 2016;534:47-54.
- 530 13. Petljak M, Alexandrov LB. Understanding mutagenesis through delineation of
531 mutational signatures in human cancer. Carcinogenesis. 2016;37:531-40.
- 532 14. Mimaki S, Totsuka Y, Suzuki Y, Nakai C, Goto M, Kojima M, et al. Hypermutation
533 and unique mutational signatures of occupational cholangiocarcinoma in printing
534 workers exposed to haloalkanes. Carcinogenesis. 2016;37:817-26.

- 535 15. Polak P, Kim J, Braunstein LZ, Karlic R, Haradhavala NJ, Tiao G, et al. A mutational
536 signature reveals alterations underlying deficient homologous recombination repair in
537 breast cancer. *Nat Genet.* 2017;49:1476-86.
- 538 16. Sawrycki P, Domagalski K, Cechowska M, Gąsior M, Jarkiewicz-Tretyn J, Tretyn A.
539 Relationship between CYP1B1 polymorphisms (c.142C > G, c.355G > T, c.1294C >
540 G) and lung cancer risk in Polish smokers. *Future Oncol.* 2018;14:1569-77.
- 541 17. Zerp SF, van Elsas A, Peltenburg LT, Schrier PI. p53 mutations in human cutaneous
542 melanoma correlate with sun exposure but are not always involved in
543 melanomagenesis. *Br J Cancer.* 1999;79:921-6.
- 544 18. Hochreiter S, Schmidhuber J. Long short-term memory. *Neural Comput.*
545 1997;9:1735-80.
- 546 19. Gers FA, Schmidhuber J, Cummins F. Learning to forget: continual prediction with
547 LSTM. *Neural Comput.* 2000;12:2451-71.
- 548 20. Tayara H, Chong KT. Improving the quantification of DNA sequences using
549 evolutionary information based on deep learning. *Cells.* 2019;8:1635.
- 550 21. Liu Q, Fang L, Yu G, Wang D, Xiao CL, Wang K. Detection of DNA base
551 modifications by deep recurrent neural network on Oxford Nanopore sequencing data.
552 *Nat Commun.* 2019;10:2449.
- 553 22. Zhou J, Lu Q, Xu R, Gui L, Wang H. EL_LSTM: prediction of DNA-binding residue
554 from protein sequence by combining long short-term memory and ensemble learning.
555 *IEEE/ACM Trans Comput Biol Bioinform.* 2018;17:124-35.
- 556 23. Markey MK, Lo JY, Tourassi GD, Floyd CE, Jr. Self-organizing map for cluster
557 analysis of a breast cancer database. *Artif Intell Med.* 2003;27:113-27.
- 558 24. Furukawa T. SOM of SOMs. *Neural Netw.* 2009;22:463-78.

- 559 25. Kolasa M, Długosz R, Pedrycz W, Szulc M. A programmable triangular
560 neighborhood function for a Kohonen self-organizing map implemented on chip.
561 Neural Netw. 2012;25:146-60.
- 562 26. Fukuoka Y, Zhou M, Vittinghoff E, Haskell W, Goldberg K, Aswani A. Objectively
563 measured baseline physical activity patterns in women in the mPED trial: cluster
564 analysis. JMIR Public Health Surveill. 2018;4:e10.
- 565 27. Edge SB, Compton CC. The American Joint Committee on Cancer: the 7th edition of
566 the AJCC cancer staging manual and the future of TNM. Ann Surg Oncol.
567 2010;17:1471-4.
- 568 28. Sahin S, Kozat S. Nonuniformly sampled data processing using LSTM networks.
569 IEEE Trans Neural Netw Learn Syst. 2018;30:1452-61.
- 570 29. Meier B, Cooke SL, Weiss J, Bailly AP, Alexandrov LB, Marshall J, et al. C. elegans
571 whole-genome sequencing reveals mutational signatures related to carcinogens and
572 DNA repair deficiency. Genome Res. 2014;24:1624-36.
- 573 30. Huang MN, Yu W, Teoh WW, Ardin M, Jusakul A, Ng AWT, et al. Genome-scale
574 mutational signatures of aflatoxin in cells, mice, and human tumors. Genome Res.
575 2017;27:1475-86.
- 576 31. Nik-Zainal S, Kucab JE, Morganella S, Glodzik D, Alexandrov LB, Arlt VM, et al.
577 The genome as a record of environmental exposure. Mutagenesis. 2015;30:763-70.
- 578 32. Kucab JE, Zou X, Morganella S, Joel M, Nanda AS, Nagy E, et al. A compendium of
579 mutational signatures of environmental agents. Cell. 2019;177:821-36.e16.
- 580 33. McCloskey K, Taly A, Monti F, Brenner MP, Colwell LJ. Using attribution to decode
581 binding mechanism in neural network models for chemistry. Proc Natl Acad Sci U S
582 A. 2019;116:11624-9.

