

# Tumor microenvironment-aware, single-transcriptome prediction of microsatellite instability in colorectal cancer using meta-analysis

Mi-Kyoung Seo

Yonsei University College of Medicine

Hyundeok Kang

Yonsei University College of Medicine

Sangwoo Kim (✉ [swkim@yuhs.ac](mailto:swkim@yuhs.ac))

Yonsei Univ. College of Medicine

---

## Research Article

**Keywords:** Microsatellite instability, Colorectal cancer, Recursive feature elimination-random forest, Tumor microenvironment, Machine learning

**Posted Date:** November 15th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-1019124/v1>

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

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Scientific Reports on April 15th, 2022. See the published version at <https://doi.org/10.1038/s41598-022-10182-3>.

1 **Tumor microenvironment-aware, single-transcriptome prediction of**  
2 **microsatellite instability in colorectal cancer using meta-analysis**

3

4 Mi-Kyoung Seo<sup>1,2</sup>, Hyundeok Kang<sup>1</sup>, Sangwoo Kim<sup>1\*</sup>

5

6 <sup>1</sup> Department of Biomedical Systems Informatics and Brain Korea 21 PLUS Project for  
7 Medical Science, Yonsei University College of Medicine, Seoul 03722, South Korea

8 <sup>2</sup> Department of Nuclear Medicine, Seoul National University Hospital, Seoul 03082, South  
9 Korea

10

11 \*Correspondence

12 Sangwoo Kim

13 Postal address: #613, Yonsei Univ. College of Medicine, 50-1 Yonsei-ro, Seodaemun-gu, Seoul  
14 03722, South Korea

15 E-mail: [swkim@yuhs.ac](mailto:swkim@yuhs.ac)

16 Telephone: +82-2-2228-2589

17

18

19

20 **Abstract**

21 Detecting microsatellite instability (MSI) in colorectal cancers (CRCs) is essential since it is  
22 therapeutic strategy determinant feature, including immunotherapy and chemotherapy. Yet, no  
23 attempt has been made to exploit transcriptomic profile and tumor microenvironment (TME)  
24 of it to unveil MSI status in CRC. Hence, we developed a novel TME-aware, single-  
25 transcriptome predictor of MSI for CRC, called MAP (Microsatellite instability Absolute single  
26 sample Predictor). MAP was developed utilizing recursive feature elimination-random forest  
27 with 466 CRC samples from The Cancer Genome Atlas, and its performance was validated in  
28 independent cohorts, including 1118 samples. MAP showed robustness and predictive power  
29 in predicting MSI status in CRC. Additional advantages for MAP were demonstrated through  
30 comparative analysis with existing MSI classifier and other cancer types. Our novel approach  
31 will provide access to untouched vast amounts of publicly available transcriptomic data and  
32 widen the door for MSI CRC research and be useful for gaining insights to help with  
33 translational medicine.

34

35 **Keywords:** Microsatellite instability; Colorectal cancer; Recursive feature elimination-random  
36 forest; Tumor microenvironment, Machine learning

37

38

## 39 **Introduction**

40 Microsatellite instability (MSI) is characterized by genetic hypermutability due to impaired  
41 DNA mismatch repair (MMR) system<sup>1</sup>. MSI is observed in many solid tumors, including  
42 gastric, and endometrium cancers, as well as in colorectal cancer (CRC, approximately 15%)<sup>2,3</sup>.  
43 Exhibiting prognostic and predictive features of a high tumor mutational burden, a high  
44 neoantigen load, and an immune-active tumor microenvironment (TME) characterized by high  
45 levels of tumor-infiltrating lymphocytes and overexpression of immune checkpoint molecules,  
46 cancers with MSI are known to be great candidates for immune checkpoint inhibitors (ICIs)  
47 treatment, such as pembrolizumab and atezolizumab (anti-PD-1 and anti-PD-L1 monoclonal  
48 antibody, respectively)<sup>4,5</sup>. In addition, MSI primary tumors face a better prognosis than patients  
49 with microsatellite stability (MSS) tumors<sup>2</sup>. MSI status is meaningful as a predictive indicator  
50 for cancer treatment and as a prognostic determinant, identifying a patient's MSI status is  
51 essential in clinical setting and research areas.

52 With recent escalation of its importance in CRC, it has been explored from publicly obtained  
53 samples, such as The Cancer Genome Atlas (TCGA) and Gene Expression Omnibus (GEO)  
54 database, resulting in numerous studies which broaden our understanding in MSI and expand  
55 therapeutic options for MSI CRC patients<sup>6</sup>. However, prior to utilization, MSI status  
56 information must be provided beforehand by quantifying the extent of genomic events in  
57 microsatellite loci, at genomic level, using the Bethesda Panel, a PCR-based five marker panel,  
58 or by examining the loss of mismatch repair proteins using immunohistochemistry (IHC) at the  
59 protein level<sup>1</sup>. Additionally, with recent advances in next-generation sequencing (NGS)  
60 technology, MSI predictors, such as MANTIS<sup>7,8</sup>, MSIsensor<sup>9</sup>, and MOSAIC<sup>10</sup>, have been  
61 developed to extract MSI status from whole exome and whole genome sequencing data.  
62 However, assigning the MSI status from expression data had not been possible until a k-Nearest  
63 Neighbors (k-NN, k=5) classifier called preMSIm using 15 gene-set for pan-cancer was

64 recently proposed<sup>11</sup> and several other attempts which had been made, although the software  
65 has not been made readily available<sup>12,13</sup>.

66 The preMSIm has constructed as a pan-cancer MSI predictor by utilizing three MSI dominant  
67 carcinomas as training data<sup>11</sup>, but it did not reflect the distinct expression profiles of its cancer  
68 of origin. Furthermore, individual MSI tumors have unique tumor microenvironment (TME)  
69 and molecular pathway characteristics. For example, immune inflamed MSI microenvironment  
70 could be characterized by higher infiltration of anti-tumorigenic immune cells, such as adaptive  
71 immune cells (T and B lymphocytes) and innate immune cells (dendritic cells, macrophages,  
72 and natural killer cells) than immune desert MSS tumors, and, in CRC, when mutations or  
73 activation of MYC and RAS pathways occur, chemokine *CCL9* is expressed and an  
74 immunosuppressive environment is established, which prevents enrichment of cytotoxic NK  
75 cells and T cells around the tumor<sup>14</sup>. Therefore, transcriptome based MSI predictor which  
76 integrates both TME and molecular pathway characteristics in CRC is needed.

77 Here, in this study, we have developed an enhanced single-sample MSI classifier called MAP  
78 (Microsatellite instability Absolute single sample Predictor) that integrates transcriptomic  
79 characteristics of TME and molecular pathways to predict MSI in CRC. Our TME and  
80 molecular pathways aware predictor will open a way to utilize CRC expression data to  
81 elucidate MSI CRC. Hence, massive amounts of publicly available expression data without  
82 MSI status will be utilized to drive valuable MSI CRC research through our novel approach,  
83 and, furthermore, to give patient benefit in clinical setting.

