

1 **Evaluating multiple next-generation sequencing derived tumor features to accurately**
2 **predict DNA mismatch repair status**

3 Romy Walker^{1,2}, Peter Georgeson^{1,2}, Khalid Mahmood^{1,2,3}, Jihoon E. Joo^{1,2}, Enes Makalic⁴, Mark
4 Clendenning^{1,2}, Julia Como^{1,2}, Susan Preston^{1,2}, Sharelle Joseland^{1,2}, Bernard J. Pope^{1,3}, Ryan
5 Hutchinson^{1,2}, Kais Kasem⁵, Michael D. Walsh⁶, Finlay A. Macrae^{7,8}, Aung K. Win^{2,4}, John L.
6 Hopper⁴, Dmitri Mouradov^{9,10}, Peter Gibbs^{9,10,11}, Oliver M. Sieber^{9,10,12,13}, Dylan E.
7 O’Sullivan^{14,15}, Darren R. Brenner^{14,15,16}, Steven Gallinger^{17,18,19}, Mark A. Jenkins^{2,4}, Christophe
8 Rosty^{1,2,20,21}, Ingrid M. Winship^{7,22}, Daniel D. Buchanan^{1,2,7#}

9
10 ¹ Colorectal Oncogenomics Group, Department of Clinical Pathology, Victorian Comprehensive
11 Cancer Centre, The University of Melbourne, Parkville, Victoria, Australia

12 ² University of Melbourne Centre for Cancer Research, Victorian Comprehensive Cancer Centre,
13 Parkville, Victoria, Australia

14 ³ Melbourne Bioinformatics, The University of Melbourne, Melbourne, Victoria, Australia

15 ⁴ Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health,
16 The University of Melbourne, Carlton, Victoria, Australia

17 ⁵ Department of Clinical Pathology, Medicine Dentistry and Health Sciences, The University of
18 Melbourne, Parkville, Victoria, Australia

19 ⁶ Sullivan Nicolaides Pathology, Bowen Hills, Queensland, Australia

20 ⁷ Genomic Medicine and Family Cancer Clinic, Royal Melbourne Hospital, Parkville, Melbourne,
21 Victoria, Australia

22 ⁸ Colorectal Medicine and Genetics, The Royal Melbourne Hospital, Parkville, Victoria,
23 Australia

24 ⁹ Personalized Oncology Division, The Walter and Eliza Hall Institute of Medial Research,
25 Parkville, Victoria, Australia

26 ¹⁰ Department of Medical Biology, The University of Melbourne, Parkville, Victoria, Australia

27 ¹¹ Department of Medical Oncology, Western Health, Victoria, Australia

28 ¹² Department of Surgery, The University of Melbourne, Parkville, Victoria, Australia

29 ¹³ Department of Biochemistry and Molecular Biology, Monash University, Clayton, Victoria,
30 Australia

31 ¹⁴ Department of Oncology, University of Calgary, Calgary, Canada

32 ¹⁵ Department of Community Health Sciences, University of Calgary, Calgary, Canada

33 ¹⁶ Department of Cancer Epidemiology and Prevention Research, Alberta Health Services,
34 Calgary, Canada

35 ¹⁷ Ontario Institute for Cancer Research, Toronto, Ontario, Canada

36 ¹⁸ Department of Laboratory Medicine and Pathobiology, University of Toronto, Toronto,
37 Ontario, Canada

38 ¹⁹ Lunenfeld Tanenbaum Research Institute, Mount Sinai Hospital, University of Toronto,
39 Toronto, Ontario, Canada

40 ²⁰ Envoi Specialist Pathologists, Brisbane, Australia

41 ²¹ University of Queensland, Brisbane, Australia

42 ²² Department of Medicine, The University of Melbourne, Parkville, Australia

43

44

45 **Running Title: NGS tumor mismatch-repair deficiency**

46

47 #To whom correspondence should be addressed:

48

49 Associate Professor Daniel D. Buchanan

50 Head, Colorectal Oncogenomics Group

51 Department of Clinical Pathology

52 The University of Melbourne

53 Victorian Comprehensive Cancer Centre

54 305 Grattan Street

55 Parkville, Victoria, 3010 Australia

56 Ph: +61 385597004

57 Email: daniel.buchanan@unimelb.edu.au

58

59

60

61 **Number of text pages:** 53 pages

62 **Figures & Tables:** 4 figures and 4 tables

63 **Running head:** 41 (with characters)

64 **Abstract:** 220 words

65 **References:** 76 references

66

67 **Abstract**

68 Identifying tumor DNA mismatch repair deficiency (dMMR) is important for precision medicine.
69 We assessed tumor features, individually and in combination, in whole-exome sequenced (WES)
70 colorectal cancers (CRCs) and in panel sequenced CRCs, endometrial cancers (ECs) and
71 sebaceous skin tumors (SSTs) for their accuracy in detecting dMMR. CRCs (n=300) with WES,
72 where MMR status was determined by immunohistochemistry, were assessed for microsatellite
73 instability (MSMuTect, MANTIS, MSIseq, MSISensor), COSMIC tumor mutational signatures
74 (TMS) and somatic mutation counts. A 10-fold cross-validation approach (100 repeats) evaluated
75 the dMMR prediction accuracy for 1) individual features, 2) Lasso statistical model and 3) an
76 additive feature combination approach. Panel sequenced tumors (29 CRCs, 22 ECs, 20 SSTs) were
77 assessed for the top performing dMMR predicting features/models using these three approaches.
78 For WES CRCs, 10 features provided >80% dMMR prediction accuracy, with MSMuTect,
79 MSIseq, and MANTIS achieving $\geq 99\%$ accuracy. The Lasso model achieved 98.3%. The additive
80 feature approach with $\geq 3/6$ of MSMuTect, MANTIS, MSIseq, MSISensor, INDEL count or TMS
81 ID2+ID7 achieved 99.7% accuracy. For the panel sequenced tumors, the additive feature
82 combination approach of $\geq 3/6$ achieved accuracies of 100%, 95.5% and 100%, for CRCs, ECs,
83 and SSTs, respectively. The microsatellite instability calling tools performed well in WES CRCs,
84 however, an approach combining tumor features may improve dMMR prediction in both WES and
85 panel sequenced data across tissue types.

86

87 **Keywords:** Colorectal cancer, DNA mismatch repair deficiency, endometrial cancer, Lynch
88 syndrome, microsatellite instability, *MLH1* promoter methylation, sebaceous skin tumor, tumor
89 mutation burden, tumor mutational signatures

90

91 **Declared conflicts of interest**

92 The authors have no conflicts of interest to declare.

93

94 **Data availability statement:** The datasets generated during and/or analyzed during the current
95 study are available from the corresponding author on reasonable request.

96

97 **Funding**

98 Funding by a National Health and Medical Research Council of Australia (NHMRC) project grant
99 GNT1125269 (PI- Daniel Buchanan), supported the design, analysis, and interpretation of data.
100 RW is supported by the Margaret and Irene Stewardson Fund Scholarship and by the Melbourne
101 Research Scholarship. DDB is supported by an NHMRC Investigator grant (GNT1194896) and
102 University of Melbourne Dame Kate Campbell Fellowship. PG is supported by the University of
103 Melbourne Research Scholarship. MAJ is supported by an NHMRC Investigator grant
104 (GNT1195099). AKW is supported by an NHMRC Investigator grant (GNT1194392). JLH is
105 supported by the University of Melbourne Dame Kate Campbell Fellowship. OMS is supported
106 by an NHMRC Senior Research Fellowship (GNT1136119). DEO is supported by a Canadian
107 Institutes of Health Research (CIHR) Post-doctoral Fellowship. BP is supported by a Victorian
108 Health and Medical Research Fellowship from the Victorian Government.

109

110 Introduction

111 DNA mismatch-repair (MMR) deficiency (dMMR) is an important molecular phenotype of
112 solid tumors characterized by the presence of microsatellite instability (MSI) and/or loss of
113 expression of one or more of the DNA MMR proteins, MLH1, MSH2, MSH6 and PMS2.
114 Identifying dMMR tumors is important for understanding disease prognosis¹, response to immune
115 checkpoint inhibition therapy² and to identify people with Lynch syndrome. Lynch syndrome is
116 the most common inherited cancer predisposition disorder and, therefore, the Evaluation of
117 Genomic Applications in Practice and Prevention Working Group recommends that all newly
118 diagnosed colorectal (CRC) and endometrial cancers (EC) are screened for dMMR to improve the
119 identification of carriers^{3,4}.

120 The dMMR mutator phenotype arises in tumors where errors occur during the DNA
121 replication process⁵. Specifically, defects in the components of the MMR system responsible for
122 the recognition of mismatches such as single nucleotide variants (SNVs) and insertion-deletions
123 (INDELs), can lead to the development of numerous frameshift mutations in coding and non-
124 coding microsatellite regions⁶. dMMR is related to biallelic inactivation of one of the MMR genes,
125 resulting from either somatic methylation of the *MLH1* gene promoter region⁷ or double somatic
126 MMR gene mutations⁸ (sporadic dMMR), or germline pathogenic variants in the MMR genes⁹ or
127 deletions in the 3' end of the *EPCAM* gene¹⁰ (inherited dMMR). CRC, EC and sebaceous skin
128 tumors (SSTs), including sebaceous adenomas, carcinomas and sebaceomas, are tissue types that
129 demonstrate the highest frequencies of dMMR where up to 26%¹¹, 31%¹¹ and 31%¹² of these tissue
130 types, respectively, present with the dMMR phenotype, followed by stomach cancer at 19%¹¹.

131 The most common approach for identifying dMMR tumors is by assessing MMR protein
132 expression through immunohistochemistry (MMR IHC)^{13,14} and/or by testing for high levels of

133 microsatellite instability using polymerase chain reactions (MSI-PCR)¹⁵. While both screening
134 methodologies are commonly used, each present advantages and limitations. The advantages of
135 performing MMR IHC include simple experimental execution, short turnaround time, low
136 associated costs as well as giving an indication of the defective gene¹⁶. However, false positive or
137 false negative MMR IHC results can occur due to technical artefacts, variable performance of
138 different MMR antibodies and inherent variability in the interpretation of the staining by different
139 pathologists¹⁶. Further challenges include the interpretation of weaker staining in less proliferative
140 tissue and heterogenous patterns of MMR protein loss¹⁷⁻²⁴.

141 While MMR IHC is more widely adopted in the clinical setting, MSI-PCR remains the gold
142 standard for detecting dMMR¹⁶; to date multiple markers have been identified to call MSI in tumor
143 samples²⁵. The limitations for MSI-PCRs include additional laboratory implementation
144 requirements related to tissue DNA extraction and increased labor costs; both can lead to a delay
145 in receiving test results¹⁶. Nonetheless, MMR IHC and MSI-PCR methodologies have proven to
146 be effective for identifying dMMR in CRC samples²⁶ with a reported concordance of 91.9%¹⁶, but
147 the accuracy for either of these tools can decrease when applied to different tissue types²⁷. As next-
148 generation sequencing (NGS) becomes more widely adopted for precision oncology, there is an
149 increasing need to accurately determine tumor MMR status using NGS data.

150 To date, several tools have been developed to assess MSI from NGS data, including
151 MSISensor²⁸, MSIseq²⁹, MANTIS³⁰ and more recently MSMuTect³¹. To the best of our
152 knowledge, the comparison of these four MSI tools on the same tumors has not yet been performed.
153 In addition to MSI, other tumor features derived from NGS have been shown to be associated with
154 dMMR, such as tumor mutational burden (TMB)³² and tumor mutational signatures (TMS)³³.

155 TMB, characterized by high SNV and INDEL counts, is a biomarker for response to immune
156 checkpoint inhibition therapy^{34,35} and is increased in dMMR tumors³⁶.

157 TMS aggregate tens to thousands of the observed somatic mutations within a tumor into
158 patterns related to the underlying mutational processes^{37,38}. The predominant TMS framework,
159 published on the COSMIC website, defines 107 different signature definitions categorized into
160 three distinct subgroups: 1) 78 single base substitutions (SBS) where seven of the SBS signatures
161 (SBS6, SBS14, SBS15, SBS20, SBS21, SBS26 and SBS44) are associated with dMMR; 2) 18
162 small (1 to 50 base pair) insertions and deletions or ID signatures where ID1, ID2, and ID7 are
163 associated with dMMR, and 3) 11 doublet base substitutions or DBS signatures where DBS7 and
164 DBS10 have both been previously associated with dMMR³³. However, DBS signatures have a
165 reported low prevalence in CRC compared with other tissue types so were excluded from our
166 study³⁸. Previously, we have shown that the combination of individual TMS can improve the
167 ability of TMS to discriminate important molecular and genetic subtypes of CRC, including
168 identifying germline biallelic carriers of pathogenic variants in the *MUTYH* gene by combining
169 SBS18 and SBS36^{39,40}. We further observed that the combination of ID2 with ID7 (TMS ID2+ID7)
170 was the most informative for differentiating dMMR from pMMR CRCs amongst all possible TMS
171 combinations³⁹. To date, the comparison of MSI calling tools, somatic mutation counts, TMB and
172 TMS tumor features for determining the dMMR status in CRC tumors has not yet been undertaken.

173 In this study, we assessed 104 tumor features derived from whole-exome sequencing
174 (WES) (**Table 1**), consisting of the MSI prediction tools (MSMuTect, MANTIS, MSIseq and
175 MSISensor), TMS (78 SBS and 18 ID signatures), TMS ID2+ID7, TMB and individual SNV and
176 INDEL somatic mutation counts for their accuracy in predicting dMMR status in 300 well-
177 characterized CRCs. Secondly, we investigated whether a combination of these tumor features,

178 using either a statistical model or a simple approach that added individual features together
179 (additive feature combination), could improve the dMMR prediction accuracy in WES CRC
180 tumors. Finally, we evaluated the effectiveness of the top performing tumor features from the WES
181 analysis, individually and in combination, in an independent set of CRC, EC and SST tumors that
182 had undergone targeted multigene panel sequencing for their dMMR prediction accuracy.