583

584 **Figure captions**

585 **Fig. 1. Training process of the LSTM-SOM model.** Details of the LSTM-SOM model are
586 described in the methods. Two classifications were used for each training period. Ten classes
587 of mutant sequences were obtained after 3 rounds and an extra round of training. Three of the
588 eight dimensions in LSTM output vectors are shown in the space rectangular coordinate
589 system.

590 **Fig. 2. Mutation type and composition of flanking bases in different MBs.** Each bar
591 except for “Reference Allele” and “Mutation Allele” represents one flanking genetic locus.
592 Bars on the left of “Reference Allele” represent bases on the 5’ end of the mutation site, and
593 bars on the right of “Mutation Allele” represent bases on the 3’ end of the mutation site.

594 **Fig. 3. Quantity and proportion of MBs in different cancers.** The left subgraph shows the
595 proportion of different MBs in all SBS mutation data points from different kinds of cancers.
596 The right subgraph shows the quantity and proportion of different MBs in patients.
597 Differences in quantity are reflected in the size of the point, and differences in proportion are
598 reflected in the color of the point.

599 **Fig. 4. Relationship between patient survival and MB in genes with high mutation**
600 **frequencies.** The top 4 most frequently mutated genes are shown (other genes with high
601 mutation frequencies are shown in Additional file 1: Figure S5). For each gene, the left
602 subgraph shows the proportion of MB in all mutation data points from different cancers; the
603 middle subgraph shows the p value of the log-rank test between groups in the whole
604 population; and the right subgraph shows the p value of the log-rank test between groups of
605 patients with different cancers with high incidence. Only p values less than 0.05 are shown in
606 the heat map.

607 **Fig. 5. MBs in patients with different clinical features.** *: $p < 0.05$ in the t test or ANOVA
608 between groups; **: $p < 0.005$ in the t test or ANOVA between groups. The proportion is
609 shown as the mean \pm standard deviation, and error bars represent standard deviation.

610 **Fig. 6. Differences in survival and clinical features between patients clustered according**
611 **to MB composition.** a: Characteristics of MB composition in patients of 7 classes clustered
612 by the K-means method; each line represents one patient. b: Survivorship curve of each class
613 of patients. c: Log-rank test between classes; differences in the p value are reflected in color.
614 d, e: Clinical features of patients in different classes (*: $p < 0.05$ ANOVA or the chi-square
615 test; **: $p < 0.005$ ANOVA or the chi-square test; error bars represent standard deviation).