84

85

## 86 **Results**

### 87 **Overview of MAP development**

88 As an MSI single sample predictor (SSP), the MAP model was developed with the following

89 four components (Fig. S1): 1) identification of the MAP signature (MAPgene model); 2)  
90 modeling based on pairwise gene expression of the MAP signature genes (MAPpairs model);  
91 3) modeling based on ssGSEA scores of cancer-, molecular-, TME-, and immune-related  
92 signatures (MAPsig model); and 4) post-refinement of the final model and prediction of MSI  
93 status. To develop an SSP of MSI status without relying on a relative approach (e.g., comparing  
94 a patient's data with other samples) and with limited platform bias, we constructed a recursive  
95 feature elimination-random forest (RFE-RF) model (MAPpairs model) with pairwise gene  
96 comparisons, leveraging an informative gene-set (MAP signature from the MAPgene model),  
97 rather than gene expression profiles, on a training dataset. In brief, RFE trains the model, ranks  
98 the features, and selects features through the process of repeatedly removing lower-ranked  
99 features<sup>15</sup>. The method of building a model by selecting features with the RFE method based  
100 on the RF algorithm is called RFE-RF<sup>15</sup>. We built another RFE-RF model (MAPsig model)  
101 based on ssGSEA scores for 101 signatures to reflect the degree of activity of cancer-, immune-  
102 and TME-related signatures of the samples. To select the best RFE-RF model from the parts  
103 mentioned above, we evaluated the area under the receiver operating characteristic curve (AUC)  
104 and confirmed the model performance in validation datasets (Table S1). Finally, at the post-  
105 refinement stage, an integrated MAPpairs and MAPsig model was used to precisely predict  
106 MSI status. We named this final model MAP and evaluated its accuracy, kappa, sensitivity,  
107 specificity, F1, and balanced accuracy in the validation datasets (Table S2).

108

### 109 **MAP signature**

110 To minimize the size of the informative gene-set utilized in the MAPgene model, we, first,  
111 identified differentially expressed genes (DEGs) between MSI and MSS samples using the  
112 Wilcoxon rank-sum test. We assessed 718 DEGs with criteria of  $P < 0.001$  and  $|\log_2 \text{fold}$   
113  $\text{change}| > 1$ , and selected a gene-set of 31 genes by performing RFE-RF modeling with an AUC

114 of 99.2%. We called this gene-set the MAP signature (11 up- and 20 down-regulated DEGs,  
115 Fig. 1a and Table S3). Among genes comprising the MAP signature, the *MLH1* gene, which is  
116 commonly downregulated and/or hypermethylated in sporadic MSI samples, ranked as the top  
117 feature gene, based on both accuracy (the importance of the features that improves  
118 classification accuracy) and Gini index values (the importance of the features that reduces the  
119 impurity of classification) (Fig. 1b and Table S3). We also noted that *LY6G6D*<sup>16</sup> and *EPM2AIP1*  
120 genes<sup>17</sup>, which share a promoter with *MLH1*, were included in the gene-set. Additionally, we  
121 found that a known predictive marker for chemotherapy, thymidylate synthase (*TYMS*)<sup>18</sup>, was  
122 included in the gene-set, and its expression was higher in the MSI samples than the MSS  
123 samples (Fig. 1a and 1b). Other genes belonging to the following pathways were also included  
124 in the MAP signature: the WNT signaling pathway (*RNF43*, *TCF7*, and *NKDI*), Hippo  
125 signaling (*TCF7*, *NKDI*, and *TGFBR2*), and MAPK signaling (*DUSP4* and *TGFBR2*). Three  
126 well-known frameshift mutated genes (*DDX27*, *TGFBR2*, and *RNF43*) in microsatellite loci in  
127 MSI CRC were also included. In terms of MMR, 718 DEGs were initially used when  
128 constructing the MAP gene model, although three MMR genes (*MSH2*, *MSH6*, and *PMS2*) were  
129 not included because their statistical significance or fold change did not meet the inclusion  
130 criteria (Fig. S2).

131 To assess the representativeness of the MAP signature (31-genes) in reflection of MSI status,  
132 the expression patterns of the genes in the gene-set were investigated in a validation dataset,  
133 and expression patterns similar to those observed in the discovery dataset were noted. To  
134 further investigate whether the MAP signature could serve as a surrogate marker of MSI status  
135 and to evaluate its generalizability, we obtained and compared ssGSEA scores for the MAP  
136 signature in MSI and MSS tumors, as well as in MSI tumors of each of the four consensus  
137 molecular subtypes (CMSs). The general comparison between MSI and MSS samples revealed  
138 significant differences in MAP signature scores ( $P = 7.6 \times 10^{-36}$ ), but not among the CMSs (Fig.

139 1c). This suggests that the composite genes of the MAP signature can capture MSI's behavior-  
140 related features and discriminate between MSS and MSI status independent of CMS context.

141

## 142 **MAP model**

143 Although the MAPgene model built based on gene expression showed high performance, to  
144 develop a true SSP with unnormalized data that does not rely on a relative comparison among  
145 multiple samples, we employed a pairwise gene comparison approach for model building: for  
146 example, if the expression of gene A was greater than that of gene B, the sample would be  
147 assigned MSI status. An RF model with 1000 trees of such rules was constructed utilizing the  
148 RFE-RF algorithm with five-fold cross-validation repeated 100 times. Finally, the MAPpairs  
149 model, comprised of 187 rules from 465 ( ${}_{31}C_2$ ) rules at a starting point, was selected (AUC of  
150 99.7%). To assess its performance and reproducibility, we applied the model to internal and  
151 external RNA-seq validation datasets and obtained accuracies of 99.1% and 95.4%,  
152 respectively, indicating it to be robust and highly accurate. In the MAPpairs model, *MLH1*-  
153 related features (*MLH1/HPSE*, *MLH1/FECH*, and *MLH1/GNLY*) were the highest ranked  
154 features (Fig. 2a), and the expression levels of these features separated the two groups well  
155 (Fig. 2b). To investigate the features of MAPpairs further, we calculated the number of MSI  
156 and MSS samples that satisfied each rule in MAPpairs (Fig. S3). Most rules were able to  
157 classify MSI and MSS samples, and as such, they were considered to be reflective of the overall  
158 characteristics of MSI, although not all samples may show similar profiles: For example, MSI  
159 samples are known to have a loss of *MLH1* and an immunity-activated characteristic<sup>2</sup>, as well  
160 as high expression of thymidylate synthase (*TYMS*) (chemotherapy response-associated gene)<sup>18</sup>.  
161 Features of the MAPpairs model, *MLH1/GNLY* and *TYMS/MLH1*, respectively, described these  
162 MSI characteristics well (Fig. 2b), but not in all tumors. This may suggest that the MAPpairs  
163 model, a random forest classifier, captures and reflects the more complexity of MSI CRC, not

164 just a simple single rule.