183

184 **Materials and Methods**

185

186 ***Study Cohort***

187 The study population included men and women retrospectively identified from five studies
188 where pMMR or dMMR status was determined by MMR IHC and where an etiology for dMMR
189 status could be defined, namely a sporadic etiology caused by tumor *MLH1* methylation or double
190 somatic MMR mutations, or an inherited etiology caused by a germline MMR gene pathogenic
191 variant (Lynch syndrome). The breakdown of participants included in this study by their dMMR
192 and pMMR status, tissue type and by WES or panel sequencing is shown in **Figure 1**:

193 **1)** the ANGELS study (*Applying Novel Genomic approaches to Early-onset and suspected Lynch*
194 *Syndrome colorectal and endometrial cancers*)³⁹ recruited participants that were diagnosed with
195 CRC or EC between 2014 – 2021 who were referred from family cancer clinics across Australia
196 (n=79). All ANGELS study participants provided informed consent and the study was approved
197 by the University of Melbourne human research ethics committee (HREC#1750748) and
198 institutional review boards at each family cancer clinic;

199 **2)** CRC- or EC-affected participants from the ACCFR (*Australasian Colorectal Cancer Family*
200 *Registry*) were selected from both population-based and clinic-based recruitment (n=139);

201 **3)** CRC-affected participants from the OFCCR (*Ontario Familial Colorectal Cancer Registry*)
202 were population-based patients (<50 years old) recruited from the Cancer Care Ontario, Toronto,
203 Canada (n=53). Study participants from both the ACCFR and OFCCR were recruited between
204 1998 and 2008, and were included according to the recruitment policy and eligibility criteria
205 previously described^{41,42}. Informed consent was obtained from all study participants and the study
206 protocol was approved by the institutional human ethics committee at both study sites;

207 **4)** CRC-affected participants from the WEHI study (*Walter and Eliza Hall Institute of Medical*
208 *Research*) were recruited from the Royal Melbourne Hospital (Parkville, VIC, Australia) and the
209 Western Hospital Footscray (Footscray, VIC, Australia), between Jan 1, 1993, and Dec 31, 2009³⁹.
210 All patients provided written informed consent. The study was approved by human research ethics
211 committees at both sites (HREC 12/19) (n = 80);

212 **5)** SST-affected participants from the MTS study (*Muir-Torre Syndrome Study*) were referred
213 between July 2016 and September 2021 following clinical diagnostic MMR IHC testing by
214 Sullivan Nicolaides Pathology service in Brisbane¹² or by family cancer clinics in Australia.
215 Informed consent was obtained from the study participants and the study protocol was approved
216 by the human research ethics committee from the University of Melbourne (HREC#1648355) and
217 by the relevant institutional human ethics committees (n = 20).

218

219 ***Tumor Categorization***

220 MMR IHC testing was performed on formalin-fixed paraffin embedded (FFPE) tissues for
221 all four MMR proteins for the ACCFR and OFCCR as previously described⁴²⁻⁴⁴, and a subset of
222 these tumors also underwent MSI-PCR testing as previously described⁴⁵. MMR IHC testing for
223 the ANGELS and MTS studies was part of routine clinical assessment in pathology laboratories

224 across Australia, reported by the duty pathologist. Fresh-frozen tissue specimens from the WEHI
225 study were assessed for MLH1, MSH2 and MSH6 MMR IHC and MSI-PCR tested using BAT25,
226 BAT26, D5S346, D2S123 and D17S250 MSI markers. Germline MMR gene testing (as described
227 in Buchanan *et al.*⁴³) and tumor *MLH1* promoter methylation testing by MethyLight (as described
228 in Buchanan *et al.*⁴⁶) were performed on all dMMR tumors showing loss of MLH1/PMS2 protein
229 expression or sole PMS2 loss by IHC. Tumors were considered to have double somatic MMR
230 mutations when they were found to have two pathogenic/likely pathogenic somatic mutations or a
231 single somatic pathogenic/likely pathogenic mutation in combination with presence of loss of
232 heterozygosity. Germline pathogenic variants and somatic MMR gene mutations were confirmed
233 in WES and targeted panel sequencing data prior to analysis. Therefore, for each of the dMMR
234 tumors included in this study we could confirm an inherited or acquired cause for their respective
235 pattern of MMR IHC protein loss. Concurrently, for the pMMR tumors, we did not find evidence
236 of a germline MMR pathogenic variant or double MMR somatic mutation in these tumor samples.

237 All tumors in the study were assigned to one of four categories based on dMMR or pMMR
238 status determined from MMR IHC and/or MSI-PCR and based on the cause for dMMR:

- 239 **1) dMMR-Lynch syndrome (dMMR-LS)** – identified carrier of a germline pathogenic variant
240 in one of the DNA MMR genes where the corresponding tumor showed commensurate loss of
241 MMR protein expression by IHC;
- 242 **2) dMMR-*MLH1* methylation (dMMR-*MLH1me*)** – tumors were positive for methylation of
243 the *MLH1* gene promoter “C region”⁴⁷ and showed loss of MLH1 and PMS2 protein expression
244 by IHC without a germline MMR gene pathogenic variant;

245 **3) dMMR-double somatic (dMMR-DS)** – tumors harbored two somatic mutations (SNVs and/or
246 loss of heterozygosity) in the same MMR gene that showed loss of protein expression by IHC with
247 no identified pathogenic germline MMR gene variant; and

248 **4) MMR-proficient (pMMR)** – tumors showed normal expression of all four MMR proteins and
249 did not show presence of double somatic MMR gene mutations or a germline MMR gene
250 pathogenic variant.

251 The three dMMR subtypes dMMR-LS, dMMR-DS and dMMR-MLH1me were combined as a
252 single dMMR tumor group in downstream analysis.

253

254 ***Whole-Exome and Targeted Panel Sequencing Capture Regions***

255 The targeted panel was based on the design described in Zaidi *et al.*⁴⁸ consisting of probes
256 targeting the following regions: 1) 298 genes incorporating key hereditary CRC^{49–51} and EC⁵² risk
257 genes and genes that are frequently mutated as identified by The Cancer Genome Atlas (TCGA)
258 data^{32,53,54}, 2) 28 microsatellite loci including the five ‘gold standard’ MSI markers (BAT25,
259 BAT26, NR-21, NR-24, and MONO-27) currently implemented in routine MSI-PCR diagnostics,
260 3) 212 homopolymer regions distributed genome-wide to assess for MSI in tumor samples and 4)
261 56 copy number variants known to be susceptible to copy number changes in CRCs. The panel
262 capture was 2.005 megabases (Mb) in size. The WES capture incorporates all exonic regions
263 within the genome and is 67.296 Mb in size. The panel additionally included capture of intronic
264 regions within the MMR genes, which the WES capture did not cover.

265 *Next-Generation Sequencing*

266 In total, 300 CRC tumors were sequenced by WES and 71 tumors (29 CRCs, 22 ECs and
267 20 SSTs) were sequenced by the targeted multigene panel (**Figure 1**). FFPE CRC, EC or SST
268 tissues were macrodissected and DNA extracted using the QIAmp DNA FFPE Tissue Kit (Qiagen,
269 Hilden, Germany) according to the manufacturer's instructions. Peripheral blood-derived DNA
270 was extracted using the DNeasy blood and tissue kit (Qiagen, Hilden, Germany) and sequenced as
271 germline references.

272 The WES capture was the Agilent Clinical Research Exome V2 kit (Agilent Technologies
273 Santa Clara, United States) with sequencing performed on an Illumina NovaSeq 6000 comprising
274 150 base pair (bp) paired-end reads performed at the Australian Genome Research Facility³⁹. For
275 the WEHI CRCs, exome-enrichment was performed using the TruSeq Exome Enrichment Kit
276 (Illumina, San Diego, United States) and 100 bp paired-end read sequencing performed on an
277 Illumina HiSeq 2000 at the Australian Genome Research Facility³⁹. The on-target coverage for the
278 300 WES samples had a median of 323.7 for the FFPE tumor DNA samples and 137.4 for blood-
279 derived DNA samples, with an interquartile range of 111.8 – 426.4 and 100.6 – 204.9, respectively.

280 Library preparation for targeted panel sequencing was performed using the SureSelect™
281 Low Input Target Enrichment System (Agilent Technologies, Santa Clara, United States) using
282 standard protocol and sequenced on an Illumina NovaSeq 6000 comprising 150 bp paired end
283 reads performed at the Australian Genome Research Facility. The on-target coverage for the 71
284 panel sequenced samples was (median and interquartile range) 919.3 and 694.6 – 1164.9 for FFPE
285 tumor DNA samples and 160.6 and 135.8 – 178.0 for blood-derived DNA samples.

286 *Bioinformatics Pipeline*

287 For both WES and targeted panel sequenced samples, adapter sequences were trimmed
288 from raw FASTQ files using trimmomatic 0.38⁵⁵ and aligned to the GRCh37 human reference
289 genome using Burrows-Wheeler Aligner v. 0.7.12. Germline variants, somatic variants (SNVs)
290 and somatic INDELS were called using Strelka (v. 2.9.2., Illumina) using the recommended
291 workflow⁵⁶. TMS were calculated using the pre-defined set of 78 SBS and 18 ID signatures
292 published on COSMIC as version 3.2 (COSMIC, <https://cancer.sanger.ac.uk/signatures/>, last
293 accessed date: June 15, 2022)³³. Variants outside the WES and panel capture regions were
294 excluded and variants with the PASS filter called from Strelka were retained. Additional variant
295 filters included were restrictions to a minimum depth of 50x for germline and tumor samples with
296 a minimum variant allele frequency of 10% as detailed previously³⁹.

297

298 *Selection of Features of Interest*

299 The 104 tumor features selected for analysis in this study are shown in **Table 1**. Several
300 tools have been developed to assess MSI from NGS data. Our analysis focused on MSMuTect³¹,
301 MANTIS³⁰, MSIseq²⁹ and MSISensor²⁸. Tumors were classified as having high levels of MSI
302 (MSI-H) or as microsatellite stable (MSS). We assessed all SBS (n=78) and ID (n=18) TMS as
303 described by COSMIC³³, but the DBS TMS were excluded due to their reported low prevalence
304 in CRCs³⁸. Previously, we have shown that combining ID2 and ID7 TMS enabled detection of
305 dMMR CRCs³⁹ and, therefore, was included as a tumor feature in this study. Somatic mutation
306 counts, namely SNVs or INDELS, as well as TMB (SNV and INDEL mutation count combined /
307 Mb) were each included, given previous associations with tumor dMMR status⁵⁷.

308

309 *Feature Performance Evaluation in WES data from CRCs*

310 We assessed the 104 tumor features calculated from WES from 209 pMMR CRCs and 91
311 dMMR CRCs (pMMR:dMMR ratio = 2.3:1) (**Figure 1**). The dMMR CRCs comprised dMMR-LS
312 tumors (n=49), dMMR-*MLH1*me tumors (n=26) and dMMR-DS tumors (n=16). All 300 CRCs
313 were randomly partitioned into a training set (80% of CRCs) and a test set (20% of CRCs), while
314 maintaining the same pMMR:dMMR ratio, using *caret* R package⁵⁸. We performed a 10-fold cross
315 validation approach on the training set (repeated 100x) to calculate the average classification
316 accuracy by fitting a generalized linear model and determining the error rate, specificity,
317 sensitivity, and the area under the curve (AUC) with corresponding 95% confidence intervals
318 (CIs). Based on the unequal distribution of dMMR and pMMR tumors in the WES dataset, the *no*
319 *information rate* was 69.5%, indicating that any feature with this prediction accuracy was
320 equivalent to selecting a dMMR sample by chance.

321

322 Tumor feature analysis of the WES CRC dataset comprised of three different approaches:

323 *A) Individual tumor feature assessment*

324 Each of the 104 tumor features were assessed individually and then ranked by their accuracy in
325 identifying dMMR tumors. Individual CRC tumor features with a prediction accuracy >80% from
326 the WES data were considered good predictors for differentiating dMMR from pMMR tumors and
327 were included in downstream analyses.

328

329 *B) Generation of a statistical model by combining tumor features*

330 We investigated whether combining tumor features using a Lasso penalized regression model⁵⁹
331 could improve the overall dMMR prediction accuracy in CRC. Lasso enables the simultaneous

332 parameter estimation and variable selection as well as having been shown to reduce overfitting
333 when compared to conventional maximum likelihood regression models. Lasso regression has a
334 tuning parameter called lambda that controls which features are included in the regression model
335 by shrinking the coefficient or “weighting” of individual features within the model towards zero,
336 helping with the exclusion of some of the features from integration into the final model via a
337 penalization process using cross-validation.

338

339 *C) Applying an additive feature combination count*

340 Our third approach investigated combining the top ranked individual tumor features in an additive
341 approach (additive feature combination). Specifically, the tumor features that achieved a mean
342 prediction accuracy >95% from the WES CRC analysis (from part A), were included in this
343 approach and added together to give an overall count. The bimodal distribution supported a
344 majority vote decision on dMMR status.

345

346 *Assessment of individual tumor features, the statistical model and additive feature combination* 347 *approaches derived from the WES analysis on panel sequenced CRCs, ECs, and SSTs*

348 The top individual tumor features determined from (A), best performing Lasso model (B)
349 and the additive feature combination approach (C) were then assessed for their dMMR prediction
350 accuracy in three independent tumor sets comprised of n=29 CRCs, n=22 ECs and n=20 SSTs
351 tested by targeted multigene panel sequencing. The *no information rate* for features analyzed from
352 the panel dataset was at 71.8%, indicating a prediction accuracy of this value was similar to
353 selecting a dMMR sample by chance.

