

1 **MSIsensor-RNA: microsatellite instability detection** 2 **for bulk and single-cell gene expression data**

3

4 Peng Jia^{1,2,†}, Xuanhao Yang^{1,2,†}, Xiaofei Yang^{1,3}, Tingjie Wang^{1,3}, Yu Xu⁴, and Kai

5 Ye^{1,2,4,5,*}

6 ¹MOE Key Lab for Intelligent Networks & Networks Security, Faculty of Electronic

7 and Information Engineering, Xi'an Jiaotong University, Xi'an, 710049, China

8 ²School of Automation Science and Engineering, Faculty of Electronic and

9 Information Engineering, Xi'an Jiaotong University, Xi'an, 710049, China

10 ³School of Computer Science, Faculty of Electronic and Information Engineering,

11 Xi'an Jiaotong University, Xi'an 710049, China

12 ⁴School of Life Science and Technology, Xi'an Jiaotong University, Xi'an, 710049,

13 China

14 ⁵Genome Institute, the First Affiliated Hospital of Xi'an Jiaotong University, Xi'an

15 710061, China

16 *Correspondence: kaiye@xjtu.edu.cn (Ye K)

17 [†]These authors contributed equally to this work.

18 **Abstract**

19 **Background:** Microsatellite instability (MSI) is an indispensable biomarker in
 20 cancer immunotherapy. Currently, MSI scoring methods by high-throughput omics
 21 methods have gained popularity and demonstrated better performance than the gold
 22 standard method for MSI detection. However, MSI detection method on expression
 23 data, especially single-cell expression data is still lacking, limiting the scope of
 24 clinical application and prohibiting investigation of MSI at single cell level.

25 **Results:** Herein, we developed MSIsensor-RNA, an accuracy, robust, adaptable, and
 26 standalone software, to detect MSI status by its associated genes' expression values.
 27 We demonstrated the favorable performance and promise of MSIsensor-RNA in both
 28 bulk and single-cell gene expression data in multiplatform technologies including
 29 RNA-seq, Microarray, and single-cell RNA-seq.

30 **Conclusions:** MSIsensor-RNA is a versatile, efficient, and robust method for MSI
 31 status detection from both bulk and single-cell gene expression data in clinical
 32 researches and applications. MSIsensor-RNA is available at
 33 <https://github.com/xjtu-omics/msisensor-rna>.

34 **Keywords:** microsatellite instability, cancer, gene expression, multiplatform,
 35 single-cell RNA-seq, RNA-seq, Microarray

36

37 **Background**

38 Microsatellite instability (MSI) refers to hypermutations of microsatellite sites due
 39 to inactivating alterations of mismatch repair (MMR) genes in malignancies [1, 2].
 40 Currently, MSI is an indispensable pan-cancer biomarker in cancer immunotherapy
 41 therapy and prognosis, and it is routinely examined in multiple cancer types,
 42 particularly in colorectal cancer (CRC), stomach adenocarcinoma (STAD), and
 43 uterine corpus endometrial carcinoma (UCEC) [2-5]. For example, MSI positive
 44 patients are often resistant to 5-fluorouracil treatment but have a better outcome for
 45 immune checkpoint blockade treatment [4, 5].

46 In clinical settings, MSI detection mainly relies on the gold-standard experimental
 47 method, MSI-PCR [6], which is laborious and time-consuming. With the
 48 advancement of next-generation-sequencing technology, numerous features of
 49 genomics, epigenomics, transcriptomics, and histology are investigated, and novel
 50 MSI computational algorithms have been developed for a variety of scenarios [7-20].
 51 Genomics-based methods quantify MSI according to genetic mutations at
 52 microsatellite sites, which achieve high accuracy and are becoming popular in
 53 clinical MSI detection. For example, MSIsensor [9] detects MSI with high
 54 concordance as 99.4% on MSK-IMPACT panel [21]. Epigenomics-based method
 55 MIRMMR [18] detects MSI using methylation levels in MMR pathway with 0.97

56 AUC. In addition, transcription levels of MSI-associated genes exhibit correlation
57 with MSI, hinting possibility of MSI detection using transcriptomics data [15-17].
58 Besides these high-throughput technologies, deep learning algorithms were also
59 applied to hematoxylin and eosin-stained slides to detect MSI [19, 20]. However, all
60 these MSI methods detected MSI at a sample level, lacking cell-level measuring of
61 MSI. Recently, single-cell RNA-seq (scRNA-seq) technology enables investigation
62 of cell specific transcriptome and sheds light on tumor heterogeneity and tumor
63 stages. In particular, the single-cell and spatial transcriptome enable the dynamic
64 analysis of MSI in the complex tumor microenvironment, such as in metastatic and
65 recurrent cancer [22]. However, current MSI detection methods designed for bulk
66 gene expression data do not perform well on scRNA-seq samples. For example, the
67 only software for gene expression data, PreMSIm [16], only provided fixed
68 signatures and a fixed model for all cancers, which limits the widely application of
69 the methods. Moreover, the normalized method in PreMSIm also leads to poor
70 performance with abnormal samples. Here, we developed MSIsensor-RNA, a robust
71 method for MSI-associated genes detection and MSI evaluation for both bulk gene
72 expression data and single-cell RNA-seq data.

