

Mutational signatures driven by epigenetic determinants stratify patients for therapeutic interventions in gastric cancer

Jaqueline Ramalho Buttura^{#1}, Monize Nakamoto Provisor^{#1,2}, Renan Valieris¹, Vinicius Fernando Calsavara³, Rodrigo Duarte Drummond¹, Alexandre Defelicibus¹, Joao Paulo Lima¹, Helano Carioca Freitas^{4,5}, Vladmir C. Cordeiro Lima⁴, Thais Fernanda Bartelli⁵, Marc Wiedner⁶, Rafael Rosales⁷, Kenneth John Gollob⁸, Joanna Loizou^{6,9}, Emmanuel Dias-Neto^{5,10}, Diana Noronha Nunes⁵, Israel Tojal da Silva¹

¹Laboratory of Bioinformatics and Computational Biology, A.C. Camargo Cancer Center, So Paulo, SP 01509-010, Brazil.

²present address: Department of Genomics, Fleury Group, So Paulo, SP, Brazil.

³Department of Statistics and Epidemiology, A.C. Camargo Cancer Center, So Paulo, SP, Brazil.

⁴Medical Oncology Department, A.C. Camargo Cancer Center, So Paulo, SP, Brazil.

⁵Laboratory of Medical Genomics, A.C. Camargo Cancer Center, So Paulo, SP 01509-010, Brazil.

⁶CeMM Research Center for Molecular Medicine of the Austrian Academy of Sciences, Lazarettgasse 14, AKH BT 25.3, 1090 Vienna, Austria.

⁷University of So Paulo. Department of Mathematics and Computer Science. 14049-900, Ribeiro Preto, SP, Brazil.

⁸Translational Immuno-oncology Group, A.C. Camargo Cancer Center, So Paulo, SP 01509-010, Brazil.

⁹Institute of Cancer Research, Department of Medicine I, Medical University of Vienna and Comprehensive Cancer Center, Vienna 1090, Austria.

¹⁰University of So Paulo. Laboratory of Neurosciences, Institute of Psychiatry. So Paulo, SP, Brazil.

- equal contribution;

Abstract

DNA mismatch repair deficiency (dMMR) leads to increased mutation load, which in turn may impact anti-tumor immune responses and treatment effectiveness. Currently, there are different mutational signatures described in primary cancers that are associated with dMMR. Whether the somatic and epigenetic changes in MMR genes precede one or more dMMR signa-

Preprint submitted to medRxiv

April 17, 2020

tures, and if so by which mechanism remains unknown. To investigate the relationship between these changes and dMMR signatures, we performed a *de novo* extraction of mutational signatures in a large cohort of 787 gastric cancer patients. We detected three dMMR-related signatures, one of which clearly discriminates tumors with *MLH1* gene silencing caused by hypermethylation within its promoter (AUC = 98%). We then demonstrate that samples with the highest exposures to signature share features related to better prognosis, encompassing clinical and molecular aspects, as well as altered immune infiltrate composition, predictive of a better response to immune checkpoint inhibitors. Overall, our analysis explored the impact of modifications in MMR-related genes on shaping specific mutational signatures and we provide evidence that patient classification based on mutational signature exposure can identify a group of patients with a good prognosis and who are potentially good candidates for immunotherapy.

Keywords: Mutational Signature, Bioinformatics, Gastric Cancer, DNA, Mismatch Repair Prognosis

1 Introduction

Cancer results from the sequential accumulation of DNA alterations, including single nucleotide mutations [1] that arise from different endogenous or exogenous processes [2]. Distinct DNA-damaging processes leave characteristic nucleotide base-change footprints known as mutational signatures [3]. Previous studies [4] have extracted distinct mutational signatures by examining a large set of human cancer genomes and some of these have been reported in the COSMIC database (denoted hereafter as CS). This pan-cancer analysis revealed significant heterogeneity of operational mutational processes, that encompass mutation-triggering events as diverse as the off-target activity AID/APOBEC family of cytidine deaminases, the exposure to ultraviolet light, tobacco-smoking and the defective DNA mismatch repair [5, 6].

Collectively, the understanding of the mechanistic basis of mutational signatures, as well as to their etiology, may provide clues for cancer diagnosis and hold prognostic value [7]. For example, six mutational signatures have been associated with BRCA1/BRCA2 dysfunction, which most likely are predictive of response to treatment with PARP inhibitors [8]. Thus, homologous recombination repair (HRR)-deficiency features based on these signatures

20 allowed the prediction of BRCAness in breast cancer patients with 98.7%
21 sensitivity [8]. Additionally, given that nucleotide excision repair (NER) de-
22 ficient tumors are more sensitive to certain treatments, somatic variations in
23 the *ERCC2* gene, which encodes a key protein of the NER pathway, have
24 also been linked with characteristic mutational signatures [9, 10]. Other mu-
25 tational processes are associated with patients harboring biallelic *MUTYH*
26 germline mutations [11], a finding that may indicate deficient base excision
27 repair (BER). Such patients are eligible for genetic counseling [12] and might
28 benefit from immunotherapy [13].

29
30 In addition to HRR, NER and BER repair pathways, another mecha-
31 nism underlying oncogenic genomic variations, with important effects on
32 anti-tumor immune responses occur in tumors with impaired DNA mismatch
33 repair (MMR), which harbor elevated frequencies of single-nucleotide vari-
34 ants (SNVs) and exceptionally high indel rates [14]. Recent studies demon-
35 strated that various MMR-deficient (dMMR) tumor types (gastrointestinal,
36 glioblastoma, endometrial and prostate) are more responsive to programmed
37 cell death protein 1 (PD1) immune checkpoint inhibitors as compared to
38 MMR-proficient tumors [15, 16, 17]. A set of four mutational signatures
39 (CS-6, CS-15, CS-20, and CS-26) have been associated with dMMR. Never-
40 theless, it is still unclear if somatic and epigenetic changes in MMR genes
41 lead to one or more dMMR signatures.

42
43 In this study we investigated the significance of molecular events in MMR-
44 genes that shape characteristic mutation signatures found in MMR-deficient
45 gastric adenocarcinomas. The presence of these signatures was evaluated for
46 their prognostic value in a cohort of 787 gastric cancer patients with pub-
47 licly available data, including 439 patients from the TCGA, and validated
48 in a second cohort composed of 170 gastric cancer patients [18]. We fur-
49 ther investigated whether local tumor immune response and prognosis varied
50 according to MMR-deficiency exposure load. The consequences of these ap-
51 pear to be predictive of the responsiveness to immune checkpoint blockade
52 and may be used to support treatment strategies in the future.

53 2. Materials and Methods

54 2.1. Clinical and genomic data from public cohort

55 The non-redundant public cohorts assessed here contained clinical and
56 molecular information of gastric adenocarcinoma samples provided by: i)
57 The Cancer Genome Atlas (TCGA, N=439), ii) cBioPortal, N=226; iii) and
58 International Cancer Genome Consortium (ICGC, N=122), totaling 787 pa-
59 tients (Supplementary Material Table S1). TCGA data was assessed on
60 October 4th, 2018 and corresponds to the MC3 variant calling project, which
61 is a comprehensive effort to detect consensus mutations and forms the basis
62 of Pan-Cancer Atlas initiative [19]. cBioPortal and ICGC cohorts' com-
63 prise Asian samples which were last assessed on January 9th, 2019. Raw
64 reads from matched non-tumor exomes from TCGA dataset, encompassing
65 the MMR genes were downloaded and used to detect the germline SNVs
66 following the Genome Analysis Toolkits (GATK)s best practice for germi-
67 line alterations calling. We also used additional filters considering mutations
68 with VAF (variant allele frequency) ≥ 0.3 and minimum depth coverage of
69 10 reads. Furthermore, *dbNSFP_MetaLR_rankscore* was used to filter out
70 (≤ 0.6) the synonymous mutations. The methylation levels in the form of
71 beta-values ranging from 0 to 1 were addressed for the TCGA cohort [20].
72 We then used the CpG sites in the promoter of MMR genes to detect those
73 that were hypermethylated or hypomethylated. The baseline clinical features
74 are summarized in Supplementary Material Table S2.

75 2.2. Clinical and genomic data from validation cohort

76 Patients in the validation cohort were prospectively enrolled in an insti-
77 tutional study to unveil the epidemiology and genomics of gastric adenocar-
78 cinomas in Brazil [18]. This study was approved by the local ethics com-
79 mittee and all participants provided written informed consent. An overview
80 of the clinical characteristics of patients in the validation cohort is provided
81 in Table S3. Genomic DNA from frozen tissues (n=165) was extracted with
82 AllPrep DNA/RNA Mini Kit (Qiagen), QIASymphony THC 400 (Qiagen)
83 or phenol/chloroform/isoamyl alcohol precipitation. gDNA from FFPE tis-
84 sue (n=4) was extracted with RecoverAll Total Nucleic Acid Isolation Kit
85 (Thermo Fisher), and there was one sample from gastric wash. Exome li-
86 braries were prepared using Agilent SureSelect V6 kit and sequenced us-
87 ing Illumina platforms (HiSeq4000, 100bp, n=33; Novaseq, 150bp, n=137 -

88 pairedend reads for both). The raw sequencing data (.fastq files) were de-
89 posited in SRA (<http://www.ncbi.nlm.nih.gov/sra>) under accession number
90 PRJNA505810.

91
92 For our local independent validation cohort, the somatic SNVs were called
93 by using an in-house pipeline following the Broad Institute GATK Best Prac-
94 tices guidelines [21] as described [6]. Briefly, the raw reads were aligned
95 using Burrows Wheeler Aligner (BWA-mem) with default settings to assem-
96 bly GRCh38. Next, alignment files in SAM format were converted to BAM
97 files, sorted and filtered to exclude reads with mapq score <15. The re-
98 tained reads were processed using SAMtools (v1.9) and Picard (v3.8) ([https :](https://broadinstitute.github.io/picard/)
99 [//broadinstitute.github.io/picard/](https://broadinstitute.github.io/picard/)) respectively, which excludes low-quality
100 reads and PCR duplicates. Finally, the somatic SNVs calling was performed
101 for the whole exome data from analysis-ready BAM files using Mutect2 (v3.8)
102 for tumor samples and further with a panel of 16 unmatched non-tumor leuko-
103 cyte samples. Extensive filtering was applied to remove low mapping quality,
104 as well as strand, position bias and OxoG oxidative artifacts. Furthermore,
105 any residual germline mutations from the database of germline mutations
106 of of gnomAD (<https://gnomad.broadinstitute.org/>) and Online Archive of
107 Brazilian Mutations (ABraOM, available at <http://abraom.ib.usp.br/>) were
108 removed.