354

355 **Statistical Analysis**

356 All statistical analyses were done using the R programming language (v.4.1.0). The *tidyverse*
357 package (v.1.3.1.)⁶⁰ was used for data import, tidying and visualization purposes and the *caret*
358 (v.6.0-9.0) package⁵⁸ was used for cross-validation. Receiving operator curves (ROC) were
359 generated using the *pROC* package (v.1.18.0)⁶¹, with the AUC being determined using the *cvAUC*
360 package (v.1.1.4)⁶². Statistical models were fitted using the Lasso (*glmnet*, v.4.1-3)⁶³ package. We
361 used the *cutpointr* (v.1.1.1) package⁶⁴ for estimation of the best “cut points” or “thresholds” which
362 maximize the Youden-index (true positive rate minus false positive rate over all possible cut
363 points), defined as the most optimal threshold in binary disease classification tasks. Here, the
364 *cutpointr* package determines a recommended threshold that best differentiates dMMR from
365 pMMR cases for each feature and validates its performance using bootstrapping. The average
366 weight for each group was calculated using the *plyr* (v.1.0.7) package⁶⁵. The *ggplot2* (v.3.3.5)
367 package⁶⁶ was used for data visualization in combination with *hrbrthemes* (v.0.8.0)⁶⁷ for histogram
368 generation and *ggrepel* (v.0.9.1)⁶⁸ for histogram annotations. Correlation scores between the
369 dMMR and pMMR groups were estimated by a *heteroscedastic two-tailed t-test*. P-values <0.05
370 were considered statistically significant. The 95% CIs for the WES data were calculated using the
371 binomial (Clopper-Pearson) “exact” method⁶⁹ and for the targeted panel data using the *binom*
372 (v.1.1-1) package⁷⁰ in R.

373

374 **Results**

375 For the initial performance evaluation of 104 tumor features we assessed 209 (69.7%) pMMR
376 CRCs and 91 (30.3%) dMMR CRCs sequenced by WES. The clinicopathological characteristics,
377 pattern of MMR IHC loss and dMMR etiology are summarized in **Supplementary Table 1**. The

378 mean age at CRC diagnosis (\pm standard deviation, SD) for the dMMR group was 51 ± 15.0 with
379 62.6% being female and 49 ± 16.3 with 55.5% being female for the pMMR group. The
380 clinicopathological characteristics, pattern of MMR IHC loss and dMMR etiology for panel
381 sequenced CRC (n=29), EC (n=22) and SST (n=20) tumors are summarized in **Supplementary**
382 **Table 2**. Within the panel sequenced tumors, the proportion of dMMR for the CRC, EC and SST
383 subsets was 72.4% (21/29), 81.8% (18/22) and 65.0% (13/20), respectively. The predominant
384 dMMR subtype across the CRC WES and targeted panel sequenced tumors was dMMR-LS
385 (53.8% and 66.7%, respectively). Within the dMMR subgroup, the most predominant pattern of
386 loss observed in CRCs and ECs was MLH1/PMS2 (WES CRCs: 65.9%, panel CRCs: 47.6% and
387 ECs: 50.0%), whereas for the SSTs tumors, this was MSH2/MSH6 loss (76.9%). Tumors showing
388 less common patterns of MMR loss including solitary loss of MSH6 or PMS2 by IHC were present
389 in both the WES CRCs (16.5%) and panel sequenced tumors (19.2%), however, sole PMS2 loss
390 cases were absent from the EC and SST cohorts.

391

392 Assessment of Tumor Features for dMMR Prediction Accuracy in WES CRCs

393 *A) Individual tumor feature assessment*

394 Twelve of the 104 tumor features derived from WES had a mean dMMR prediction
395 accuracy $>80\%$ on the test dataset (**Table 2**). The mean accuracy for the remaining 92 features is
396 shown in **Supplementary Table 3**. The four MSI tools were among the best predictors, with
397 MSMuTect, MSIseq and MANTIS each achieving a mean prediction accuracy of $\geq 99.0\%$ with
398 MSMuTect achieving the highest accuracy (99.3%, 95% CI: 99.1%-99.5%) (**Table 2**). The
399 combination of TMS ID2+ID7 achieved an accuracy of 96.8% (95% CI: 96.4%-97.2%), and
400 outperformed these signatures individually (**Table 2**). To avoid collinearity issues between the

401 combined TMS ID2+ID7 variable with the individual TMS ID2 and TMS ID7 features, the latter
402 were excluded from downstream analysis as they provided a lower prediction score. Therefore, the
403 remaining 10 features were considered as the top 10 dMMR predictors and included in subsequent
404 analyses (**Figure 1**).

405 The mean, SD, and range of values for each of these top 10 dMMR predictive features by
406 MMR status and by dMMR subtype for the 300 WES CRCs are shown in **Supplementary Table**
407 **4**. For each of these features, the mean values were significantly different between the dMMR and
408 pMMR CRCs (all $p < 1 \times 10^{-12}$ from a *two-tailed t-test*), with TMS ID2+ID7 showing the most
409 significant difference ($p\text{-value} = 7.775 \times 10^{-98}$), although MSISensor presented with the highest
410 Cohen's *d* effect size of 4.5, indicating that the means of the pMMR and dMMR groups differed
411 by more than four times the SD (**Supplementary Table 4**). The variation in proportion or counts
412 was larger in the dMMR tumors than in the pMMR tumors for all but one of these top 10 features
413 where TMS ID2+ID7 demonstrated a broad range of values in the pMMR CRCs compared with
414 the dMMR CRCs (**Figure 2, Supplementary Table 4**).

415 The AUCs for the top 10 features when taking all possible thresholds into account are
416 shown in **Supplementary Figure 1**. The MSI prediction tools MSMuTect, MSIseq, and MANTIS
417 as well as INDEL count demonstrated the best AUCs. In addition, we calculated recommended
418 thresholds for each feature for differentiating dMMR from pMMR CRCs using the methodology
419 described in the methods (**Supplementary Table 5**). When applying these thresholds, it was not
420 possible to achieve a complete separation between the dMMR and pMMR tumors for each of the
421 tumor features (**Figure 3**).

422 Investigation of the CRCs misclassified based on the individual tumor feature analysis
423 demonstrated that the misclassification rate (error rate) for the MSI tools was low with MSMuTect

424 (2/300), MANTIS (1/300), MSIseq (1/300) and MSISensor (5/300) calling ≤ 5 incorrectly out of
425 300 tumors ($\leq 1.7\%$ error rate). Of the CRCs misclassified by the MSI tools, only two tumors were
426 misclassified by more than one MSI tool, both were dMMR-MLH1me CRCs classified as pMMR.
427 Of note, one of these dMMR-MLH1me CRCs was misclassified as a pMMR tumor by 9 out of the
428 top 10 tumor features. The second misclassified dMMR-MLH1me CRC was classified as pMMR
429 by MSMuTect and MSISensor but classified as dMMR by MSIseq and MANTIS (overall 6/10
430 features classified this CRC as dMMR). For INDEL count, 3/300 were incorrectly classified,
431 where two pMMR CRCs were classified as dMMR. TMS ID2+ID7 had 10/300 incorrect
432 classifications with seven pMMR tumors incorrectly called as dMMR. The remaining features
433 from the top 10 prediction accuracy list demonstrated the following incorrect classifications:
434 SBS20 (34/300), SBS54 (55/300), SBS15 (44/300) and TMB (19/300) encompassing incorrect
435 calls in both directions (dMMR to pMMR and vice versa).

436

437 *B) Generation of a statistical model by combining tumor features*

438 We assessed whether a combination of features within a statistical model could improve
439 dMMR prediction accuracy. For this, we performed a Lasso penalized logistic regression. Here,
440 after calculating the best lambda value, we found that the combination of TMS ID2+ID7
441 (coefficient = 5.29), MANTIS (coefficient = 1.70), MSISensor (coefficient = 0.09) with SBS15
442 (coefficient = 2.25) provided the best prediction accuracy from all possible feature combinations,
443 demonstrating a mean accuracy of 98.3% (95% CI: 0.981-0.986), sensitivity of 0.973 (95% CI:
444 0.966-0.980) and specificity of 1.000 (95% CI: 1.000-1.000) on the test set.

445 *C) Assessing an additive feature combination count for dMMR prediction*

446 Based on the observation that the top performing tumor features from the individual feature
447 analysis did not all misclassify the same CRCs lead us to explore a novel approach of combining
448 tumor features together to increase the overall accuracy i.e., an additive tumor feature combination
449 approach. This approach used a majority count of individual tumor features to overcome the small
450 inaccuracies that each of the top tumor features displayed individually i.e., if one of these top
451 dMMR predictive tumor features misclassified a CRC then the other top dMMR predictive tumor
452 features would correctly classify the same CRC and, thereby, achieve the correct classification
453 overall. Six of the top 10 features from the 10-fold cross-validation analysis demonstrated a mean
454 prediction accuracy of >95% and thus had the least number of incorrect CRC tumor classifications,
455 consisting of MSMuTect, MANTIS, MSISEq, MSISensor, INDEL count, and TMS ID2+ID7. We
456 applied the recommended threshold for determining dMMR status determined previously for each
457 tumor feature (**Figure 3, Supplementary Table 5**) to derive a count out of these six selected
458 features, in which each feature is weighted equally. The results show a bimodal distribution across
459 the 300 CRCs (**Figure 4**) where 0/6 to 2/6 features correctly classified all the pMMR CRCs and
460 4/6 to 6/6 correctly classified all but one of the dMMR tumors with an accuracy of 99.7%. The
461 only exception was the previously mentioned dMMR-MLH1me tumor, which did not meet the
462 recommended thresholds for all six features and thus received a count of 0/6 features suggestive
463 the CRC is pMMR rather than its initial dMMR status.

464

465 A summary of the results from the WES CRC analysis for the three approaches is shown
466 in **Table 3** and **Figure 1**.

467

468 Assessment of individual tumor features, Lasso statistical model and additive feature combination
469 approaches derived from the WES analysis on panel sequenced CRCs, ECs, and SSTs

470 To determine the generalizability of the findings from the three approaches performed on the
471 WES CRCs, we tested 71 tumors with targeted panel sequencing data to evaluate performance on
472 both a smaller capture and across different tissue types known to have a high prevalence of dMMR.

473

474 *A) Evaluation of the top performing individual features from WES analysis on the panel sequenced*
475 *CRC, EC, and SST tumors*

476 Out of the top 10 dMMR tumor features from the WES CRC analysis, only four achieved a
477 mean dMMR prediction accuracy of >80% in the panel sequenced CRC tumors (**Table 4**). For EC
478 and SST tumors only one feature (MANTIS) and two features (MANTIS and TMS ID2+ID7),
479 respectively, of the top 10 tumor features achieved a mean dMMR prediction accuracy of >80%
480 (**Table 4**). Across the three tissue types, MANTIS demonstrated the highest mean accuracy,
481 achieving 100% (95% CI: 88.1%-100.0%) accuracy in the panel sequenced CRCs, 86.4% accuracy
482 in ECs (95% CI: 65.1%-97.1%) and 85% accuracy in SSTs (95% CI: 62.1%-96.8%) (**Table 4**).
483 MSMuTect and INDEL count performed poorly in all three panel sequenced tissue types compared
484 with their accuracy in the WES CRCs. MSMuTect and INDEL count are features that provide
485 absolute counts that in our data were two orders of magnitude smaller in the panel sequenced
486 tumors compared with the WES CRCs. The reduction in discriminatory ability is likely related to
487 differences in the size (WES: 67.7 Mb and panel: 2.0 Mb) and location (additional coverage of
488 intronic regions of the MMR genes in the panel capture) of the regions covered by the WES and
489 panel captures resulting in a lower somatic mutation count.

490 The mean, SD, and range of values for each of these top 10 dMMR predicting features by
491 MMR status and by dMMR subtype for each of CRC, EC and SST tissue types are shown in
492 **Supplementary Tables 6A, 6B, 6C** and in **Supplementary Figure 2, Supplementary Figure 3,**
493 **and Supplementary Figure 4**, respectively. The mean values of each of the top 10 predictors were
494 significantly different between the dMMR and pMMR tumors in all three tissue types except for
495 TMS SBS15 in CRCs, MSISensor in ECs, TMB in ECs and SSTs and, TMS SBS20 and TMS
496 SBS54 in SSTs. MSMuTect consistently had the highest Cohen's d effect size of all top 10 tumor
497 features for each tissue type with the highest effect size observed in CRCs (3.2), indicating the
498 mean of the dMMR and pMMR subgroups for this feature differ by approximately three SDs.

499

500 *B) Evaluation of the Lasso statistical model on the panel sequenced CRC, EC, and SST tumors*

501 From WES analysis, the Lasso statistical model comprised of TMS ID2+ID7, MANTIS,
502 MSISensor and SBS15 achieved a mean prediction accuracy of 98.3%. When this model was
503 applied, with the coefficients determined from the WES analysis, on these three independent panel
504 sequenced tissue types, the prediction accuracies were lower (CRC: 89.7%, EC: 68.2% and SST:
505 85.0%) (**Table 3**).

506

507 *C) Evaluation of the additive tumor feature combination approach on the panel sequenced CRC,*
508 *EC, and SST tumors*

509 For each of the top 10 dMMR predictive tumor features we determined the optimal thresholds
510 for the panel sequenced CRCs, ECs, and SSTs (**Supplementary Table 5**) and plotted them by
511 tissue type (CRC - **Supplementary Figure 5**), (EC - **Supplementary Figure 6**), (SST -
512 **Supplementary Figure 7**). The determined thresholds for MANTIS were consistent across both

513 WES and panel captures as well as across tissue types while the calculated thresholds for MSIseq
514 were consistent for CRC across WES and panel captures but different to the thresholds determined
515 for EC and SST. The remaining eight tumor features showed variability in their determined
516 thresholds across both capture type and tissue type (**Supplementary Table 5**). As such, we applied
517 the thresholds determined for each tissue type for the panel sequenced data in the additive feature
518 combination approach below.

519

520 The additive feature combination approach incorporates a count of MSMuTect, MANTIS,
521 MSIseq, MSISensor, INDEL count and TMS ID2+ID7 tumor features to classify a tumor as
522 dMMR. The distribution of the counts of these six tumor features determined for each tumor are
523 shown for CRC (**Supplementary Figure 8**), EC (**Supplementary Figure 9**) and SSTs
524 (**Supplementary Figure 10**). For each tissue type, all the dMMR tumors had $\geq 3/6$ tumor features
525 classify them as dMMR, except for a single dMMR-MLH1me EC (1/71, 1.4%) which scored 0/6
526 and, therefore, was suggestive of pMMR status. This approach achieved accuracy scores of 100%,
527 95.5% and 100%, for CRC, EC and SST, respectively (**Table 3**).

