

Assessing Genotype-Phenotype Correlations with Deep Learning in Colorectal Cancer: A Multi-Centric Study

Marco Gustav (1), Marko van Treeck (1), Nic G. Reitsam (2), Zunamys I. Carrero (1), Chiara M. L. Loeffler (1,3), Asier Rabasco Meneghetti (1), Bruno Märkl (2), Lisa A. Boardman (4), Amy J. French (5), Ellen L. Goode (6), Andrea Gsur (7), Stefanie Brezina (7), Marc J. Gunter (8,9), Neil Murphy (8), Pia Hönscheid (10,11,12), Christian Sperling (10), Sebastian Foersch (13), Robert Steinfeldler (14), Tabitha Harrison (14,15), Ulrike Peters (14,15), Amanda Phipps (14,15), Jakob Nikolas Kather⁺ (1,3,16,17)

1. Else Kroener Fresenius Center for Digital Health, Faculty of Medicine and University Hospital Carl Gustav Carus, TUD Dresden University of Technology, 01307 Dresden, Germany.
2. Pathology, Faculty of Medicine, University of Augsburg, Augsburg, Germany
3. Department of Medicine I, Faculty of Medicine and University Hospital Carl Gustav Carus, TUD Dresden University of Technology, 01307 Dresden, Germany.
4. Division of Gastroenterology and Hepatology, Mayo Clinic, Rochester, Minnesota, USA.
5. Division of Laboratory Genetics, Department of Laboratory Medicine and Pathology, Mayo Clinic, Rochester, Minnesota, USA.
6. Department of Quantitative Health Sciences, Division of Epidemiology, Mayo Clinic, Rochester, Minnesota, USA.
7. Center for Cancer Research, Medical University of Vienna, Vienna, Austria.
8. Nutrition and Metabolism Branch, International Agency for Research on Cancer, World Health Organization, Lyon, France.
9. Cancer Epidemiology and Prevention Research Unit, School of Public Health, Imperial College London, London, United Kingdom.
10. Institute of Pathology, University Hospital Carl Gustav Carus (UKD), Technical University Dresden (TUD), Dresden, Germany.
11. National Center for Tumor Diseases (NCT), Partner Site Dresden, German Cancer Research Center Heidelberg, Dresden, Germany.
12. German Cancer Consortium (DKTK) and German Cancer Research Center (DKFZ), Heidelberg, Germany.
13. Institute of Pathology, University Medical Center Mainz, Mainz, Germany.
14. Division of Public Health Sciences, Fred Hutchinson Cancer Center, Seattle, WA, USA.
15. Department of Epidemiology, University of Washington, Seattle, WA, USA.
16. Medical Oncology, National Center for Tumor Diseases (NCT), University Hospital Heidelberg, Heidelberg, Germany.
17. Pathology & Data Analytics, Leeds Institute of Medical Research at St James's, University of Leeds, Leeds, United Kingdom.

+ Correspondence to:
Jakob Nikolas Kather, MD, MSc
Professor of Clinical Artificial Intelligence
Else Kroener Fresenius Center for Digital Health
Technical University Dresden
Fetscherstrasse 74
01307 Dresden, Germany
jakob-nikolas.kather@alumni.dkfz.de
<http://www.kather.ai>

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

51 Abstract

52 **Background:** Deep Learning (DL) has emerged as a powerful tool to predict genetic biomarkers
53 directly from digitized Hematoxylin and Eosin (H&E) slides in colorectal cancer (CRC). However,
54 few studies have systematically investigated the predictability of biomarkers beyond routinely
55 available alterations such as microsatellite instability (MSI), and *BRAF* and *KRAS* mutations.

56
57 **Methods:** Our primary dataset comprised H&E slides of CRC tumors across five cohorts totaling
58 1,376 patients who underwent comprehensive panel sequencing, with an additional 536 patients
59 from two public datasets for validation. We developed a DL model using a single transformer model
60 to predict multiple genetic alterations directly from the slides. The model's performance was
61 compared against conventional single-target models, and potential confounders were analyzed.