109 *2.3. Mutational signatures estimation*

110 All somatic SNVs of the six classes (C>A, C>G, C>T, T>A, T>C and
111 T>G) were mapped onto trinucleotide sequences by including the 5' and 3'
112 neighboring base-contexts. Next, the SNV spectrum with 96 trinucleotide
113 mutations types for all samples were loaded into signeR [22] to estimate
114 the optimal number of mutational signatures, which is based on the me-
115 dian Bayesian Information Criterion (BIC) value. We next used the cosine
116 similarity to compare the extracted *de novo* mutational signatures to those
117 described in the COSMIC signatures (v2), considering cosine similarity > 0.7
118 as a measure of closeness to COSMIC signatures. Patients with higher ex-
119 posure for a given signature (Exposure value greater or equal to the third
120 quartile) were named as *high* and those with lower exposure values (Exposure
121 value less than third quartile) were named as *low*.

122 2.4. Molecular features

123 We used the MSIseq [23] software for microsatellite instability (MSI) sta-
124 tus prediction (MSI-H and Non-MSI-H) from whole exome data. Briefly, this
125 software is based on four machine-learning frameworks, which requires a cat-
126 alog of somatic SNVs and microindels of samples, a file containing the exact
127 locations of mononucleotides ($\text{length} \geq 5$) and microsatellites consisting of di,
128 tri, and tetranucleotide repeats, as annotated in the simpleRepeats track and
129 available at [http : //hgdownload.cse.ucsc.edu/goldenpath/hg19/database/](http://hgdownload.cse.ucsc.edu/goldenpath/hg19/database/).
130 MSIseq is available at The Comprehensive R Archive Network (CRAN).

131
132 Consistent with a method previously proposed by Chalmers et al. [24],
133 the Tumor Mutational Burden (TMB) was calculated as the total number of
134 mutations divided by the length of the target region in megabases.

135
136 The tumor heterogeneity estimation was performed using math.score (MATH)
137 function from package maftools version 3.8 [25]. A higher MATH score indi-
138 cates increased tumor heterogeneity.

139
140 Neoantigen count: The list of neoantigens available for 77 TCGA-STAD
141 samples was extracted from The Cancer Immunome Atlas (TCIA) [26], as-
142 sessed on June 17th, 2017 at [https : //tcia.at/neoantigens](https://tcia.at/neoantigens).

143 2.5. Statistical analyses

144 The baseline patient characteristics are expressed as absolute and relative
145 frequencies for qualitative variables and as the mean \pm standard deviation
146 (SD) for quantitative variables. Mutational signature exposure and TMB
147 were considered as continuous variables. The association between qualita-
148 tive variables was evaluated by chi-squared test or Fishers exact test, as
149 appropriate.

150
151 Overall survival functions were estimated by the Kaplan-Meier estima-
152 tor and the log-rank test was used to compare the survival functions among
153 groups (eg, patients with higher mutational signature exposure (S^{high}) ver-
154 sus other (S^{low})). The Cox semiparametric proportional hazards model was
155 fitted to the dataset to describe the relationship between overall survival and
156 the main clinical features. Hazard ratio (HR) and 95% confidence intervals
157 (95%CI) were calculated for all variables. A backward stepwise selection al-
158 gorithm was applied, with different significance levels to enter ($p=0.10$) and

159 remain ($p=0.05$) in the model. Variables were removed from the model if they
160 were non-significant or acted as confounders (change in coefficient $>20\%$).
161 The proportional hazards assumption was assessed based on the Schoenfeld
162 residuals [27]. There was evidence that covariates had a constant effect over
163 time in all cases.

164

165 Multivariate analyses were performed considering the main clinical fea-
166 tures (such as age, pathological stage, Lauren tumor subtype and ethnicity),
167 previously associated with overall survival, and with exposures of mutational
168 signatures associated with dMMR, besides molecular features TMB and MSI
169 status. Forest plots were created based on the final multiple Cox regression
170 model. Metastatic patients were excluded from these analyses. In addition,
171 we fitted simple and multiple logistic regression models in order to assess the
172 effect of S2, S4 and S5 exposures in the *MLH1* methylation. Overall perfor-
173 mance, calibration, and the discriminatory power of the final multiple logistic
174 regression model were assessed using the Brier score, the HosmerLemeshow
175 goodness-of-fit test, and the area under the receiver operating characteristic
176 (ROC) curve (AUC), respectively [28]. Besides, we assessed the goodness-of-
177 fit through a Q-Q plot. The significance level was fixed at 5% for all tests
178 (two-sided). Statistical analysis was performed using R software (v3.5).

179 *2.6. Mutational signatures in cell lines*

180 The CRISPR-Cas9 knockout clones for *MLH1* were generated in human
181 HAP1 cells using the following guide RNA (gRNA) sequence: 5 - AAGA-
182 CAATGGCACCGGGATC - 3. Clonal populations with a frameshift muta-
183 tion within *MLH1* were subsequently cultured for three months to allow for
184 the accumulation of mutations during cellular division [29]. To identify muta-
185 tions, genomic DNA was submitted to whole genome sequencing (WGS). *De*
186 *nov*o somatic mutations including substitutions, indels and rearrangements
187 in subclones were obtained by removing all mutations seen in parental clones.
188 Next, SNVs were mapped onto trinucleotide sequences by including the 5'
189 and 3' neighboring base-contexts and then the level of samples' exposure to
190 previously found mutational signatures was estimated [22].

191 *2.7. Significantly mutated genes and pathway analysis*

192 To assess the impact of dMMR pathway on the genes throughout the
193 genome, we searched for genes more frequently mutated than would be ex-
194 pected by chance [30]. The proper gene symbol annotation in MAF (Mu-

195 tation Annotation Format) files was addressed by maftools (v.8) [25] (pre-
196 pareMutSig function) and then loaded into online MutSigCV server (v.1.3.4)
197 (<https://cloud.genepattern.org/gp/pages/index.jsf>). The oncoplots were
198 built by using the significantly mutated genes from MutSigCV analysis. The
199 significantly mutated genes associated to $S4^{high}$ and $S4^{low}$ group were entered
200 into Gene Set Enrichment Analysis (GSEA) according to the Investigate gene
201 sets function available at MSigDB (Molecular Signatures Database (v7.0),
202 <http://software.broadinstitute.org/gsea/msigdb/index.jsp>). KEGG, RE-
203 ACTOME, GO biological process, oncogenic signatures (module C6) and im-
204 munologic signatures (module C7) were also considered to compute overlaps,
205 and the top 20 gene sets with FDR q-value <0.05 were used to summarize
206 these analysis.

207 *2.8. Inflammatory infiltrate and immune aspects*

208 We estimated the cellular composition from the bulk expression datasets
209 (TCGA) by using two complementary approaches. For both analyses, FPKM
210 (Fragments Per Kilobase Million) from 380 tumor TCGA-STAD samples
211 were used as normalized gene expression profile, retrieved on January 22th,
212 2018. First, the CIBERSORT software based on the deconvolution method
213 for characterizing cell composition of complex tissues from their gene ex-
214 pression profiles, was used [31]. CIBESORT takes advantage of a validated
215 leukocyte gene signature matrix, termed LM22. This gene signature contains
216 547 genes that distinguish 22 human hematopoietic cell phenotypes, includ-
217 ing seven T cell types, nave and memory B cells, plasma cells, natural killer
218 (NK) cells, and myeloid subsets. Simultaneously, a recent technique based on
219 gene set enrichment analysis (GSEA) termed as xCell [32] was used to infer
220 34 immune cell types. Herein, we used this method to confirm the findings
221 by CIBERSORT.

222
223 CIBERSORT analysis was performed online using a public server (<http://cibersort.stanford.edu/>) for characterizing absolute and relative immune
224 cell composition with 1000 permutations and disabled quantile normaliza-
225 tion as set parameters. From the 380 TCGA-STAD tumor samples, 215
226 215 samples (56%) yielded data on infiltrating immune cells (p-value <0.05),
227 which were considered for further analysis (50 samples as $S4^{high}$ and 165 as
228 $S4^{low}$). We also used the second approach known as xCell to reinforce the
229 findings when comparing $S4^{high}$ and $S4^{low}$ samples. xCell analysis was per-
230 formed using the R package with default parameters (available at <https://>

232 //github.com/dviraran/xCell). In order to verify the immune effector re-
233 sponse present in $S4^{high}$ and $S4^{low}$ samples, differential expression of key im-
234 munoregulatory/inflammatory or cytotoxic markers was also performed. The
235 comparison of the groups in this section was performed by *Mann-Whitney U*
236 Test with statistical significance set at $p\text{-value} < 0.05$.

237

238 We also used the pre-processed immune subtypes previously described by
239 Thorsson et al. [33] for TCGA samples (available for 103 in $S4^{high}$ samples
240 and 285 in $S4^{low}$ which can be assessed in Table S2).