528

529 A summary of the WES CRC and CRC, EC, and SST panel sequencing results for all three
530 approaches is provided in **Table 3**.

531

532 **Discussion**

533 In this study, we compared tumor features calculated from next generation sequencing data
534 for their accuracy in predicting dMMR status in 300 CRCs, 91 of which were dMMR determined
535 by immunohistochemistry or MSI-PCR and with an established sporadic or inherited etiology for

536 their dMMR status. Ten features achieved >80% dMMR prediction accuracy from the WES CRC
537 tumors, with the highest accuracy predictors being the MSI tools MSMuTect, MSIseq and
538 MANTIS, all of which achieved $\geq 99\%$ accuracy. The combination of TMS ID2+ID7 achieved the
539 highest mean accuracy for dMMR prediction out of the 97 TMS features assessed. When applied
540 to the targeted multi-gene panel setting, the performance of these 10 features was reduced not only
541 in CRC but also for the EC and SST tumors. In addition, we investigated two approaches that
542 combined these top 10 performing tumor features to improve the overall prediction accuracy. The
543 Lasso generated model achieved 98.3% accuracy in WES CRCs although the performance of the
544 model was reduced in the panel sequenced CRC, EC, and SST tumors. For both the WES CRCs
545 and panel sequencing across tissue types, the additive tumor feature combination approach, where
546 having ≥ 3 of the top 6 tumor features classify a tumor as dMMR, achieved the highest prediction
547 accuracies of the three approaches tested.

548
549 To date, multiple tools to detect MSI from NGS data have been developed⁷¹. NGS based MSI
550 tool development has been constantly evolving since the introduction of MSISensor²⁸ and
551 mSINGS⁷², which were followed by MSIseq²⁹, MANTIS³⁰ and MSMuTect³¹. However, to the best
552 of our knowledge, neither a comparison of more than three MSI detection tools on the same tumor
553 sample nor the effectiveness of these MSI tools specifically on SST tumors has been performed to
554 date. Previously, MANTIS has been compared to MSISensor with the former showing superior
555 sensitivity (97.18% vs. 96.48%) and specificity (99.68% vs. 98.73%)³⁰. This was supported by our
556 findings, and we additionally showed that across the WES and panel tested CRCs, MANTIS
557 provided the highest dMMR prediction accuracy and was shown to be the top performing feature
558 in the EC and SST tumors as well. Recently, the United States Food and Drug Administration

559 (FDA) approved MSISensor for detecting MSI in metastatic CRCs for selecting patients for
560 immune checkpoint inhibition therapy⁷³. In our study, MSISensor had the lowest accuracy (97.7%)
561 in WES CRCs of the four MSI tools tested, incorrectly classifying 5/300 CRCs. Seeking FDA
562 approval for other MSI tools in addition to MSISensor is warranted based on our findings.

563
564 MSMuTect has been trained on 20 different tissue types using WES data and, therefore, it
565 was not surprising it had the highest mean accuracy of the top performing tumor features in our
566 WES CRC analysis. MSMuTect has been designed to accurately detect somatic MSI indels using
567 a count of indels from the captured sequencing region³¹. Thus, the MSI indel count from WES data
568 (67.7 Mb) could be up to ~34x larger than that from panel data (2.0 Mb), which likely explains the
569 poor performance of this tool observed in our panel sequencing data test sets. When we adjusted
570 the MSMuTect threshold for calling dMMR for panel data, MSMuTect showed improved
571 discrimination of dMMR from pMMR tumors. This increase in prediction accuracy was also
572 observed for the INDEL count where adjusting the threshold for panel data improved the overall
573 performance. Adjusting the threshold for panel sequencing data enabled the inclusion of
574 MSMuTect and INDEL count as two of the six tumor features in our additive feature combination
575 approach that ultimately performed well on panel sequenced tumors. Tumor features that calculate
576 a percentage rather than raw counts such as MANTIS, MSISensor, SBS TMS and ID TMS are
577 more adaptable to changes in capture size. For example, our results showed that the calculated
578 thresholds for differentiating dMMR from pMMR for MANTIS were consistent across both WES
579 and panel captures as well as across tissue types. Therefore, we recommend training features that
580 incorporate a count of genomic variants, such as INDELS, SNVs and MSMuTect on the capture
581 size to improve dMMR prediction accuracy.

582

583 While three ID TMS (ID1, ID2 and ID7) are reported to be associated with dMMR³³, our
584 results showed that the combination of ID2 and ID7 TMS achieved the highest dMMR prediction
585 accuracy of any of the TMS features in WES CRC tumors, outperforming ID2 or ID7 alone. Of
586 the seven SBS TMS that are associated with dMMR (SBS6, SBS14, SBS16, SBS20, SBS21,
587 SBS26 and SBS44)³³, only two, TMS SBS15 and TMS SBS20, showed >80% dMMR prediction
588 accuracy in WES CRC tumors, but were shown to be poor predictors in the panel sequenced
589 tumors. Interestingly, TMS SBS54 was one of the top 10 dMMR predictors from the WES CRC
590 analysis, although currently its proposed etiology in COSMIC is related to a “possible sequencing
591 artefact and/or a possible contamination with germline variants”³³. Another study has shown that
592 SBS15, SBS20 and SBS54 are observed in CRCs with a high immune cytolytic activity (CYT)
593 compared with CYT-low CRCs⁷⁴. CYT-high CRCs have been shown to correlate with an increased
594 somatic mutation load and high levels of MSI⁷⁵, this may explain the observation of TMS SBS15,
595 TMS SBS20 and TMS SBS54 demonstrating >80% dMMR prediction accuracy in our WES CRC
596 analysis.

597

598 The combination of tumor features via the Lasso regression model achieved similar mean
599 accuracy as the four MSI tools individually in the WES CRC analysis. The Lasso calculated final
600 model that best distinguished dMMR from pMMR tumors in the WES CRC cohort consisted of
601 TMS ID2+ID7, MANTIS, MSISensor and TMS SBS15. The statistical approach used to determine
602 the final model assigns a ‘weight’ (coefficient value) or confidence of how well each feature
603 detects dMMR. As per generalized linear modelling methodology, the weight of any given feature
604 is reduced as the model incorporates additional features. Hence, with MANTIS being one of the

605 best predictors, its weighting was reduced when other features were added to the final model. This
606 resulted in the Lasso model prediction accuracy being lower than MANTIS alone. Of note, since
607 most of the approaches taken (i.e., assessing features individually or in combination) already
608 achieved a very high prediction accuracy of ~99%, alternate modelling approaches such as
609 Random Forest would not result in a significant improvement in dMMR prediction accuracy.

610
611 Strengths of our study were a large sample of tumors including dMMR tumors with
612 confirmed sporadic or inherited etiology concordant with MMR IHC and MSI-PCR results for
613 both the WES and panel sequenced datasets. Tumor MMR status combined with identified
614 etiology provided a more reliable reference group of CRCs than would a group based on MMR
615 IHC test results without etiological confirmation given the known challenges that can lead to false
616 positive and negative MMR IHC results¹⁶. We assessed many tumor features that can be readily
617 derived from NGS data ensuring that our findings have potential to be easily implemented in
618 clinical diagnostics. We applied our findings from WES to panel data to determine the
619 generalizability of our findings to smaller panel captures such as those that are currently used in
620 clinical diagnostics. We showed the applicability of our findings on tissue types that display a high
621 proportion of dMMR phenotype. Our dMMR tumor samples included those with the frequent
622 pattern of MMR IHC namely MLH1/PMS2 loss and MSH2/MSH6 loss but also tumors with
623 solitary MSH6 loss or solitary PMS2 loss, ensuring we covered the spectrum of dMMR tissue
624 types which is particularly relevant given the identified challenges associated with interpretation
625 of solitary MSH6 loss⁷⁶.

626

627 There were several limitations of our study including testing of only three tissue types.
628 Testing of these tumor features and approaches in other tissue types such as stomach cancer, which
629 also has a high prevalence of dMMR overall and dMMR related to Lynch syndrome, would
630 determine the suitability of these tumor features for inclusion in an additive feature combination
631 approach in a pan-cancer setting. In addition, the sample size for the panel sequenced tumors was
632 limited for all three tissue types, however, there was a high proportion of dMMR in the tumors
633 tested (72.4% for CRC, 81.8% for EC and 65.0% for SST). No tumor feature or approach achieved
634 100% accuracy in the CRC WES analysis. This was largely related to a single tumor (dMMR-
635 MLH1me) from the WES CRC analysis that was called incorrectly by 9/10 top individual tumor
636 features suggesting the CRC was pMMR. Therefore, we repeated the *MLH1* methylation testing
637 for this tumor using both MethyLight and MS-HRM assays. Both assays found no evidence of
638 *MLH1* methylation in the tumor. These new *MLH1* methylation results and the pMMR
639 classification from our analysis suggest the initial dMMR classification was a false positive. If this
640 CRC would initially have been categorized as a pMMR tumor, then MANTIS and MSIseq would
641 have achieved 100% accuracy in the WES CRC analysis. Furthermore, the identification of an
642 initial tumor misclassification provides strong support for evaluating multiple dMMR prediction
643 tumor features and highlights the advantage of combining these features through an additive
644 feature combination approach.

645

646 **Conclusion**

647 Our findings provide an important comparison of tumor features for dMMR prediction,
648 highlighting performance differences between capture size and tissue types. Our results
649 demonstrate the high accuracy of multiple individual tumor features including the MSI calling

650 tools MSMuTect, MSIseq, MANTIS and MSISensor, as well as INDEL count and the combination
651 of TMS ID2+ID7 for predicting dMMR status using WES CRCs. Moreover, our findings highlight
652 the benefit of combining these six tumor features in a simple additive feature combination
653 approach to improve dMMR prediction accuracy, particularly in targeted panel sequencing data
654 from CRC, EC, or SST tumors. With the reported inaccuracies of MMR IHC and the increasing
655 application of clinical NGS testing of tumor tissue, accurately deriving dMMR status from this
656 NGS data will have important implications for diagnostics and targeted therapy and likely improve
657 patient outcomes and cancer prevention.

658

659 **Acknowledgements:** We thank members of the Colorectal Oncogenomics Group and members
660 from the Genomic Medicine and Family Cancer Clinic for their support of this manuscript. We
661 thank the participants and staff from the Australasian and Ontario Colorectal Cancer Family
662 Registries (ACCFR/OFCCR) and the ANGELS, Muir-Torre and WEHI studies. We especially
663 thank Maggie Angelakos, Samantha Fox, Allyson Templeton for supporting this study. We thank
664 the Australian Genome Research Facility for their collaboration on this project. We thank A/Prof
665 Sue Finch of the Melbourne Statistical Consulting Platform and Statistical Consulting Centre at
666 the University of Melbourne for guidance with the statistical aspects of this study.

667 **References**

- 668 1. Gryfe R, Kim H, Hsieh ET, Aronson MD, Holowaty EJ, Bull SB, Redston M, Gallinger S.
669 Tumor microsatellite instability and clinical outcome in young patients with colorectal
670 cancer. *N Engl J Med*, 2000, 342:69–77
- 671 2. Le DT, Uram JN, Wang H, Bartlett BR, Kemberling H, Eyring AD, Skora AD, Luber BS,
672 Azad NS, Laheru D, Biedrzycki B, Donehower RC, Zaheer A, Fisher GA, Crocenzi TS, Lee
673 JJ, Duffy SM, Goldberg RM, de la Chapelle A, Koshiji M, Bhajjee F, Huebner T, Hruban
674 RH, Wood LD, Cuka N, Pardoll DM, Papadopoulos N, Kinzler KW, Zhou S, Cornish TC,
675 Taube JM, Anders RA, Eshleman JR, Vogelstein B, Diaz LA. PD-1 Blockade in Tumors
676 with Mismatch-Repair Deficiency. *N Engl J Med*, 2015, 372:2509–20
- 677 3. Evaluation of Genomic Applications in Practice and Prevention (EGAPP) Working Group.
678 Recommendations from the EGAPP Working Group: genetic testing strategies in newly
679 diagnosed individuals with colorectal cancer aimed at reducing morbidity and mortality from
680 Lynch syndrome in relatives. *Genet Med*, 2009, 11:35–41
- 681 4. Green RF, Ari M, Kolor K, Dotson WD, Bowen S, Habarta N, Rodriguez JL, Richardson
682 LC, Khoury MJ. Evaluating the role of public health in implementation of genomics-related
683 recommendations: a case study of hereditary cancers using the CDC Science Impact
684 Framework. *Genet Med*, 2019, 21:28–37
- 685 5. Baretta M, Le DT. DNA mismatch repair in cancer. *Pharmacol Ther*, 2018, 189:45–62
- 686 6. Eshleman JR, Markowitz SD. Mismatch repair defects in human carcinogenesis. *Hum Mol*
687 *Genet*, 1996, 5 Spec No:1489–94
- 688 7. Young J, Simms LA, Biden KG, Wynter C, Whitehall V, Karamatic R, George J, Goldblatt
689 J, Walpole I, Robin S-A, Borten MM, Stitz R, Searle J, McKeone D, Fraser L, Purdie DR,