Figures

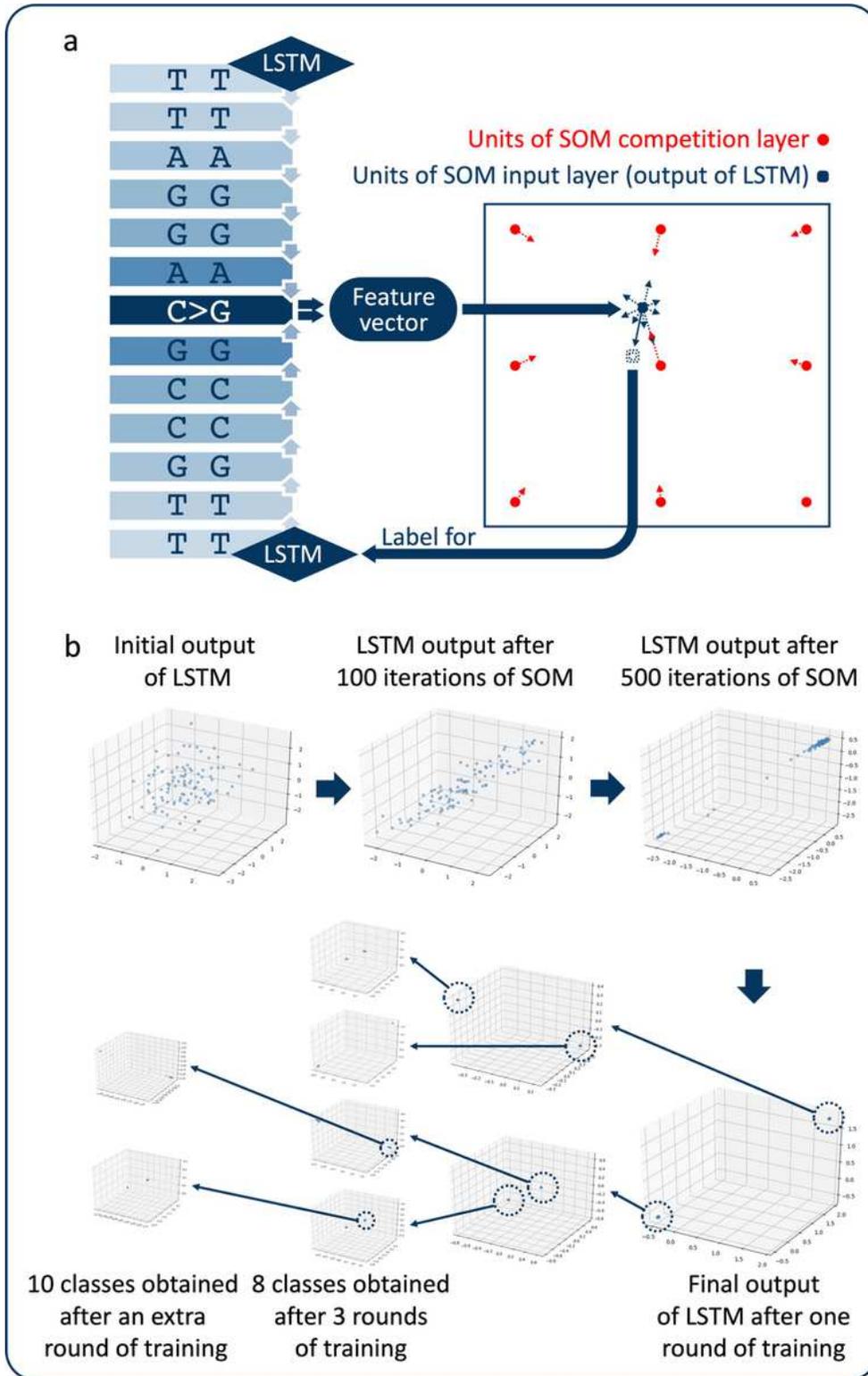


Figure 1

Training process of the LSTM-SOM model. Details of the LSTM-SOM model are described in the methods. Two classifications were used for each training period. Ten classes of mutant sequences were

obtained after 3 rounds and an extra round of training. Three of the eight dimensions in LSTM output vectors are shown in the space rectangular coordinate system.