241 3. Results

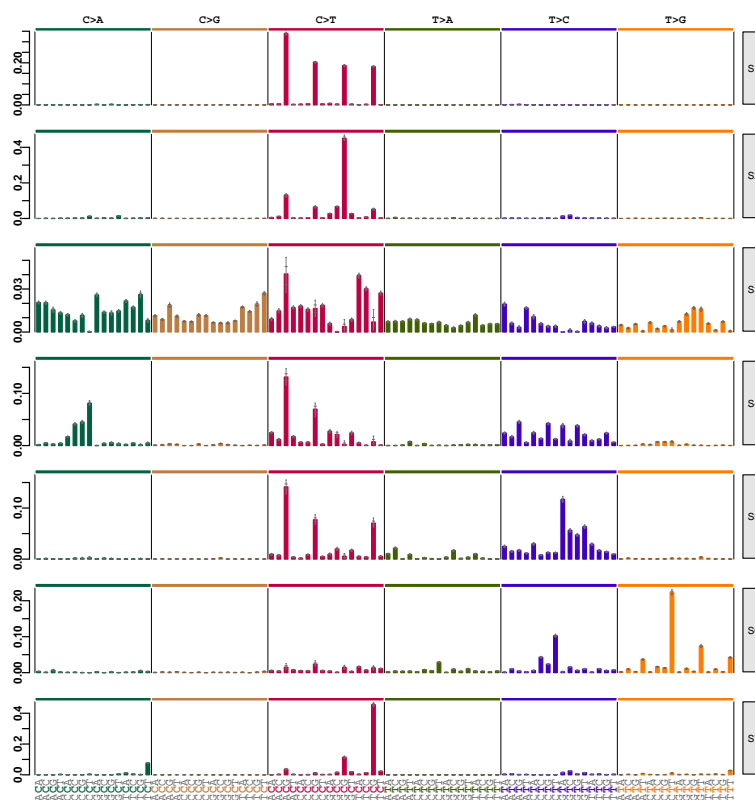
242 3.1. Mutational signatures

243 Using *signeR* [22] analysis to estimate *de novo* mutation signatures across
244 three gastric cancer cohorts, we identified seven (denoted hereafter as S[1-7])
245 mutational signatures (Figure 1A) which are related to signatures described
246 in the COSMIC database by cosine similarity scores (Figure 1B). Signature
247 1 (S1) is associated with endogenous mutational processes initiated by spon-
248 taneous deamination of 5-methylcytosine (CS-1); Signatures S2, S4 and S5
249 are associated with defective DNA mismatch repair and/or microsatellite in-
250 stability (CS-6/CS-15, CS-20 and CS-21/CS-26 respectively); Signature S3
251 is related with failure of DNA-double strand break repair by homologous re-
252 combination (CS-3); Signature S6 related to CS-17 with unknown etiology;
253 and Signature S7 is associated with error-prone polymerase activity (POLE
254 (DNA Polymerase Epsilon, Catalytic Subunit), CS-10).

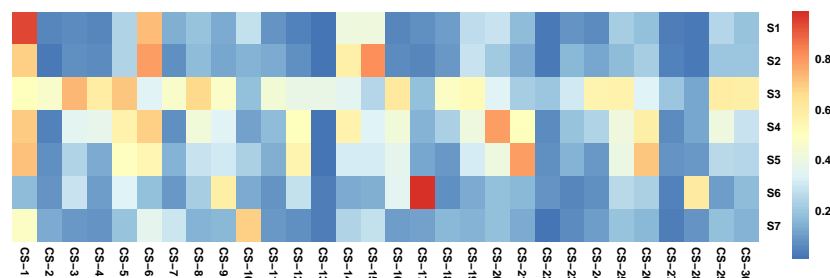
255

256 CS-3 (S3 in Figure 1 and Figure S1) was the predominant signature found
257 here, supporting previous results which have characterized this signature in
258 gastric cancer samples with a very high prevalence of small indels and base
259 substitutions due to failure of DNA double-strand break repair by homo-
260 logous recombination [34]. This finding suggests that 7-12% of gastric cancers
261 may benefit from either platinum therapy or PARP inhibitors. However, no-
262 tably, another group of patients not exposed to signature CS-3 was found to
263 be highly exposed to signatures associated with dMMR (S2, S4 and S5 in
264 Figure S1). Thus, our analysis identified a distinct group of gastric cancer
265 patients harboring features that might have therapeutic relevance and which
266 are further investigated here.

Figure 1: *De novo* mutational signatures in gastric cancer. (a) Mutational signatures in gastric cancer from 787 patients from TCGA, ICGC and cBioPortal cohorts. (b) Heatmap with cosine similarities between *de novo* mutational signatures and COSMIC signatures.



(a) Mutational signatures called by *signer*.



(b) Heatmap with cosine similarity between *de novo* signatures and COSMIC signatures.

267 *3.2. dMMR signatures and prognostic features*

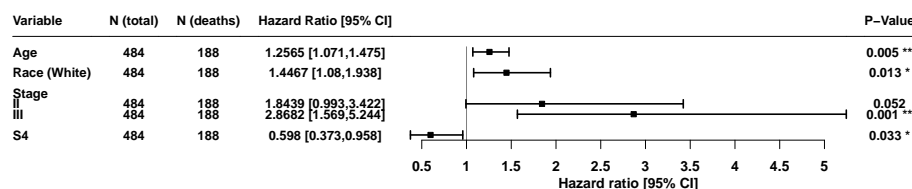
268 We reasoned that dMMR signature exposure could hold prognostic value
269 in gastric cancer. Therefore, we first evaluated the influence of each dMMR
270 signature exposure and the main possible clinical and molecular prognostic
271 features such as age at diagnosis, ethnicity, tumor pathological stage, Lauren
272 classification, anatomic site, TMB and microsatellite instability (MSI) sta-
273 tus on overall survival (OS) fitting simple Cox regression model (Figure S2).
274 Data from 584 gastric cancer patients with available vital status information
275 (Alive/Dead) and without metastasis at diagnosis were included in simple
276 and multiple Cox regression models. The median follow-up time for these
277 patients was 28.9 months (with a 95% confidence interval: 95%CI 25.8-32.1)
278 and the mean follow up time was 36.2 months (95%CI 32.9-39.5).

279
280 We then fitted a multiple Cox regression model to the dataset using prog-
281 nostic features (variables with significant p-value are shown at Figure S2),
282 and observed that S4 exposure burden was associated with improved OS
283 compared with other dMMR signatures (hazard ratio [HR] 0.59 with 95%CI
284 0.37-0.96) (Figure 2A, Figure S3). Thus, we focused on signature S4, which
285 has the potential to offer important clinically actionable information for treat-
286 ment selection.

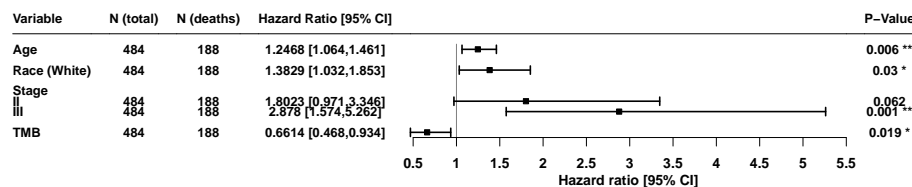
287

Figure 2: Forest plots showing the hazard ratio estimated according to multiple Cox regression models for overall survival.

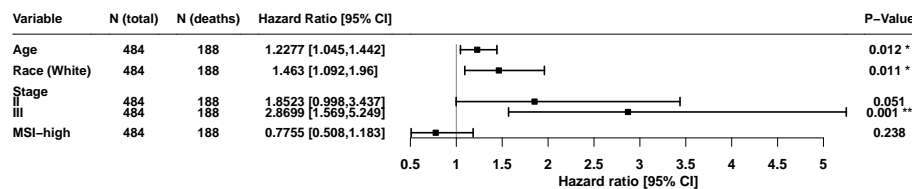
(a) Mutational signature S4



(b) Tumor Mutational Burden (TMB)



(c) MSI status



288 Our analysis also revealed that a higher TMB was also associated with
 289 improved OS (HR=0.66; 95%CI 0.46-0.93) (Figure 2B), consistent with pre-
 290 vious studies [35]. Distinctly from TMB, we found no association between
 291 predicted MSI-H status and improved OS (HR=0.78; 95%CI 0.51-1.18) (Fig-
 292 ure 2C). The calibration curves for these models considering overall survival
 293 in 2 years is given in Figure S4, which indicates that all models are adequate.

294
 295 We used the *maxstat* function (available in R language) to define groups
 296 according to the S4 exposure. The optimal cutpoint was within the range of

297 highest quartile (Q3), and thus patients with S4 exposure \geq Q3 were labeled
298 as S4^{high} otherwise S4^{low} ($<$ Q3). The survival curves from patients in the
299 S4^{high} and S4^{low} groups were statistically different (p-value $<$ 0.03) with a median OS of 72 months (95%CI 48.0- ∞) in the S4^{high} group as compared to
300 37 months (95%CI 28.0-68.0) in the S4^{low} group (Figure S5). Next, we used
301 an independent gastric cancer cohort to validate that S4^{high} has a survival
302 benefit. Kaplan-Meier was performed to analyze patients in the validation
303 cohort grouped based on samples' exposure level of signature S4. By comparing patients in the S4^{high} group whose median OS was not reached at the
304 time of 5 years (95%CI 38.2- ∞) and the S4^{low} group with a median OS of
305 48 months (95%CI 21.3- ∞), the data indicated a trend toward a survival
306 benefit for S4^{high} group, supporting our previous findings (Figure S6).
307
308

309 3.3. dMMR signatures associated with MLH1 hypermethylation

310 Although the genes associated with dMMR are known, the underlying
311 gene modifications that lead to each of the dMMR signatures still remain
312 poorly characterized.
313

314 To improve our understanding of the determinant changes that influence the different dMMR signatures detected in this study, we first looked
315 for somatic and germline SNVs and indels in MMR genes (*LIG1*, *POLE*,
316 *EXO1*, *MLH1*, *MLH3*, *MSH2*, *MSH3*, *MSH5*, *MSH6*, *PCNA*, *PMS1*, *PMS2*,
317 *PMS2L3*, *PMS2L4*, *POLD1*, *POLD2*, *POLD3*, *POLD4* and *SSBP1*). We
318 observed that only 6% of patients (12/197) harbored somatic variations in
319 the *MLH1* gene within either S2^{high} or S4^{high} groups and no mutated patients
320 within S^{low} groups. Likewise, our results showed that 9% of the patients in
321 the S2^{high} group harbored somatic variations in the *MLH3* gene. We also
322 found that only 8% of patients in the S5^{high} group (8/100 considering TCGA
323 cohort) harbor germline mutations in the *MSH5* gene. Altogether, we could
324 observe only few cases harboring mutated MMR genes with some association
325 with the S^{high} groups.
326