- 690 Podger K, Price R, Buttenshaw R, Walsh MD, Barker M, Leggett BA, Jass JR. Features of
691 Colorectal Cancers with High-Level Microsatellite Instability Occurring in Familial and
692 Sporadic Settings. *Am J Pathol*, 2001, 159:2107–16
- 693 8. Haraldsdottir S, Hampel H, Tomsic J, Frankel WL, Pearlman R, de la Chapelle A, Pritchard
694 CC. Colon and endometrial cancers with mismatch repair deficiency can arise from somatic,
695 rather than germline, mutations. *Gastroenterology*, 2014, 147:1308-1316.e1
- 696 9. Lynch HT, Lynch PM, Lanspa SJ, Snyder CL, Lynch JF, Boland CR. Review of the Lynch
697 syndrome: history, molecular genetics, screening, differential diagnosis, and medicolegal
698 ramifications. *Clin Genet*, 2009, 76:1–18
- 699 10. Ligtenberg MJ, Kuiper RP, Chan TL, Goossens M, Hebeda KM, Voorendt M, Lee TY,
700 Bodmer D, Hoenselaar E, Hendriks-Cornelissen SJ, Tsui WY, Kong CK, Brunner HG, van
701 Kessel AG, Yuen ST, van Krieken JH, Leung SY, Hoogerbrugge N. Heritable somatic
702 methylation and inactivation of MSH2 in families with Lynch syndrome due to deletion of
703 the 3' exons of TACSTD1. *Nat Genet*, 2009, 41:112–7
- 704 11. Bonneville R, Krook MA, Kautto EA, Miya J, Wing MR, Chen H-Z, Reeser JW, Yu L,
705 Roychowdhury S. Landscape of Microsatellite Instability Across 39 Cancer Types. *JCO*
706 *Precis Oncol*, 2017, 2017
- 707 12. Walsh MD, Jayasekara H, Huang A, Winship IM, Buchanan DD. Clinico-pathological
708 predictors of mismatch repair deficiency in sebaceous neoplasia: A large case series from a
709 single Australian private pathology service. *Australas J Dermatol*, 2019, 60:126–33
- 710 13. Mascarenhas L, Shanley S, Mitchell G, Spurdle AB, Macrae F, Pachter N, Buchanan DD,
711 Ward RL, Fox S, Duxbury E, Driessen R, Boussioutas A. Current mismatch repair

- 712 deficiency tumor testing practices and capabilities: A survey of Australian pathology
713 providers. *Asia Pac J Clin Oncol*, 2018, 14:417–25
- 714 14. Shia J. Immunohistochemistry versus microsatellite instability testing for screening
715 colorectal cancer patients at risk for hereditary nonpolyposis colorectal cancer syndrome.
716 Part I. The utility of immunohistochemistry. *J Mol Diagn*, 2008, 10:293–300
- 717 15. Boland CR, Thibodeau SN, Hamilton SR, Sidransky D, Eshleman JR, Burt RW, Meltzer SJ,
718 Rodriguez-Bigas MA, Fodde R, Ranzani GN, Srivastava S. A National Cancer Institute
719 Workshop on Microsatellite Instability for cancer detection and familial predisposition:
720 development of international criteria for the determination of microsatellite instability in
721 colorectal cancer. *Cancer Res*, 1998, 58:5248–57
- 722 16. Chen M-L, Chen J-Y, Hu J, Chen Q, Yu L-X, Liu B-R, Qian X-P, Yang M. Comparison of
723 microsatellite status detection methods in colorectal carcinoma. *Int J Clin Exp Pathol*, 2018,
724 11:1431–8
- 725 17. Chapusot C, Martin L, Bouvier AM, Bonithon-Kopp C, Ecarnot-Laubriet A, Rageot D,
726 Ponnelle T, Laurent Puig P, Faivre J, Piard F. Microsatellite instability and intratumoural
727 heterogeneity in 100 right-sided sporadic colon carcinomas. *Br J Cancer*, 2002, 87:400–4
- 728 18. Graham RP, Kerr SE, Butz ML, Thibodeau SN, Halling KC, Smyrk TC, Dina MA, Waugh
729 VM, Rumilla KM. Heterogenous MSH6 loss is a result of microsatellite instability within
730 MSH6 and occurs in sporadic and hereditary colorectal and endometrial carcinomas. *Am J*
731 *Surg Pathol*, 2015, 39:1370–6
- 732 19. Joost P, Veurink N, Holck S, Klarskov L, Bojesen A, Harbo M, Baldetorp B, Rambech E,
733 Nilbert M. Heterogenous mismatch-repair status in colorectal cancer. *Diagn Pathol*, 2014,
734 9:126

- 735 20. McCarthy AJ, Capo-Chichi J-M, Spence T, Grenier S, Stockley T, Kamel-Reid S, Serra S,
736 Sabatini P, Chetty R. Heterogenous loss of mismatch repair (MMR) protein expression: a
737 challenge for immunohistochemical interpretation and microsatellite instability (MSI)
738 evaluation. *J Pathol Clin Res*, 2019, 5:115–29
- 739 21. Pai RK, Plesec TP, Abdul-Karim FW, Yang B, Marquard J, Shadrach B, Roma AR. Abrupt
740 loss of MLH1 and PMS2 expression in endometrial carcinoma: molecular and morphologic
741 analysis of 6 cases. *Am J Surg Pathol*, 2015, 39:993–9
- 742 22. Shia J, Zhang L, Shike M, Guo M, Stadler Z, Xiong X, Tang LH, Vakiani E, Katabi N,
743 Wang H, Bacares R, Ruggeri J, Boland CR, Ladanyi M, Klimstra DS. Secondary mutation in
744 a coding mononucleotide tract in MSH6 causes loss of immunoexpression of MSH6 in
745 colorectal carcinomas with MLH1/PMS2 deficiency. *Mod Pathol*, 2013, 26:131–8
- 746 23. Watkins JC, Nucci MR, Ritterhouse LL, Howitt BE, Sholl LM. Unusual Mismatch Repair
747 Immunohistochemical Patterns in Endometrial Carcinoma. *Am J Surg Pathol*, 2016, 40:909–
748 16
- 749 24. Watson N, Grieu F, Morris M, Harvey J, Stewart C, Schofield L, Goldblatt J, Iacopetta B.
750 Heterogeneous staining for mismatch repair proteins during population-based prescreening
751 for hereditary nonpolyposis colorectal cancer. *J Mol Diagn*, 2007, 9:472–8
- 752 25. Baudrin LG, Deleuze J-F, How-Kit A. Molecular and computational methods for the
753 detection of microsatellite instability in cancer. *Frontiers in Oncology*, 2018, 8
- 754 26. Vasen HFA, Hendriks Y, de Jong AE, van Puijenbroek M, Tops C, Bröcker-Vriens AHJT,
755 Wijnen JTh, Morreau H. Identification of HNPCC by Molecular Analysis of Colorectal and
756 Endometrial Tumors. *Dis Markers*, 2004, 20:207–13

- 757 27. Siemanowski J, Schömig-Markiefka B, Buhl T, Haak A, Siebolts U, Dietmaier W, Arens N,
758 Pauly N, Ataseven B, Büttner R, Merkelbach-Bruse S. Managing Difficulties of
759 Microsatellite Instability Testing in Endometrial Cancer-Limitations and Advantages of Four
760 Different PCR-Based Approaches. *Cancers (Basel)*, 2021, 13:1268
- 761 28. Niu B, Ye K, Zhang Q, Lu C, Xie M, McLellan MD, Wendl MC, Ding L. MSIsensor:
762 microsatellite instability detection using paired tumor-normal sequence data. *Bioinformatics*,
763 2014, 30:1015–6
- 764 29. Ni Huang M, McPherson JR, Cutcutache I, Teh BT, Tan P, Rozen SG. MSIseq: Software for
765 Assessing Microsatellite Instability from Catalogs of Somatic Mutations. *Sci Rep*, 2015,
766 5:13321
- 767 30. Kautto EA, Bonneville R, Miya J, Yu L, Krook MA, Reeser JW, Roychowdhury S.
768 Performance evaluation for rapid detection of pan-cancer microsatellite instability with
769 MANTIS. *Oncotarget*, 2017, 8:7452–63
- 770 31. Maruvka YE, Mouw KW, Karlic R, Parasuraman P, Kamburov A, Polak P, Haradhvala NJ,
771 Hess JM, Rheinbay E, Brody Y, Koren A, Braunstein LZ, D’Andrea A, Lawrence MS, Bass
772 A, Bernards A, Michor F, Getz G. Analysis of somatic microsatellite indels identifies driver
773 events in human tumors. *Nat Biotechnol*, 2017, 35:951–9
- 774 32. Muzny DM, Bainbridge MN, Chang K, Dinh HH, Drummond JA, Fowler G, Kovar CL,
775 Lewis LR, Morgan MB, Newsham IF, Reid JG, Santibanez J, Shinbrot E, Trevino LR, Wu
776 Y-Q, Wang M, Gunaratne P, Donehower LA, Creighton CJ, Wheeler DA, Gibbs RA,
777 Lawrence MS, Voet D, Jing R, Cibulskis K, Sivachenko A, Stojanov P, McKenna A, Lander
778 ES, Gabriel S, Getz G, Ding L, Fulton RS, Koboldt DC, Wylie T, et al. Comprehensive
779 molecular characterization of human colon and rectal cancer. *Nature*, 2012, 487:330–7

- 780 33. Tate JG, Bamford S, Jubb HC, Sondka Z, Beare DM, Bindal N, Boutselakis H, Cole CG,
781 Creatore C, Dawson E, Fish P, Harsha B, Hathaway C, Jupe SC, Kok CY, Noble K, Ponting
782 L, Ramshaw CC, Rye CE, Speedy HE, Stefancsik R, Thompson SL, Wang S, Ward S,
783 Campbell PJ, Forbes SA. COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic
784 Acids Res*, 2019, 47:D941–7
- 785 34. Panda A, Betigeri A, Subramanian K, Ross JS, Pavlick DC, Ali S, Markowski P, Silk A,
786 Kaufman HL, Lattime E, Mehnert JM, Sullivan R, Lovly CM, Sosman J, Johnson DB,
787 Bhanot G, Ganesan S. Identifying a Clinically Applicable Mutational Burden Threshold as a
788 Potential Biomarker of Response to Immune Checkpoint Therapy in Solid Tumors. *JCO
789 Precis Oncol*, 2017, 2017
- 790 35. Zheng M. Tumor mutation burden for predicting immune checkpoint blockade response: the
791 more, the better. *J Immunother Cancer*, 2022, 10:e003087
- 792 36. Chang H, Sasson A, Srinivasan S, Golhar R, Greenawalt DM, Geese WJ, Green G, Zerba K,
793 Kirov S, Szustakowski J. Bioinformatic Methods and Bridging of Assay Results for Reliable
794 Tumor Mutational Burden Assessment in Non-Small-Cell Lung Cancer. *Mol Diagn Ther*,
795 2019, 23:507–20
- 796 37. Alexandrov LB, Nik-Zainal S, Wedge DC, Aparicio SAJR, Behjati S, Biankin AV, Bignell
797 GR, Bolli N, Borg A, Børresen-Dale A-L, Boyault S, Burkhardt B, Butler AP, Caldas C,
798 Davies HR, Desmedt C, Eils R, Eyfjörd JE, Foekens JA, Greaves M, Hosoda F, Hutter B,
799 Ilicic T, Imbeaud S, Imielinski M, Jäger N, Jones DTW, Jones D, Knappskog S, Kool M,
800 Lakhani SR, López-Otín C, Martin S, Munshi NC, Nakamura H, et al. Signatures of
801 mutational processes in human cancer. *Nature*, 2013, 500:415–21

- 802 38. Alexandrov LB, Kim J, Haradhvala NJ, Huang MN, Tian Ng AW, Wu Y, Boot A,
803 Covington KR, Gordenin DA, Bergstrom EN, Islam SMA, Lopez-Bigas N, Klimczak LJ,
804 McPherson JR, Morganella S, Sabarinathan R, Wheeler DA, Mustonen V, Getz G, Rozen
805 SG, Stratton MR. The repertoire of mutational signatures in human cancer. *Nature*, 2020,
806 578:94–101
- 807 39. Georgeson P, Pope BJ, Rosty C, Clendenning M, Mahmood K, Joo JE, Walker R,
808 Hutchinson RA, Preston S, Como J, Joseland S, Win AK, Macrae FA, Hopper JL, Mouradov
809 D, Gibbs P, Sieber OM, O’Sullivan DE, Brenner DR, Gallinger S, Jenkins MA, Winship IM,
810 Buchanan DD. Evaluating the utility of tumour mutational signatures for identifying
811 hereditary colorectal cancer and polyposis syndrome carriers. *Gut*, 2021, 70:2138–49
- 812 40. Georgeson P, Harrison TA, Pope BJ, Zaidi SH, Qu C, Steinfeldt RS, Lin Y, Joo JE,
813 Mahmood K, Clendenning M, Walker R, Amitay EL, Berndt SI, Brenner H, Campbell PT,
814 Cao Y, Chan AT, Chang-Claude J, Doheny KF, Drew DA, Figueiredo JC, French AJ,
815 Gallinger S, Giannakis M, Giles GG, Gsur A, Gunter MJ, Hoffmeister M, Hsu L, Huang W-
816 Y, Limburg P, Manson JE, Moreno V, Nassir R, Nowak JA, et al. Identifying colorectal
817 cancer caused by biallelic MUTYH pathogenic variants using tumor mutational signatures.
818 *Nat Commun*, 2022, 13:3254
- 819 41. Jenkins MA, Win AK, Templeton AS, Angelakos MS, Buchanan DD, Cotterchio M,
820 Figueiredo JC, Thibodeau SN, Baron JA, Potter JD, Hopper JL, Casey G, Gallinger S, Le
821 Marchand L, Lindor NM, Newcomb PA, Haile RW, Colon Cancer Family Registry Cohort I.
822 Cohort Profile: The Colon Cancer Family Registry Cohort (CCFRC). *Int J Epidemiol*, 2018,
823 47:387–388i