165 In order to complement the MAPpairs model with the characteristics of immune and TME  
166 profiles, as well as the transcriptomic profile and tumor's characteristics of MSI, we built the  
167 MAPsig model based on molecular, cancer, immune, TME, and MAP signature scores inferred  
168 by single-sample gene set enrichment analysis (ssGSEA). The top signatures used in the final  
169 MAPsig model (44 signatures) included the MAP signature, antitumorogenic immune  
170 lymphocytes (effector memory CD8 T cell, Teff (CD8 T effector), Th2 cells, activated CD4 T  
171 cell), complement, INF- $\gamma$  signatures, Wnt- $\beta$ /catenin signaling, glycolysis, and cell cycle  
172 signaling (Fig. 2c). To find out the degree of activation of 44 signatures, we investigated the  
173 heatmap based on the inferred ssGSEA score. Compared to MSS, antitumorogenic immune  
174 lymphocytes, complement, glycolysis, cell cycle, and INF- $\gamma$  related signatures were up-  
175 regulated in MSI, whereas MAP signature, Notch, angiogenesis, epithelial signature, and Wnt-  
176  $\beta$ /catenin signaling were down-regulated (Fig. S4). The final MAP model was established after  
177 integrating the MAPpairs and the MAPsig models, and post-refinement processing was done  
178 by utilizing probability. Next, we applied the final MAP model to validation datasets to evaluate  
179 any potential overfitting and its applicability across multiple platforms. A total of 1118 samples  
180 (240 MSI and 878 MSS tumors) were tested, and MAP exhibited an average accuracy of 96.1%  
181 (95% confidence interval (CI) 94.3-98.9), a sensitivity of 93.1%, a specificity of 97.5%, and  
182 an F1 score of 92.0% (Fig. 2d and 2e), indicating outstanding performance and feasibility as  
183 an MSI predictor.

184

### 185 **MSI signatures in other cancer types**

186 Using TCGA-STAD and TCGA-UCEC RNA-seq datasets, we evaluated whether MAP, which  
187 was developed for CRC, could be applied to other cancers. In the stomach adenocarcinoma  
188 (STAD) and uterine corpus endometrial carcinoma (UCEC) data, accuracies of 80.2% and 75.4%

189 were observed, respectively. To investigate why the MSI classifier of CRC is not suitable for  
190 other cancers, the same method used to construct the MSI signature (MAP signature) in CRC  
191 was applied to examine MSI signatures in gastric cancer and uterine cancer, and then the  
192 differences in expression patterns were investigated. Uterine cancer showed an accuracy of  
193 90.9%, with only nine genes (*CXCL13*, *EPM2AIP1*, *H2AFJ*, *HOXA9*, *MLH1*, *RNLS*, *SDR42E1*,  
194 *TNFSF9*, and *ZNF300*), whereas gastric cancer reached an accuracy of 83.4% using 78 genes  
195 (Table S4). We further probed how cancer-specific MSI signatures are expressed in each cancer  
196 and observed that individual MSI signatures tend to correspond to DEGs not statistically  
197 significant in other cancers (Fig. 3a and 3b). *MLH1* and *EPM2AIP1* were differentially  
198 expressed in all three cancers, *RPL22L1* was included in the MAP signature of CRC and STAD,  
199 and *H2AFJ* was observed in both CRC and UCEC. In addition, comparing the MAP signature  
200 and MSI signature from the recently developed preMSIm, five genes (*MLH1*, *RPL22L1*,  
201 *EPM2AIP1*, *DDX27*, and *SHROOM4*) were observed in both signatures in CRC (Fig. 3c). It is  
202 also worth mentioning that all of the genes used in preMSIm are down-regulated in MSI, except  
203 *RPL22L1*, whereas the MAP signature additionally includes both up- and down-regulated  
204 genes in CRC. In addition, the expression pattern of the signature of preMSIm did not appear  
205 to suitably reflect genes important in gastric cancer, such as *DDX27*, *SMAP1*, and *ZSWIM3*,  
206 and in uterine cancer, such as *DDX27*, *SHROOM4*, *SMAP1*, and *ZSWIM3*, thereby making it  
207 unable to efficiently differentiate MSI and MSS tumors ( $|\log_2 \text{fold change}| < 0.5$ ) in these cancer  
208 types (Fig. 3d and 3e).

209

## 210 **Discussion**

211 Not all MSI status information is available in publicly available colorectal cancer expression  
212 data, such as NCBI Gene Expression Omnibus (NCBI GEO), thus such data cannot be utilized  
213 in MSI CRC research. For example, it hampers studies determining why most MSS samples

214 belong to the immune desert type or the mechanism by which immune evasion occurred in a  
215 subset of MSI tumors by utilizing molecular or immunological characterization of MSI and  
216 MSS. Although at the research level, if these studies are conducted, this may give clues to  
217 convert the immune-inactivated tumors into immune-activated types or to discover drugs  
218 targeting abnormally activated oncogenic pathways or suppressed TMEs which can be  
219 combined with ICIs. Additionally, since MSI samples are rare, with the difficulty of producing  
220 expression data due to RNA degradation, meta-analysis using multiple cohorts is required, but  
221 the use is hindered due to the absence of MSI status information in them. Furthermore, MSI  
222 research analysis can be performed after cross-validation of the MSI status of RNA-seq data of  
223 the tumor identified as the MSI sample at the DNA level. Here, we present MAP, a tumor  
224 microenvironment-aware, single-transcriptome predictor of MSI in CRC, with robust accuracy  
225 validated using large samples from multiple cohorts of primary tumors. ( $N = 1118$ ). We expect  
226 that the MAP will open the door to make such datasets of use in future MSI studies. MAP has  
227 the advantage of not requiring a matched normal sample as a control and sufficiently predicts  
228 MSI status with a single-sample transcriptome profile.

229 Attempts to create an absolute predictor for subtype classification of cancer and stratification  
230 of patients by applying relationships or ratios between two genes, not the expression value of  
231 the gene itself, are ongoing<sup>19,20</sup>. MAP is an absolute classifier, not relative, and was developed  
232 to reflect tumor molecular characteristics, immune-related signatures, and tumor-infiltrating  
233 immune cells in TME of CRC. Also, since MAP is an RF classifier, one feature does not  
234 represent all MSI in common, but the MSI status is determined through the complex reflection  
235 of various features. Therefore, it may be difficult to interpret clinical and biological  
236 significance of features, and it might be considered to be included technical as well as  
237 biological rules to improve the accuracy of classification.