327
328 We next searched for epigenetic changes in the MMR genes. In line with
329 previous studies [36, 37], we observed downregulation of *MLH1* gene expression driven by hypermethylation of its promoter (Figure S7). To further assess how the mutational exposure is associated with epigenetic changes in the
330 *MLH1* gene, simple and multiple logistic regression models were fitted to the
331
332

333 dataset (Table 1). This analysis revealed that S4 exposure burden was associ-
 334 ated with an increased chance of *MLH1* promoter being methylated (odds ra-
 335 tio [OR]=22.561; 95%CI 7.909-64.353). On the other hand, S5 exposure bur-
 336 den was associated with a decreased chance of *MLH1* promoter methylation
 337 (OR=0.107; 95%CI 0.048-0.238). Finally, no difference was observed in S2
 338 exposure burden (OR=3.682; 95%CI 0.881-15.386). The performance of this
 339 model was adequate (HosmerLemeshow goodness-of-fit test $\chi^2(8)=10.257$;
 340 p-value=0.247) (Figure 3A), with a good performance observed (Brier score
 341 0.0364), and an excellent power of discrimination (AUC=0.982; 95%CI 0.971-
 342 0.994) (Figure 3B). Using the Youden index, the best cutoff value (threshold)
 343 was 0.125, which had a sensitivity of 95.45% and specificity of 95.82% (Fig-
 344 ure 3B).

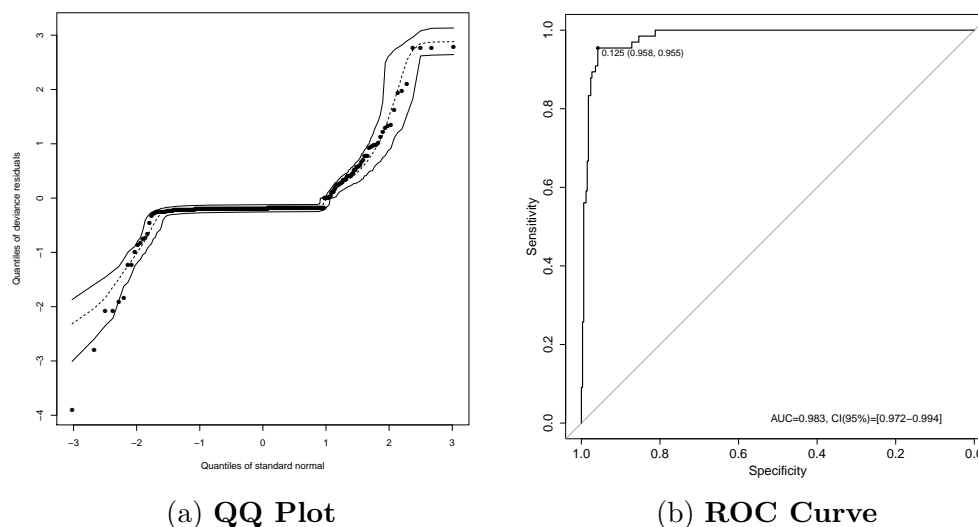
345

Table 1: Simple and multiple logistic regression models for *MLH1* methylation and dMMR mutational signatures

Variable	Simple logistic regression model							
	Coefficient	Standard Error	CI(95%) for coefficient		p-value	OR	CI(95%) for OR	
			Lower	Upper			Lower	Upper
ExpS2	4.772	0.5226	3.748	5.796	< 0.0001	118.155	42.424	329.078
ExpS4	3.2424	0.3469	2.562	3.922	< 0.0001	25.595	12.968	50.518
ExpS5	0.2725	0.1674	-0.056	0.601	0.104	1.313	0.946	1.823

Variable	Multiple logistic regression model							
	Coefficient	Standard Error	CI(95%) for coefficient		p-value	OR	CI(95%) for OR	
			Lower	Upper			Lower	Upper
Intercept	-2.640	0.300	-3.227	-2.052	< 0.0001			
ExpS2	1.304	0.730	-0.126	2.733	0.074	3.682	0.881	15.386
ExpS4	3.1162	0.5348	2.068	4.164	< 0.0001	22.561	7.909	64.353
ExpS5	-2.231	0.4067	-3.028	-1.434	< 0.0001	0.107	0.048	0.238

Figure 3: **Performance and power discrimination of logistic model.** (a) QQ Plot showing that the propensity score model had an adequate level of calibration according to the Hosmer-Lemeshow test. (b) Summary of ROC curve showing the power of discrimination for *MLH1* methylation at signature S4.



346 It should be noted that these observations suggest that neither germline
347 SNVs, somatic SNVs or indels are the major modifications affecting gene
348 expression levels, moreover, in fact, *MLH1* promoter is hypermethylated in
349 almost 60% of individuals in the S4^{high} group (beta-value ≥ 0.3). Taken
350 together, we conclude that the main mechanism of impaired MMR associated
351 with the signature S4 (CS-20) in gastric cancer samples is driven by *MLH1*
352 promoter hypermethylation. Moreover, by using genomic sequencing data
353 from three HAP1 cells samples (2 *MLH1*^{KO} and 1 *MLH1*^{WT} cell), we have
354 observed that *MLH1*^{KO} cell lines have higher exposures of the S4 signature,
355 while the parental cell line has higher exposure of the S5 signature (Figure
356 S8), which identifies the absence of MLH1 as the cause of the S4 mutational
357 signature.

358 3.4. Clinical and molecular features

359 Given our observation that the dMMR signature S4 was associated with
360 better prognosis, and possibly related to an epigenetic causative mechanism,
361 we next tried to further characterize the clinical and molecular features within

362 $S4^{high}$ and $S4^{low}$ groups.

363

364 Previously defined clinical features that were associated with improved
365 prognosis in gastric cancer were also enriched in the $S4^{high}$ group (Table 2),
366 such as distal anatomic site and intestinal histology [38]. On the other hand,
367 known clinical variables associated with a worse prognosis in gastric cancer,
368 such as cardia/proximal anatomic site, diffuse histology, positive lymph node
369 metastasis (stage N+) and advanced pathological stages (stage III and IV)
370 [38] were significantly higher in the $S4^{low}$ group (Table 2). In addition, the
371 predicted MSI-H status, MSI and POLE molecular subtypes were also en-
372 riched in the $S4^{high}$ group, while genomically stable (GS) and chromosomal
373 instability (CIN) molecular subtypes were enriched in the $S4^{low}$ group (Table
374 2). Our data also reveal that most cases of MSI-H (n=119/160, 74%, Table 2)
375 were grouped within the $S4^{high}$ group, however, it has not escaped from our
376 attention that a smaller fraction of MSI-H cases were unexpectedly grouped
377 in the $S4^{low}$ group. Similarly, we have also found Non-MSI-H patients in
378 the $S4^{high}$ group. Comparing the survival curves of these groups, we found
379 MSI-H within $S4^{low}$ group trends to have a worse prognosis with 9.07 months
380 as median OS (95%CI 9.0-∞) than Non-MSI-H within $S4^{high}$ group with 53
381 months as median OS (95%CI 20.0-∞). Similarly, diffuse histologic subtype
382 grouped into $S4^{high}$ (median OS not reached, 95%CI 24.0-∞) trends to have
383 a better prognosis than intestinal histologic subtype grouped into $S4^{low}$ (me-
384 dian OS of 43.1 months, 95%CI 28.0-∞) (Figure S9). Thus, we conclude that
385 the mutational signature classification was able to improve the stratification
386 of patients within the prognostic groups, independent of their previous clin-
387 ical or molecular classification.

388

389 To further understand tumor heterogeneity in the $S4^{high}$ and $S4^{low}$ pa-
390 tients, we examined the spread of allele frequencies according to the quan-
391 titative measure of the degree of heterogeneity [39]. We then performed a
392 correlation analysis based on this score, S4 exposure and TMB (Figure 4).
393 We noted that the correlation of the tumor heterogeneity score (MATH) with
394 either TMB or S4 exposure are opposite in the $S4^{high}$ and $S4^{low}$ groups. In
395 the $S4^{high}$ group, there was a negative correlation of S4 exposure or TMB
396 with MATH and, in the $S4^{low}$ group there was a positive correlation. We also
397 observed that the MATH score is higher in the $S4^{low}$ than the $S4^{high}$ group
398 (p-value=3.711x10⁻¹²). Lastly, the TMB and neoantigen load (by TCIA)
399 have shown a positive correlation with signature S4 exposure in both groups

Table 2: The clinical-pathological features of gastric cancer according to S4 dMMR mutational signature groups

