Title: Intratumoral heterogeneity in microsatellite instability status at single cell resolution

Authors and affiliations: Harrison Anthony[1,2] and Cathal Seoighe[1,2]

[1] School of Mathematical and Statistical Sciences, University of Galway, Galway, Ireland

[2] The Research Ireland Centre for Research Training in Genomics Data Science, Galway, Ireland


Author for correspondence: Cathal Seoighe (cathal.seoighe@universityofgalway.ie), School of Mathematical and Statistical Sciences, University Road, University of Galway, Ireland.

The authors declare no potential conflicts of interest.

**Abstract**

Subclonal diversity within a tumor is highly relevant for tumor evolution and treatment. This diversity is often referred to as intratumoral heterogeneity and is known to complicate the interpretation of single-test biomarkers. Microsatellite instability (MSI) is one such biomarker, which is used to help guide immune checkpoint inhibitor treatment through the classification of samples as either having high microsatellite instability (MSI-H) or as being microsatellite stable (MSS). One area that has yet to be addressed in depth is whether MSI itself is a heterogeneous phenomenon, such that only some subclones are MSI-H. To investigate heterogeneity in MSI status, we curated and analyzed data from several single-cell RNA sequencing studies that had paired clinical MSI status and developed a computational pipeline to infer MSI-H cells and quantify heterogeneity In MSI status. We found evidence of heterogeneity in MSI status both in individuals originally classified as MSI-H and MSS. Out of 49 individuals, 15 showed evidence of divergence in MSI status between distinct clusters of cancer cells and most had distinct MSI-H and MSS subclones. These results raise questions about the current practice of treating MSI as a binary biomarker, ignoring heterogeneity between cancer subclones. Accounting for heterogeneity may lead to improved biomarker performance and, potentially, help explain reports of intrinsic treatment resistance and low overall responder rate in MSI-H cancers. Further studies are warranted to determine the frequency of heterogeneity in this biomarker at the population level, and whether the presence of both MSI-H and MSS subclones can have clinical implications.

## 1. Introduction

Subclonal diversity within a tumor is a critical consideration in cancer research and treatment. The overall diversity found in a single neoplasia is called intratumoral heterogeneity (ITH). While ITH was first conceptualized to be genetic in nature[1], it is now used to describe genetic, epigenetic, and phenotypic differences between subclones[2]. The diversity within a tumor is important because ITH has been linked to poor patient outcomes, therapy resistance, and relapse[3,4]. Furthermore, biomarkers that rely on single-sample tests can be susceptible to

49    sampling bias when ITH is present[5,6]. While its origins are still debated[7], one well known driver

50    of ITH is genomic instability[8].

51    Genomic instability is a hallmark of cancer, characterized by a higher rate of accumulation of

52    mutations during replication, typically due to deficiencies in DNA repair genes[9]. The two most

53    common forms of genomic instability are at the chromosomal level, where instability is

54    characterized by aneuploidy and chromosomal aberrations[10], and at the microsatellite level,

55    where short tandem repeats expand and contract in a mutator phenotype manner[11]. The latter,

56    referred to as microsatellite instability (MSI), is hypothesized to be the result of a deficient

57    mismatch repair (dMMR) pathway and is commonly used as a biomarker to help guide immune

58    checkpoint inhibitor treatment. This is done by classifying cancers as either having high

59    microsatellite instability (MSI-H) or as being microsatellite stable (MSS)[12]. The classification is

60    normally carried out using a single-sample test that compares five microsatellite markers

61    between a tumor and paired-normal sample[13,14]. While the interplay between chromosomal

62    instability and ITH is well defined and explained[15–17], the relationship between MSI and ITH is

63    less clear.

64    Up to this point, most research on MSI and ITH has been framed around how MSI can impact

65    and shape the variation present within a tumor. Most notable have been studies focusing on

66    specific mutations[18,19] and the immune cell types present in the tumor microenvironment[20,21].

67    Although, there have been reported cases of ITH in MSI status[22–27], the question of whether

68    MSI itself is frequently a heterogeneous phenomenon, with some subclones displaying MSI

69    while others do not, has yet to be examined in detail. This warrants further investigation as it

70    may help to explain limitations of MSI-H as a biomarker for precision medicine, such as low

71    treatment response rates and intrinsic treatment resistance[28–30]. Taking heterogeneity into

72    account may, ultimately, lead to improved biomarker performance.

73    The current literature suggests that subclonality of MSI status is relatively rare or entirely

74    absent[31,32], but that is not always the case. There are many examples of individuals not only

75    with discordant MSI statuses between the primary tumor site and metastases[22–25] but also

76    between multiple sites in the primary tumor[26,27]. While these are small case studies, they

77 provide anecdotal evidence for cancers comprising MSI-H and MSS subclones. However, there

78 has, as yet, been no attempt to evaluate the frequency with which this occurs. Furthermore, a

79 detailed examination of heterogeneity requires an assessment of MSI at the single-cell level

80 with next-generation sequencing, not with the traditional methods of PCR and IHC used in

81 these case studies, as these are limited to detecting clear spatial heterogeneity.

82 Here we aimed to address these gaps through an analysis of published single-cell datasets that

83 include paired clinical MSI status. To do this, we developed a custom Snakemake[33] pipeline that

84 identifies MSI-H cells and uses novel methods to assess levels of heterogeneity and have made

85 this pipeline available as an open-source, scalable resource to the scientific community. We

86 evaluated the pipeline by mixing varying numbers of MSI-H and MSS cells from different

87 samples. Applying this framework, we show evidence of heterogeneity in MSI status at the

88 single-cell level and estimate its prevalence in the curated data. We also examine the nature of

89 MSI heterogeneity through a detailed investigation of single-cell data from two individuals –

90 one classified as MSI-H and the other MSS through PCR/IHC tests.

91 ## 2. Materials and Methods

92 **2.1 Datasets**

93 We used single-cell RNA sequencing data that was generated as part of three previous

94 studies[34–36] (Table 1). Raw FASTQ files were downloaded from either the European Genome-

95 Phenome Archive (RRID:SCR 004944) or from the Sequence Read Archive (RRID:SCR 001370).

96 All other data was downloaded in matrix format from the Gene Expression Omnibus (RRID:SCR

97 005012). The data consists of 134 samples from 49 individuals with metastatic or non-

98 metastatic colorectal cancer. Individuals were grouped into MSI-H and MSS categories based on

99 the original PCR/IHC clinical status reported in previous studies. In total there were 29 deemed

100 MSI-H, 18 MSS, and two did not have a reported MSI status (Table 2). Each sample was created

101 with either Single Cell 3' v2, 3' v3, or 5' Reagent Kit from 10X Genomics and was sequenced

102 either on an Illumina NextSeq 500, NovaSeq 6000, BGISEQ DNBSEQ-T7, or HiSeq X Ten machine.

103 Complete sequencing and library preparation information can be found by referencing the

104 Dataset ID in Table 1. Individuals from datasets EGAD00001008555, EGAD00001008584,

105　EGAD00001008585 had multi-regional samples from the same tumor and multi-site samples

106　from metastatic tissue and lymph nodes. Even though variation in the multi-site samples could

107　be considered intra-individual heterogeneity rather than ITH, we kept them in the analysis to

108　retain as many cancer cells and as much heterogeneity as possible. The other two datasets

109　GSE205506 and PRJNA932556 include individuals that had treatment for MSI (anti-PD-1 and

110　celecoxib).  We excluded the following samples because we did not identify any cancer cells:

111　XHC080-SI-GA-B11, XHC082-SI-GA-C1, XHC127-SI-GA-F10, EXT129, EXT051, and EXT097.

112　**2.2 Data processing**

113　We aligned FASTQ files to the GRCh38 human reference genome and converted them to a gene

114　count matrix using the 10x Genomics Cell Ranger v7.2.0 software suite. From there all matrix

115　files were processed following the Seurat best practices tutorials (https://satijalab.org/seurat/).

116　Briefly, Seurat objects were created and only genes detected in a minimum of 3 cells were used

117　for downstream analysis. We further filtered out cells with fewer than 100 features, more than

118　3000 features, and if a cell had more than 35 percent of all genes labeled as mitochondrial.