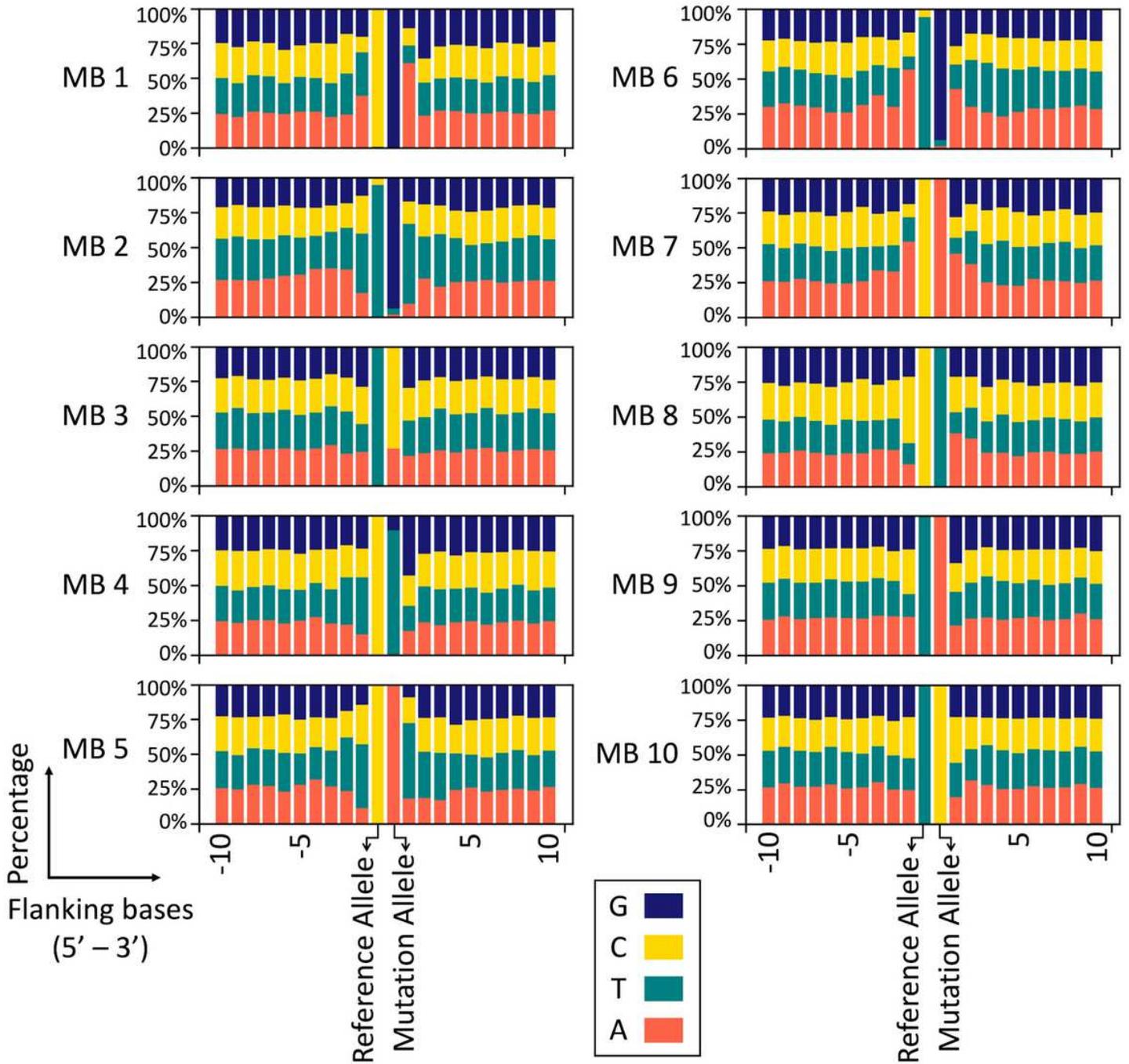


Figure 2

Mutation type and composition of flanking bases in different MBs. Each bar except for “Reference Allele” and “Mutation Allele” represents one flanking genetic locus. Bars on the left of “Reference Allele” represent bases on the 5’ end of the mutation site, and bars on the right of “Mutation Allele” represent bases on the 3’ end of the mutation site.