62
63 **Findings:** The multi-target model was able to predict numerous biomarkers from pathology slides,
64 matching and partly exceeding single-target transformers. The Area Under the Receiver Operating
65 Characteristic curve (AUROC, mean \pm std) on the primary external validation cohorts was: *BRAF*
66 (0.78 \pm 0.01), hypermutation (0.88 \pm 0.01), MSI (0.93 \pm 0.01), *RNF43* (0.86 \pm 0.01); this biomarker
67 predictability was mirrored across metrics and co-occurrence analyses. However, biomarkers with
68 high AUROCs largely correlated with MSI, with model predictions depending considerably on MSI-
69 associated morphology upon pathological examination.

70
71 **Interpretation:** Our study demonstrates that multi-target transformers can predict the biomarker
72 status for numerous genetic alterations in CRC directly from H&E slides. However, their pre-
73 dictability is mainly associated with MSI phenotype, despite indications of slight biomarker-inherent
74 contributions to a phenotype. Our findings underscore the need to analyze confounders in AI-based
75 oncology biomarkers. To enable this, we developed a validated model applicable to other cancers
76 and larger, diverse datasets.

77
78 **Funding:** The German Federal Ministry of Health, the Max-Eder-Programme of German Cancer
79 Aid, the German Federal Ministry of Education and Research, the German Academic Exchange
80 Service, and the EU.

81 Introduction

82 Exome sequencing, including targeted panel sequencing, plays a key role in precision oncology of
83 colorectal cancer (CRC)^{1,2} but remains inaccessible to many CRC patients worldwide due to the
84 need for cost-intensive laboratory equipment and complex data analysis.³ In contrast, Hematoxylin
85 and Eosin (H&E) stained histopathology slides are a standard diagnostic tool available for almost
86 every cancer patient globally. The advent of deep learning (DL) has unlocked these slides as a
87 quantifiable data resource. Recent studies have extensively suggested that DL can predict
88 molecular biomarkers directly from digitized H&E slides, including microsatellite instability (MSI)⁴⁻
89 ¹⁰, hypermutation status, and gene mutations such as *TP53*, *BRAF*, and *KRAS*.^{7,8,10-13} When ap-
90 plied as pre-screening tools, such DL systems can streamline the diagnostic workflow^{9,14}, identify-
91 ing those cases that need further testing and ruling out others.¹⁵

92
93 Previous DL studies in CRC have predominantly focused on specific genetic alterations as poten-
94 tial biomarkers, referred to as 'prediction targets' in computational pathology, with only few studies
95 adopting a pan-molecular alteration approach.^{11,16} Such efforts have often been constrained by the
96 challenge of obtaining diverse datasets with extensive sequencing, and even then, traditional
97 methods require training separate models for individual prediction targets^{5-8,12,13}, making them
98 labor- and resource-intensive. Addressing both challenges, this study introduces the first pan-
99 biomarker DL approach for CRC, utilizing a single model to predict multiple molecular targets. The
100 transformer-based model was trained and validated on a unified dataset from the Genetics and
101 Epidemiology of Colorectal Cancer Consortium (GECCO), which consolidates sequencing data
102 from diverse cohorts.¹⁷ Furthermore, its generalizability was evaluated using two extensively
103 studied CRC cohorts.

104 As DL has demonstrated robust capability in linking phenotype to genotype, particularly for pre-
105 dicting MSI⁵⁻⁹, we leverage the strengths of our dataset and model to extend this analysis to mul-
106 tiple prediction targets. Specifically, we analyzed the co-occurrence of genetic alterations with MSI,
107 evaluated the ability of DL to predict these alterations, and examined the corresponding slide
108 morphology, relating it to features known to be relevant for MSI detection. Our study included
109 genetic alterations investigated in prior DL-based studies and clinically relevant mutations, such as
110 *BMPR2*, *RNF43* and *BRAF* as well as MSI and hypermutation status.