- 824 42. Newcomb PA, Baron J, Cotterchio M, Gallinger S, Grove J, Haile R, Hall D, Hopper JL,
825 Jass J, Le Marchand L, Limburg P, Lindor N, Potter JD, Templeton AS, Thibodeau S,
826 Seminara D, Colon Cancer Family R. Colon Cancer Family Registry: an international
827 resource for studies of the genetic epidemiology of colon cancer. *Cancer Epidemiol*
828 *Biomarkers Prev*, 2007, 16:2331–43
- 829 43. Buchanan DD, Clendenning M, Rosty C, Eriksen SV, Walsh MD, Walters RJ, Thibodeau
830 SN, Stewart J, Preston S, Win AK, Flander L, Ouakrim DA, Macrae FA, Boussioutas A,
831 Winship IM, Giles GG, Hopper JL, Southey MC, English D, Jenkins MA. Tumour testing to
832 identify Lynch syndrome in two Australian colorectal cancer cohorts. *J Gastroenterol*
833 *Hepatol*, 2017, 32:427–38
- 834 44. Walsh MD, Buchanan DD, Pearson S-A, Clendenning M, Jenkins MA, Win AK, Walters RJ,
835 Spring KJ, Nagler B, Pavluk E, Arnold ST, Goldblatt J, George J, Suthers GK, Phillips K,
836 Hopper JL, Jass JR, Baron JA, Ahnen DJ, Thibodeau SN, Lindor N, Parry S, Walker NI,
837 Rosty C, Young JP. Immunohistochemical testing of conventional adenomas for loss of
838 expression of mismatch repair proteins in Lynch syndrome mutation carriers: a case series
839 from the Australasian site of the colon cancer family registry. *Mod Pathol*, 2012, 25:722–30
- 840 45. Cicek MS, Lindor NM, Gallinger S, Bapat B, Hopper JL, Jenkins MA, Young J, Buchanan
841 D, Walsh MD, Le Marchand L, Burnett T, Newcomb PA, Grady WM, Haile RW, Casey G,
842 Plummer SJ, Krumroy LA, Baron JA, Thibodeau SN. Quality assessment and correlation of
843 microsatellite instability and immunohistochemical markers among population- and clinic-
844 based colorectal tumors results from the Colon Cancer Family Registry. *J Mol Diagn*, 2011,
845 13:271–81

- 846 46. Buchanan DD, Tan YY, Walsh MD, Clendenning M, Metcalf AM, Ferguson K, Arnold ST,
847 Thompson BA, Lose FA, Parsons MT, Walters RJ, Pearson SA, Cummings M, Oehler MK,
848 Blomfield PB, Quinn MA, Kirk JA, Stewart CJ, Obermair A, Young JP, Webb PM, Spurdle
849 AB. Tumor mismatch repair immunohistochemistry and DNA MLH1 methylation testing of
850 patients with endometrial cancer diagnosed at age younger than 60 years optimizes triage for
851 population-level germline mismatch repair gene mutation testing. *J Clin Oncol*, 2014, 32:90–
852 100
- 853 47. Weisenberger DJ, Siegmund KD, Campan M, Young J, Long TI, Faasse MA, Kang GH,
854 Widschwendter M, Weener D, Buchanan D, Koh H, Simms L, Barker M, Leggett B, Levine
855 J, Kim M, French AJ, Thibodeau SN, Jass J, Haile R, Laird PW. CpG island methylator
856 phenotype underlies sporadic microsatellite instability and is tightly associated with BRAF
857 mutation in colorectal cancer. *Nat Genet*, 2006, 38:787–93
- 858 48. Zaidi SH, Harrison TA, Phipps AI, Steinfeld R, Trinh QM, Qu C, Banbury BL, Georgeson
859 P, Grasso CS, Giannakis M, Adams JB, Alwers E, Amitay EL, Barfield RT, Berndt SI,
860 Borozan I, Brenner H, Brezina S, Buchanan DD, Cao Y, Chan AT, Chang-Claude J,
861 Connolly CM, Drew DA, Farris AB, Figueiredo JC, French AJ, Fuchs CS, Garraway LA,
862 Gruber S, Guinte MA, Hamilton SR, Harlid S, Heisler LE, Hidaka A, et al. Landscape of
863 somatic single nucleotide variants and indels in colorectal cancer and impact on survival. *Nat*
864 *Commun*, 2020, 11:3644
- 865 49. Belhadj S, Terradas M, Munoz-Torres PM, Aiza G, Navarro M, Capellá G, Valle L.
866 Candidate genes for hereditary colorectal cancer: Mutational screening and systematic
867 review. *Hum Mutat*, 2020, 41:1563–76

- 868 50. Seifert BA, McGlaughon JL, Jackson SA, Ritter DI, Roberts ME, Schmidt RJ, Thompson
869 BA, Jimenez S, Trapp M, Lee K, Plon SE, Offit K, Stadler ZK, Zhang L, Greenblatt MS,
870 Ferber MJ. Determining the clinical validity of hereditary colorectal cancer and polyposis
871 susceptibility genes using the Clinical Genome Resource Clinical Validity Framework.
872 *Genet Med*, 2019, 21:1507–16
- 873 51. Weren RDA, Ligtenberg MJL, Kets CM, de Voer RM, Verwiel ETP, Spruijt L, van Zelst-
874 Stams WAG, Jongmans MC, Gilissen C, Hehir-Kwa JY, Hoischen A, Shendure J, Boyle EA,
875 Kamping EJ, Nagtegaal ID, Tops BBJ, Nagengast FM, Geurts van Kessel A, van Krieken
876 JHJM, Kuiper RP, Hoogerbrugge N. A germline homozygous mutation in the base-excision
877 repair gene NTHL1 causes adenomatous polyposis and colorectal cancer. *Nat Genet*, 2015,
878 47:668–71
- 879 52. Spurdle AB, Bowman MA, Shamsani J, Kirk J. Endometrial cancer gene panels: clinical
880 diagnostic vs research germline DNA testing. *Mod Pathol*, 2017, 30:1048–68
- 881 53. Cancer Genome Atlas Research Network, Kandoth C, Schultz N, Cherniack AD, Akbani R,
882 Liu Y, Shen H, Robertson AG, Pashtan I, Shen R, Benz CC, Yau C, Laird PW, Ding L,
883 Zhang W, Mills GB, Kucherlapati R, Mardis ER, Levine DA. Integrated genomic
884 characterization of endometrial carcinoma. *Nature*, 2013, 497:67–73
- 885 54. Cherniack AD, Shen H, Walter V, Stewart C, Murray BA, Bowlby R, Hu X, Ling S, Soslow
886 RA, Broaddus RR, Zuna RE, Robertson G, Laird PW, Kucherlapati R, Mills GB, Akbani R,
887 Ally A, Auman JT, Balasundaram M, Balu S, Baylin SB, Beroukheim R, Bodenheimer T,
888 Bogomolny F, Boice L, Bootwalla MS, Bowen J, Bowlby R, Broaddus R, Brooks D,
889 Carlsen R, Cherniack AD, Cho J, Chuah E, Chudamani S, et al. Integrated Molecular
890 Characterization of Uterine Carcinosarcoma. *Cancer Cell*, 2017, 31:411–23

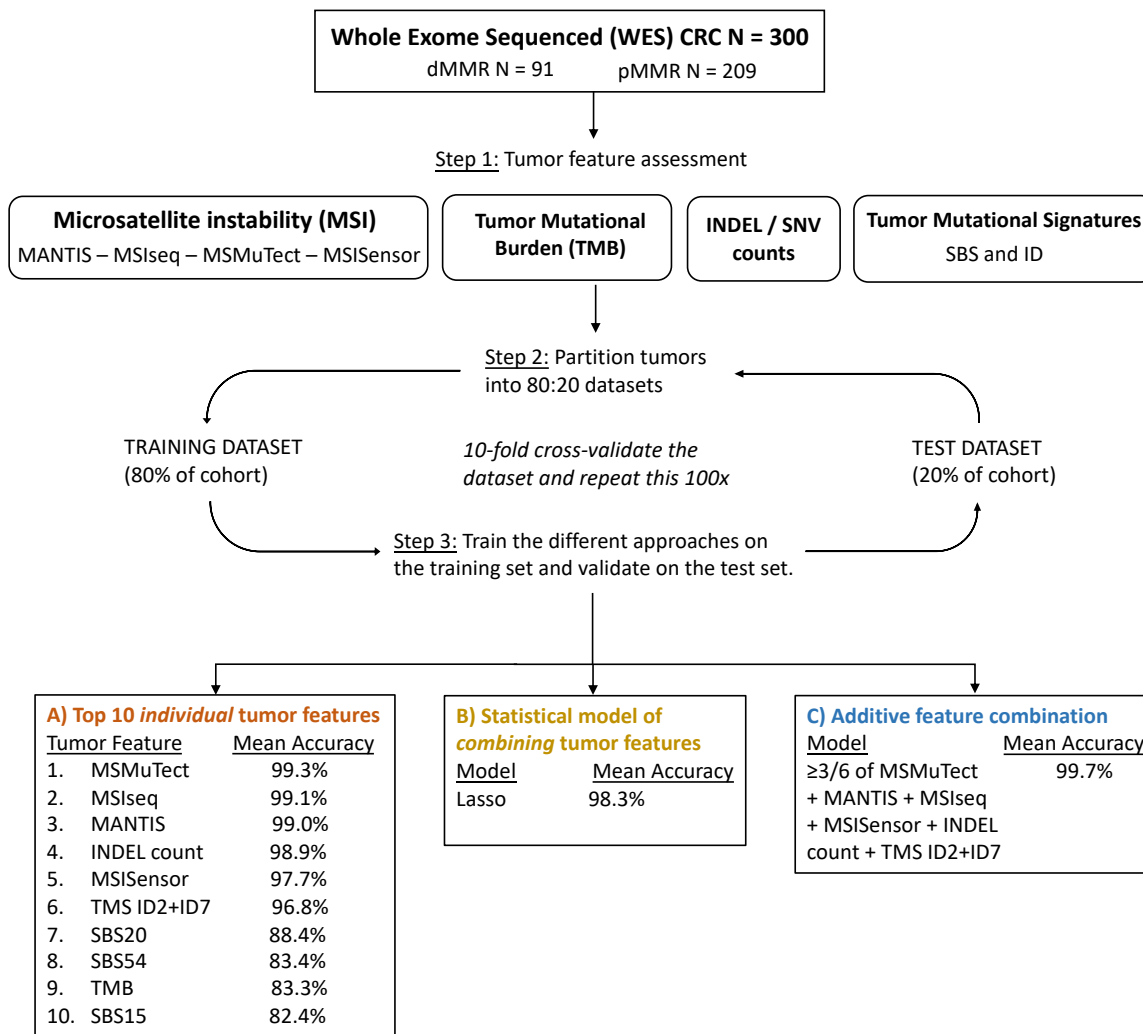
- 891 55. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
892 data. *Bioinformatics*, 2014, 30:2114–20
- 893 56. Saunders CT, Wong WSW, Swamy S, Becq J, Murray LJ, Cheetham RK. Strelka: accurate
894 somatic small-variant calling from sequenced tumor–normal sample pairs. *Bioinformatics*,
895 2012, 28:1811–7
- 896 57. Sha D, Jin Z, Budzcies J, Kluck K, Stenzinger A, Sinicrope FA. Tumor Mutational Burden
897 (TMB) as a Predictive Biomarker in Solid Tumors. *Cancer Discov*, 2020, 10:1808–25
- 898 58. Kuhn M, cre, Wing J, Weston S, Williams A, Keefer C, Engelhardt A, Cooper T, Mayer Z,
899 Kenkel B, R Core Team, Benesty M, Lescarbeau R, Ziem A, Scrucca L, Tang Y, Candan C,
900 Hunt T. caret: Classification and Regression Training. 2022
- 901 59. Tibshirani R. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal*
902 *Statistical Society Series B (Methodological)*, 1996, 58:267–88
- 903 60. Wickham H, Averick M, Bryan J, Chang W, McGowan LD, François R, Golemund G,
904 Hayes A, Henry L, Hester J, Kuhn M, Pedersen TL, Miller E, Bache SM, Müller K, Ooms J,
905 Robinson D, Seidel DP, Spinu V, Takahashi K, Vaughan D, Wilke C, Woo K, Yutani H.
906 Welcome to the Tidyverse. *Journal of Open Source Software*, 2019, 4:1686
- 907 61. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez J-C, Müller M. pROC: an
908 open-source package for R and S+ to analyze and compare ROC curves. *BMC*
909 *Bioinformatics*, 2011, 12:77
- 910 62. LeDell E, Petersen M, Laan M van der. cvAUC: Cross-Validated Area Under the ROC
911 Curve Confidence Intervals. 2022
- 912 63. Friedman JH, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear Models
913 via Coordinate Descent. *Journal of Statistical Software*, 2010, 33:1–22

- 914 64. Thiele C, Hirschfeld G. cutpointr: Improved Estimation and Validation of Optimal Cutpoints
915 in R. Journal of Statistical Software, 2021, 98:1–27
- 916 65. Wickham H. The Split-Apply-Combine Strategy for Data Analysis. Journal of Statistical
917 Software, 2011, 40:1–29
- 918 66. Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York, 2016
- 919 67. Rudis B, cre, Kennedy P, Reiner P, support) DW (Secondary axis, Adam X, Fonts) G
920 (Roboto C& TW, Font) I (Plex S, Font) IT (Public S, Barnett J, Leeper TJ, Meys J.
921 hrbrthemes: Additional Themes, Theme Components and Utilities for “ggplot2.” 2020
- 922 68. Slowikowski K, Schep A, Hughes S, Dang TK, Lukauskas S, Irisson J-O, Kamvar ZN, Ryan
923 T, Christophe D, Hiroaki Y, Gramme P, Abdol AM, Barrett M, Cannoodt R, Krassowski M,
924 Chirico M, Aphalo P. ggrepel: Automatically Position Non-Overlapping Text Labels with
925 “ggplot2.” 2021
- 926 69. Clopper CJ, Pearson ES. THE USE OF CONFIDENCE OR FIDUCIAL LIMITS
927 ILLUSTRATED IN THE CASE OF THE BINOMIAL. Biometrika, 1934, 26:404–13
- 928 70. Dorai-Raj S. binom: Binomial Confidence Intervals for Several Parameterizations. 2022
- 929 71. Renault V, Tubacher E, How-Kit A. Assessment of Microsatellite Instability from Next-
930 Generation Sequencing Data. Edited by Laganà A. Computational Methods for Precision
931 Oncology, Cham, Springer International Publishing, 2022, pp. 75–100
- 932 72. Salipante SJ, Scroggins SM, Hampel HL, Turner EH, Pritchard CC. Microsatellite Instability
933 Detection by Next Generation Sequencing. Clinical Chemistry, 2014, 60:1192–9
- 934 73. Ratovomanana T, Cohen R, Svrcek M, Renaud F, Cervera P, Siret A, Letourneur Q, Buhard
935 O, Bourgoin P, Guillerme E, Dorard C, Nicolle R, Ayadi M, Touat M, Bielle F, Sanson M,
936 Rouzic PL, Buisine M-P, Piessen G, Collura A, Fléjou J-F, Reyniès A de, Coulet F,