238 During the development, it showed an accuracy of 99.1% (1/115) in the correct identification

239 of MSI in the internal validation TCGA dataset. Only one sample (TCGA-DC-6154) with MSI  
240 status was incorrectly predicted as being MSS by the MAP model, and it was also marked as  
241 MSS with the MOSAIC program, a tool which predicts MSI status at genomic level<sup>10</sup>. We  
242 speculated such discrepancy may stem from the different tissue sampling locations (MSI typing  
243 vs. DNA and RNA sequencing) or MSI intratumor heterogeneity, rather than MAP  
244 misinforming. We also encountered misclassification of a 11CO070 (MSS) hypermutated  
245 sample from an external RNA-seq validation dataset and five MSS samples from the  
246 GSE39582 dataset as MSI by MAP. Using the clinical information available, we further  
247 investigated the five MSS samples from the GSE39582 dataset and they all carried BRAF  
248 mutation and high CpG island methylator phenotype (CIMP). In sporadic MSI CRC, the  
249 accompanying characteristics of BRAF mutation and high CIMP are known to be strongly  
250 correlated with MSI<sup>21</sup>, but it was not possible to determine misassignment or tumor  
251 heterogeneity characteristics in detail due to the absence of lynch syndrome status or mutation  
252 information of other MMR genes of the samples. Additionally, in research on CMS reported  
253 by the Colorectal Cancer Subtyping Consortium, the distribution of CMS2 (known as immune-  
254 desert type) samples with MSI status was exceedingly rare (10 of 270, 2.7%)<sup>2</sup>, whereas eight  
255 out of the 10 CMS2-MSI samples belonged to one cohort (GSE13294 dataset). This particular  
256 cohort carried a slightly dissimilar CMS2-MSI population distribution from the other datasets,  
257 and out of these eight samples, five were classified as MSS by MAP.

258 MAP showed accuracies of 98.6% (95% CI 97.6-99.6) in RNA-seq and 95.1% (95% CI 91.6-  
259 98.7) in microarray data, all primary tumor and MSI detected based on PCR panel, showing a  
260 slight difference depending on the platform. Although MAP is an absolute SSP with a  
261 specificity of approximately 97% and a high accuracy of 96.1%, it may be due to the inherent  
262 characteristics derived from development based on RNA-seq, or a rare MSI subgroup (eg.  
263 immune-desert CMS2-MSI) that exists in a specific cohort (GSE13294). Due to the paucity of

264 clinical information, we were unable to thoroughly characterize the samples that were not  
265 accurately predicted.

266 The recently developed preMSIm, a pan-cancer MSI predictor, is a k-NN classifier using 15  
267 genes identified by using only three frequent cancer types (COAD, STAD, and UCEC) as  
268 training data. However, due to the limitations mentioned by the author of preMSIm<sup>11</sup> and based  
269 on our findings, these 15 genes are not enough to predict MSI in pan-cancer. This is because  
270 tumor biology and tumor microenvironment are distinct for individual cancer origins,  
271 suggesting diverse tumor-intrinsic gene expression patterns. In this context, MAP is superior  
272 when predicting the MSI status in CRC as it was designed to reflect both the molecular  
273 characteristics of CRC and the complexity of its surroundings.

274 As the MAP model includes the *MLH1* gene, sporadic CRC, characterized by *MLH1* promoter  
275 hypermethylation or *MLH1* loss, can be classified well, whereas Lynch syndrome, a familial  
276 syndrome, due to germline mutations of MMR or *EPCAM* gene<sup>1</sup>, may not be reflected. In  
277 addition, due to the lack of IHC and clinical information (e.g., *KRAS*, *BRAF* mutations, and  
278 Lynch syndrome status) in the validation datasets, the characteristics of samples with  
279 incorrectly predicted MSI status (e.g., *MSH2/MSH6*-negative CRC) could not be thoroughly  
280 assessed. Although MAP reflects the characteristics of sporadic MSI CRC well, *MSH2/MSH6*-  
281 negative MSI CRC reflection is somewhat limited because the expression patterns of *MSH2*  
282 and *MSH6* among MMR genes are not distinctly distinguished from MSS and MSI in TCGA  
283 and external validation dataset.

284 In conclusion, we provided MAP, an MSI predictor for CRC that is robust and accurate.  
285 Although MSI prediction based on IHC and PCR is well established and available at a low cost  
286 for clinical application, MAP will find use in MSI-related research seeking to employ the large  
287 amounts of publicly available CRC expression data and will be useful for gaining insights to  
288 help with translational medicine.

289

290

## 291 **Methods**

### 292 **Dataset acquisition**

293 This meta-analysis was performed in accordance with the PRISMA guidelines (Fig. S5). For  
294 the discovery cohort, 581 RNA-seq data (rsem.norm.expression) from TCGA-COADREAD  
295 were downloaded from the TCGA data portal (<https://portal.gdc.cancer.gov/>). Matching data  
296 on MSI status (82 MSI and 499 MSS) was downloaded from The Cancer Imaging Archive  
297 (TCIA) (<https://tcia.at/>). For the validation cohort, 106 RNA-seq data from 24 MSI and 82  
298 MSS samples (rsem.norm.expression) from an independent study<sup>22</sup> were downloaded. MSI-  
299 low tumors were grouped with MSS tumors as in previous studies<sup>23,24</sup>. The gene expression  
300 values of RNA-seq were log-transformed (with base 2) for analysis. Five independent  
301 microarray-based cohorts were used as an additional validation dataset, particularly to test  
302 platform compatibility<sup>22,25-29</sup>. Detailed information on the datasets is available in Table S1.  
303 Information on consensus molecular subtype (CMS) classification was obtained from the  
304 Colorectal Cancer Subtyping Consortium for all array datasets<sup>2</sup>. For cases with missing CMS  
305 information, CMS labels were inferred by using the random forest (RF) method provided by  
306 the CMSclassifier R package<sup>2</sup>. Genes covered in both of the discovery and validation datasets  
307 were used for further analysis.

308

### 309 **Development of the MAP predictor**

#### 310 *Development of a gene-based predictor (MAPgene)*

311 A schematic drawing of the MAP development process is provided in Fig. S1. To select  
312 informative genes for MSI prediction, we first identified differentially expressed genes (DEGs)  
313 between MSI and MSS tumors using the Wilcoxon rank-sum test in the discovery cohort. To

314 construct and train a prediction model, RNA-seq data were divided into training and internal  
315 validation datasets at a ratio of 4 to 1. To extract the most discriminative genes from the DEGs,  
316 the recursive feature elimination-random forest (RFE-RF) algorithm was used on the 466  
317 training dataset. Briefly, feature selection was conducted by the backward selection method,  
318 wherein the RFE-RF repeatedly constructed an RF model by eliminating features with the least  
319 importance. The selection process was repeated 100 times, applying an upsampling approach  
320 to the MSI group (due to the small group size) using caret<sup>30</sup> and randomForest R package. The  
321 final model (MAPgene) was then selected based on that with the best area under a receiver  
322 operating characteristic curve (AUC) for 31 genes.