111 **Materials and Methods**

112 **Patient Samples**

113 GECCO is an international collaboration utilizing genotype and sequencing data from a growing
114 resource of over 60 studies.¹⁷ Our experiments were conducted on five cohorts within GECCO,
115 comprising N=1,376 patients with complete data and digitized whole-slide images (WSIs): Euro-
116 pean Prospective Investigation into Cancer (EPIC, N=183)¹⁸, Colorectal Cancer Study of Austria
117 (CORSA, N=158)¹⁹, Iowa Women's Health Study (IWHS, N=390)²⁰, Cancer Risk Assessment study
118 (CRA, N=321) and Women's Health Initiative (WHI, N=324)²¹ (Fig. 1A-C). These studies provided
119 digitized WSIs with unified demographic, clinical, and lifestyle data (relevant data for this study in
120 Tab. S1), establishing them as the primary dataset for this study. Centralized targeted tumor
121 sequencing of up to 356 genes was conducted on all samples, with MSI and hypermutation status
122 also assessed. The analysis targeted non-silent mutations and mutational signatures using panel
123 sequencing of 1.8 Megabases (Mb), with tumor and normal sequencing coverage of 975x and
124 273x, respectively. To ensure generalizability of our model and comparability with existing
125 literature, we utilized a secondary external dataset with two well-characterized publicly available
126 CRC cohorts from The Cancer Genome Atlas (TCGA, N=426)²² and the Clinical Proteomic Tumor
127 Analysis Consortium (CPTAC, N=110, fresh-frozen tissue samples)²³ (Fig. 1C). We included
128 patients with WSIs and ground truth data on the primary genetic alterations analyzed in this study.
129 The MSI status for TCGA was defined per Liu et al.²⁴, with high grade MSI categorized as MSI and
130 low grade MSI and microsatellite stable (MSS) grouped as MSS²⁵. Ground truth labels were derived
131 from the clinical data accessible at <https://portal.gdc.cancer.gov/> (accessed on Nov 20, 2024).

132 **Experimental Design**

133 Our study included genetic alterations investigated in prior DL-based studies, clinically relevant
134 mutations, and genetic alterations strongly associated with MSI in the TCGA CRC cohort²⁶, in-
135 cluding *APC*, *BMP2R*, *BRAF*, *KRAS*, *RNF43*, *TP53*, *ZNRF3*.^{7,8,11-13,26,27} Together with MSI status
136 and hypermutation, these exemplify the main prediction targets among the broad target set utilized
137 from the GECCO cohorts (Tab. S2). Continuous prediction targets were discretized into classes
138 using predefined thresholds to achieve balanced case distribution. Only prediction targets with at
139 least 20 samples per class were included in model training to ensure robust analysis.

140 Model training involved 731 patients from three GECCO cohorts (EPIC, CORSA, IWHS; Fig. 1C)
141 using seven-fold cross-validation to balance training and validation sets, ensure class distribution,
142 and identify a median-performing model. The resulting seven models were deployed on external
143 test sets. The primary test set consisted of 645 patients from two GECCO cohorts (CRA, WHI; Fig.
144 1C), ensuring a diverse and representative split while addressing missing cases (Tab. S1).
145 Exclusively female cohorts were included in both datasets (IWHS for training, WHI for testing). To

146 assess generalizability and enable comparison to the literature, selected models were validated on
147 a secondary test set comprising TCGA and CPTAC cohorts. 'External validation' and 'test set'
148 typically refer to the primary test set unless stated otherwise; cohort-wise analyses are explicitly
149 noted. In total, we trained eleven models using the same training set: a primary model incorporating
150 the broad target set (Tab. S2) and a secondary model excluding MSI as a target to evaluate its
151 impact on other predictions—both tested on the primary and secondary test sets. Additionally, we
152 trained nine single-target models, each dedicated to one of the main prediction targets, and
153 evaluated them on the primary test set. All models were evaluated in direct comparison to the
154 primary multi-target model and compared to the literature.