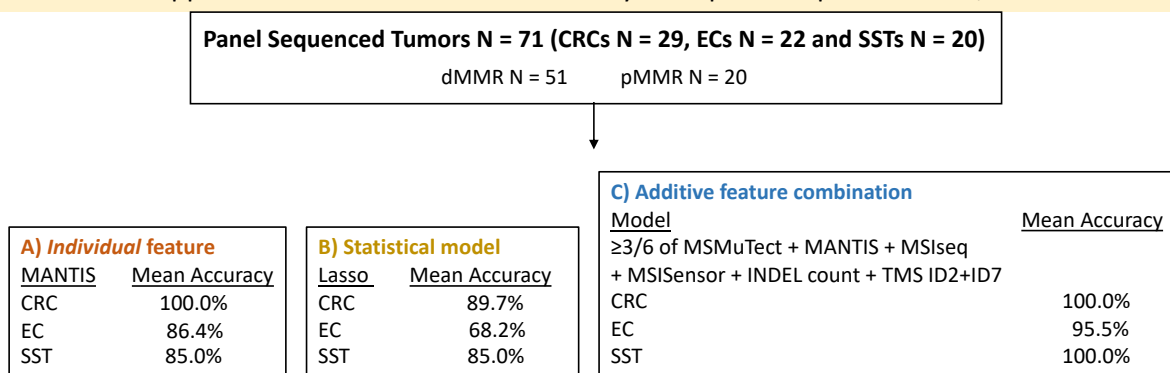
- 937 Ghiringhelli F, André T, Jonchère V, Duval A. Performance of Next-Generation Sequencing
938 for the Detection of Microsatellite Instability in Colorectal Cancer With Deficient DNA
939 Mismatch Repair. *Gastroenterology*, 2021, 161:814-826.e7
- 940 74. Roufas C, Georgakopoulos-Soares I, Zaravinos A. Molecular correlates of immune cytolytic
941 subgroups in colorectal cancer by integrated genomics analysis. *NAR Cancer*, 2021,
942 3:zcab005
- 943 75. Zaravinos A, Roufas C, Nagara M, de Lucas Moreno B, Oblovatskaya M, Efstathiades C,
944 Dimopoulos C, Ayiomamitis GD. Cytolytic activity correlates with the mutational burden
945 and deregulated expression of immune checkpoints in colorectal cancer. *J Exp Clin Cancer*
946 *Res*, 2019, 38:364
- 947 76. Chen W, Pearlman R, Hampel H, Pritchard CC, Markow M, Arnold C, Knight D, Frankel
948 WL. MSH6 immunohistochemical heterogeneity in colorectal cancer: comparative
949 sequencing from different tumor areas. *Hum Pathol*, 2020, 96:104–11
950

951 **Figures**

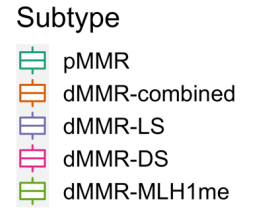
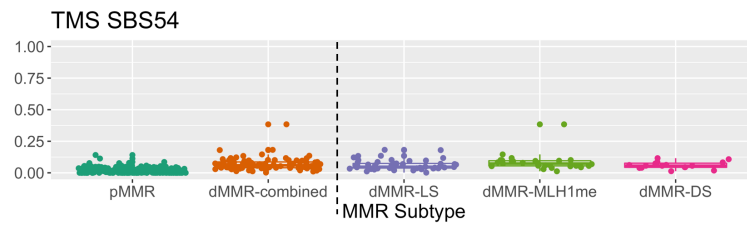
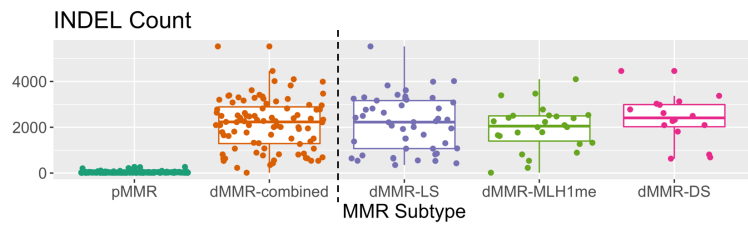
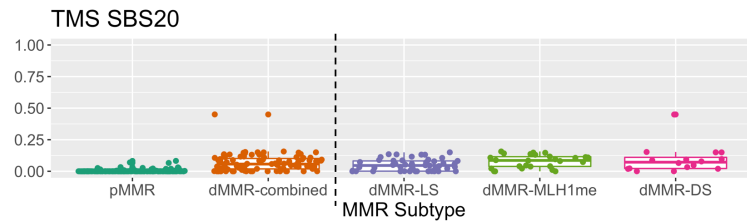
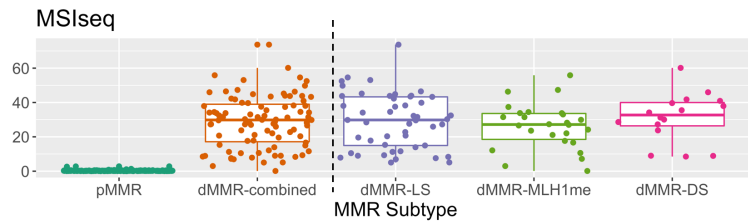
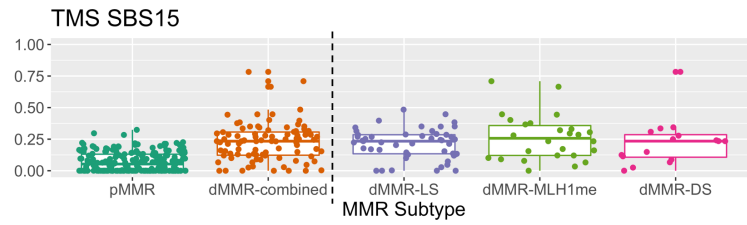
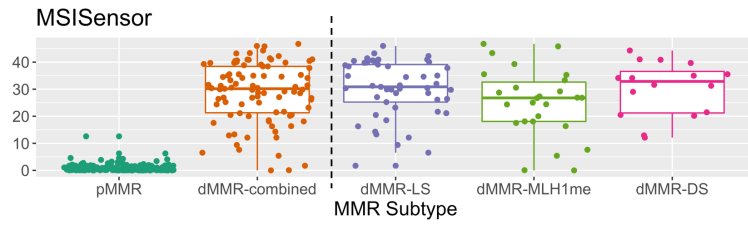
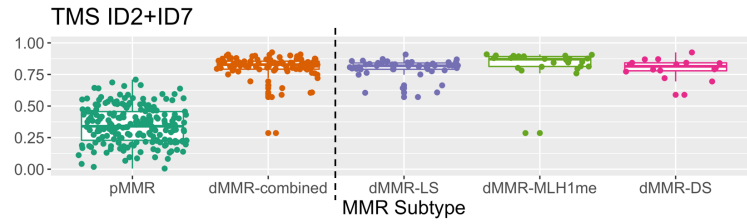
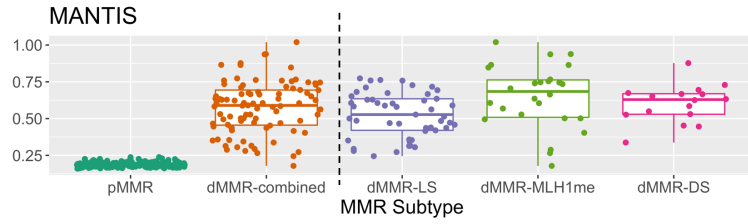
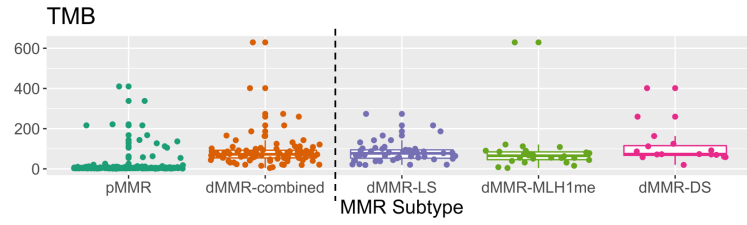
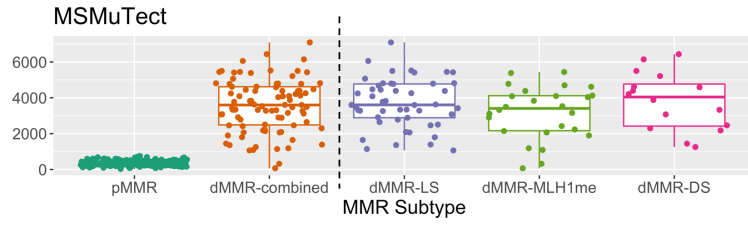
Analysis 1. Assessment of tumor features for dMMR prediction accuracy in WES CRCs



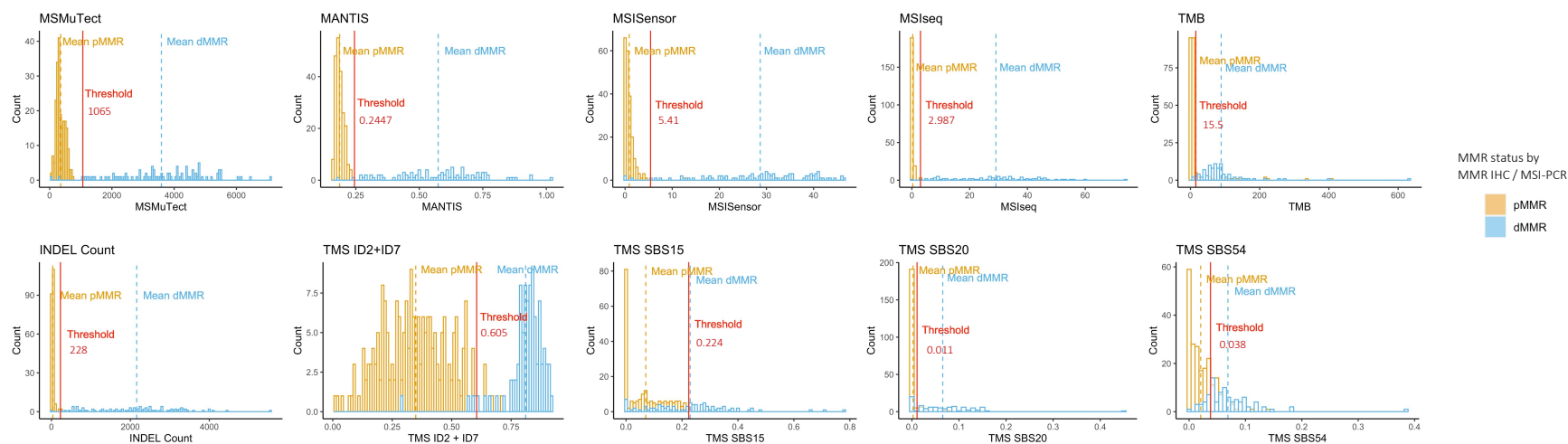
Analysis 2. Assessment of A) individual tumor features, B) statistical model and C) additive feature combination approaches derived from the WES analysis on panel sequenced CRCs, ECs and SSTs



953 **Figure 1.** Overview of the study design. In total, 300 whole-exome sequenced (WES) colorectal
954 cancers (CRCs) consisting of 91 DNA mismatch repair deficient (dMMR) and 209 DNA mismatch
955 repair proficient (pMMR) tumors were analyzed. We investigated 104 tumor features for their
956 ability to distinguish dMMR from pMMR tumors consisting of four MSI tools, 97 tumor
957 mutational signature definitions (TMS), tumor mutation burden (TMB) calculated as mutations
958 per mega base, somatic insertion / deletion (INDEL) and somatic single nucleotide variant (SNV)
959 counts. We performed a 10-fold cross-validation approach with 100 repeats to calculate the mean
960 accuracy on the test dataset. (A) The top 10 ranked individual tumor features, (B) a Lasso
961 regression model and (C) an additive feature combination approach was tested to determine the
962 benefit of combining tumor features to improve dMMR prediction. The findings from these three
963 approaches were tested on an independent set of targeted panel sequenced tumors of CRC,
964 endometrial cancer (EC) and sebaceous skin tumor (SST) tissue types with reported mean
965 accuracies.