323

#### 324 *Development of an absolute, gene-pair-based predictor (MAPpairs)*

325 To make the MAPgene model absolute (i.e., to predict MSI status from a single patient without  
326 comparison to a reference cohort or sample-wise normalization), a new model (MAPpairs) was  
327 developed using pairwise gene expression values instead of single gene expression values. A  
328 total of 465 ( ${}_{31}C_2$ ) gene-expression pairs were generated for the selected 31 genes. These gene  
329 expression pairs were then converted to rules that indicated the relative over- or under-  
330 expression between two genes. For example, if the expression of gene A was higher than that  
331 of gene B, the rule (gene A > gene B) was generated. Another RFE-RF model was then  
332 constructed using the 465 rules and trained with a five-fold cross-validation. Similar to the  
333 feature selection procedure, RFE-RF was applied with a five-fold cross-validation and repeated  
334 100 times. The final absolute model was selected according to its AUC.

335

#### 336 *Development of a tumor microenvironment-integrated model (MAPsig)*

337 To construct a more sophisticated model, we exploited the molecular differences in cancer-,  
338 immune-, and TME-related signatures between MSI and MSS tumors. We collected 101

339 signatures, including immune and stromal cells (TCIA and MCP-counter)<sup>31,32</sup>, cancer  
340 hallmarks from MSigDB<sup>33</sup>, immune-related signatures, such as epithelial and mesenchymal  
341 signatures<sup>34</sup>; stromal and immune signatures<sup>35</sup>; immunoinhibitory signatures and  
342 immunostimulatory signatures<sup>31</sup>; T-cell-inflamed gene expression profile (GEP)<sup>36</sup> and IFN- $\gamma$   
343 expanded signatures<sup>36</sup>; cell cycle signature<sup>37</sup>; cell cycle regulator<sup>38</sup>; mismatch repair (KEGG),  
344 C-ECM signature<sup>39</sup>; angiogenesis, HLA class I and II family signature<sup>40</sup>; pro-inflammatory  
345 cytokines and chemokines<sup>40</sup>; CD8 T cells (Teff)<sup>41</sup>; and the MAP signature. To obtain signature  
346 scores for each individual sample, single sample gene-set enrichment analysis (ssGSEA), with  
347  $ssgsea.norm = F$ , was applied for the signatures above. Additionally, for cross-platform  
348 comparability, the acquired score was adjusted to a value between 1 to 10. We used the same  
349 modeling method as that for the MAPgene and MAPpairs models, although with different input  
350 values. Finally, the MAPsig model and features were selected for inclusion in the final  
351 according to those that provided the best AUC.

352

### 353 *Model refinement*

354 When applying the MAPpairs model, we noted that true MSI samples tended to be classified  
355 with MSI at a probability much higher than 70%. Thus, only samples with a probability of  
356 having MSI that was more than 70% were assigned MSI status. Samples with a predicted  
357 probability of MSI that was lower than 70% were further examined by applying the MAPsig  
358 model to determine final MSI status, as it showed high overall AUC, accuracy, and specificity,  
359 but low sensitivity, making it of use in only verifying a true MSS sample. The software is  
360 available at <https://sourceforge.net/p/mapmsi/wiki/MAP/>.

361

### 362 **Validation dataset**

363 To evaluate the predictive performance of the MAP model, we employed RNA-seq data for

364 CRCs ( $N=106$ ) with  $\log_2$ -transformed rsem.norm data. Also, to assess platform independency  
365 and the applicability of MAP on different array datasets, we collected data for five cohorts. In  
366 the microarray datasets, the probes per gene were selected using Jetset  
367 (<http://www.cbs.dtu.dk/biotools/jetset/>)<sup>42</sup>. The array datasets were processed using fRMA R  
368 package per sample<sup>43</sup>. A total of five datasets were evaluated for the following: accuracy,  
369 sensitivity, specificity, F1 score, and balanced accuracy. All information on the datasets is  
370 provided in Table S1. For RNA-seq of stomach adenocarcinoma (STAD), and uterine corpus  
371 endometrial carcinoma (UCEC) were downloaded from the TCGA data portal  
372 (<https://portal.gdc.cancer.gov/>).

373

### 374 **Consistency of genes in a microsatellite instability classifier model based on gene** 375 **expression**

376 To verify the consistency of feature genes with discriminative value in an MSI classifier model  
377 using gene expression, the Wilcoxon rank-sum test was used to analyze the external RNA-seq  
378 validation dataset. In addition, to assess the utility of MAP for MSI prediction, we calculated  
379 MAP signature scores (31-gene-set signature) using ssGSEA and compared them between MSI  
380 and MSS groups, as well as among MSI CMSs, using the Wilcoxon rank-sum test and Kruskal-  
381 Wallis test.

382

### 383 **MSI signature construction at UCEC and STAD**

384 To investigate the MSI signature that can distinguish MSS and MSI in each cancer types, the  
385 same method was applied when constructing the MAP signature, except that the  $P < 0.02$  and  
386  $|\log_2 \text{ fold change}| > 1$  criteria was applied to identify sufficient number of DEG from two types  
387 of cancer. TCGA-UCEC and STAD expression dataset were download TCGA-STAD and  
388 UCEC RNA-seq data were downloaded from

389 EBPlusPlusAdjustPANCAN\_IlluminaHiSeq\_RNASeqV2.geneExp.txt file at  
390 <https://gdc.cancer.gov/about-data/publications/panimmune>. In this file, only 12 of 15  
391 signatures of preMSIm existed. The missing genes were *HENMT1*, *NOL4L*, and *RTF2*.

392

### 393 **Statistical analysis**

394 Comparisons of two groups were conducted using the Wilcoxon rank-sum test, while  
395 comparisons of multiple groups were performed using the Kruskal-Wallis test. All statistical  
396 analyses were conducted using R language software (<https://www.r-project.org/>).

397

398

399 **Acknowledgements**

400 This research was supported by a grant of the Korea Health Technology R&D Project through  
401 the Korea Health Industry Development Institute (KHIDI), funded by the Ministry of Health  
402 & Welfare, Republic of Korea (grant number HI14C1324).

403

404 **Author Information**

405 Affiliations

406 **Department of Biomedical Systems Informatics and Brain Korea 21 PLUS Project for**  
407 **Medical Science, Yonsei University College of Medicine, Seoul 03722, South Korea**

408 Mi-Kyoung Seo, Hyundeok Kang, Sangwoo Kim

409 **Department of Nuclear Medicine, Seoul National University Hospital, Seoul 03082, South**  
410 **Korea**

411 Mi-Kyoung Seo

412

413 **Author Contributions**

414 Conceptualization, M-K.S.; methodology, M-K.S.; software, M-K.S.; validation, M-K.S.;  
415 formal analysis, M-K.S.; investigation, M-K.S.; resources, M-K.S., and S.K.; data curation,  
416 M.S.; writing-original draft preparation, M-K.S., and H.K.; writing-review and editing, M-K.S.,  
417 and H.K.; visualization, M-K.S.; supervision, S.K.; project administration, M-K.S.; funding  
418 acquisition, S.K. All authors have read and agreed to the published version of the manuscript.