	<i>All</i> (<i>n</i> = 787)		<i>S4^{low}</i> (<i>n</i> = 590)		<i>S4^{high}</i> (<i>n</i> = 197)		<i>P</i> value
	<i>N</i>	%	<i>N</i>	%	<i>N</i>	%	
Age (mean ± SD)	64.17 ± 11.73		63.44 ± 12		66.25 ± 10.68		0.0024
Gender	767	97	571	91	196	99	
<i>Female</i>	274	36	192	34	82	42	0.0473
<i>Male</i>	493	64	379	66	114	58	
Race	726	92	546	93	180	91	
<i>White</i>	275	28	203	37	72	40	0.7427
<i>Black</i>	13	2	11	2	2	1	
<i>Asian</i>	437	60	331	61	106	59	
<i>Other</i>	1	0	1	0	0	0	
Anatomic Site	627	80	495	84	132	67	
<i>Cardia/Proximal</i>	168	27	148	30	20	15	0.0015
<i>Fundus/Body</i>	212	34	166	34	46	35	
<i>Antrum/Distal</i>	242	39	178	36	64	48	
<i>Other</i>	5	1	3	1	2	2	
Histology Lauren	467	59	383	65	84	43	
<i>Diffuse</i>	150	32	132	34	18	21	0.0041
<i>Intestinal</i>	301	64	235	61	66	79	
<i>Mixed</i>	16	3	116	4	0	0	
Stage T	712	90	526	89	186	94	
<i>T1 – T2</i>	181	25	132	25	49	26	0.8116
<i>T3 – T4</i>	531	75	394	75	137	74	
Stage N	712	90	526	89	186	94	
<i>N0</i>	173	24	115	22	58	31	0.0144
<i>N+</i>	539	76	411	78	128	69	
Stage M	707	90	524	89	183	93	
<i>M0</i>	623	88	461	88	162	89	0.9422
<i>M1</i>	62	1	47	9	15	8	
<i>MX</i>	22	3	16	3	6	3	
Pathological Stage	715	91	546	93	169	86	
<i>I</i>	85	12	58	11	27	16	0.0386
<i>II</i>	220	31	160	29	60	36	
<i>III</i>	289	40	228	42	61	36	
<i>IV</i>	121	15	100	18	21	12	
Molecular Subtype	403	51	289	49	114	58	
<i>CIN</i>	223	55	206	71	17	15	<0.0001
<i>GS</i>	50	12	47	16	3	3	
<i>EBV</i>	38	9	33	11	5	4	
<i>MSI</i>	85	21	0	0	85	75	
<i>POLE</i>	7	2	3	1	4	4	
MSIseq Status	787	100	590	100	197	100	
<i>MSI – H</i>	160	20	41	7	119	60	<0.0001
<i>Non – MSI – H</i>	627	80	549	93	78	40	
Immune Subtype	388	49	285	48	103	52	
<i>C1</i>	128	33	107	35	27	26	<0.0001
<i>C2</i>	209	54	135	47	74	72	
<i>C3</i>	35	9	34	12	1	1	
<i>C4</i>	9	2	8	3	1	1	
<i>C6</i>	7	2	7	2	0	0	

400 (Figure 4).

401

402 These findings suggest that tumors highly exposed to signature S4 are
403 more homogeneous in the $S4^{high}$ group and, together with high TMB and
404 high neoantigen load, a reduced tumor heterogeneity appears to be determi-
405 nant of a good prognosis. In this sense, we speculate that the methylation of
406 *MLH1* promoter associated with the $S4^{high}$ signature may be an early event
407 in tumorigenesis.

408

409 In order to check if signature S4 is represented equally across the three
410 cohorts, we compared their samples exposures. To avoid performing statisti-
411 cal tests with different numbers of samples, a subsampling procedure was
412 applied, randomly selecting 24 samples from each cohort and then perform-
413 ing the KruskalWallis test. This was repeated 1000 times, always generating
414 p-values above 0.05, leading us to conclude that the S4 exposure is similar
415 for all cohorts.

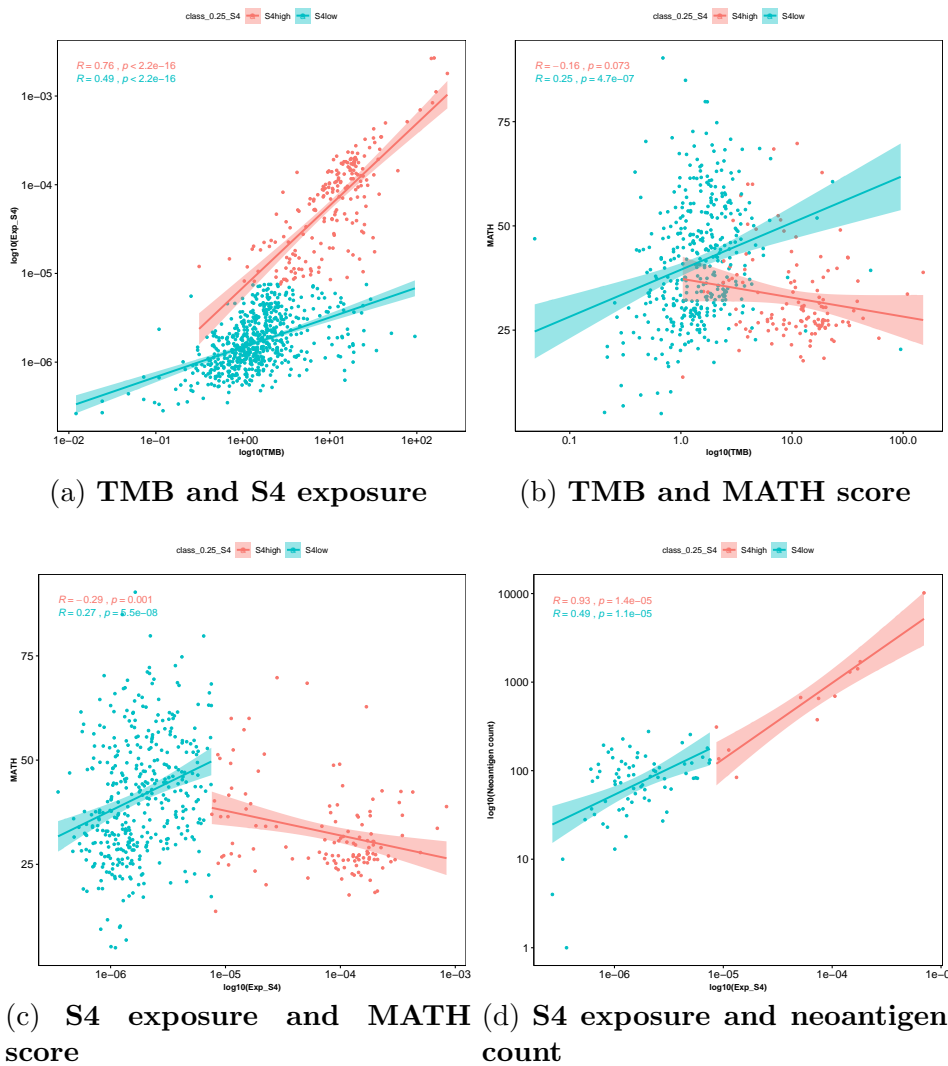
416 3.5. Significantly mutated genes and related pathways in $S4$ dMMR groups

417 MMR-deficiency leads to an elevated frequency of mutations in the genome
418 [14] and the consequences of MMR-deficiency may be derived from func-
419 tional alterations in many distinct genes. In order to verify the existence
420 of a common set of genes commonly mutated and their main related path-
421 ways between $S4^{high}$ and $S4^{low}$ groups (Supplementary Material Table S4)
422 we investigated the presence of consistent SNVs differentiating these groups.
423 At least one somatic mutation, including SNVs and indels, was detected for
424 83.25% of $S4^{high}$ patients, and 78.64% within the $S4^{low}$ group. We observed
425 an increased number of deletions in the $S4^{high}$ group, while within the $S4^{low}$
426 group the mutations basically consisted of SNVs. These results are expected
427 when considering that MSI/dMMR would lead to a higher number of dele-
428 tions [14].

429

430 The gene set found as significantly mutated in the $S4^{high}$ group is com-
431 posed of 102 genes. The most commonly mutated genes in this group are
432 *ARID1A* (42%), *KMT2D* (35%) and *TP53* (31%). In addition, there are an-
433 other 56 genes presenting mutations in at least 10% of patients (Table S4A).
434 The enrichment analysis of these mutated genes identified pathways related
435 to immune cell differentiation, protein and RNA metabolism, gene expres-
436 sion regulation, cell differentiation and embryogenesis (Table S4B). It was

Figure 4: **Scatter plots showing the Spearman correlation between molecular features.** Blue dots represent patients in $S4^{low}$ group and red those in $S4^{high}$ group.



437 previously suggested that somatic mutations in chromatinregulating genes
438 such as *KMT2D* (also known as *MLL2*) and *ARID1A* are associated with
439 improved survival [37].

440

441 In the gene set found as significantly mutated in the $S4^{low}$ group, 12 out
442 of 24 genes are known oncogenes, associated with tumor progression, or tu-
443 mor suppressor genes. These 12 genes are *PIK3CA*, *KRAS*, *RHOA*, *CDH1*,
444 *CTNNB1*, *ITGAV*, *SMAD4*, *TP53*, *CDKN2A*, *APC*, *PTEN* and *PIK3R1*
445 (details in the Table S4C). The most frequently mutated genes in the $S4^{low}$
446 group are *TP53* (47%), *ARID1A* (13%) and *CDH1* (9%) (Table S4C). The
447 other 21 significantly mutated genes for this group were mutated in up to
448 8% of patients. In addition to the common pathways related to cancer,
449 we also found pathways associated with regulation of cell death, phospho-
450 rus metabolism, regulation of transferase activity, morphogenesis pathways
451 (gland development and anatomical structure of a tube) and VEGF and neu-
452 rotrophin signaling pathways (Table S4D).

453

454 These findings were in accordance with some genes found previously in
455 215 non-hypermuted tumors from the TCGA cohort as *APC*, *CTNNB1*,
456 *SMAD4* and *SMAD2*, with somatic mutations in *CDH1* and *RHOA* enriched
457 in the genomically stable and/or diffuse histology [40], subtypes enriched in
458 $S4^{low}$ group.