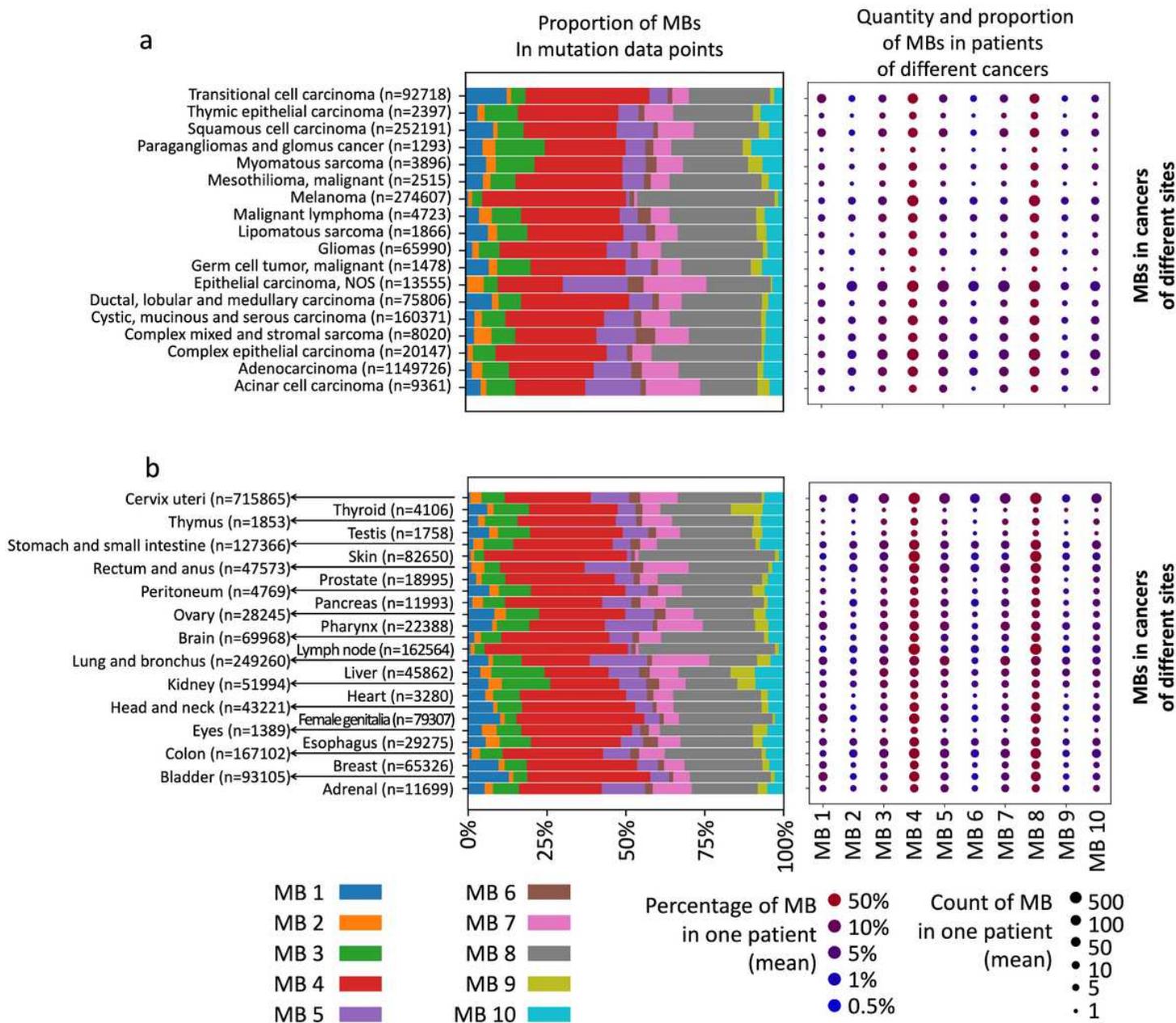


Figure 3

Quantity and proportion of MBs in different cancers. The left subgraph shows the proportion of different MBs in all SBS mutation data points from different kinds of cancers. The right subgraph shows the quantity and proportion of different MBs in patients. Differences in quantity are reflected in the size of the point, and differences in proportion are reflected in the color of the point.