419

420 **Corresponding authors**

421 Correspondence to Sangwoo Kim

422

423 **Ethics declarations**

424 **Ethics approval and consent to participate**

425 Not applicable. No permissions were required to use any of the repository data. All methods  
426 were performed in accordance with the PRISMA guidelines.

427

428 **Conflicts of Interest**

429 The authors declare no conflict of interest.

430

431 **Data Availability**

432 All relevant datasets used in the current study are available in the TCGA  
433 (<https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga>) and  
434 GEO (<https://www.ncbi.nlm.nih.gov/geo/>). This study analysis used all publicly available  
435 datasets, and the dataset accession numbers included in Table S1. The software is available at  
436 <https://sourceforge.net/p/mapmsi/wiki/MAP/>.

437 **References**

- 438 1 Evrard, C., Tachon, G., Randrian, V., Karayan-Tapon, L. & Tougeron, D. Microsatellite  
439 Instability: Diagnosis, Heterogeneity, Discordance, and Clinical Impact in Colorectal  
440 Cancer. *Cancers (Basel)* **11**, doi:10.3390/cancers11101567 (2019).
- 441 2 Guinney, J. *et al.* The consensus molecular subtypes of colorectal cancer. *Nat Med* **21**,  
442 1350-1356, doi:10.1038/nm.3967 (2015).
- 443 3 Cortes-Ciriano, I., Lee, S., Park, W. Y., Kim, T. M. & Park, P. J. A molecular portrait of  
444 microsatellite instability across multiple cancers. *Nat Commun* **8**, 15180,  
445 doi:10.1038/ncomms15180 (2017).
- 446 4 Le, D. T. *et al.* PD-1 Blockade in Tumors with Mismatch-Repair Deficiency. *N Engl J Med*  
447 **372**, 2509-2520, doi:10.1056/NEJMoa1500596 (2015).
- 448 5 Sveen, A. *et al.* Colorectal Cancer Consensus Molecular Subtypes Translated to Preclinical  
449 Models Uncover Potentially Targetable Cancer Cell Dependencies. *Clin Cancer Res* **24**,  
450 794-806, doi:10.1158/1078-0432.CCR-17-1234 (2018).
- 451 6 Danaher, P. *et al.* Gene expression markers of Tumor Infiltrating Leukocytes. *J*  
452 *Immunother Cancer* **5**, 18, doi:10.1186/s40425-017-0215-8 (2017).
- 453 7 Kautto, E. A. *et al.* Performance evaluation for rapid detection of pan-cancer microsatellite  
454 instability with MANTIS. *Oncotarget* **8**, 7452-7463, doi:10.18632/oncotarget.13918  
455 (2017).
- 456 8 Bonneville, R. *et al.* Landscape of Microsatellite Instability Across 39 Cancer Types. *JCO*  
457 *Precis Oncol* **2017**, doi:10.1200/PO.17.00073 (2017).
- 458 9 Niu, B. *et al.* MSIsensor: microsatellite instability detection using paired tumor-normal  
459 sequence data. *Bioinformatics* **30**, 1015-1016, doi:10.1093/bioinformatics/btt755 (2014).
- 460 10 Hause, R. J., Pritchard, C. C., Shendure, J. & Salipante, S. J. Classification and

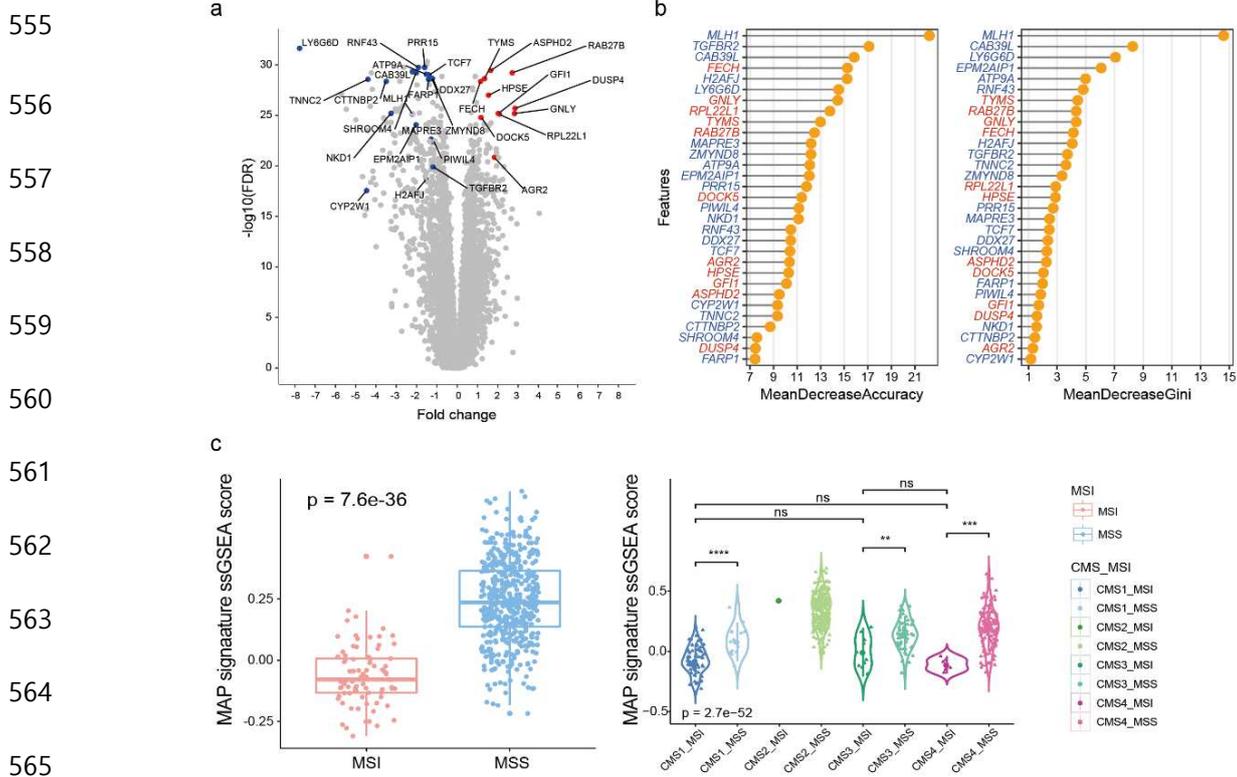
- 461 characterization of microsatellite instability across 18 cancer types. *Nat Med* **22**, 1342-  
462 1350, doi:10.1038/nm.4191 (2016).
- 463 11 Li, L., Feng, Q. & Wang, X. PreMSIm: An R package for predicting microsatellite  
464 instability from the expression profiling of a gene panel in cancer. *Comput Struct*  
465 *Biotechnol J* **18**, 668-675, doi:10.1016/j.csbj.2020.03.007 (2020).
- 466 12 Danaher, P. *et al.* A gene expression assay for simultaneous measurement of microsatellite  
467 instability and anti-tumor immune activity. *J Immunother Cancer* **7**, 15,  
468 doi:10.1186/s40425-018-0472-1 (2019).
- 469 13 Pacinkova, A. & Popovici, V. Cross-platform Data Analysis Reveals a Generic Gene  
470 Expression Signature for Microsatellite Instability in Colorectal Cancer. *Biomed Res Int*  
471 **2019**, 6763596, doi:10.1155/2019/6763596 (2019).
- 472 14 Giraldo, N. A. *et al.* The clinical role of the TME in solid cancer. *Br J Cancer* **120**, 45-53,  
473 doi:10.1038/s41416-018-0327-z (2019).
- 474 15 Darst, B. F., Malecki, K. C. & Engelman, C. D. Using recursive feature elimination in  
475 random forest to account for correlated variables in high dimensional data. *BMC Genet* **19**,  
476 65, doi:10.1186/s12863-018-0633-8 (2018).
- 477 16 Giordano, G. *et al.* JAK/Stat5-mediated subtype-specific lymphocyte antigen 6 complex,  
478 locus G6D (LY6G6D) expression drives mismatch repair proficient colorectal cancer. *J*  
479 *Exp Clin Cancer Res* **38**, 28, doi:10.1186/s13046-018-1019-5 (2019).
- 480 17 Hesson, L. B. *et al.* Lynch syndrome associated with two MLH1 promoter variants and  
481 allelic imbalance of MLH1 expression. *Hum Mutat* **36**, 622-630, doi:10.1002/humu.22785  
482 (2015).
- 483 18 Klingbiel, D. *et al.* Prognosis of stage II and III colon cancer treated with adjuvant 5-  
484 fluorouracil or FOLFIRI in relation to microsatellite status: results of the PETACC-3 trial.  
485 *Ann Oncol* **26**, 126-132, doi:10.1093/annonc/mdu499 (2015).