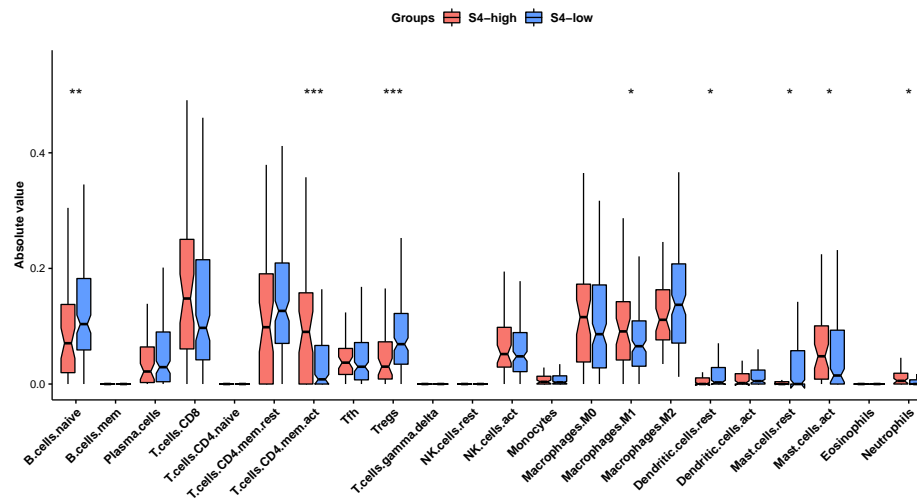
459 3.6. Immune diversity in $S4$ dMMR groups

460 To investigate a possible role of the immune system being associated with
461 the improved clinical outcomes seen in the $S4^{high}$ as compared to $S4^{low}$, we
462 performed a series of analysis to determine the immune cell infiltrate com-
463 position in each group. These analyses used two different analytical method-
464 ologies (see Materials and Methods), and demonstrated a significantly higher
465 proportion of infiltrating cytotoxic and pro-inflammatory immune cells in the
466 group $S4^{high}$, as exemplified by increased CD8+ central and effector memory
467 T cells, CD4+ memory T cells, Th1 cells, gamma/delta T cells, NK cells,
468 M1 macrophages and plasmacytoid dendritic cells (pDC), as compared to
469 the $S4^{low}$ group (Figure 5A and S10). In contrast, immature and immune
470 regulatory dendritic cells were higher in the $S4^{low}$ group (Figure 5A and S10).

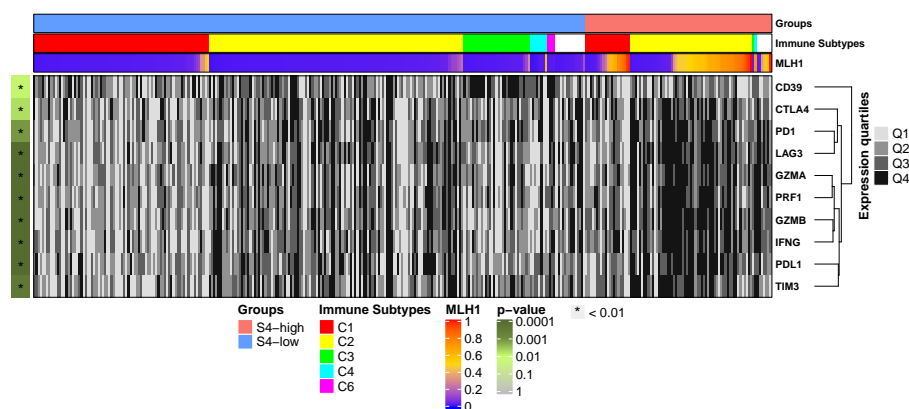
471

472 To further characterize the immune regulatory environment in the $S4^{high}$
473 and $S4^{low}$ groups respectively associated with good or poor clinical outcomes,

Figure 5: The main immunological features associated with $S4^{high}$ and $S4^{low}$ groups.



(a) Boxplots showing the absolute quantification of immune infiltrate cells estimated by CIBERSORT (Newman et al., 2015). Blue boxes represent patients in $S4^{low}$ group and red those within $S4^{high}$ group. Resulting p-values: * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$ by *Mann-Whitney U test*.



(b) Heatmap showing the immune effector genes (cytotoxic and immune checkpoints) expression. Normalized gene expression level (FPKM) for each marker gene was classified in quartiles. Q1 means the range up to first quartile (25% lowest values); Q2: of 25% to 50% (median) expression value; Q3: of 50% to 75% expression values and Q4 means the range of 75% (third quartile) to highest expression values. Left to heatmap, the p value labels consider the *Mann-Whitney U test* by comparing gene expression between $S4^{high}$ (red samples) versus $S4^{low}$ (blue samples) groups. These comparisons are also shown in Figure S9 as boxplots.

474 we performed comparisons between gene expression levels of key genes cod-
475 ing for immunoregulatory and effector molecules proven to be important for
476 tumor control in many cancers [41, 42]. First, the genes for CD8+ T cell
477 related cytolytic molecules Granzyme A/B and Perforin-1 (*GZMA*, *GZMB*
478 and *PRF1* genes, respectively) all displayed higher expression in the $S4^{high}$
479 group (Figures 5B and S11). Moreover, the inflammatory T cell response
480 related cytokine, IFN-gamma (*IFNG* gene) and other proinflammatory cy-
481 tokines (*IL1B*, *IL6* and *IL8*), T cell activation marker genes (*IL2RA* and
482 *ICOS*) and NK cell KIR family receptors were also higher in the $S4^{high}$ than
483 the $S4^{low}$ group (Figures 5B and S11). Second, a series immunosuppression
484 related genes (*TGFB1*, *IL10*, and *FOXP3*) showed no differences between
485 the groups, while *ENTPD1* (CD39 gene), a protein associated with Treg
486 immunosuppression activity [43], was higher in $S4^{low}$ vs. $S4^{high}$ (Figures 5B
487 and S11).

488

489 Importantly, the expression of the immune checkpoint inhibitor genes
490 (*PDCD1* for PD1 receptor, *CD274* for PD-L1 ligand, *PDCD1LG2* for PD-
491 L2, *HAVCR2* for TIM3, *LAG3* and *CTLA4*) were also higher in the $S4^{high}$
492 group (Figures 5B and S11) possibly indicating a relationship with a more
493 immunologically activated tumor microenvironment [44]. Expression of *HLA*,
494 antigen processing and presentation-related genes (such as *CD86*, *B2M*, vari-
495 ous *HLA* class II genes, *HLA-E*, *HLA-C*, *TAP1* and *TAP2*) were also higher
496 in the $S4^{high}$ group (Figure S11). Together these findings indicate a highly
497 activated immune microenvironment in the $S4^{high}$ group as compared to the
498 $S4^{low}$ group.

499

500 The immune subtypes previously characterized by Thorsson et al. [33]
501 reinforce the finding that the $S4^{high}$ group primarily presents a more immuno-
502 logically active tumor microenvironment which is composed predominantly
503 of C2 (interferon-gamma) immune subtype, with a significantly higher pro-
504 portion of this subtype as compared to $S4^{low}$ (Table 2 and Figure 5B). Im-
505 portantly, the C2 immune subtype has been associated with highly mutated
506 tumors [33]. On the other hand, the $S4^{low}$ group displayed a higher propor-
507 tion of C3 (inflammatory) and C1 (wound healing) immune subtypes [33]
508 (Table 2 and Figure 5B). Lastly, the C2 immune subtype in the $S4^{low}$ group
509 seems to be less activated than in the $S4^{high}$ group with reduced relative gene
510 expression of immune effector molecules (Figure 5B).

511 *3.7. Discussion*

512 This study provides a comprehensive and integrated analysis of the im-
513 pact of MMR-related gene alterations in shaping specific mutational signa-
514 tures associated with gastric cancer. We present evidence that the determi-
515 nation of MMR-deficiency can be used not only for MSI-phenotype classifi-
516 cation, but also as a potential indicator of prognosis and to select potential
517 candidates for treatment with checkpoint inhibitors.

518
519 We performed a *de novo* extraction of mutational signatures based on
520 somatic SNVs across four WES cohorts, spanning 787 gastric cancer samples
521 derived mainly from populations with European and Asian descentance. We
522 found 7 different mutational signatures, with three related to MMR-deficiency.

523
524 Next, we examined the prognostic value of these dMMR signatures in mul-
525 tivariate survival analysis by employing a Cox proportional hazards model.
526 This analysis revealed signature S4, related to the previously described CS-
527 20, as the only dMMR signature with significant prognostic value. This prog-
528 nostic value was validated using our local cohort of gastric cancer patients,
529 distinct in terms of molecular ancestry as well as some clinical and molecular
530 features such as Lauren's histology and tumor heterogeneity. This cohort
531 was predominantly composed of diffuse/mixed histology samples, while pub-
532 lic cohorts were enriched for the intestinal subtype. Furthermore, the S4^{low}
533 group in this independent cohort was less heterogeneous than the S4^{low} group
534 from the public cohorts and even than the S4^{high} groups from both cohorts.
535 Nevertheless, we observed a better prognosis for the patients of the S4^{high}
536 group, also for this cohort.

537
538 Interestingly, after performing a comprehensive analysis of patients ex-
539 posed to signature S4 by an in depth evaluation of molecular and immune fea-
540 tures, we observed that the main mechanism associated with impaired MMR
541 seems to be the hypermethylation of the *MLH1* gene promoter (*hMLH1*).
542 Moreover, we show that disruption of *MLH1 in vitro* using CRISPR/Cas9
543 assay reproduces the CS-20 signature [45] that resembles the S4 signature.
544 Here, we have shown an endogenous epigenetic mechanism for this signature
545 in gastric cancer patients. Remarkably, we also reproduce the S4 signature
546 in an isogenic cell model in which the *MLH1*^{KO} cells had a high exposure of
547 the S4 signature. Thus importantly, we conclude that independently of the
548 primary mechanism that leads to the loss of *MLH1* gene expression - due to

549 promoter hypermethylation or loss of function mutagenesis - it results in the
550 same mutational signature.

551

552 It has been well documented that CpG island methylator (CIMP) phe-
553 notype is an early event in tumorigenesis, preceding the *hMLH1* in solid
554 tumors, which in turn drives the microsatellite instability high (MSI-H) phe-
555 notype [36, 40, 37]. In contrast, MSI-low (MSI-L) and microsatellite stable
556 (MSS) gastric carcinoma subtypes have unmethylated *MLH1* promoters and
557 regular *MLH1* activity [36]. Here, we classified samples as MSI-H and non-
558 MSI-H (MSI-L and MSS) and observed that most cases of MSI-H fall within
559 the $S4^{high}$ group, however, a smaller fraction of MSI-H cases did not show
560 high exposure of this mutational signature and were grouped in the $S4^{low}$
561 group. This is in line with previous studies showing about one-quarter of
562 the MSI-H cases, despite being MSI-H, present distinct molecular features
563 and poor prognosis [46]. Similarly, we have also found Non-MSI-H patients
564 in the $S4^{high}$ group, showing that mutational signature exposure is capable
565 of clustering samples independently of their MSI-status. Furthermore, we
566 also identified a few cases (4%) in the $S4^{high}$ group that instead of presenting
567 *hMLH1*, carried somatic mutations in the *MLH1* gene that apparently lead
568 to loss-of-function of the encoded protein. Thus, for about 70% of $S4^{high}$
569 cases we found a clear genetic or epigenetic cause.