119　While we followed the Seurat best practices closely, the default settings aim to maximize

120　immune cell type identification and filtered out the majority of cancer cells. The filter settings

121　described here were designed to maximize the number of tumor cells retained while still

122　removing cells that were poor-quality or likely necrotic. These filter settings, specifically

123　mitochondrial gene percentage, are supported by a recent study that showed that filter settings

124　that are too strict remove viable cancer cells from single-cell sequencing data[37]. After filtering,

125　all gene count matrices then were normalized using the LogNormalize option with the scale

126　setting set to 10,000. The 2,000 most variable genes were found using the "vst" selection

127　method and were used to cluster together groups of cells with the RunPCA function. The first

128　15 PCA dimensions were used to run the following functions: FindNeighbors, FindClusters

129　(resolution set to 0.5), and RunUMAP. If an individual had multiple samples, they were

130　integrated together by using the IntegrateLayers function, with the method set to

131　CCAIntegration and k.weight set to 50. Each integrated sample then underwent re-clustering

132　with the settings previously mentioned. The integration step after subsetting down to only

133     cancer cells used the same settings except the FindClusters resolution was set to 0.8, and only

134     the first 10 principal components were used.

135     **2.3 Cell classification and measuring ITH**

136     After each sample is processed, all cells are classified as either cancer or normal and then MSI-H

137     or MSS. These classification steps are built upon two machine learning based programs trained

138     on large pan-cancer datasets. The first, scATOMIC[38], was used to distinguish tumor cells from

139     normal ones, and the second, MSIsensor-RNA[39] determined MSI status. Both tools were run

140     with default settings, but to get an MSI score for each cell, we had to transform the prebuilt

141     MSIsensor-RNA baseline file. This was done by filtering both the count matrix and baseline to

142     only include gene names common to both. The filtered baseline and count matrix files were

143     then used to get MSI scores for all cells within a sample. From there cells were classified as MSI-

144     H if they also were labeled as cancer by scATOMIC and if they had an MSI score of .75 or more

145     (75% probability the cell is MSI-H).

146     Levels of ITH were assessed with two methods. First, we measured ITH by testing for

147     differences in mean MSI score between cancer cell clusters with a one-way ANOVA test. The

148     ANOVA F-statistic was used to describe levels of heterogeneity in the biomarker, with a large

149     value of the F-statistic indicating greater heterogeneity in MSI. Secondly, we identified

150     subclones within each individual by comparing CNVs between MSI-H, MSS, and normal cells.

151     This was done by passing the relevant cell classification for each unique barcode to InferCNV[40]

152     (DOI: 10.18129/B9.bioc.infercnv). InferCNV was used with default settings except in the case of

153     CRC2821, which had many more cancer cells than the other samples. We increased the k_nn

154     setting from the default of 20 to 50 to take into account the larger dataset. Lastly, we ran

155     differential expression using a Wilcoxon Rank Sum Test (the default for Seurat) between

156     clusters of cancer cells and between MSI-H and MSS cancer cells for each individual.

157     We verified how well our pipeline captured heterogeneity in MSI status by mixing together

158     randomly sampled tumor and normal cells in varying proportions using a custom R script. We

159     simulated varying levels of heterogeneity by mixing the cells of one sample that had

160     homogenous MSI-H cancer cells (GSM6213995 from individual P33) and one sample with

161    homogenous MSS cancer cells (XHC118-SI-GA-F1 from individual CRC2811). In total, we had

162    eleven different mixes, with the proportion of MSI-H cells ranging from 0 to 1 (in increments of

163    0.1) and the remainder being MSS (Supplementary Table 1). The results of these mixing

164    experiments were replicated 100 times, except for the pure MSS and MSI-H cases, for which all

165    cancer cells were included.

166    Although MSIsensor-RNA has been shown to classify single-cell RNA sequencing samples

167    accurately[39], we checked its ability to distinguish between MSI-H and MSS samples in our

168    datasets at the individual level. This was done by scoring individuals with MSIsensor-RNA using

169    all available cells and again with just the cancer cells. We used the AggregateExpression

170    function in Seurat to create the two different scenarios, and measured MSIsensor-RNA

171    performance with ROC-AUC using the MLeval and caret packages in R[41,42].

172    **2.4 Statistical analysis**

173    All statistical analyses were carried out in R (R version 4.1.1; https://www.R-project.org/)[43] and

174    all plots were created with ggplot2[44]. Two statistical tests were performed as part of our

175    computational pipeline. The first is a one-way ANOVA test that we used to measure ITH by

176    comparing the difference in means between clusters of cancer cells. This was done with the aov

177    function and was followed with Tukey's Honestly Significant Difference test using the TukeyHSD

178    function, both of which are from the "stats" package[43].

179    **2.5 Data availability**

180    All data used in this study are available either publicly through the SRA (RRID:SCR 001370) and

181    Gene Expression Omnibus (RRID:SCR 005012) or through the European Genome-Phenome

182    Archive (RRID:SCR 004944) with data access requests. The associated dataset IDs and metadata

183    for each dataset are in Table 1.

184    **2.6 Code availability**

185    A distributable version of the computational pipeline, SC-MSI, used in this study is available on

186    GitHub (https://github.com/harrison-anth/sc_msi). We have written the entire workflow in

187    Snakemake to ensure reproducibility and scalability. All original results and code used in this

188    study as well as the R script used to mix together single-cell sequencing samples are stored in

189    another GitHub repository (https://github.com/harrison-anth/sc_msi_legacy).

190                                             **3.   Results**

191    **3.1 Computational pipeline distinguishes MSI-H and MSS individuals and captures ITH**

192    To determine whether MSIsensor-RNA could distinguish between MSI-H and MSS individuals,

193    we ran it on the aggregate expression of all cells and again on only the cancer cells for each

194    individual. As expected, MSI-H individuals generally had higher MSI scores than MSS individuals,

195    and MSIsensor-RNA was able to broadly distinguish between the two groups (Figure 1A,B and

196    Supplementary Figure 1A,B). These results were seen for both aggregated expression of all cells

197    and only cancer cells, but subsetting down to only cancer cells yielded lower MSI scores. There

198    were also disagreements between PCR/IHC MSI status and MSIsensor-RNA score with several

199    MSS individuals having relatively high MSI scores and several MSI-H individuals having low MSI

200    scores (Figure 1A,B).

201    Next, we simulated different levels of heterogeneity to determine how well our pipeline

202    captured ITH in MSI status. For this purpose, we simulated different levels of heterogeneity,

203    ranging progressively from pure MSS cells to pure MSI-H cells by mixing together samples from

204    two individuals comprised of homogeneously MSI-H cancer cells and homogeneously MSS

205    cancer cells (Supplementary Tables 1 and 2). As expected, the more homogeneous samples

206    (MSS, mix M1, mix M9 and MSI-H in Figure 1C) had low F-statistic values while mixtures with

207    more equal proportions of MSI-H and MSS cells (mixes M3-M5 in Figure 1C) had high F-statistic

208    values. Increasing the proportion of MSI-H cells until mix M7 resulted in an overall reduction in

209    the number of MSS subclones identified (Figure 2D and Supplementary Tables 1 and 2), and

210    there was one MSI-H subclone that was consistently detected after the proportion of MSI-H

211    cells was 0.1 (mixes M1-M9). Together, these results show that the F-statistic is sensitive to ITH,

212    and that the number of subclones can be consistently identified across replicates, providing

213    useful context to the heterogeneity.

**3.2 MSI-H and MSS individuals have evidence of ITH in MSI status**

In order to assess heterogeneity in MSI status, we first calculated F-statistics (see Materials and Methods) based on clusters of cancer cells and identified subclones based on CNV patterns. We found that MSI-H and MSS individuals both had evidence of heterogeneity in MSI. In total, 15 of 49 individuals showed evidence of divergence in MSI status between distinct clusters of cancer cells (F > 25; Figure 1 and Figure *2*). Several individuals had very large estimates of heterogeneity based on F-statistics (75.20 to 116.10) with most of these individuals being originally deemed to be MSS, and one originally deemed to be MSI-H from a PCR or IHC test. In contrast, the lowest F-statistics (1.30 to 1.68) were found in MSI-H and MSS individuals, and the ANOVA tests were not statistically significant in either case (P > 0.05; Supplementary Table 3). This was also seen in most other individuals with fewer than three cancer cell clusters (Table 2 and Supplementary Table 3).

In general, MSI-H and MSS individuals had similar distributions of F-statistic values but with several outliers having large F-statistic values among the MSS individuals (Figure 2A,B). Interestingly, nearly every individual in the analysis had both MSI-H and MSS subclones, and a larger proportion of MSS subclones (Figure 2A,B and Table 2). The exceptions were two individuals who each had a subclone proportion of 0.5 (Figure 2B), but they had very few cancer cells and too few cancer clusters to calculate an F-statistic for comparison (Supplementary Table 3). Those with the most MSI-H subclones, six and eight, were originally determined to be MSI-H, but one MSS individual also had four MSI-H subclones. Independent of MSI status, the distribution in number of clusters across individuals was relatively even, and most individuals had two or fewer samples used in the clustering process (Figure 2C,D).