73 **Implementation**

74 *Dataset.* We downloaded RNA-seq data of 1,428 TCGA samples across CRC, STAD,

75 and UCEC from TCGA Research Network (<https://portal.gdc.cancer.gov>) and
76 obtained their MSI status determined by gold standards (**Table S1**). We obtained 141
77 RNA-seq samples of ICGC from ICGC data portal (<https://dcc.icgc.org>), and their
78 MSI status reported by MIMcall [23]. Another 106 RNA-seq samples with the
79 matched MSI status were downloaded from public publication of Clinical Proteomic
80 Tumor Analysis Consortium (CPTAC) [24]. We also downloaded Microarray data
81 and their MSI status of 1,468 samples across CRC and STAD from GEO dataset
82 (<https://www.ncbi.nlm.nih.gov/geo>). For scRNA-seq data, we got the gene
83 expression data and their MSI status from 133 CRC samples in two recent
84 publications [25, 26].

85 *Overall design.* The pipeline of MSIsensor-RNA consists of data preprocessing,
86 informative genes selection, model training, and model testing (**Fig. 1 and Fig. S1**).
87 First, we preprocess the expression values of samples from Microarray, bulk
88 RNA-seq, and scRNA-seq. Next, we select an informative gene set for MSI
89 detection from 1,428 TCGA samples. Then we used these TCGA samples to train a
90 machine learning model for each cancer type for MSI scoring. Finally, we applied
91 the trained model to independent databases to test the performance of the
92 MSIsensor-RNA for each cancer type.

93 *Data preprocessing.* In MSIsensor-RNA, we accept Microarray expression value,

94 FPKM, TPM, and RESM read count as input. All values of expression matrix were
95 added 1 and followed by log2 transformed. Then, for each sample or cell, expression
96 values were normalized as a Gaussian distribution with 0 mean and 1 standard
97 deviation. For scRNA-seq sample, to obtain accurate MSI status, we only included
98 high-quality cells with at least 20% genes detected for MSI detection. If the number
99 of high-quality cells was less than 20, we sort all cells by the ratio of detected genes
100 in descending order, and the top 20 cells would be utilized for MSI detection. To
101 solve the dropout problem of scRNA-seq, we imputed zero values by the average of
102 the gene expression value in the given sample.

103 *Selection of informative genes.* We select informative genes for MSI classification
104 in terms of stability, discrimination, and generalization. Firstly, we remove
105 ribosomal genes, mitochondrial genes, and genes with low FPKM in TCGA dataset.
106 Secondly, we selected genes with discriminative gene expression signatures between
107 MSI samples and MSS samples. We perform rank-sum tests for expression values
108 between MSI samples and MSS samples for each gene, and only genes with P value
109 < 0.01 are included for the following analysis. Furthermore, we compute the fold of
110 i th gene by:

$$F^i = \left| \log_2 \left(\frac{\frac{1}{n} \sum_{j=1}^n G_j^i}{\frac{1}{m-n} \sum_{k=n+1}^m G_k^i} \right) \right|$$

111 where m is the sample number for informative genes selection, n is the MSI sample

number, G_j^i is the gene expression value of i th for j sample. We only select genes with fold > 0.5 for candidate informative genes. Finally, we keep genes with more generalization ability for MSI detection. We calculate the area under the receiver operating characteristic curve (AUC) of the gene expression value and only genes with AUC > 0.65 are kept for next step. We also calculate the 10-fold cross validation score of SVM and random forest, and only first quartet genes are included the final informative gene set (**Fig. S2**).

Machine learning model training and testing. We build a support vector machine (SVM) model to classify the MSI status for CRC, STAD, and UCEC in TCGA dataset. Firstly, we utilized SOMTE [27] to correct the imbalance between MSI and MSS in each cancer type by amplifying the MSI samples. Then, we utilized the expression values from correct data as input to train SVM model for MSI classification. To evaluate the performance of MSI sensor-RNA, we tested the trained model with 1,848 independent samples of multiplatform including 247 RNA-seq, 1,468 Microarray, and 133 scRNA-seq samples. For a scRNA-seq sample, we calculated the MSI score with SVM model for each high-quality cell. Then the average cell MSI score is used to evaluate the MSI status of a scRNA-seq sample.

PreMSIm running. To compare performance of MSI sensor-RNA with the only standalone software PreMSIm, we also apply the data of Microarray, RNA-seq, and

131 scRNA-seq from 1,848 independent samples to PreMSIm. For Microarray and
 132 RNA-seq samples, we test PreMSIm with two modes: PreMSIm-all and
 133 PreMSIm-split. In PreMSIm-all, we integrate all input samples to PreMSIm
 134 normalized module and predicted module. PreMSIm-split referred to input samples
 135 one database for each run.

136 *Performance comparison of MSIsensor-RNA and PreMSIm.* In MSIsensor-RNA,
 137 the predicted MSI probability by the SVM model was used to score the MSI status.
 138 The probabilities were further transformed to MSI status by the Youden index [28].
 139 We first compared the MSIsensor-RNA score between MSI and MSS samples to test
 140 the performance of MSIsensor-RNA in multiplatform by rank sum test. To further
 141 evaluate the performance of two MSI detected methods, we calculated AUC,
 142 accuracy, F-score, precision, sensitivity, and specificity of MSIsensor-RNA and
 143 PreMSIm in different sequencing technologies.