570

571 We also demonstrated a strong correlation between S4 exposure and
572 TMB, which showed significant prognostic value upon multivariate survival
573 analysis. Hypermutated tumors have been associated with better prognos-
574 is and a good response to immunotherapy apparently due to neoantigen
575 enrichment and intrinsic antitumor immune responses [47, 48]. However, a
576 threshold for classifying TMB-high samples usually varies with tumor type
577 [49] and in some cases may not predict a better response [48] due to intratu-
578 moral heterogeneity [50]. In this sense, it is important to highlight that most
579 mutations in the $S4^{high}$ signature are clonal, which is an important feature
580 to predict response to immune checkpoint inhibitors therapy.

581

582 High intratumoral heterogeneity has been associated with an incomplete
583 response to therapy, higher relapse rates, and poor clinical outcomes [51, 52].
584 The increased genomic instability observed in MSI/dMMR and CIN (chromo-
585 somal instability) tumors is the major driver of high intratumoral hetero-
586 geneity [53, 51]. However, the most unstable tumors (with the highest burden

587 of somatic SNVs or copy number alterations) are not the most intrinsically
588 heterogeneous[53]. Furthermore, the greatest intratumor heterogeneity was
589 found in tumors exhibiting relatively high numbers of both somatic muta-
590 tions and copy number alterations, which can be associated with exogenous
591 mutagens, including viral infection and tobacco smoking. These tumors have
592 high number of sub-clonal mutations related to late events and exhibit fre-
593 quent chromosomal instability associated with CIN subtype, *TP53* muta-
594 tions, and APOBEC-related mutational signatures (previous related to EBV
595 gastric cancer subtype [6, 53]. Likewise, here we noted that patients with
596 higher S4 exposure harbor more homogeneous tumors as compared to S4^{low}
597 group. Taken together, the S4^{high} group encompassed patients with interme-
598 diate to high TMB, in addition to the MSI and CIN molecular phenotypes
599 associated with lower tumor heterogeneity, which might allow for a more ef-
600 fective antitumor immune response in this subset of gastric cancer patients.
601 In this sense, a recent meta-analysis discussed the importance of MSI-status
602 for the treatment response in gastric cancer patients, suggesting that MSI-H
603 patients may not benefit from perioperative or adjuvant therapy and could
604 go straight to surgery [54].

605
606 Finally, several studies have shown that the tumor microenvironment
607 context, at diagnosis, is capable of predicting treatment response and clini-
608 cal outcome [55, 56]. The balance of inflammatory/cytotoxic immune cells,
609 with elements of an effective antitumor response, including regulatory cells
610 and suppressor signals, may indicate which patients have an intrinsically
611 effective antitumor response, and thus, a better prognosis. EBV and MSI
612 subtypes in gastric cancer have already been associated with higher immune
613 infiltrate and responsiveness to immunotherapy, as well as better prognosis
614 [55]. Here we found many elements indicating that the tumor microenviron-
615 ment in the S4^{high} group is more active as compared to S4^{low} patients. In
616 general, the absolute quantification by CIBERSORT [31] or GSEA scores
617 by xCell [32] of immune cell subtypes and the differential gene expression
618 pointed to higher activity of proinflammatory and cytotoxic cells, as well as
619 antigen processing and presentation in S4^{high}. In contrast, although there
620 were some immunogenic tumors in S4^{low} group, the predominant environ-
621 ment was enriched in Treg lymphocytes and M2 macrophages, both related
622 to worse prognosis [56, 55].

623
624 In conclusion, while past studies have aimed to identify patients using

625 molecular and clinical features such as MSI status, TMB load and *MLH1*
626 gene expression levels, our study provides evidence that classification based
627 on mutational signature exposure may identify groups of patients with com-
628 mon clinical, immunological and mutational features that are directly related
629 to a better prognosis, and who might benefit from immunotherapy-based
630 treatments.

631 3.8. Acknowledgements

632 This project received financial support from FAPESP (14-26897-0 and
633 16/11791-7); ED-N and KJG are research fellows from Conselho Nacional
634 de Desenvolvimento Científico e Tecnológico (CNPq, Brazil). ED-N acknowl-
635 edges the support given by Associação Beneficente Alzira Denise Hertzog
636 Silva (ABADHS).

637 References

- 638 [1] M. R. Stratton, P. J. Campbell, P. A. Futreal, The cancer genome,
639 Nature 458 (2009) 719–724.
- 640 [2] R. Hakem, DNA-damage repair; the good, the bad, and the ugly, EMBO
641 J. 27 (2008) 589–605.
- 642 [3] T. Helleday, S. Eshtad, S. Nik-Zainal, Mechanisms underlying muta-
643 tional signatures in human cancers, Nat. Rev. Genet. 15 (2014) 585–598.
- 644 [4] L. B. Alexandrov, S. Nik-Zainal, D. C. Wedge, S. A. Aparicio, S. Behjati,
645 A. V. Biankin, et al., Signatures of mutational processes in human
646 cancer, Nature 500 (2013) 415–421.
- 647 [5] M. Hollstein, L. B. Alexandrov, C. P. Wild, M. Ardin, J. Zavadil, Base
648 changes in tumour DNA have the power to reveal the causes and evolu-
649 tion of cancer, Oncogene 36 (2017) 158–167.
- 650 [6] I. Bobrovnitchaia, R. Valieris, R. D. Drummond, J. P. Lima, H. C.
651 Freitas, T. F. Bartelli, et al., APOBEC-mediated DNA alterations:
652 A possible new mechanism of carcinogenesis in EBV-positive gastric
653 cancer, Int. J. Cancer 146 (2020) 181–191.
- 654 [7] A. Van Hoeck, N. H. Tjoonk, R. van Boxtel, E. Cuppen, Portrait of
655 a cancer: mutational signature analyses for cancer diagnostics, BMC
656 Cancer 19 (2019) 457.

- 657 [8] H. Davies, D. Glodzik, S. Morganella, L. R. Yates, J. Staaf, X. Zou,
658 et al., HRDetect is a predictor of BRCA1 and BRCA2 deficiency based
659 on mutational signatures, *Nat. Med.* 23 (2017) 517–525.
- 660 [9] J. Kim, K. W. Mouw, P. Polak, L. Z. Braunstein, A. Kamburov, D. J.
661 Kwiatkowski, et al., Somatic ERCC2 mutations are associated with a
662 distinct genomic signature in urothelial tumors, *Nat. Genet.* 48 (2016)
663 600–606.
- 664 [10] M. Jager, F. Blokzijl, E. Kuijk, J. Bertl, M. Vougioukalaki, R. Janssen,
665 et al., Deficiency of nucleotide excision repair is associated with muta-
666 tional signature observed in cancer, *Genome Res.* 29 (2019) 1067–1077.
- 667 [11] M. L. Thibodeau, E. Y. Zhao, C. Reisle, C. Ch’ng, H. L. Wong, Y. Shen,
668 et al., Base excision repair deficiency signatures implicate germline and
669 somatic MUTYH aberrations in pancreatic ductal adenocarcinoma and
670 breast cancer oncogenesis, *Cold Spring Harb Mol Case Stud* 5 (2019).
- 671 [12] M. D. Giraldez, F. Balaguer, L. Bujanda, M. Cuatrecasas, J. Mu?oz,
672 V. Alonso-Espinaco, et al., MSH6 and MUTYH deficiency is a frequent
673 event in early-onset colorectal cancer, *Clin. Cancer Res.* 16 (2010) 5402–
674 5413.
- 675 [13] N. M. Volkov, G. A. Yanus, A. O. Ivantsov, F. V. Moiseenko, O. G.
676 Matorina, I. V. Bizin, et al., Efficacy of immune checkpoint blockade
677 in MUTYH-associated hereditary colorectal cancer, *Invest New Drugs*
678 (2019).
- 679 [14] R. Mandal, R. M. Samstein, K. W. Lee, J. J. Havel, H. Wang, C. Kr-
680 ishna, et al., Genetic diversity of tumors with mismatch repair deficiency
681 influences anti-PD-1 immunotherapy response, *Science* 364 (2019) 485–
682 491.
- 683 [15] D. T. Le, J. N. Durham, K. N. Smith, H. Wang, B. R. Bartlett, L. K.
684 Aulakh, et al., Mismatch repair deficiency predicts response of solid
685 tumors to PD-1 blockade, *Science* 357 (2017) 409–413.
- 686 [16] W. Abida, M. L. Cheng, J. Armenia, S. Middha, K. A. Autio, H. A.
687 Vargas, et al., Analysis of the Prevalence of Microsatellite Instability in
688 Prostate Cancer and Response to Immune Checkpoint Blockade, *JAMA*
689 *Oncol* 5 (2019) 471–478.