**3.3 Single-cell level resolution of heterogeneity in one MSI-H and one MSS individual**

We selected two individuals (P24 and CRC2786) with relatively high F-statistics and many MSI-H subclones to illustrate the heterogeneity in MSI that is evident from single cell RNA-Seq data (Figure 4 and Figure *5*). The MSI-H individual, P24, had good overlap in cells classified as cancer with scATOMIC (Figure 4A) and those with high MSI scores determined by MSIsensor-RNA

241  (Figure 4B). The re-clustered cancer cells appear to cluster by MSI score; notably, clusters two

242  and three (Figure 4C). Those larger differences in the clusters were also seen in the pseudobulk

243  analysis of the re-clustered cancer cells (Figure 4D).

244  The MSS individual, CRC2786, also had good overlap between the cells determined to be

245  cancerous and those with a high MSI score; however, there was less separation of cancer and

246  normal cells in this individual (Figure 5A,B). Similarly, in the re-clustered cancer cells, cells with

247  higher and lower MSI scores were somewhat more intermingled than for the MSI-H individual,

248  although clusters three and four seem to be predominantly MSS and MSI-H, respectively (Figure

249  5C). This result is recapitulated in the pseudobulk analysis with cluster three showing a low MSI

250  score and cluster four a much higher one (Figure 5D).

251  **3.4 Significant differences in MSI score and gene expression between clusters of cancer cells**

252  We examined MSI ITH in CRC2786 and P24 further by assessing differences between clusters of

253  cancer cells. We found both individuals to have many clusters with significantly different MSI

254  scores (Figure 6), and several genes showed differences in expression between clusters and

255  between cells identified as MSI-H and MSS (Supplementary Figures 2 and 3). Both individuals

256  had clusters with high and low MSI scores (Figure 6A,B). These differences were found to be

257  statistically significant ($P < 0.05$ using a Tukey HSD test; Supplementary Tables 4 and 5). We

258  found that 35 cluster pairs for CRC2786 had significantly different MSI scores and 17 cluster

259  pairs for P24 (Figure 6C,D) were significantly different.

260  Within the clusters of cancer cells gene expression was also significantly different between

261  clusters (Supplementary Figures 2A and 3A; Supplementary Tables 6 and 7), and between the

262  MSI-H and MSS cells in those clusters (Supplementary Figures 2B and 3B; Supplementary Tables

263  8 and 9). The top five differentially expressed genes for each cluster of cancer cells for each

264  individual were retained for analysis as well as the top 50 differentially expressed genes

265  between MSI-H and MSS cells. Individual CRC2786 had three genes: *MALAT1, EEF1A1,* and

266  *SH3BGRL3* in common between those differentially expressed between clusters and between

267  cells with different MSI status; Supplementary Figure 2A,B). P24, on the other hand had one

268  gene, *PCLAF* that was differentially expressed between clusters and between MSI-H and MSS

269    cells. When comparing the differential expression analyses for both individuals, we found three

270    genes: *BMX*, *LRMP*, and *SH2D6* that were differentially expressed between clusters and two

271    genes (*TYMS*, *OXCT1*) in common that differentiated MSI-H and MSS cells.

272                                        **4.   Discussion**

273    In our study, we showed that MSI status can be heterogeneous at the single-cell level and

274    provide a pipeline to measure that heterogeneity with the clustering of cancer cells and CNV

275    based subclone analysis. These results contrast with the assumption that is commonly made,

276    both in research and in clinical practice, that MSI is dichotomous. While this assumption has

277    proven useful, enabling MSI-H to be applied as a biomarker for immune checkpoint inhibitor

278    treatment[45], overall responder rate has been reported to be as low as 31%[30]. This could be

279    explained, at least in part, by the heterogeneity in MSI-H individuals, which a binary

280    classification fails to take into account. Furthermore, it is well known that single-sample tests

281    (like the ones used to assign MSI status) are susceptible to under-sampling bias when ITH is

282    present[3] and multi-sample, multi-regional tests may be needed to improve classification. This is

283    supported by a recent study [46] that reported a higher accuracy in predicting immunotherapy

284    effectiveness over traditional PCR/IHC tests by incorporating MSI cell type proportion into an

285    MSI score. Similarly to our study, this also found that both MSI-H and MSS individuals had a

286    mixture of MSI-H and MSS cell types in single-cell sequencing data; however, their

287    methodology, which involved clustering cells based on gene-set enrichment of MSI-H and MSS

288    signatures, did not identify any MSS individuals with only MSS cells. Our pipeline was able to

289    find examples of MSS individuals comprising MSS cells only, which would make sense given that

290    microsatellite instability is a relatively rare trait, and it would be unlikely to be present in every

291    MSS individual in a study cohort. This is likely due to the main difference between our methods,

292    as we test individual cells for microsatellite instability, whereas Zhao et al. [46] labelled cells as

293    MSI-H at the pseudo-bulk level with gene-set enrichment guided cell clustering. Our study is

294    also different as we aimed to measure ITH and provide our pipeline in an open access format.

295

296  Our finding that nearly every MSI-H and MSS individual had MSI-H and MSS subclones has not

297  yet been reported in other studies; however, two case studies that infer subclonality of dMMR

298  status from discordant IHC test results have been reported[23,47]. Even though dMMR and MSI-H

299  technically refer to different phenomena, MSI is considered to be the byproduct of dMMR and

300  both are predictive of immune checkpoint inhibitor treatment efficacy. Combined with our

301  findings, these case studies provide insights that could help explain reports of 30% or more of

302  MSI-H cancers having primary resistance to single-agent immune checkpoint inhibitor

303  treatment[28,48]. A treatment regime for an MSI-H cancer would potentially miss one or more

304  MSS subclones, leaving behind a population of cells that would not respond in the same way to

305  immunotherapy. Although this would need to be demonstrated with a clinical experiment in

306  which treatment results are measured longitudinally, our results provide a plausible mechanism

307  for treatment resistance which is not currently given adequate consideration[48].

308  Our computational pipeline is the first to identify and quantify heterogeneity in MSI status at

309  the single-cell level. We built the pipeline around MSIsensor-RNA and scATOMIC, two pan-

310  cancer, machine learning based approaches. The combination of these programs may give rise

311  to some potential issues. Naturally, as both approaches are trained on gene expression data,

312  there will be overlap in genes used to train both classifiers and consequently overlap in cell type

313  prediction. Yet, we found different genes to be differentially expressed between cancer cell

314  clusters and MSI-H and MSS cells. This is likely because there is no overlap in training data

315  between the two tools. One other caveat is the loss of microsatellite instability signal in MSI-H

316  individuals after subsetting down to the cancer cells. Despite being necessary at the single-cell

317  level to only label cells as MSI-H if they were also determined to be cancerous by scATOMIC,

318  there were likely instances where MSIsensor-RNA correctly identified MSI-H cells and

319  scATOMIC did not. Going forward, it would be beneficial for a benchmarking study to be done

320  to determine if MSIsensor-RNA could also better identify cancer cells in MSI-H individuals.

321  Another factor to consider is that there can be an overlap between the genes used in clustering

322  of cells and the genes used to generate an MSI score. Whether one or more of the 100 genes

323  used in the MSIsensor-RNA baseline are included in the 2,000 most variable genes used in

324  clustering steps of pipeline will change from individual to individual. While not included in this

325    study, we have checked clustering of cancer cells with and without the 100 genes used by

326    MSIsensor-RNA and found it did not appear to affect the clustering results.

327    MSI is typically detected in next-generation sequencing data by comparing the distribution of

328    indels in microsatellites between a paired-normal and tumor sample. Because we only had

329    access to single-cell RNA sequencing data, we chose to use MSIsensor-RNA, one of the only

330    methods reported to accurately classify single-cell RNA sequencing samples. We found that it

331    could broadly distinguish between the individuals deemed MSI-H and MSS with PCR/IHC tests

332    (Figure 1A,B and Supplementary Figure 1A,B). However, it is worth noting that this method

333    does not directly detect MSI with microsatellites but uses machine learning models trained on

334    gene expression patterns from MSI-H and MSS individuals. This technique is more suited to

335    detecting dMMR, which is traditionally measured with differences in gene expression.