- 486 19 Seo, M. K., Paik, S. & Kim, S. An Improved, Assay Platform Agnostic, Absolute Single  
487 Sample Breast Cancer Subtype Classifier. *Cancers (Basel)* **12**,  
488 doi:10.3390/cancers12123506 (2020).
- 489 20 Auslander, N. *et al.* Robust prediction of response to immune checkpoint blockade therapy  
490 in metastatic melanoma. *Nat Med* **24**, 1545-1549, doi:10.1038/s41591-018-0157-9 (2018).
- 491 21 Chang, S. C. *et al.* Clinicopathological and Molecular Profiles of Sporadic Microsatellite  
492 Unstable Colorectal Cancer with or without the CpG Island Methylator Phenotype (CIMP).  
493 *Cancers (Basel)* **12**, doi:10.3390/cancers12113487 (2020).
- 494 22 Vasaikar, S. *et al.* Proteogenomic Analysis of Human Colon Cancer Reveals New  
495 Therapeutic Opportunities. *Cell* **177**, 1035-1049 e1019, doi:10.1016/j.cell.2019.03.030  
496 (2019).
- 497 23 Popat, S., Hubner, R. & Houlston, R. S. Systematic review of microsatellite instability and  
498 colorectal cancer prognosis. *J Clin Oncol* **23**, 609-618, doi:10.1200/JCO.2005.01.086  
499 (2005).
- 500 24 Kim, C. G. *et al.* Effects of microsatellite instability on recurrence patterns and outcomes  
501 in colorectal cancers. *Br J Cancer* **115**, 25-33, doi:10.1038/bjc.2016.161 (2016).
- 502 25 Marisa, L. *et al.* Gene expression classification of colon cancer into molecular subtypes:  
503 characterization, validation, and prognostic value. *PLoS Med* **10**, e1001453,  
504 doi:10.1371/journal.pmed.1001453 (2013).
- 505 26 de Sousa, E. M. F. *et al.* Methylation of cancer-stem-cell-associated Wnt target genes  
506 predicts poor prognosis in colorectal cancer patients. *Cell Stem Cell* **9**, 476-485,  
507 doi:10.1016/j.stem.2011.10.008 (2011).
- 508 27 Barras, D. *et al.* BRAF V600E Mutant Colorectal Cancer Subtypes Based on Gene  
509 Expression. *Clin Cancer Res* **23**, 104-115, doi:10.1158/1078-0432.CCR-16-0140 (2017).
- 510 28 Jorissen, R. N. *et al.* DNA copy-number alterations underlie gene expression differences

- 511 between microsatellite stable and unstable colorectal cancers. *Clin Cancer Res* **14**, 8061-  
512 8069, doi:10.1158/1078-0432.CCR-08-1431 (2008).
- 513 29 Cancer Genome Atlas, N. Comprehensive molecular characterization of human colon and  
514 rectal cancer. *Nature* **487**, 330-337, doi:10.1038/nature11252 (2012).
- 515 30 Kuhn, M. Building Predictive Models in R Using the caret Package. *Journal of Statistical*  
516 *Software* **28** (2008).
- 517 31 Charoentong, P. *et al.* Pan-cancer Immunogenomic Analyses Reveal Genotype-  
518 Immunophenotype Relationships and Predictors of Response to Checkpoint Blockade.  
519 *Cell Rep* **18**, 248-262, doi:10.1016/j.celrep.2016.12.019 (2017).
- 520 32 Becht, E. *et al.* Estimating the population abundance of tissue-infiltrating immune and  
521 stromal cell populations using gene expression. *Genome Biol* **17**, 218,  
522 doi:10.1186/s13059-016-1070-5 (2016).
- 523 33 Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set  
524 collection. *Cell Syst* **1**, 417-425, doi:10.1016/j.cels.2015.12.004 (2015).
- 525 34 Linnekamp, J. F. *et al.* Consensus molecular subtypes of colorectal cancer are recapitulated  
526 in in vitro and in vivo models. *Cell Death Differ* **25**, 616-633, doi:10.1038/s41418-017-  
527 0011-5 (2018).
- 528 35 Yoshihara, K. *et al.* Inferring tumour purity and stromal and immune cell admixture from  
529 expression data. *Nat Commun* **4**, 2612, doi:10.1038/ncomms3612 (2013).
- 530 36 Ayers, M. *et al.* IFN-gamma-related mRNA profile predicts clinical response to PD-1  
531 blockade. *J Clin Invest* **127**, 2930-2940, doi:10.1172/JCI91190 (2017).
- 532 37 Davoli, T., Uno, H., Wooten, E. C. & Elledge, S. J. Tumor aneuploidy correlates with  
533 markers of immune evasion and with reduced response to immunotherapy. *Science* **355**,  
534 doi:10.1126/science.aaf8399 (2017).
- 535 38 Cancer Genome Atlas Research, N. Comprehensive molecular characterization of