- 690 [17] Z. R. Reichert, J. Urrutia, J. J. Alumkal, Microsatellite Instability as
691 an Emerging Biomarker for Checkpoint Inhibitor Response in Advanced
692 Prostate Cancer, *JAMA Oncol* 5 (2019) 478–479.
- 693 [18] GE4GACgroup, T. Bartelli, et al., Genomics and epidemiology for gas-
694 tric adenocarcinomas (GE4GAC): a Brazilian initiative to study gastric
695 cancer, *Appl. Cancer Res.* 39 (2019).
- 696 [19] K. Ellrott, M. Bailey, G. Saksena, K. Covington, C. Kandoth, C. Stew-
697 art, et al., Scalable open science approach for mutation calling of tumor
698 exomes using multiple genomic pipelines, *Cell Syst.* 28 (2018) 271–281.
- 699 [20] G. Ciriello, M. L. Miller, B. A. Aksoy, Y. Senbabaoglu, N. Schultz,
700 C. Sander, Emerging landscape of oncogenic signatures across human
701 cancers, *Nat. Genet.* 45 (2013) 1127–1133.
- 702 [21] M. A. DePristo, E. Banks, R. Poplin, K. V. Garimella, J. R. Maguire,
703 C. Hartl, et al., A framework for variation discovery and genotyping
704 using next-generation DNA sequencing data, *Nat. Genet.* 43 (2011)
705 491–498.
- 706 [22] R. A. Rosales, R. D. Drummond, R. Valieris, E. Dias-Neto, I. T. da Silva,
707 signeR: an empirical Bayesian approach to mutational signature discov-
708 ery, *Bioinformatics* 33 (2017) 8–16.
- 709 [23] M. N. Huang, J. R. McPherson, I. Cutcutache, B. T. Teh, P. Tan, S. G.
710 Rozen, MSIsq: Software for Assessing Microsatellite Instability from
711 Catalogs of Somatic Mutations, *Sci Rep* 5 (2015) 13321.
- 712 [24] Z. R. Chalmers, C. F. Connelly, D. Fabrizio, L. Gay, S. M. Ali, R. Ennis,
713 et al., Analysis of 100,000 human cancer genomes reveals the landscape
714 of tumor mutational burden, *Genome Med* 9 (2017) 34.
- 715 [25] A. Mayakonda, D. C. Lin, Y. Assenov, C. Plass, H. P. Koeffler, Maftools:
716 efficient and comprehensive analysis of somatic variants in cancer,
717 *Genome Res.* 28 (2018) 1747–1756.
- 718 [26] P. Charoentong, F. Finotello, M. Angelova, C. Mayer, M. Efre-
719 mova, D. Rieder, et al., Pan-cancer Immunogenomic Analyses Reveal
720 Genotype-Immunophenotype Relationships and Predictors of Response
721 to Checkpoint Blockade, *Cell Rep* 18 (2017) 248–262.

- 722 [27] D. Schoenfeld, Partial residuals for the proportional hazards regression
723 model., *Biometrika* 69 (1982) 239–241.
- 724 [28] D. Hosmer, S. Lemeshow, S. RX., *Applied Logistic Regression.*, Wiley
725 3 edition (2013).
- 726 [29] X. Zou, M. Owusu, R. Harris, S. P. Jackson, J. I. Loizou, S. Nik-Zainal,
727 Validating the concept of mutational signatures with isogenic cell mod-
728 els, *Nat Commun* 9 (2018) 1744.
- 729 [30] M. S. Lawrence, P. Stojanov, P. Polak, G. V. Kryukov, K. Cibulskis,
730 A. Sivachenko, et al., Mutational heterogeneity in cancer and the search
731 for new cancer-associated genes, *Nature* 499 (2013) 214–218.
- 732 [31] A. M. Newman, C. L. Liu, M. R. Green, A. J. Gentles, W. Feng, Y. Xu,
733 et al., Robust enumeration of cell subsets from tissue expression profiles,
734 *Nat. Methods* 12 (2015) 453–457.
- 735 [32] D. Aran, Z. Hu, A. J. Butte, xCell: digitally portraying the tissue
736 cellular heterogeneity landscape, *Genome Biol.* 18 (2017) 220.
- 737 [33] V. Thorsson, D. Gibbs, S. Brown, D. Wolf, D. Bortone, T. Ou Yang,
738 et al., The Immune Landscape of Cancer, *Immunity* 51 (2019) 411–412.
- 739 [34] L. B. Alexandrov, S. Nik-Zainal, H. C. Siu, S. Y. Leung, M. R. Stratton,
740 A mutational signature in gastric cancer suggests therapeutic strategies,
741 *Nat Commun* 6 (2015) 8683.
- 742 [35] R. Buttner, J. W. Longshore, F. Lopez-Rios, S. Merkelbach-Bruse,
743 N. Normanno, E. e. a. Rouleau, Implementing TMB measurement in
744 clinical practice: considerations on assay requirements, *ESMO Open* 4
745 (2019) e000442.
- 746 [36] S. Y. Leung, S. T. Yuen, L. P. Chung, K. M. Chu, A. S. Chan, J. C. Ho,
747 hMLH1 promoter methylation and lack of hMLH1 expression in sporadic
748 gastric carcinomas with high-frequency microsatellite instability, *Cancer*
749 *Res.* 59 (1999) 159–164.
- 750 [37] W. Hu, Y. Yang, L. Qi, J. Chen, W. Ge, S. Zheng, Subtyping of mi-
751 crosatellite instability-high colorectal cancer, *Cell Commun. Signal* 17
752 (2019) 79.

- 753 [38] G. Cammerer, A. Formentini, M. Karletshofer, D. Henne-Bruns, M. Ko-
754 rnmann, Evaluation of important prognostic clinical and pathological
755 factors in gastric cancer, *Anticancer Res.* 32 (2012) 1839–1842.
- 756 [39] A. Rajput, T. Bocklage, A. Greenbaum, J. H. Lee, S. A. Ness, Mutant-
757 Allele Tumor Heterogeneity Scores Correlate With Risk of Metastases
758 in Colon Cancer, *Clin Colorectal Cancer* 16 (2017) e165–e170.
- 759 [40] A. J. Bass, V. Thorsson, I. Shmulevich, S. M. Reynolds, M. Miller,
760 B. Bernard, T. Hinoue, et al., Comprehensive molecular characterization
761 of gastric adenocarcinoma, *Nature* 513 (2014) 202–209.
- 762 [41] G. Landskron, M. De la Fuente, P. Thuwajit, C. Thuwajit, M. A. Her-
763 moso, Chronic inflammation and cytokines in the tumor microenviron-
764 ment, *J Immunol Res* 2014 (2014) 149185.
- 765 [42] S. A. Fuertes Marraco, N. J. Neubert, G. Verdeil, D. E. Speiser, In-
766 hibitory Receptors Beyond T Cell Exhaustion, *Front Immunol* 6 (2015)
767 310.
- 768 [43] J. Bastid, A. Cottalorda-Regairaz, G. Alberici, N. Bonnefoy, J. F.
769 Eliaou, A. Bensussan, ENTPD1/CD39 is a promising therapeutic target
770 in oncology, *Oncogene* 32 (2013) 1743–1751.
- 771 [44] M. Binnewies, E. W. Roberts, K. Kersten, V. Chan, D. F. Fearon,
772 M. Merad, et al., Understanding the tumor immune microenvironment
773 (TIME) for effective therapy, *Nat. Med.* 24 (2018) 541–550.
- 774 [45] J. Drost, R. van Boxtel, F. Blokzijl, T. Mizutani, N. Sasaki, V. Sasselli,
775 et al., Use of CRISPR-modified human stem cell organoids to study the
776 origin of mutational signatures in cancer, *Science* 358 (2017) 234–238.
- 777 [46] W. Hu, Y. Yang, L. Qi, J. Chen, W. Ge, S. Zheng, Subtyping of mi-
778 crosatellite instability-high colorectal cancer, *Cell Commun. Signal* 17
779 (2019) 79.
- 780 [47] A. M. Goodman, S. Kato, L. Bazhenova, S. P. Patel, G. M. Frampton,
781 V. Miller, et al., Tumor Mutational Burden as an Independent Predictor
782 of Response to Immunotherapy in Diverse Cancers, *Mol. Cancer Ther.*
783 16 (2017) 2598–2608.

- 784 [48] S. Mishima, A. Kawazoe, Y. Nakamura, A. Sasaki, D. Kotani,
785 Y. Kuboki, et al., Clinicopathological and molecular features of re-
786 sponders to nivolumab for patients with advanced gastric cancer, *J*
787 *Immunother Cancer* 7 (2019) 24.
- 788 [49] R. M. Samstein, C. H. Lee, A. N. Shoushtari, M. D. Hellmann, R. Shen,
789 Y. Y. Janjigian, et al., Tumor mutational load predicts survival after
790 immunotherapy across multiple cancer types, *Nat. Genet.* 51 (2019)
791 202–206.
- 792 [50] Q. Jia, W. Wu, Y. Wang, P. B. Alexander, C. Sun, Z. Gong, et al.,
793 Local mutational diversity drives intratumoral immune heterogeneity in
794 non-small cell lung cancer, *Nat Commun* 9 (2018) 5361.
- 795 [51] G. Stanta, S. Bonin, Overview on Clinical Relevance of Intra-Tumor
796 Heterogeneity, *Front Med (Lausanne)* 5 (2018) 85.
- 797 [52] A. C. F. Bolhaqueiro, B. Ponsioen, B. Bakker, S. J. Klaasen, E. Ku-
798 cukkose, R. H. van Jaarsveld, et al., Ongoing chromosomal instability
799 and karyotype evolution in human colorectal cancer organoids, *Nat.*
800 *Genet.* 51 (2019) 824–834.
- 801 [53] F. Raynaud, M. Mina, D. Tavernari, G. Ciriello, Pan-cancer inference
802 of intra-tumor heterogeneity reveals associations with different forms of
803 genomic instability, *PLoS Genet.* 14 (2018) e1007669.
- 804 [54] F. Pietrantonio, R. Miceli, A. Raimondi, Y. W. Kim, W. K. Kang, R. E.
805 Langley, et al., Individual Patient Data Meta-Analysis of the Value of
806 Microsatellite Instability As a Biomarker in Gastric Cancer, *J. Clin.*
807 *Oncol.* 37 (2019) 3392–3400.
- 808 [55] M. Wang, R. A. Busuttil, S. Pattison, P. J. Neeson, A. Boussioutas,
809 Immunological battlefield in gastric cancer and role of immunotherapies,
810 *World J. Gastroenterol.* 22 (2016) 6373–6384.
- 811 [56] W. H. Fridman, L. Zitvogel, C. Saut?S-Fridman, G. Kroemer, The im-
812 mune contexture in cancer prognosis and treatment, *Nat Rev Clin Oncol*
813 14 (2017) 717–734.