336    However, the two are inherently related as both are used as predictive biomarkers for the same

337    immunotherapy, and MSI is hypothesized to be the downstream result of dMMR.  Based on

338    these differences, it would be worthwhile to reproduce our results with data generated from

339    other single-cell sequencing technologies to ensure the 3' and 5' sequencing bias does not

340    factor into the expression-based differences we found in our results. Replicating these findings

341    with DNA based assays (like whole-genome amplification and sequencing) would permit the use

342    of NGS tools that measure differences in microsatellite repeats.

343    Altogether, we found that heterogeneity in microsatellite instability is more common than

344    previously reported and we found it both in MSI-H and MSS individuals. These results could

345    help to explain why there are reports of treatment resistance and low response rates in MSI-H

346    cancers treated with immune checkpoint inhibitors; however, our study only analyzed data

347    from 49 individuals that underwent 3' and 5' single-cell RNA sequencing. Further studies are

348    warranted to determine the frequency of heterogeneity in this biomarker at the population

349    level and whether the presence of MSI-H and MSS subclones can have clinical impacts,

350    including the capacity for rapid evolution of resistance to treatments for which MSI-H is used as

351    a biomarker.

352

## Acknowledgements

**References**

376

377   1.    Nowell PC. The clonal evolution of tumor cell populations. *Science (80- )*.
378         1976;194(4260):23-28.

379   2.    Marusyk A, Almendro V, Polyak K. Intra-tumour heterogeneity: A looking glass for
380         cancer? *Nat Rev Cancer*. 2012;12(5):323-334.

381   3.    Marusyk A, Janiszewska M, Polyak K. Intratumor Heterogeneity: The Rosetta Stone of
382         Therapy Resistance. *Cancer Cell*. 2020;37(4):471-484.

383   4.    Qazi MA, Vora P, Venugopal C, et al. Intratumoral heterogeneity: Pathways to treatment
384         resistance and relapse in human glioblastoma. *Ann Oncol*. 2017;28(7):1448-1456.

385   5.    Gilson P, Merlin JL, Harlé A. Deciphering Tumour Heterogeneity: From Tissue to Liquid
386         Biopsy. *Cancers (Basel)*. 2022;14(6).

387   6.    McGranahan N, Swanton C. Biological and Therapeutic Impact of Intratumor
388         Heterogeneity in Cancer Evolution. *Cancer Cell*. 2015;28(1):141.

389   7.    Sun R, Hu Z, Curtis C. Big bang tumor growth and clonal evolution. *Cold Spring Harb
390         Perspect Med*. 2018;8(5):1-14.

391   8.    Burrell RA, McGranahan N, Bartek J, Swanton C. The causes and consequences of genetic
392         heterogeneity in cancer evolution. *Nature*. 2013;501(7467):338-345.

393   9.    Negrini S, Gorgoulis VG, Halazonetis TD. Genomic instability an evolving hallmark of
394         cancer. *Nat Rev Mol Cell Biol*. 2010;11(3):220-228.

395   10.   Thompson SL, Bakhoum SF, Compton DA. Mechanisms of Chromosomal Instability. *Curr
396         Biol*. 2010;20(6):R285-R295.

397   11.   Yamamoto H, Watanabe Y, Maehata T, Imai K, Itoh F. Microsatellite instability in cancer:
398         a novel landscape for diagnostic and therapeutic approach. *Arch Toxicol*.
399         2020;94(10):3349-3357. doi:10.1007/s00204-020-02833-z

400  12.  Lee V, Murphy A, Le DT, Diaz LA. Mismatch Repair Deficiency and Response to Immune
401      Checkpoint Blockade. *Oncologist*. 2016;21(10):1200-1211.
402      doi:10.1634/theoncologist.2016-0046

403  13.  Murphy KM, Zhang S, Geiger T, et al. Comparison of the microsatellite instability analysis
404      system and the Bethesda panel for the determination of microsatellite instability in
405      colorectal cancers. *J Mol Diagnostics*. 2006;8(3):305-311.
406      doi:10.2353/jmoldx.2006.050092

407  14.  Berg KD, Glaser CL, Thompson RE, Hamilton SR, Griffin CA, Eshleman JR. Detection of
408      microsatellite instability by fluorescence multiplex polymerase chain reaction. *J Mol
409      Diagnostics*. 2000;2(1):20-28. doi:10.1016/S1525-1578(10)60611-3

410  15.  Bakhoum SF, Landau DA. Chromosomal instability as a driver of tumor heterogeneity and
411      evolution. *Cold Spring Harb Perspect Med*. 2017;7(6):1-14.
412      doi:10.1101/cshperspect.a029611

413  16.  Furuya T, Uchiyama T, Murakami T, et al. Relationship between chromosomal instability
414      and intratumoral regional DNA ploidy heterogeneity in primary gastric cancers. *Clin
415      Cancer Res*. 2000;6(7):2815-2820.

416  17.  van den Bosch T, Derks S, Miedema DM. Chromosomal Instability, Selection and
417      Competition: Factors That Shape the Level of Karyotype Intra-Tumor Heterogeneity.
418      *Cancers (Basel)*. 2022;14(20):1-17.

419  18.  Choi EJ, Kim MS, Song SY, Yoo NJ, Lee SH. Intratumoral Heterogeneity of Frameshift
420      Mutations in MECOM Gene is Frequent in Colorectal Cancers with High Microsatellite
421      Instability. *Pathol Oncol Res*. 2017;23(1):145-149.

422  19.  Jo YS, Kim MS, Yoo NJ, Lee SH. Somatic Mutations and Intratumoral Heterogeneity of
423      MYH11 Gene in Gastric and Colorectal Cancers. *Appl Immunohistochem Mol Morphol*.
424      2018;26(8):562-566.

425  20.  Jung M, Lee JA, Yoo SY, Bae JM, Kang GH, Kim JH. Intratumoral spatial heterogeneity of

426      tumor-infiltrating lymphocytes is a significant factor for precisely stratifying prognostic

427      immune subgroups of microsatellite instability-high colorectal carcinomas. *Mod Pathol*.

428      2022;35(12):2011-2022.

429  21.    Wu W, Liu Y, Zeng S, Han Y, Shen H. Intratumor heterogeneity: the hidden barrier to

430      immunotherapy against MSI tumors from the perspective of IFN-γ signaling and tumor-

431      infiltrating lymphocytes. *J Hematol Oncol*. 2021;14(1):1-28.

432  22.    Chapusot C, Martin L, Bouvier AM, et al. Microsatellite instability and intratumoural

433      heterogeneity in 100 right-sided sporadic colon carcinomas. *Br J Cancer 2002 874*.

434      2002;87(4):400-404.

435  23.    Evrard C, Messina S, Sefrioui D, et al. Heterogeneity of Mismatch Repair Status and

436      Microsatellite Instability between Primary Tumour and Metastasis and Its Implications

437      for Immunotherapy in Colorectal Cancers. *Int J Mol Sci*. 2022;23(8).

438  24.    Huang Q, Yu T, Li L, et al. Intraindividual Tumor Heterogeneity of Mismatch Repair Status

439      in Metastatic Colorectal Cancer. *Appl Immunohistochem Mol Morphol*. 2023;31(2):84-93.

440  25.    Luchini C, Mafficini A, Chatterjee D, et al. Histo-molecular characterization of pancreatic

441      cancer with microsatellite instability: intra-tumor heterogeneity, B2M inactivation, and

442      the importance of metastatic sites. *Virchows Arch*. 2022;480(6):1261-1268.

443  26.    Riedinger CJ, Esnakula A, Haight PJ, et al. Characterization of mismatch-

444      repair/microsatellite instability-discordant endometrial cancers. *Cancer*.

445      2024;130(3):385-399.

446  27.    Tachon G, Frouin E, Karayan-Tapon L, et al. Heterogeneity of mismatch repair defect in

447      colorectal cancer and its implications in clinical practice. *Eur J Cancer*. 2018;95:112-116.

448  28.    Heregger R, Huemer F, Steiner M, Gonzalez-Martinez A, Greil R, Weiss L. Unraveling

449      Resistance to Immunotherapy in MSI-High Colorectal Cancer. *Cancers (Basel)*.

450      2023;15(20):1-18.

451    29.    Battaglin F, Naseem M, Lenz HJ, Salem ME. Microsatellite Instability in Colorectal Cancer:
452         Overview of Its Clinical Significance and Novel Perspectives. *Clin Adv Hematol Oncol*.
453         2018;16(11):735.