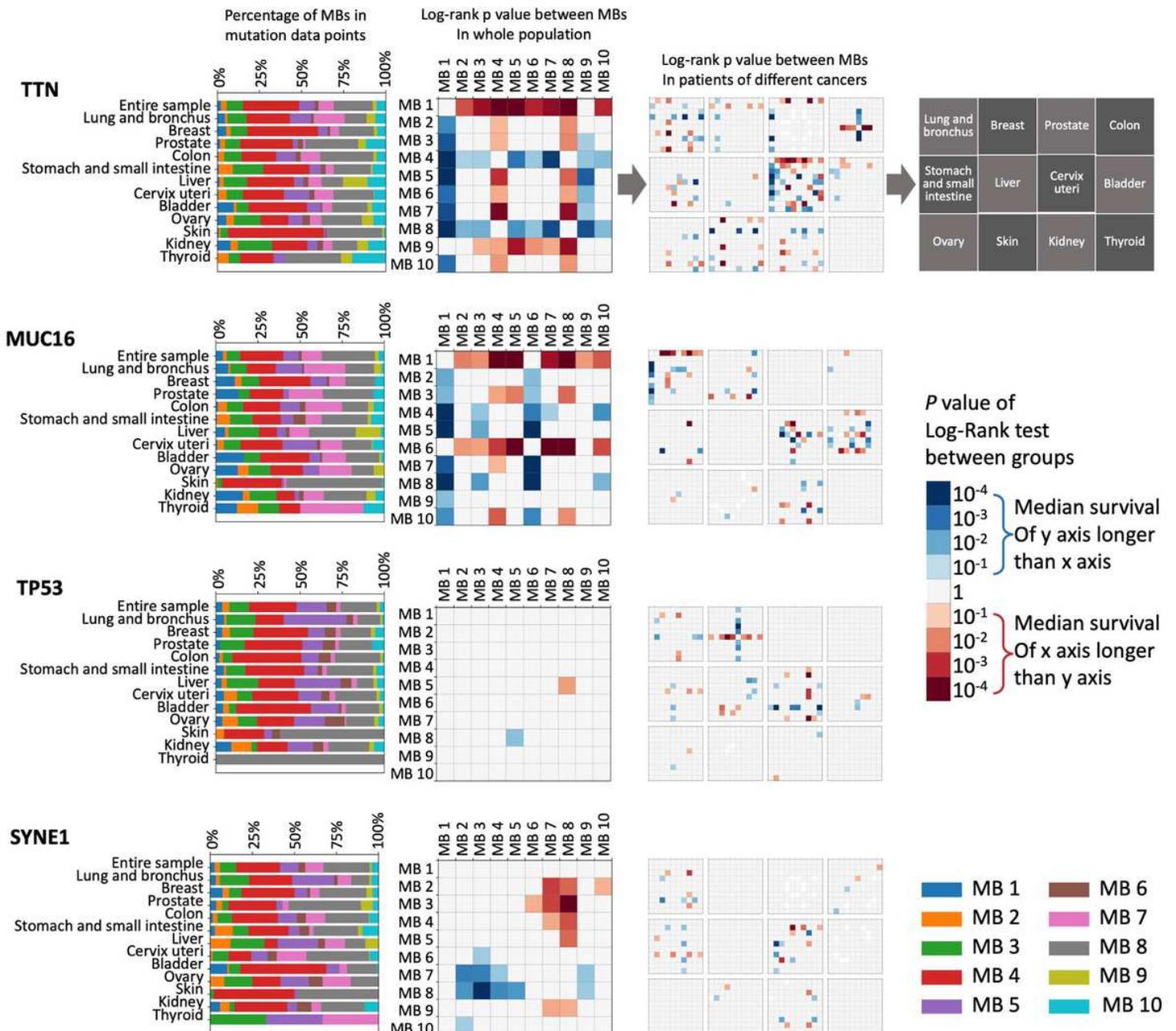


Figure 4

Relationship between patient survival and MB in genes with high mutation frequencies. The top 4 most frequently mutated genes are shown (other genes with high mutation frequencies are shown in Additional file 1: Figure S5). For each gene, the left subgraph shows the proportion of MB in all mutation data points from different cancers; the middle subgraph shows the p value of the log-rank test between groups in the whole population; and the right subgraph shows the p value of the log-rank test between groups of patients with different cancers with high incidence. Only p values less than 0.05 are shown in the heatmap.