144 *Robustness testing of MSIsensor-RNA and PreMSIm.* To test the performance of
 145 MSIsensor-RNA and PreMSIm at different normalized methods, we tested these two
 146 methods with FPKM, TPM, and read counts format of TCGA samples and
 147 calculated the AUC, F1-score, accuracy, precision, sensitivity, and specificity of
 148 each normalized method. To overcome the bias of different normalized methods and
 149 sequencing technology, we normalized the input data of each sample to a Gaussian

distribution with 0 mean and 1 standard deviation. However, in PreMSIm, the normalization process was performed by genes, which means the normalized input data of a sample would be influenced by other samples in the bulk. Here, we tested the PreMSIm in two ways. Firstly, we input TCGA samples by three cancer types and calculated the performance of predicted MSI result. Secondly, we input all TCGA samples together to evaluate its performance. We further compared the MSI result and performance of these two ways and found that the performance of PreMSIm was affected by the way input was provided.

Results

The workflow of MSIsensor-RNA includes four modules (**Fig. 1 and Fig. S1**). First, we preprocess the expression value of Microarray, bulk RNA-seq, and scRNA-seq data. Then, we select a set of informative genes for MSI detection. Next, we train a support vector machine (SVM) model to estimate MSI scores using gene expression values of the selected informative genes. Finally, we apply the trained model to predict MSI score for either one clinical sample or a single cell (**Table S1**). For a given scRNA-seq sample, we also developed a model to report MSI status of this sample by integrating MSI scores of cells within.

MSIsensor-RNA accepts a variety of expression data including FPKM, RESM normalized read count, TPM, or microarray expression format as input. Input

expression values were added 1 and then log2 transformed following Z-score normalization per sample or cell. In particular, for single cell module of MSIsensor-RNA, we only included high-quality cell in following steps, and the missing values of each gene in high-quality cells were imputed by the average of the gene expression value in this sample.

The informative gene selection module consists of three key steps (**Fig. S2**): (i) removing mitochondrial genes and ribosomal genes; (ii) filtering of genes, of which expression values do not differ significantly between MSI and MSS samples; (iii) keeping genes, of which expression values have high generalized scores for MSI detection (online methods). We applied the gene selection module to 1,428 samples based on the gene expressions (FPKM values) from three MSI-popular cancer types (CRC, STAD, and UCEC) in TCGA dataset and finally obtained 109 informative genes for MSI classification. We also performed this step for each type of CRC, STAD, and UCEC, yielding 397, 206, and 86 informative genes, respectively (**Fig. S4 and Table S2-S5**). We found that only eight informative genes are detected in all three cancer types. Of which, we found that *MLH1* was the most important informative gene for MSI detection, as confirmed by previous reports [15-17] (**Fig. S5**).

To assess the performance of MSIsensor-RNA in bulk sample data, we first

188 trained tumor-specific models for CRC, STAD, and UCEC, as well as a model for
189 all three MSI-popular cancer types in the TCGA dataset. Then we compared the two
190 kinds of models (tumor-specific and MSI-popular) with the standalone software,
191 PreMSIm, in terms of the area under the curve (AUC) of the receiver operating
192 characteristic (ROC), accuracy, sensitivity, and specificity in 1,715 (1468
193 Microarray and 247 bulk RNA-seq samples) independent samples. Notably,
194 MSIsensor-RNA normalizes the expression value of informative genes for each
195 sample independently, while PreMSIm must normalize each gene for multiple
196 samples at the same time. Thus, we examined PreMSIm with all samples normalized
197 together (PreMSIm-all) or by database (PreMSIm-split).

198 For Microarray data, we computed MSI status by MSIsensor-RNA and PreMSIm
199 in 1,468 samples from 12 GEO accessions. The result showed that MSIsensor-RNA
200 predicted MSI with 0.952 AUC, while PreMSIm only performed 0.628 AUC in
201 PreMSIm-split and 0.912 AUC in PreMSIm-all mode (**Fig. 2A, S6, S7; Table S6**
202 **and S7**). Meanwhile, MSIsensor-RNA achieved much higher sensitivities than
203 PreMSIm-split, and preMSI-all (MSIsensor-RNA: 0.968, PreMSIm-split: 0.912,
204 PreMSIm-all: 0.384) and comparable specificities with PreMSIm-split, and
205 preMSI-all (MSIsensor-RNA: 0.843, PreMSIm-split: 0.912, PreMSIm-all: 0.873).

206 To evaluate the performance using bulk RNA-seq data, we compared

MSIsensor-RNA and two modes of PreMSIm on 247 independent samples from ICGC and CPTAC. We noticed that MSIsensor-RNA achieved 0.997 AUC in tumor-specific model and 0.985 AUC in MSI-popular model, which were significantly greater than PreMSIm-all (0.5) and PreMSIm-split (0.870) (**Fig. 2B, S8, S9; Table S8 and S9**). In addition, MSIsensor-RNA performed much better than PreMSIm for both sensitivity (MSIsensor-RNA with tumor-specific model: 0.951, MSIsensor-RNA with MSI-popular model: 0.973, PreMSIm-split: 0.834, PreMSIm-all: 0.25) and specificity (MSIsensor-RNA with tumor-specific mode: 1, MSIsensor-RNA with MSI-popular model: 0.923, PreMSIm-split: 0.906, PreMSIm-all: 0.75). To further investigate the robustness of MSIsensor-RNA for different input data types, we evaluated the performance of MSIsensor-RNA and PreMSIm with FPKM, read count, and TPM normalized samples in TCGA as input. We found that MSIsensor-RNA achieved 0.982 ± 0.040 AUC indicating the robustness of MSIsensor-RNA regardless of the measurements of gene expression (**Table S10**).