454    30.    Wang R, Lian J, Wang X, et al. Intrinsic resistance and efficacy of immunotherapy in
455         microsatellite instability-high colorectal cancer: A systematic review and meta-analysis.
456         *Biomol Biomed*. 2023;23(2):198.

457    31.    Evrard C, Tachon G, Randrian V, Karayan-tapon L, Tougeron D. Discordance , and Clinical
458         Impact in Colorectal Cancer. *Cancers (Basel)*. 2019;1-25.

459    32.    Georgiades IB, Curtis LJ, Morris RM, Bird CC, Wyllie AH. Heterogeneity studies identify a
460         subset of sporadic colorectal cancers without evidence for chromosomal or
461         microsatellite instability. *Oncogene*. 1999;18(56):7933-7940.

462    33.    Köster J, Mölder F, Jablonski KP, et al. Sustainable data analysis with Snakemake.
463         *F1000Research*. 2021;10:33.

464    34.    Joanito I, Wirapati P, Zhao N, et al. Single-cell and bulk transcriptome sequencing
465         identifies two epithelial tumor cell states and refines the consensus molecular
466         classification of colorectal cancer. *Nat Genet*. 2022;54(7):963-975.

467    35.    Wu T, Zhang X, Liu X, et al. Single-cell sequencing reveals the immune microenvironment
468         landscape related to anti-PD-1 resistance in metastatic colorectal cancer with high
469         microsatellite instability. *BMC Med*. 2023;21(1):1-18.

470    36.    Li J, Wu C, Hu H, et al. Remodeling of the immune and stromal cell compartment by PD-1
471         blockade in mismatch repair-deficient colorectal cancer. Cancer Cell.

472    37.    Yates J, Kraft A, Boeva V. Filtering cells with high mitochondrial content depletes viable
473         metabolically altered malignant cell populations in cancer single-cell studies. *Genome
474         Biol*. 2025;26(1):1-26.

475    38.    Nofech-Mozes I, Soave D, Awadalla P, Abelson S. Pan-cancer classification of single cells

476        in the tumour microenvironment. *Nat Commun*. 2023;14(1):1-14.

477    39.    Jia P, Yang X, Yang X, Wang T, Xu Y, Ye K. MSIsensor-RNA: Microsatellite Instability

478        Detection for Bulk and Single-cell Gene Expression Data. *Genomics Proteomics*

479        *Bioinformatics*. Published online 2024.

480    40.    Tickle T, Tirosh I, Georgescu C, Brown M, Haas B. inferCNV of the Trinity CTAT Project.

481        Published online 2019. https://github.com/broadinstitute/inferCNV

482    41.    Christopher M, John R. Package "MLeval" Machine Learning Model Evaluation. Published

483        online 2022. https://cran.r-project.org/package=MLeval

484    42.    Kuhn M. Building predictive models in R using the caret package. *J Stat Softw*.

485        2008;28(5):1-26.

486    43.    R Studio Team. A language and environment for statistical computing. *R Found Stat*

487        *Comput*. 2021;3:https://www.R-project.org. http://www.r-project.org

488    44.    Wickham, Hadley, Navarro D. *Ggplot2: Elegant Graphics for Data Analysis*. Vol 35.

489        Springer-Verlag New York; 2010.

490    45.    Zhao P, Li L, Jiang X, Li Q. Mismatch repair deficiency/microsatellite instability-high as a

491        predictor for anti-PD-1/PD-L1 immunotherapy efficacy. *J Hematol Oncol*. 2019;12(1):1-

492        14.

493    46.    Zhao F, Wang S, Bai Y, et al. Cellular MSI-H score: a robust predictive biomarker for

494        immunotherapy response and survival in gastrointestinal cancer. *Am J Cancer Res*.

495        2024;14(11):5551-5567. doi:10.62347/AIWP6518

496    47.    Amemiya K, Hirotsu Y, Nagakubo Y, et al. Simple IHC reveals complex MMR alternations

497        than PCR assays: Validation by LCM and next-generation sequencing. *Cancer Med*.

498        2022;11(23):4479-4490.

499    48.    Zhang Q, Li J, Shen L, Li Y, Wang X. Opportunities and challenges of immunotherapy for

500        dMMR/ MSI-H colorectal cancer. *Cancer Biol Med*. 2023;20(10):1-7.

501    Table 1: Single-cell sequencing datasets

| Dataset ID | Cancer type | Individuals | Samples | Sequencing |
|---|---|---|---|---|
| EGAD00001008555 | Colorectal/metastatic | 15 | 77 | Illumina HiSeq 4000 |
| EGAD00001008584 | Colorectal/metasetatic | 3 | 6 | Illumina HiSeq 4000 |
| EGAD00001008585 | Colorectal/metastatic | 6 | 18 | Illumina NextSeq 500/NovaSeq6000 |
| GSE205506 | Colorectal | 19 | 27 | Illumina NovaSeq 6000/DNBSEQ-T |
| PRJNA932556 | Colorectal | 6 | 6 | HiSeq X Ten |

502    Table 1 legend: The metadata for all single-cell RNA sequencing data used in the study. The

503    dataset column contains the project ID for the European Genome-phenome archive (EGA

504    prefix), GEO Expression Omnibus (GSE prefix), or Sequence Read Archive (PRJ prefix).

505

506

507

508

509

510

511

512

513

514

515

516      Table 2: Individual summary statistics and subclone information

| Individual | Clusters | F | Samples | MSI-H cells | MSS cells | MSS | MSI-H | PCR/IHC | Treatment |
|---|---|---|---|---|---|---|---|---|---|
| CRC2783 | 4 | 17.15 | 3 | 40 | 211 | 5 | 2 | MSI-H | untreated |
| CRC2786 | 10 | 75.67 | 4 | 101 | 2182 | 36 | 4 | MSS | untreated |
| CRC2787 | 2 | 1.62 | 2 | 4 | 121 | 5 | 1 | MSS | untreated |
| CRC2794 | 6 | 9.51 | 4 | 3 | 939 | 28 | 1 | MSS | untreated |
| CRC2795 | 6 | 8.17 | 4 | 0 | 367 | 7 | 0 | MSS | untreated |
| CRC2801 | 7 | 4.66 | 5 | 1 | 1019 | 7 | 0 | MSS | untreated |
| CRC2803 | 8 | 15.91 | 4 | 5 | 619 | 19 | 1 | MSS | untreated |
| CRC2810 | 3 | 1.8 | 4 | 1 | 155 | 7 | 0 | MSS | untreated |
| CRC2811 | 8 | 19.63 | 4 | 1 | 1255 | 7 | 0 | MSS | untreated |
| CRC2816 | 7 | 8.68 | 4 | 2 | 585 | 17 | 1 | MSS | untreated |
| CRC2817 | 7 | 10.62 | 10 | 56 | 442 | 12 | 2 | MSI-H | untreated |
| CRC2821 | 12 | 104.26 | 9 | 240 | 9698 | 41 | 1 | MSS | untreated |
| CRC2829 | 9 | 10.04 | 2 | 0 | 1147 | 7 | 0 | Unknown | untreated |
| CRC2841 | 8 | 93.35 | 11 | 5 | 1392 | 29 | 1 | MSS | untreated |
| CRC2899 | 9 | 20.28 | 4 | 3 | 1444 | 28 | 1 | MSS | untreated |
| P11 | 3 | 6.94 | 1 | 0 | 119 | 7 | 0 | MSI-H | anti-PD-1+celecoxib |
| P12 | 3 | 3.02 | 1 | 11 | 92 | 3 | 1 | MSI-H | anti-PD-1 |
| P14 | 1 | NA | 1 | 1 | 20 | 7 | 0 | MSI-H | anti-PD-1+celecoxib |
| P15 | 3 | 2.25 | 1 | 0 | 139 | 7 | 0 | MSI-H | anti-PD-1 |
| P17 | 1 | NA | 1 | 0 | 36 | 7 | 0 | MSI-H | anti-PD-1 |