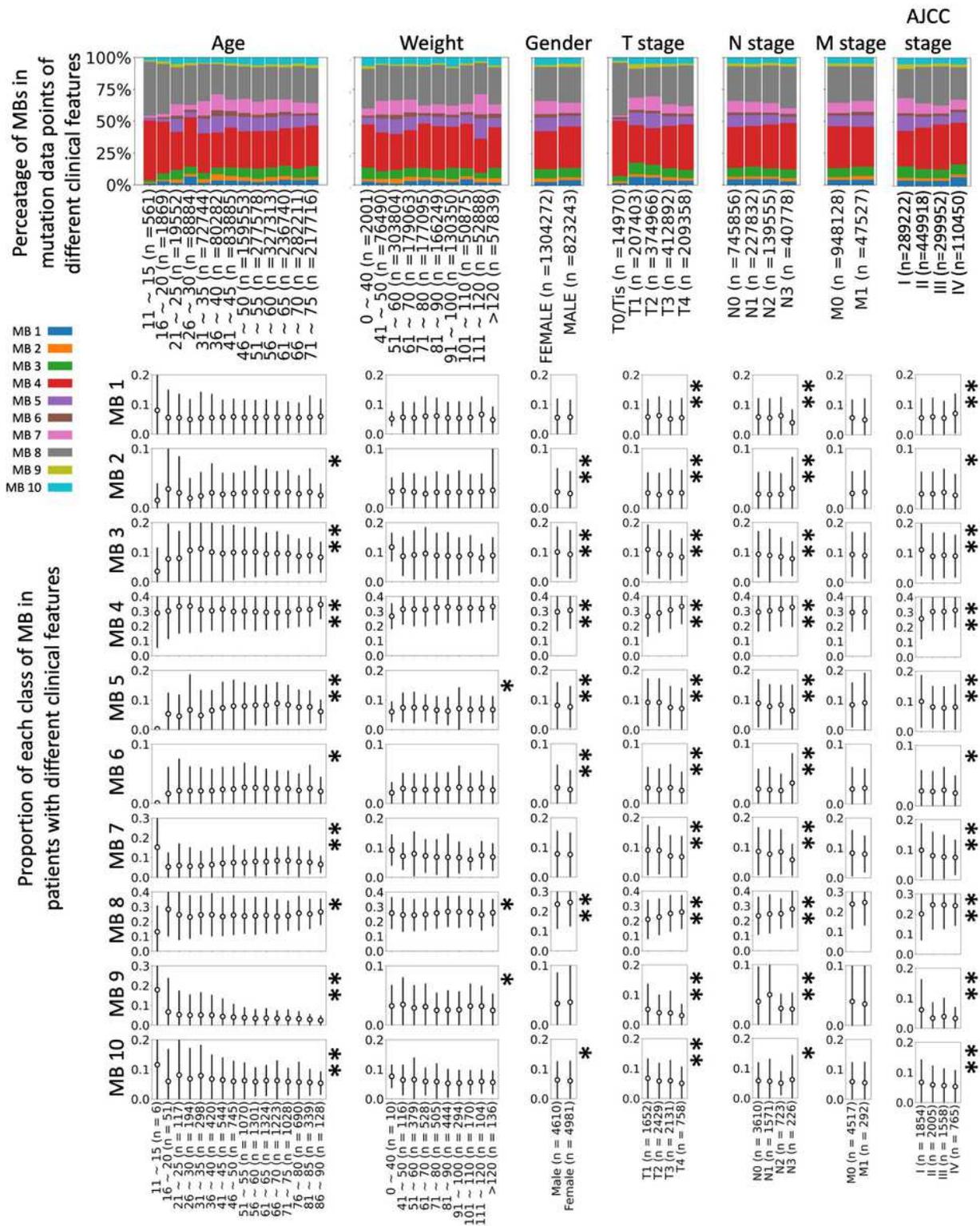


Figure 5

MBs in patients with different clinical features. *: $p < 0.05$ in the t test or ANOVA between groups; **: $p < 0.005$ in the t test or ANOVA between groups. The proportion is shown as the mean \pm standard deviation, and error bars represent standard deviation.