To assess the performance of MSIsensor-RNA and PreMSIm in scRNA-seq samples, we applied the trained model of MSIsensor-RNA to 23,902 high-quality cells from 133 samples to obtain sample specific MSI status and compared to the ratio of cells labeled as MSI by PreMSIm. The result showed MSIsensor-RNA detected MSI for scRNA-seq samples with 0.958 AUC, 0.9231 sensitivity, and

0.9362 specificity, while PreMSIm with 0.4969 AUC, 1 sensitivity, and 0.0319 specificity (**Fig. 2A, S10; Table S11 and S12**). The sample level MSI scores based on scRNA-seq was significantly different between MSI and MSS samples by MSIsensor-RNA (rank-sum test, $P = 1.01 \times 10^{-16}$) while no significant difference was detected for PreMSIm (rank-sum test, $P = 0.9547$) (**Fig. 2B**). Having established the effectiveness of MSIsensor-RNA on scRNA-seq sample, we investigated cell-level MSI. We computed the MSI scores of 21,438 high-quality cells from 100 samples (GSE178341) and found cell-type dependent MSI scores. For example, MSI scores of epithelial and immune cells in MSI samples were greater than that in MSS samples while no significant difference was detected between MSI and MSS for stromal cells (**Fig. 2C, S11 and Table S13**). This indicated the potential of MSIsensor-RNA to assess MSI at the single-cell level, providing a novel measurement for the investigation of tumorigenic process.

Discussion

Microsatellite instability is important for the prognosis assessment of both 5-FU chemotherapy [4] and immunotherapy [5]. In addition to gold-standard experimental methods [6], MSI status is also evaluated according to genomic sequencing data [7-14], gene expression data [15-17], methylation data [18], and H&E-stained slides [19, 20]. Compared to variants in microsatellite regions, gene expression values are

more directly reflective of the features of MSI and easier to obtain. In this study, we developed a robust method, MSIsensor-RNA, for MSI detection with gene expression data. MSIsensor-RNA provided informative gene selections, model training, and MSI detection modules. MSIsensor-RNA is able to process data from multiple platforms, including Microarray, RNA-seq, and single cell RNA-seq. Compared to the standalone method PreMSIm, MSIsensor-RNA also provided modules for informative gene selection and model training so that users could apply MSIsensor-RNA for different cancer types. MSIsensor-RNA also improved the normalization method of the data, yielding a more robust result than PreMSIm (Fig. 2). In addition, MSIsensor-RNA facilitates the evaluation of MSI status at the single cell level, which will be critical to better understanding the mechanism of MSI in cancer immunotherapy in the future.

In most MSI detection methods, such as MSIsensor [10] and MSIsensor-pro [11], MSI is quantified according to genetic mutations at microsatellite sites, the consequence of MSI rather than the deficiency of the MMR system, the direct cause of MSI. In this study, a set of MSI-associated genes was identified, and their expression values were used for MSI evaluation. We found that *MLH1* is the most important gene in all tested cancer types. In addition, unexpected expression of *MLH1* is commonly seen in Lynch syndrome [29]. Thus, we test the performance of MSIsensor-RNA for samples with abnormal *MLH1* expression. We train a model

266 based on all informative genes and tested it by samples with simulated abnormal
267 *MLH1* gene expression (**Table S14**). We found that the model achieved 0.974 and
268 0.972 AUCs when we set the *MLH1* expression value as the maximum and
269 minimum of all gene expression values, respectively. Furthermore, when *MLH1* was
270 excluded from the informative gene set, MSI-sensor-RNA also achieved a 0.977
271 AUC, indicating the robustness of MSI-sensor-RNA for MSI detection.

272 We demonstrate that MSI-sensor-RNA achieved higher performance than other
273 methods based on gene expression and comparable performance compared to
274 DNA-based methods (**Table S15**). In our study design, MSI-sensor-RNA detects MSI
275 according to the gene expression signature of genes on MSI associated pathways,
276 while MSI-sensor evaluates MSI by computing the ratio of somatic microsatellite
277 mutations. Although MSI-sensor achieved slightly higher performance than
278 MSI-sensor-RNA, it cannot replace the applications of MSI-sensor-RNA in gene
279 expression data. Currently, MSI-sensor-RNA reports favorable performance in all
280 three MSI-popular cancers, including colorectal cancer, stomach adenocarcinoma,
281 and uterine corpus endometrial carcinoma. The MSI features are different in
282 different cancer types. Thus, the model obtained low performance when the testing
283 samples were inconsistent with training samples in cancer types (**Table S16-S18**).
284 Therefore, the performance of MSI-sensor-RNA in other cancer types needs further
285 validation in the future.

286 **Conclusions**

287 MSI sensor-RNA is a cross-platform, efficient, and robust method for MSI status
 288 determination from both bulk and single-cell gene expression data. We demonstrated
 289 the effectiveness and robustness of MSI sensor-RNA across different platforms,
 290 hinting its potential in clinical research. Moreover, MSI sensor-RNA enables
 291 single-cell level MSI evaluation, providing a new tool to discover the role of MSI in
 292 tumorigenic process and to monitor cell-level dynamic changes during
 293 immunotherapy.