- 536 urothelial bladder carcinoma. *Nature* **507**, 315-322, doi:10.1038/nature12965 (2014).
- 537 39 Chakravarthy, A., Khan, L., Bensler, N. P., Bose, P. & De Carvalho, D. D. TGF-beta-  
538 associated extracellular matrix genes link cancer-associated fibroblasts to immune evasion  
539 and immunotherapy failure. *Nat Commun* **9**, 4692, doi:10.1038/s41467-018-06654-8  
540 (2018).
- 541 40 Tamborero, D. *et al.* A Pan-cancer Landscape of Interactions between Solid Tumors and  
542 Infiltrating Immune Cell Populations. *Clin Cancer Res* **24**, 3717-3728, doi:10.1158/1078-  
543 0432.CCR-17-3509 (2018).
- 544 41 Rosenberg, J. E. *et al.* Atezolizumab in patients with locally advanced and metastatic  
545 urothelial carcinoma who have progressed following treatment with platinum-based  
546 chemotherapy: a single-arm, multicentre, phase 2 trial. *Lancet* **387**, 1909-1920,  
547 doi:10.1016/S0140-6736(16)00561-4 (2016).
- 548 42 Li, Q., Birkbak, N. J., Gyorffy, B., Szallasi, Z. & Eklund, A. C. Jetset: selecting the optimal  
549 microarray probe set to represent a gene. *BMC Bioinformatics* **12**, 474, doi:10.1186/1471-  
550 2105-12-474 (2011).
- 551 43 McCall, M. N., Bolstad, B. M. & Irizarry, R. A. Frozen robust multiarray analysis (fRMA).  
552 *Biostatistics* **11**, 242-253, doi:10.1093/biostatistics/kxp059 (2010).

554 **Figures**

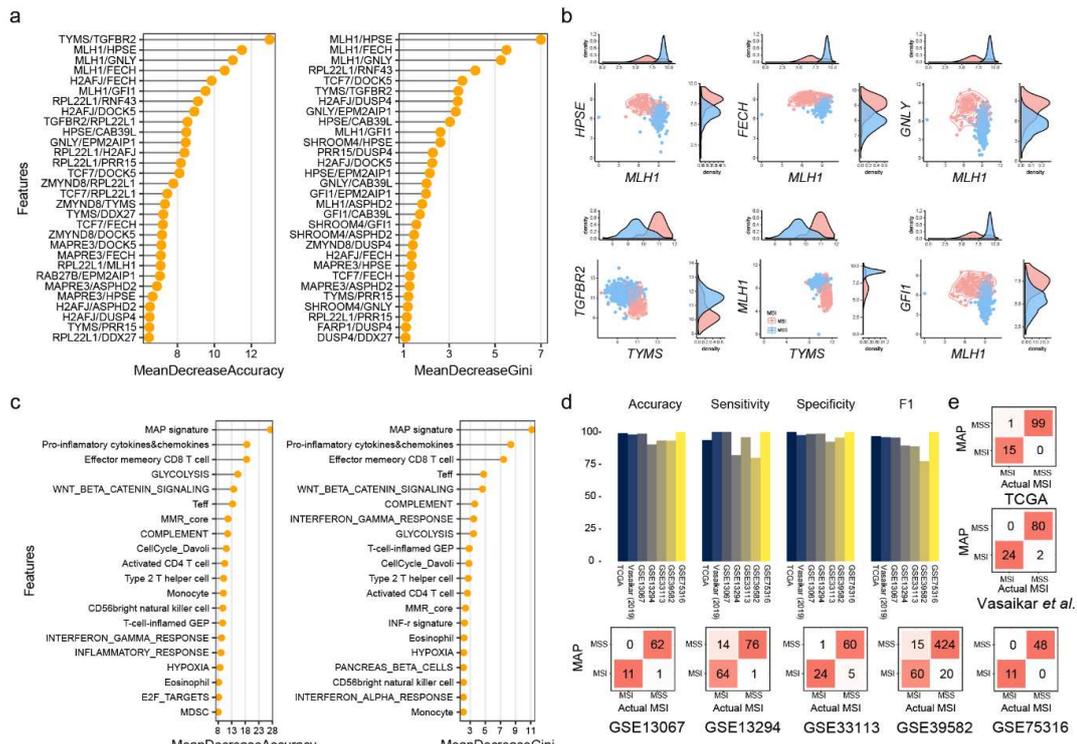


566 **Figure 1.** Designing the MAP signature from RFE-RF analysis of gene expression data. (a) A  
 567 volcano plot for DEGs between MSI and MSS samples. The x axis represents  $\log_2$  fold changes  
 568 in gene expression data for MSI versus MSS samples. Colored dots are significant DEGs in  
 569 MAP signature; red and blue indicate up- and downregulated genes, respectively. (b) The  
 570 importance of 31 features is based on accuracy and Gini index scores. The mean decrease in  
 571 accuracy is a measure of how much influence it has in improving classification accuracy. The  
 572 mean decrease in Gini is a measure of how impurity can be reduced by features used when  
 573 separating nodes. The genes with red and blue colors indicate up- and downregulated genes in  
 574 MSI, compared with MSS, respectively. (c) MAP signature. A box-plot of MAP signature  
 575 ssGSEA scores according to MSI status (left) and CMS-MSI and MSS subtypes (right). The  
 576 dots represent samples. MAP signature scores differ significantly between MSI and MSS  
 577 samples independent of CMS subtypes. CMS2-MSI did not confirm statistical significance

578 because the number of samples was small. \*  $P < 0.05$ , \*\*  $P < 0.01$ , \*\*\*  $P < 0.005$ . DEG ;  
579 differentially expressed gene, CMS ; consensus molecular subtype, ssGSEA ; single-sample  
580 gene set enrichment analysis, FDR ; false discovery rate.

581

582



584 **Figure 2.** MAP model. (a) Top 30 important features of the MAPpairs model. The mean  
 585 decrease in accuracy (left) is a measure of how much influence a feature has in improving  
 586 classification accuracy. The mean decrease in Gini (right) is a measure of how impurity can be  
 587 reduced by features used when separating nodes. (b) A scatter plot and histogram of the gene  
 588 pairs. *MLH1*-related rules and *TFGFR2/TYMS* rule are shown. (c) Top 20 important features  
 589 (signatures) of the MAPsig model. (d) Performance (accuracy, sensitive, specificity, and F1)  
 590 of the MAP model. (e) Confusion matrices of the validation dataset. The actual MSI means  
 591 MSI status provided in the dataset study. The red color-scale reflects percentages of class  
 592 predictions against the actual class.



619 the genes of the corresponding signatures marked in each panel. The blue dotted line on the x-  
620 axis means -1 and 1 of the  $\log_2$  fold change scale, and 2 ( $-\log_{10}(0.01)$ ) on the y-axis.

## Supplementary Files

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

- [MAPsupplementaryv1.pdf](#)