| Individual | Clusters | F | Samples | MSI-H cells | MSS cells | MSS | MSI-H | PCR/IHC | Treatment |
|---|---|---|---|---|---|---|---|---|---|
| P18 | 10 | 30.89 | 1 | 58 | 1633 | 35 | 3 | MSI-H | anti-PD-1 |
| P19 | 2 | 1.3 | 1 | 0 | 41 | 7 | 0 | MSI-H | anti-PD-1+celecoxib |
| P21 | 1 | NA | 2 | 3 | 27 | 1 | 1 | MSI-H | anti-PD-1+celecoxib |
| P23 | 10 | 29.07 | 1 | 200 | 1785 | 39 | 8 | MSI-H | untreated |
| P24 | 7 | 75.2 | 2 | 22 | 732 | 17 | 1 | MSI-H | anti-PD-1+celecoxib |
| P25 | 7 | 34.94 | 2 | 13 | 630 | 16 | 1 | MSI-H | anti-PD-1+celecoxib |
| P26 | 6 | 5.09 | 1 | 12 | 427 | 12 | 1 | MSI-H | anti-PD-1+celecoxib |
| P27 | 5 | 11.26 | 2 | 4 | 385 | 10 | 1 | MSI-H | anti-PD-1 |
| P28 | 5 | 51.15 | 2 | 4 | 192 | 5 | 1 | MSI-H | anti-PD-1 |
| P29 | 5 | 10.99 | 1 | 3 | 380 | 12 | 1 | MSI-H | anti-PD-1 |
| P30 | 6 | 26.69 | 2 | 62 | 434 | 12 | 2 | MSI-H | anti-PD-1 |
| P31 | 10 | 18.71 | 2 | 149 | 1833 | 39 | 6 | MSI-H | anti-PD-1 |
| P32 | 7 | 19.41 | 2 | 3 | 370 | 13 | 1 | MSI-H | anti-PD-1+celecoxib |
| P33 | 6 | 15.71 | 1 | 24 | 301 | 8 | 1 | MSI-H | untreated |
| SC024 | 5 | 21.77 | 2 | 4 | 338 | 11 | 1 | MSS | untreated |
| SC027 | 8 | 20.24 | 2 | 3 | 769 | 19 | 1 | MSS | untreated |
| SC029 | 4 | 5.78 | 2 | 0 | 300 | 7 | 0 | MSS | untreated |
| SC035 | 6 | 47.22 | 2 | 37 | 415 | 13 | 1 | MSI-H | untreated |

| Individual | Clusters | F | Samples | MSI-H cells | MSS cells | MSS | MSI-H | PCR/IHC | Treatment |
|---|---|---|---|---|---|---|---|---|---|
| SC040 | 9 | 116.1 | 4 | 82 | 1067 | 23 | 3 | MSS | untreated |
| SC041 | 5 | 10.01 | 2 | 32 | 279 | 13 | 1 | MSS | untreated |
| SC042 | 3 | 31.24 | 2 | 0 | 141 | 7 | 0 | Unknown | untreated |
| SC043 | 7 | 6.75 | 2 | 0 | 813 | 7 | 0 | MSS | untreated |
| SC044 | 8 | 32.87 | 3 | 58 | 412 | 9 | 2 | MSI-H | untreated |
| SRR23490337 | 1 | NA | 1 | 3 | 31 | 1 | 1 | MSI-H | tislelizumab |
| SRR23490338 | 2 | 7.66 | 1 | 11 | 70 | 3 | 1 | MSI-H | tislelizumab |
| SRR23490339 | 1 | NA | 1 | 0 | 18 | 7 | 0 | MSI-H | tislelizumab |
| SRR23490340 | 4 | 27.42 | 1 | 24 | 190 | 5 | 1 | MSI-H | tislelizumab |
| SRR23490341 | 3 | 48.02 | 1 | 14 | 74 | 3 | 1 | MSI-H | tislelizumab |
| SRR23490342 | 1 | NA | 1 | 0 | 48 | 7 | 0 | MSI-H | tislelizumab |

517 Table 2 legend: This table contains the summary statistics on each individual included in the

518 analysis. The clusters column refers to the number of unique cancer clusters, and F is the

519 ANOVA F-statistic used to measure heterogeneity in those clusters. The samples column

520 describes the number of samples each individual had for the analysis, and the MSI-H and MSS

521 cells column is the number of cell type for that individual. The MSS and MSI-H columns refer to

522 the number of microsatellite stable and microsatellite instability high subclones for each

523 individual. The original IHC/PCR status and treatment metadata for each individual is included

524 and were established in previous studies (Table 1). Any NA values represent F-statistics that

525 could not be calculated due to fewer than 2 clusters of cancer cells being present.
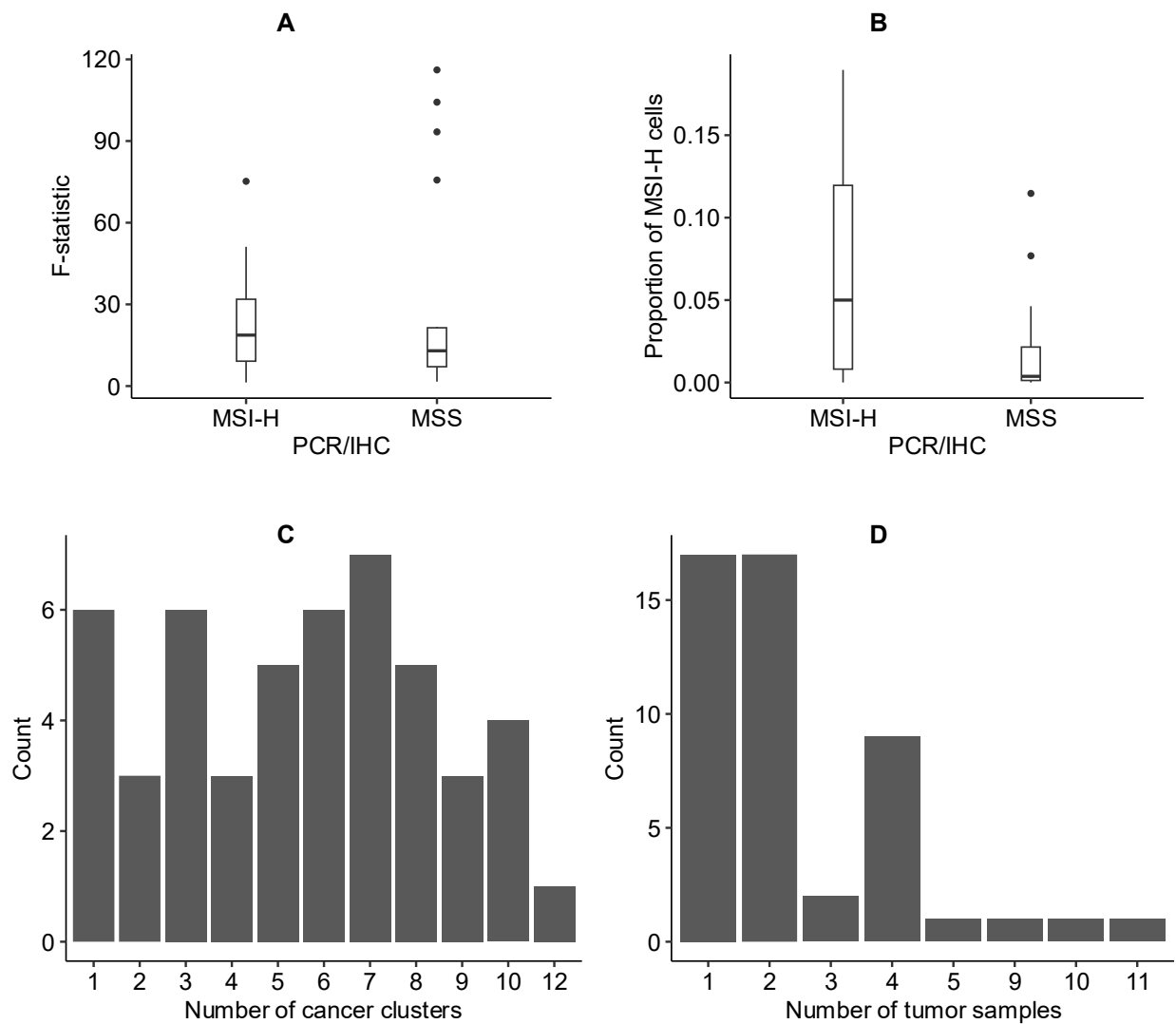
526