294

295 **Availability and requirements**

296 Project name: msisensor-rna

297 Project home page: <https://github.com/xjtu-omics/msisensor-rna>

298 Operating system(s): Unix System or Docker

299 Programming language: python

300 Other requirements: python packages including numpy, pandas and scikit-learn.

301 License: Custom License (see at homepage)

302 Any restrictions to use by non-academics: MSI sensor-RNA is free for

303 non-commercial use by academic, government, and non-profit/not-for-profit
 304 institutions. A commercial version of the software is available and licensed through
 305 Xi'an Jiaotong University. For more information, please contact with
 306 pengjia@stu.xjtu.edu.cn or kaiye@xjtu.edu.cn.

307 **Abbreviations**

308 MSI: Microsatellite Instability
 309 MMR: Mismatch Repair
 310 CRC: Colorectal Cancer
 311 STAD: Stomach Adenocarcinoma
 312 UCEC: Uterine Corpus Endometrial Carcinoma
 313 NGS: Next Generation Sequencing
 314 ROC: Receiver Operating Characteristic
 315 AUC: Area Under the Curve
 316 scRNA-seq single-cell RNA sequencing

317 **Declarations**

318 **Ethics approval and consent to participate**

319 Not applicable.

320 **Consent for publication**

321 Not applicable.

322 **Availability of data and materials**

323 MSIsensor-RNA is a python program which is available at

324 <https://github.com/xjtu-omics/msisensor-rna>; Supplementary information:

325 Supplementary data are available at xxxx online.

326 **Competing interests**

327 The authors declare that they have no competing interests.

328 **Funding**

329 K.Y. and X.Y. are supported by the National Natural Science Foundation of China

330 (32125009, 32070663 and 62172325), by the National Key R&D Program of China

331 (2022YFC3400300), the Key Construction Program of the National “985” Project,

332 the Fundamental Research Funds for the Central Universities (xzy012020012), and

333 the Natural Science Basic Research Program of Shaanxi (2021GXLH-Z-098).

334 **Authors' contributions**

335 KY conceived and designed the study. PJ collected the data. PJ, XUY, XIY, TW, and YX

336 discussed and developed the method. PJ and XUY implemented the method. PJ and KY

337 wrote the manuscript. The authors read and approved the final manuscript.

338 Acknowledgments

339 We thank everyone in the Ye lab at Xi'an Jiaotong University for helpful comments
340 and discussion. We also thank the producers of data we used.

341 References

- 342 1. Yamamoto H, Imai K: **Microsatellite instability: an update.** *Arch Toxicol*
343 2015, **89**(6):899-921.
- 344 2. Hause RJ, Pritchard CC, Shendure J, Salipante SJ: **Classification and**
345 **characterization of microsatellite instability across 18 cancer types.** *Nat*
346 *Med* 2016, **22**(11):1342-1350.
- 347 3. Rizvi NA, Hellmann MD, Snyder A, Kvistborg P, Makarov V, Havel JJ, Lee W,
348 Yuan J, Wong P, Ho TS *et al*: **Cancer immunology. Mutational landscape**
349 **determines sensitivity to PD-1 blockade in non-small cell lung cancer.**
350 *Science* 2015, **348**(6230):124-128.
- 351 4. Ribic CM, Sargent DJ, Moore MJ, Thibodeau SN, French AJ, Goldberg RM,
352 Hamilton SR, Laurent-Puig P, Gryfe R, Shepherd LE *et al*: **Tumor**
353 **microsatellite-instability status as a predictor of benefit from**
354 **fluorouracil-based adjuvant chemotherapy for colon cancer.** *N Engl J*
355 *Med* 2003, **349**(3):247-257.
- 356 5. Le DT, Uram JN, Wang H, Bartlett BR, Kemberling H, Eyring AD, Skora AD,
357 Lubner BS, Azad NS, Laheru D *et al*: **PD-1 Blockade in Tumors with**
358 **Mismatch-Repair Deficiency.** *N Engl J Med* 2015, **372**(26):2509-2520.
- 359 6. Baudrin LG, Deleuze JF, How-Kit A: **Molecular and Computational**
360 **Methods for the Detection of Microsatellite Instability in Cancer.** *Front*
361 *Oncol* 2018, **8**:621.
- 362 7. Janavicius R, Matiukaite D, Jakubauskas A, Griskevicius L: **Microsatellite**
363 **instability detection by high-resolution melting analysis.** *Clin Chem* 2010,
364 **56**(11):1750-1757.
- 365 8. Boland CR, Thibodeau SN, Hamilton SR, Sidransky D, Eshleman JR, Burt
366 RW, Meltzer SJ, Rodriguez-Bigas MA, Fodde R, Ranzani GN *et al*: **A**
367 **National Cancer Institute Workshop on Microsatellite Instability for**
368 **cancer detection and familial predisposition: development of**