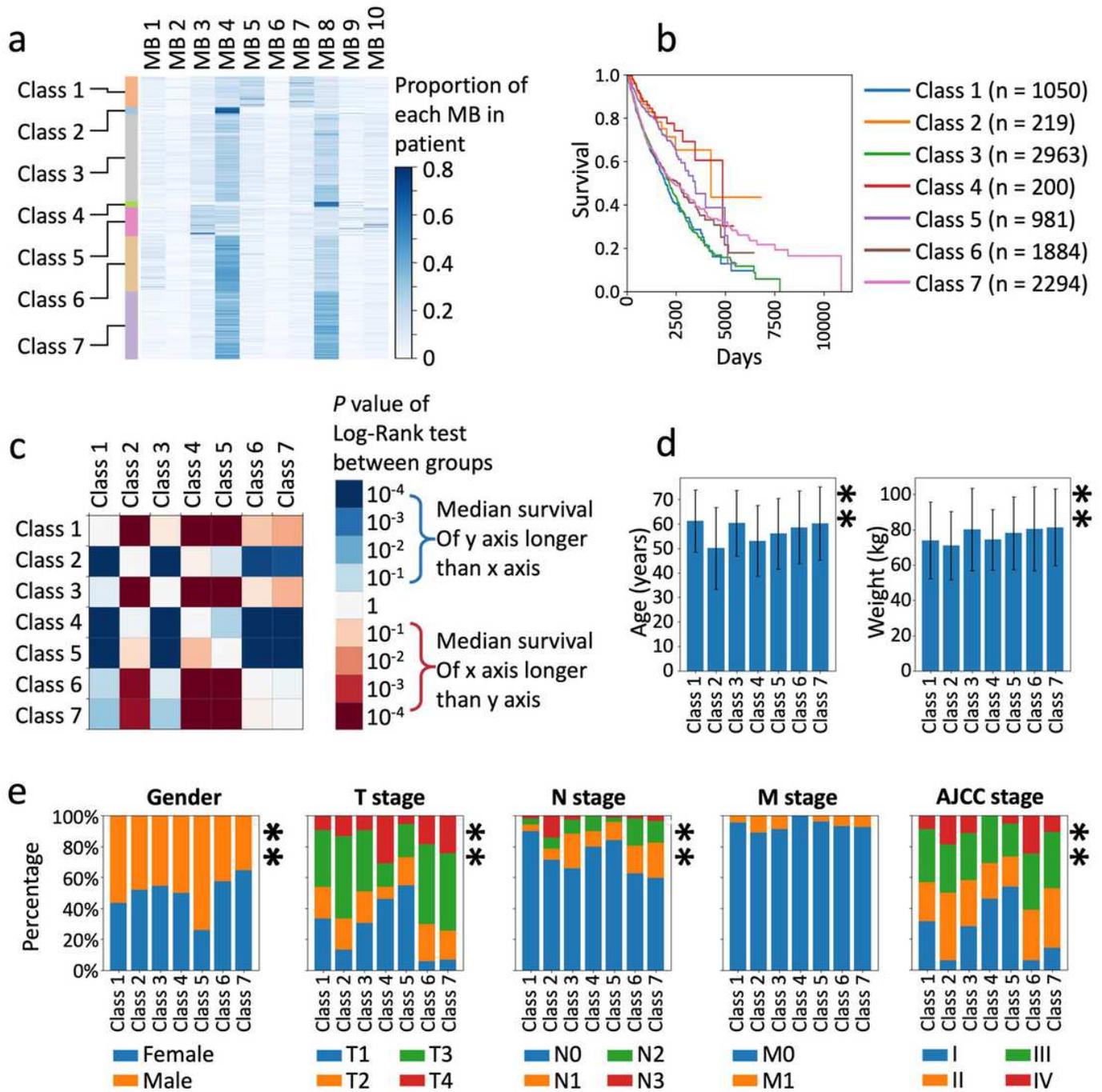


Figure 6

Differences in survival and clinical features between patients clustered according to MB composition. a: Characteristics of MB composition in patients of 7 classes clustered by the K-means method; each line represents one patient. b: Survivorship curve of each class of patients. c: Log-rank test between classes; differences in the p value are reflected in color. d, e: Clinical features of patients in different classes (*: $p < 0.05$ ANOVA or the chi-square test; **: $p < 0.005$ ANOVA or the chi-square test; error bars represent standard deviation).

Supplementary Files

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

- [Additionalfile1.docx](#)
- [Additionalfile2.docx](#)