527

528

529    Figure 1



530    Figure 1 legend: Box plots showing the distribution of (A) MSI score for individuals calculated
531    using the aggregate expression of all cells, and (B) MSI score for individuals calculated using the
532    aggregate expression of only cancer cells. Also shown are the mean values of (C) the F-statistic
533    and (D) the number of subclones for the different cell mixes shown on the x-axes (with
534    increasing proportions of MSI-H cells ranging from 0.1 in mix M1 to 0.9 in mix M9). The error
535    bars in (C) and (D) correspond to plus/minus twice the standard error around the mean. The
536    MSS and MSI-H samples in panels C and D are the obtained values for all cells in those samples
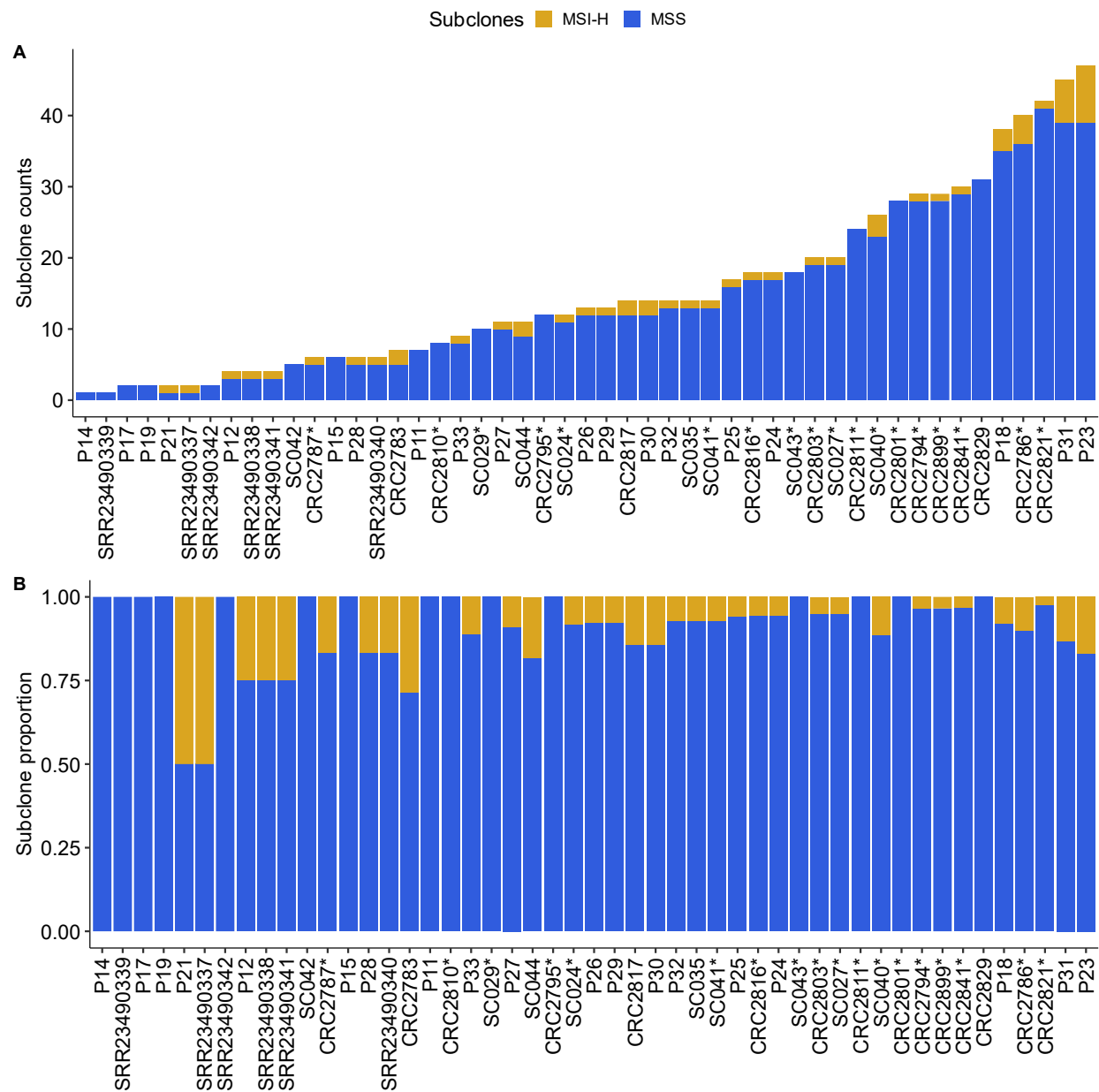537    and do not represent an average.

538

539    Figure 2



540

541    Figure 2 legend: Box plots showing the distribution of (A) F-statistics grouped by PCR/IHC MSI

542    status, (B) the proportion of MSI-H to MSS cells grouped by PCR/IHC MSI status. Also shown are

543    histograms displaying the frequency of (C) the number of cancer cell clusters and (D) the

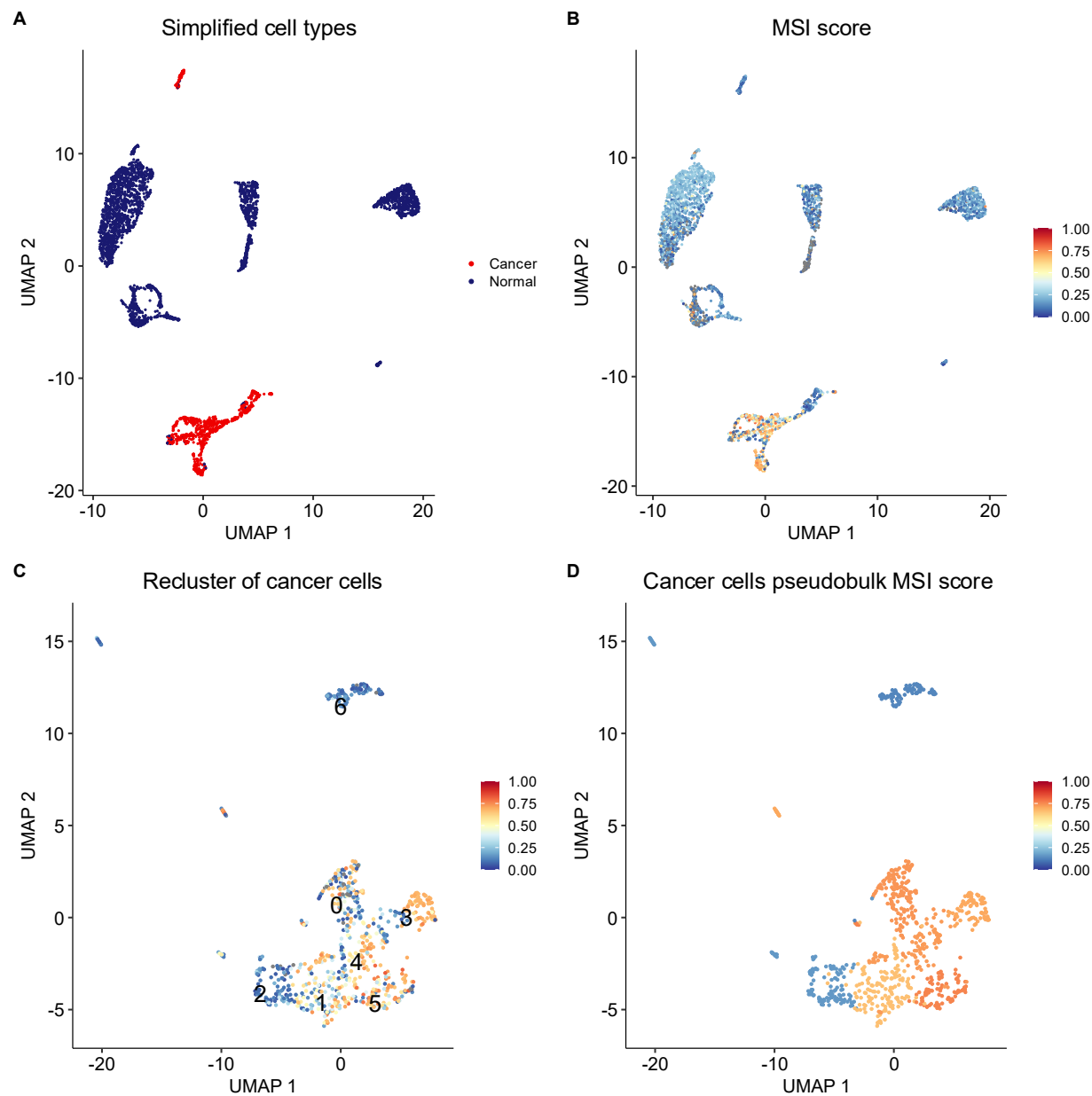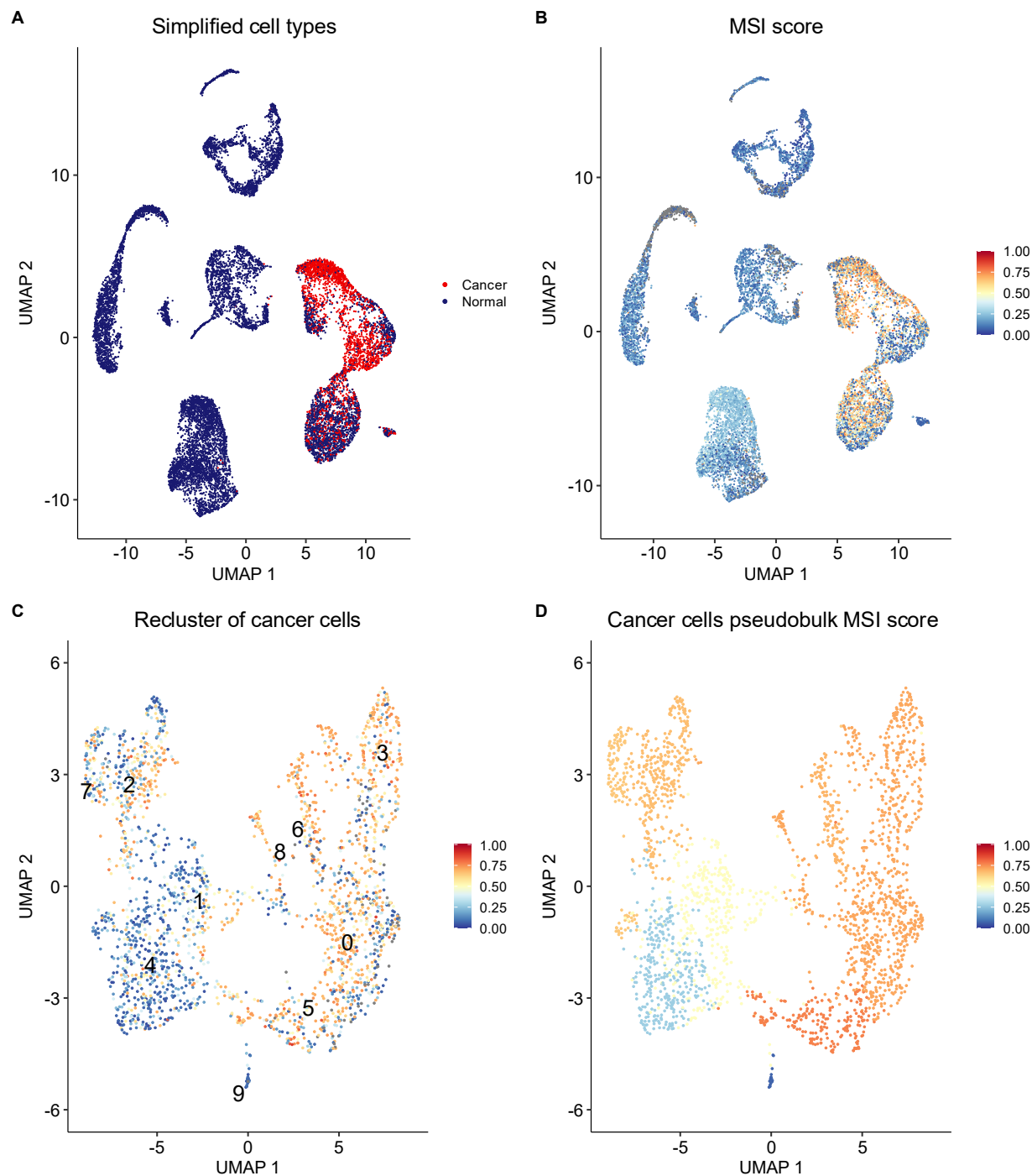544    number of tumor samples for all individuals.