-
- 369 **international criteria for the determination of microsatellite instability in**
370 **colorectal cancer.** *Cancer Res* 1998, **58**(22):5248-5257.
- 371 9. Kautto EA, Bonneville R, Miya J, Yu L, Krook MA, Reeser JW,
372 Roychowdhury S: **Performance evaluation for rapid detection of**
373 **pan-cancer microsatellite instability with MANTIS.** *Oncotarget* 2016,
374 **8**(5).
- 375 10. Huang MN, McPherson JR, Cutcutache I, Teh BT, Tan P, Rozen SG: **MSIseq:**
376 **Software for Assessing Microsatellite Instability from Catalogs of**
377 **Somatic Mutations.** *Sci Rep* 2015, **5**:13321.
- 378 11. Niu B, Ye K, Zhang Q, Lu C, Xie M, McLellan MD, Wendl MC, Ding L:
379 **MSIsensor: microsatellite instability detection using paired**
380 **tumor-normal sequence data.** *Bioinformatics* 2014, **30**(7):1015-1016.
- 381 12. Jia P, Yang X, Guo L, Liu B, Lin J, Liang H, Sun J, Zhang C, Ye K:
382 **MSIsensor-pro: Fast, Accurate, and Matched-normal-sample-free**
383 **Detection of Microsatellite Instability.** *Genomics Proteomics*
384 *Bioinformatics* 2020, **18**(1):65-71.
- 385 13. Hempelmann JA, Scroggins SM, Pritchard CC, Salipante SJ: **MSIplus for**
386 **Integrated Colorectal Cancer Molecular Testing by Next-Generation**
387 **Sequencing.** *J Mol Diagn* 2015, **17**(6):705-714.
- 388 14. Salipante SJ, Scroggins SM, Hampel HL, Turner EH, Pritchard CC:
389 **Microsatellite instability detection by next generation sequencing.** *Clin*
390 *Chem* 2014, **60**(9):1192-1199.
- 391 15. Maruvka YE, Mouw KW, Karlic R, Parasuraman P, Kamburov A, Polak P,
392 Haradhvala NJ, Hess JM, Rheinbay E, Brody Y *et al*: **Analysis of somatic**
393 **microsatellite indels identifies driver events in human tumors.** *Nat*
394 *Biotechnol* 2017, **35**(10):951-959.
- 395 16. Ratovomanana T, Cohen R, Svrcek M, Renaud F, Cervera P, Siret A,
396 Letourneur Q, Buhard O, Bourgoin P, Guillermin E *et al*: **Performance of**
397 **Next-Generation Sequencing for the Detection of Microsatellite**
398 **Instability in Colorectal Cancer With Deficient DNA Mismatch Repair.**
399 *Gastroenterology* 2021, **161**(3):814-826 e817.
- 400 17. Pacinkova A, Popovici V: **Cross-platform Data Analysis Reveals a Generic**
401 **Gene Expression Signature for Microsatellite Instability in Colorectal**
402 **Cancer.** *Biomed Res Int* 2019, **2019**:6763596.
- 403 18. Li L, Feng Q, Wang X: **PreMSIm: An R package for predicting**

-
- 404 **microsatellite instability from the expression profiling of a gene panel in**
405 **cancer.** *Comput Struct Biotechnol J* 2020, **18**:668-675.
- 406 19. Danaher P, Warren S, Ong S, Elliott N, Cesano A, Ferree S: **A gene**
407 **expression assay for simultaneous measurement of microsatellite**
408 **instability and anti-tumor immune activity.** *Journal for immunotherapy of*
409 *cancer* 2019, **7**(1):15.
- 410 20. Foltz SM, Liang WW, Xie M, Ding L: **MIRMMR: binary classification of**
411 **microsatellite instability using methylation and mutations.** *Bioinformatics*
412 2017, **33**(23):3799-3801.
- 413 21. Echle A, Grabsch HI, Quirke P, van den Brandt PA, West NP, Hutchins GGA,
414 Heij LR, Tan X, Richman SD, Krause J *et al*: **Clinical-Grade Detection of**
415 **Microsatellite Instability in Colorectal Tumors by Deep Learning.**
416 *Gastroenterology* 2020, **159**(4):1406-1416 e1411.
- 417 22. Kather JN, Pearson AT, Halama N, Jager D, Krause J, Loosen SH, Marx A,
418 Boor P, Tacke F, Neumann UP *et al*: **Deep learning can predict**
419 **microsatellite instability directly from histology in gastrointestinal cancer.**
420 *Nat Med* 2019, **25**(7):1054-1056.
- 421 23. Middha S, Zhang L, Nafa K, Jayakumaran G, Wong D, Kim HR, Sadowska J,
422 Berger MF, Delair DF, Shia J *et al*: **Reliable Pan-Cancer Microsatellite**
423 **Instability Assessment by Using Targeted Next-Generation Sequencing**
424 **Data.** *JCO Precis Oncol* 2017, **2017**(1):1-17.
- 425 24. Longo SK, Guo MG, Ji AL, Khavari PA: **Integrating single-cell and spatial**
426 **transcriptomics to elucidate intercellular tissue dynamics.** *Nat Rev Genet*
427 2021, **22**(10):627-644.
- 428 25. Fujimoto A, Fujita M, Hasegawa T, Wong JH, Maejima K, Oku-Sasaki A,
429 Nakano K, Shiraishi Y, Miyano S, Yamamoto G *et al*: **Comprehensive**
430 **analysis of indels in whole-genome microsatellite regions and**
431 **microsatellite instability across 21 cancer types.** *Genome Res* 2020.
- 432 26. Vasaikar S, Huang C, Wang X, Petyuk VA, Savage SR, Wen B, Dou Y, Zhang
433 Y, Shi Z, Arshad OA *et al*: **Proteogenomic Analysis of Human Colon**
434 **Cancer Reveals New Therapeutic Opportunities.** *Cell* 2019,
435 **177**(4):1035-1049 e1019.
- 436 27. Pelka K, Hofree M, Chen JH, Sarkizova S, Pirl JD, Jorgji V, Bejnood A,
437 Dionne D, Ge WH, Xu KH *et al*: **Spatially organized multicellular immune**
438 **hubs in human colorectal cancer.** *Cell* 2021, **184**(18):4734-4752 e4720.