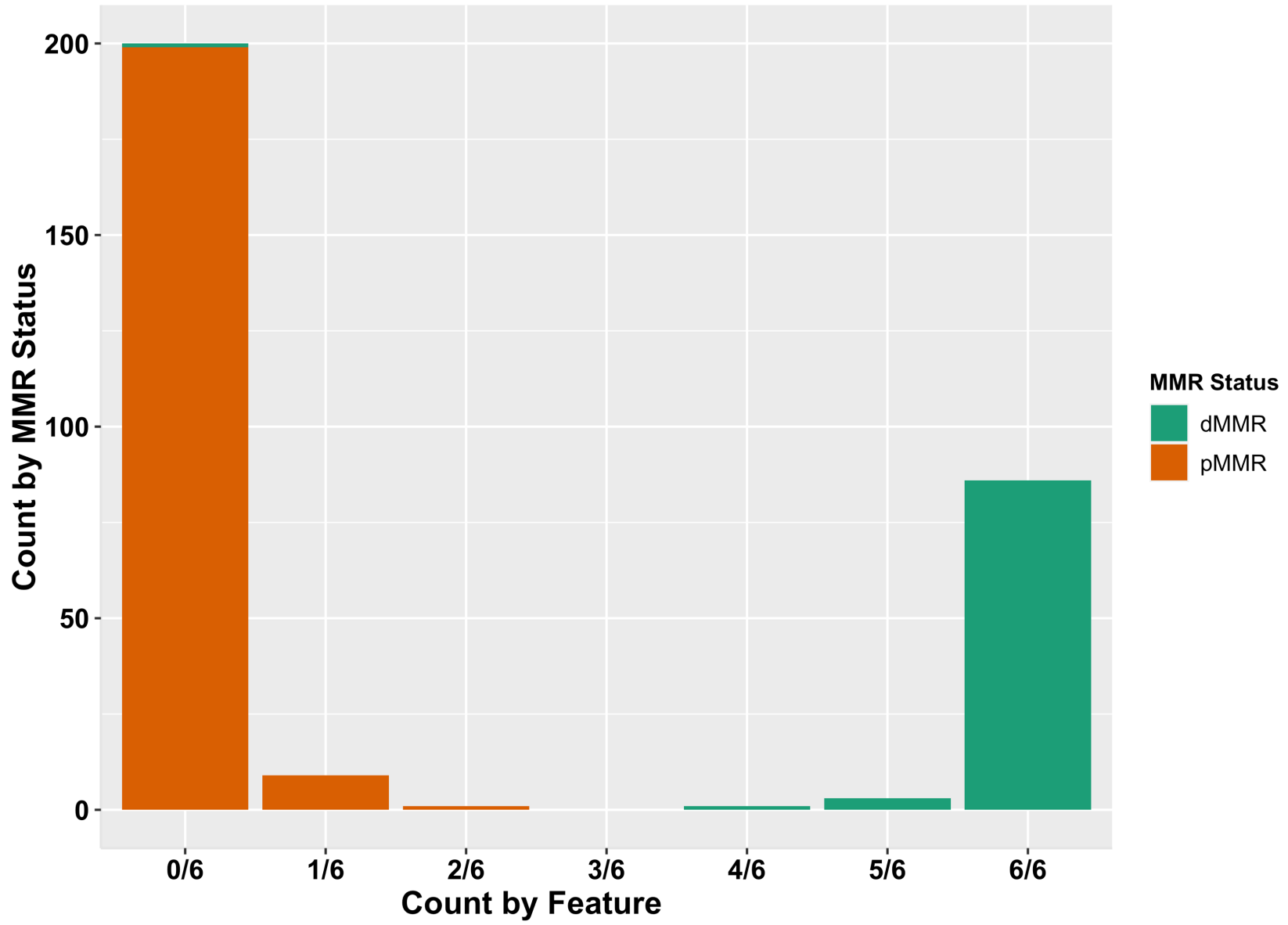
967 **Figure 2.** Tumor distribution of the top 10 DNA mismatch repair (MMR) deficient (dMMR) predicting features in the whole exome
968 sequenced (WES) colorectal cancers (CRCs) by MMR subtype. Boxplots showing the distribution of tumors by MMR status (MMR-
969 proficient (pMMR) versus dMMR) as well as stratified by dMMR subtype - dMMR-LS (Lynch syndrome), dMMR-DS (double somatic
970 MMR gene mutations) and dMMR-MLH1me (*MLH1* promoter methylation) for each of the top 10 predicting features MSMuTect,
971 MANTIS, MSISensor, MSISEq, INDEL (insertion / deletion) count, TMB (tumor mutation burden calculated as mutations / mega base),
972 TMS (tumor mutational signature) ID2+ID7, TMS SBS15, TMS SBS20 and TMS SBS54 as determined from the WES CRC analysis.
973 ID, small insertions / deletions; SBS, single base substitution.



974

975 **Figure 3.** Determination of thresholds for differentiating DNA mismatch repair (MMR) deficient (dMMR) from MMR-proficient
 976 (pMMR) colorectal cancers (CRCs) using whole exome sequencing (WES) data for each of the top 10 performing tumor features. Bar
 977 graphs presenting the distribution of tumors after applying the recommended thresholds (red line) for each of the top 10 predicting tumor
 978 features MSMuTect, MANTIS, MSISensor, MSIseq, INDEL count, TMB, TMS ID2+ID7, TMS SBS15, TMS SBS20 and TMS SBS54
 979 as determined from the WES CRC analysis. Orange coloring indicates pMMR and blue coloring represents dMMR status. ID, small
 980 insertions / deletions; SBS, single base substitution.

Feature Counting by MMR Status for WES Tumors



982 **Figure 4.** The additive tumor feature combination approach demonstrating the distribution of counts of the top six tumor features by the
983 DNA mismatch repair (MMR) status of the 300 colorectal cancers (CRCs) with whole exome sequencing (WES). Bar graphs presenting
984 the distribution of tumors after applying the additive tumor feature combination approach with the recommended thresholds from the
985 WES CRC analysis using a count of ≥ 3 out of the top six predictors from the WES CRC analysis, consisting of MSMuTect, MANTIS,
986 MSIseq, MSISensor, INDEL (insertion / deletion) count and TMS (tumor mutational signature) ID2+ID7 (small insertions / deletions)
987 for MMR status calling: MMR-deficient (dMMR) versus MMR-proficient (pMMR).

988 **Tables**

989 **Table 1.** The breakdown of the 104 tumor features calculated from next generation sequencing analysis included in this study.

<i>Feature Type</i>	<i>Count</i>	<i>Name</i>	<i>Reference</i>
<i>Total</i>	<i>N = 104</i>		
Microsatellite instability (MSI) Tools	N = 4	MSISensor	Niu <i>et al.</i> , 2014
		MSIseq	Huang <i>et al.</i> , 2015
		MANTIS	Kautto <i>et al.</i> , 2017
		MSMuTect	Maruvka <i>et al.</i> , 2017
Tumor mutational signatures (TMS)	N = 97	SBS (N = 78)	Tate <i>et al.</i> , 2018
		ID (N = 18)	Tate <i>et al.</i> , 2018
		ID2+ID7	Georgeson <i>et al.</i> , 2021
Somatic mutation counts	N = 3	INDELs	
		SNVs	
		TMB (SNVs + INDELs/ MB)	Muzny <i>et al.</i> , 2012

990 The 104 tumor features can be categorized into three distinct groups: microsatellite instability (MSI) tools, tumor mutational signatures
 991 (TMS) and somatic mutation counts. These features have previously been shown to be associated with MSI / DNA mismatch repair
 992 status as indicated by the provided references. The MSI group consists of four MSI tools namely MSISensor, MSIseq, MANTIS and
 993 MSMuTect. TMS consisted of 78 single base substitutions (SBS), 18 small insertions / deletions (IDs) and TMS ID2+ID7. The somatic

994 mutation count consisted of the single nucleotide variant count, larger insertions / deletions count and the tumor mutation burden (TMB),

995 which was calculated as the combination of SNVs and INDELs counts per megabase.

996

997 **Table 2.** Performance of the top tumor features demonstrating a prediction accuracy >80% ranked by highest mean accuracy from
 998 whole-exome sequenced (WES) colorectal cancers (CRCs).

Tumor Feature	Mean Accuracy	Error Rate	95% CI: (Accuracy)	Mean Sensitivity	95% CI: (Sensitivity)	Mean Specificity	95% CI: (Specificity)	Mean AUC	95% CI: (AUC)
MSMuTect	99.3%	0.7%	99.1% - 99.5%	97.6%	96.9% - 98.3%	100.0%	-	98.8%	98.5% - 99.1%
MSIseq	99.1%	0.9%	98.9% - 99.4%	97.7%	97.0% - 98.3%	99.8%	99.6% - 100.0%	98.7%	98.4% - 99.1%
MANTIS	99.0%	1.0%	98.8% - 99.2%	97.1%	96.4% - 97.7%	99.9%	99.8% - 100.0%	98.5%	98.1% - 98.8%
INDEL count	98.9%	1.1%	98.7% - 99.2%	97.7%	97.0% - 98.3%	99.5%	99.2% - 99.8%	98.6%	98.2% - 98.9%
MSISensor	97.7%	2.3%	97.3% - 98.0%	93.4%	92.4% - 94.5%	99.5%	99.3% - 99.7%	96.5%	96.0% - 97.0%
TMS ID2+ID7	96.8%	3.2%	96.4% - 97.2%	94.2%	93.2% - 95.2%	97.9%	97.5% - 98.4%	96.0%	95.5% - 96.6%

TMS ID2	93.3%	6.7%	92.8% - 93.8%	90.7%	89.5% - 91.9%	94.4%	93.7% - 95.1%	92.6%	92.0% - 93.1%
TMS SBS20	88.4%	11.6%	87.6% - 89.2%	68.9%	66.6% - 71.2%	97.0%	96.4% - 97.6%	82.9%	81.8% - 84.1%
TMS ID7	87.6%	12.4%	87.0% - 88.3%	74.2%	72.6% - 75.9%	93.5%	92.8% - 94.2%	83.9%	83.0% - 84.7%
TMS SBS54	83.4%	16.6%	82.6% - 84.2%	59.4%	57.5% - 61.4%	93.9%	93.1% - 94.7%	76.7%	75.6% - 77.7%
TMB	83.3%	16.7%	82.6% - 83.9%	57.8%	55.2% - 60.4%	94.5%	93.7% - 95.2%	76.1%	75.0% - 77.3%
TMS SBS15	82.4%	17.6%	81.5% - 83.3%	58.8%	56.5% - 61.1%	92.8%	91.9% - 93.7%	75.8%	74.6% - 77.0%

999 The mean accuracy values after 10-fold cross-validation with 100 repeats, error rate, mean sensitivity, mean specificity, and mean area
1000 under the curves (AUCs) with corresponding 95% confidence intervals (CIs) are shown for each of the top 10 predicting tumor features
1001 MSMuTect, MSIseq, MANTIS, INDEL (insertion / deletion) count, MSISensor, TMS (tumor mutational signature) ID2+ID7, TMS
1002 ID2, TMS SBS20, TMS ID7, TMS SBS54, TMB (tumor mutation burden) and TMS SBS15 from the WES CRC analysis. ID, small
1003 insertions, and deletions; SBS, single base substitutions.

1004

1005 **Table 3.** Summary of the best dMMR prediction results by individual tumor feature, Lasso regression model and the additive feature
 1006 combination approach for the whole-exome sequencing (WES) colorectal cancers (CRCs) and the panel sequenced CRCs, endometrial
 1007 cancers (ECs) and sebaceous skin tumors (SST).

	Performance of best Individual Feature		Performance of Statistical Model		Performance of Additive Feature Combination Approach	
WES	Feature	Mean Accuracy	Lasso	Mean Accuracy	Feature Combination	Mean Accuracy
CRC	MSMuTect	99.3%	MANTIS + TMS ID2+ID7 + MSISensor + TMS SBS15	98.3%	MSMuTect + MANTIS + MSIseq + MSISensor + INDEL count + TMS ID2+ID7	99.7%
PANEL	Feature	Accuracy	Lasso	Accuracy	Feature Combination	Accuracy
CRC	MANTIS	100.0%	MANTIS + TMS ID2+ID7 + MSISensor + TMS SBS15	89.7%	MSMuTect + MANTIS + MSIseq + MSISensor + INDEL count + TMS ID2+ID7	100.0%
EC	MANTIS	86.4%	MANTIS + TMS ID2+ID7 + MSISensor + TMS SBS15	68.2%	MSMuTect + MANTIS + MSIseq + MSISensor + INDEL count + TMS ID2+ID7	95.5%
SST	MANTIS	85.0%	MANTIS + TMS ID2+ID7 + MSISensor + TMS SBS15	85.0%	MSMuTect + MANTIS + MSIseq + MSISensor + INDEL count + TMS ID2+ID7	100.0%

1008 This table provides the top performing results from A) individual tumor feature, B) statistical model application (Lasso) and C) additive
1009 feature combination approach assessments for WES CRCs as well as targeted panel sequenced CRCs, ECs and SSTs.
1010 TMS, tumor mutational signature; ID, small insertions, and deletions; SBS, single base substitution; INDEL count, insertions / deletions.
1011

1012 **Table 4.** Assessment of top performing tumor features from whole-exome sequenced (WES) colorectal cancers (CRCs) in panel
 1013 sequenced CRC, endometrial cancer (EC) and sebaceous skin tumor (SST) test sets.

Tumor Feature	CRC			EC			SST		
	Mean Accuracy	95% CI	Error Rate	Mean Accuracy	95% CI	Error Rate	Mean Accuracy	95% CI	Error Rate
MSMuTect	27.6%	12.7% - 47.2%	72.4%	18.2%	5.2% - 40.3%	81.8%	35.0%	15.4% - 59.2%	65.0%
MSIseq	82.8%	64.2% - 94.2%	17.2%	68.2%	45.1% - 86.1%	31.8%	65.0%	40.8% - 84.6%	35.0%
MANTIS	100.0%	88.1% - 100.0%	0.0%	86.4%	65.1% - 97.1%	13.6%	85.0%	62.1% - 96.8%	15.0%
INDEL count	27.6%	12.7% - 47.2%	72.4%	18.2%	5.2% - 40.3%	81.8%	35.0%	15.4% - 59.2%	65.0%
MSISensor	96.6%	82.2% - 99.9%	3.4%	77.3%	54.6% - 92.2%	22.7%	75.0%	50.9% - 91.3%	25.0%
TMS ID2+ID7	82.8%	64.2% - 94.2%	17.2%	63.6%	40.7% - 82.8%	36.4%	85.0%	62.1% - 96.8%	15.0%

TMS SBS20	69.0%	49.2% - 84.7%	31.0%	50.0%	28.2%	-	50.0%	40.0%	19.1%	-	60.0%
					71.8%				63.9%		
TMS SBS54	51.7%	32.5% - 70.6%	48.3%	36.4%	17.2%	-	63.6%	40.0%	19.1%	-	60.0%
					59.3%				63.9%		
TMB	44.8%	26.4% - 64.3%	55.2%	31.8%	13.9%	-	68.2%	35.0%	15.4%	-	65.0%
					54.9%				59.2%		
TMS SBS15	44.8%	26.4% - 64.3%	55.2%	27.3%	10.7%	-	72.7%	60.0%	36.1%	-	40.0%
					50.2%				80.9%		

1014 Table presents the prediction accuracies, error rates and corresponding 95% confidence intervals (CIs) for panel sequenced CRCs, ECs
1015 and SSTs for the top 10 predicting tumor features MSMuTect, MSIseq, MANTIS, INDEL (insertions / deletions count), MSISensor,
1016 TMS (tumor mutational signature) ID2+ID7, TMS SBS20, TMS SBS54, TMB (tumor mutation burden, mutations / mega base) and
1017 TMS SBS15 from WES CRC analysis applied on panel sequenced CRCs, ECs and SSTs. ID, small insertions, and deletions; SBS, single
1018 base substitution.