545

546

Figure 3



Figure 3 legend: Stacked bar plots of (A) the number of subclones for each individual in the analysis and (B) the proportion of subclone types for each individual. Individuals that had a PCHR/IHC test result of MSS are indicated with an asterisk.

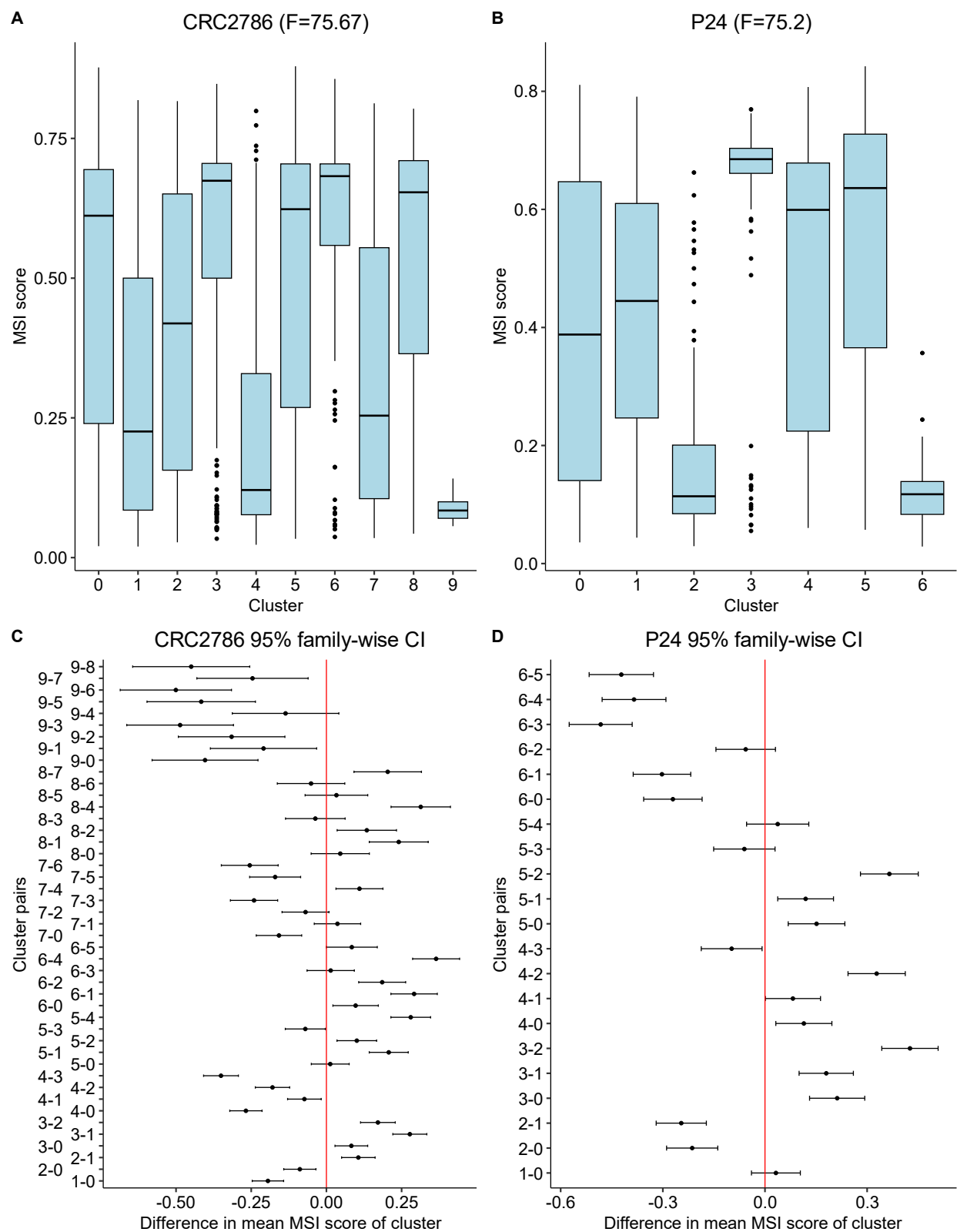554    Figure 4



555

556    Figure 4 legend: UMAP plots for MSI-H individual P24 showing (A) tumor versus normal cell

557    classification, (B) MSI scores for each cell, (C) MSI scores for re-clustered cancer cells, and (D)

558    MSI score for aggregated pseudobulk expression of each cancer cell cluster.

559

560    Figure 5



561

562    Figure 5 legend: UMAP plots for MSS individual CRC2786 showing (A) tumor versus normal cell

563    classification, (B) MSI scores for each cell, (C) MSI scores for re-clustered cancer cells, and (D)

564    MSI score for aggregated pseudobulk expression of each cancer cell cluster.

565    Figure 6



Figure 6 legend: Box plots showing the distribution of MSI scores for each cluster of cancer cells

567    in (A) individual CRC2786 and (B) individual P24. Also shown are the 95% confidence intervals

568    for the difference in mean MSI scores between each cluster pair for (C) individual CRC2786 and

569    (D) individual P24.