28. Lee HO, Hong Y, Etlioglu HE, Cho YB, Pomella V, Van den Bosch B, Vanhecke J, Verbandt S, Hong H, Min JW *et al*: **Lineage-dependent gene expression programs influence the immune landscape of colorectal cancer**. *Nat Genet* 2020, **52**(6):594-603.
29. Chawla NV, Bowyer KW, Hall LO, Kegelmeyer WP: **SMOTE: Synthetic Minority Over-sampling Technique**. *Journal of Artificial Intelligence Research* 2002, **16**:321-357.
30. Schisterman EF, Perkins NJ, Liu A, Bondell H: **Optimal cut-point and its corresponding Youden Index to discriminate individuals using pooled blood samples**. *Epidemiology* 2005, **16**(1):73-81.
31. Chen W, Hampel H, Pearlman R, Jones D, Zhao W, Alsomali M, Knight D, Frankel WL: **Unexpected expression of mismatch repair protein is more commonly seen with pathogenic missense than with other mutations in Lynch syndrome**. *Hum Pathol* 2020, **103**:34-41.

Figure legends

Fig. 1. Workflow of MSIsensor-RNA. MSIsensor-RNA includes four modules: data preprocessing, informative gene selection, SVM model training, and testing. MSIsensor-RNA selects informative genes and trains SVM model by RNA-seq samples from TCGA. MSI scores are predicted by the trained model for Microarray, RNA-seq, and scRNA-seq samples.

Fig. 2. Performance of MSIsensor-pro. A-C. AUC of MSIsensor-RNA and PreMSIm in Microarray (A), RNA-seq (B), and scRNA-seq (C) samples. Tumor-specific: MSI results with tumor specific model; MSI-popular: MSI results with three MSI-popular cancer types. **D.** Boxplot of MSIsensor-RNA score in scRNA-seq samples. **E.** Violin plot of MSIsensor-RNA score of different cell types

464 in scRNA-seq samples. Epithelial, stromal, and immune cell types are defined in

465 Pelka et al.[25]

466 **Supplementary tables**

467 Table S1. Overview of samples in this study.

468 Table S2. Details of informative genes in CRC.

469 Table S3. Details of informative genes in STAD.

470 Table S4. Details of informative genes in UCEC.

471 Table S5. Details of informative genes in three MSI-popular cancers.

472 Table S6. MSI results of Microarray samples by MSIsensor-RNA and PreMSIm.

473 Table S7. MSI detection performance of MSIsensor-RNA and PreMSIm in

474 Microarray samples.

475 Table S8. MSI results of RNA-seq samples by MSIsensor-RNA and PreMSIm.

476 Table S9. MSI detection performance of MSIsensor-RNA and PreMSIm in RNA-seq

477 samples.

478 Table S10. MSI detection performance of MSIsensor-RNA and PreMSIm in

479 different normalized samples.

480 Table S11. MSI results of scRNA-seq samples by MSIsensor-RNA.

481 Table S12. MSI detection performance of MSIsensor-RNA and preMSIm in
482 scRNA-seq samples.

483 Table S13. MSI results of scRNA-seq cells by MSIsensor-RNA.

484 Table S14. Performance of MSIsensor-RNA with abnormal MLH1 expression
485 values.

486 Table S15. Performance of MSIsensor-RNA and MSIsensor in TCGA dataset.

487 Table S16. AUC of MSIsensor-RNA with inconsistent training and testing samples.

488 Table S17: Performance of train models for cancer with low frequency MSI.

489 Table S18: Performance of MSIsensor-RNA for cancer with low frequency MSI by
490 5-fold cross validation.

491 **Supplementary figures:**

492 See in supplementary materials.

493

494 **Figures:**

495 **Fig. 1**

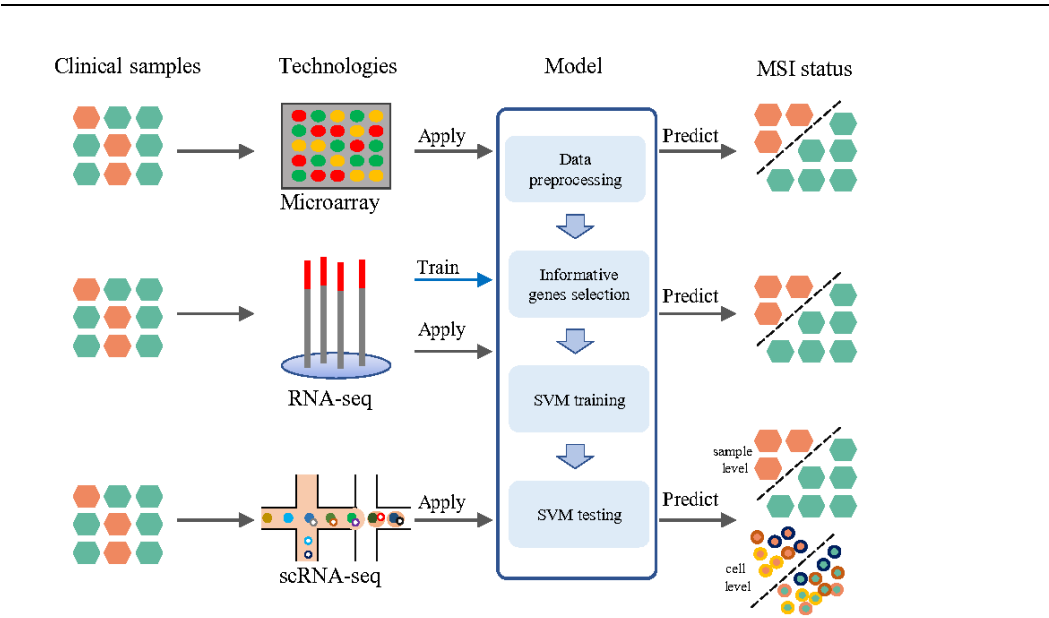
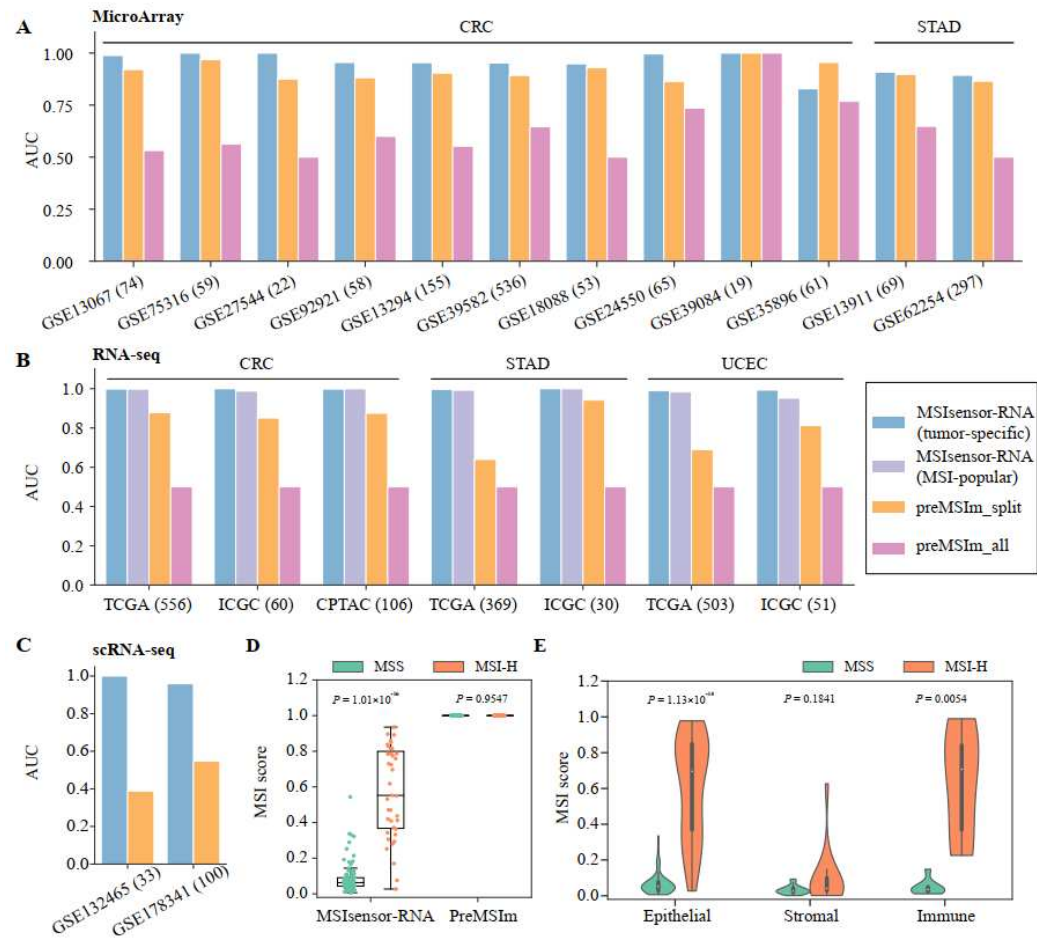


Fig. 2

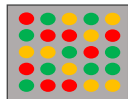


Clinical samples

Technologies

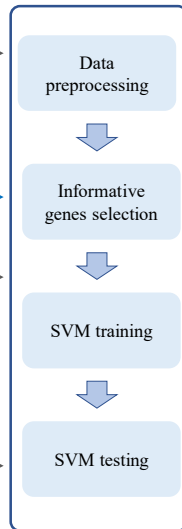
Model

MSI status

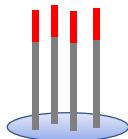


Microarray

Apply



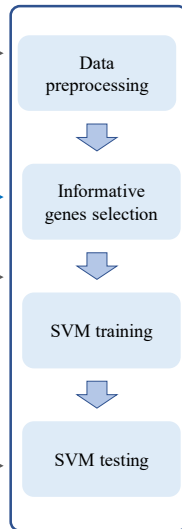
Predict



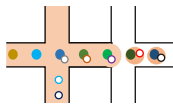
RNA-seq

Train

Apply

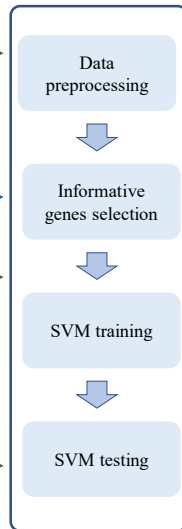


Predict



scRNA-seq

Apply



Predict