155 **Data Analysis of Genetic Alterations**

156 In CRC, there is a pronounced co-occurrence of genetic alterations with MSI.^{26,27} We employed
157 two distinct methods to investigate the interplay and dependencies between various genetic alter-
158 ations within the five GECCO cohorts: First, we used hierarchical clustering to group molecular
159 alterations in our dataset using the 'Euclidean' metric and 'Ward' procedure²⁸, using only alter-
160 ations with complete data (Tab. S1). Subsequently, we applied association rule mining²⁹ to reveal
161 the co-occurrence of gene mutations and conditions such as MSI, considering an initiating genetic
162 alteration ('antecedent') and a potentially resulting alteration ('consequent'). The presence of an
163 initiating alteration statistically increases the probability of observing a resulting alteration, without
164 implying a cause-effect biological sequence. We used six complementary metrics to assess co-
165 occurrence: 'Support', 'Confidence', 'Lift', 'Leverage', 'Conviction', 'Zhang's Metric'.

166 **Image Processing and Deep Learning Techniques**

167 According to the framework of Wagner et al.²⁵, the digitized WSIs were tessellated into tiles of
168 224×224 pixels, corresponding to 256×256 micrometers (Fig. 1B). Tiles predominantly containing
169 background (brightness value ≥ 224) and blurred regions, identified using Canny edge detection³⁰
170 (thresholds: 40, 100) with fewer than two percent of the tile's pixels identified as edges were
171 excluded. We used the pre-trained CTransPath feature extractor^{25,31} to extract a 768-dimensional
172 feature representation for each tile, which we used as an input for a subsequent transformer model
173 (Fig. 1B). Advancing upon previous work²⁵ we developed a model with an encoder-decoder
174 architecture to simultaneously predict multiple targets from tile embeddings (Fig. 1B). Specifically,
175 the 768-dimensional tile features are projected into a 512-dimensional space using a fully
176 connected layer to reduce complexity and enhance computational efficiency for the subsequent
177 model. The encoded tokens are decoded in class tokens³², each with a dimension of 1×512, with
178 one token dedicated for each prediction target. The decoded tokens are fed through a fully
179 connected layer to generate target-specific predictions with scores ranging from zero (negative
180 prediction) to one (positive prediction) for each class and patient. To address class imbalance
181 during training, cross-entropy loss is calculated per target, weighted by the inverse mutation

182 frequency, and summed to ensure proportional importance of rare mutations. The training and
183 deployment procedures were performed on a NVIDIA RTX A6000 featuring 48 GB of GPU memory.
184 All source codes for preprocessing and DL are available in our open source repositories (Tab. S3).
185 TRIPOD reporting guidelines were used.

186 **Explainability**

187 To investigate the model's detection of relevant regions for predicting alterations in CRC, we gen-
188 erated heatmaps using Grad-CAM³³, focusing on the fold with the median AUROC for MSI detec-
189 tion to ensure consistent comparisons. These heatmaps highlight the contribution of each tile to
190 patient-level predictions, enabling analysis of whether the model relies on distinct patterns and
191 areas in WSIs for different targets or similar features across targets. Notably, a high number of
192 positively contributing tiles does not necessarily correspond to a high final score due to the non-
193 linear aggregation process, and the model may also incorporate global cues not captured by the
194 heatmap. To examine morphology in greater detail, we identified highly predictive tiles ('top tiles')
195 for representative cases, selected based on their scores and attention values assigned by the
196 model. For direct comparability, we included top tiles for each prediction target alongside those for
197 MSI from the same slide. As an additional mode of explainability, we analyzed class token in-
198 teractions in the decoder during deployment of the primary model to assess overlap for the main
199 prediction targets. Grad-CAM was used to compute activations, capturing class token-score inter-
200 actions, which were aggregated into a cross-correlation matrix, averaged across patients, and vi-
201 sualized as a heatmap.

202 **Statistical Analysis**

203 We used descriptive statistics to summarize sociodemographic and clinicopathologic characteris-
204 tics across cohorts (Tab. S1). Model performance was evaluated using AUROC during training and
205 testing, with additional metrics and target-specific prediction scores (ranging from 0.00 to 1.00)
206 analyzed. For the main prediction targets we applied a two-sided DeLong test to compare AUROCs
207 of single-target models and the primary multi-target model on the primary external dataset, utilizing
208 mean prediction scores across seven folds per target. Mean and median AUROC values for each
209 target were calculated over the seven folds. DeLong tests were also applied to compare the primary
210 model (including MSI as a target) with the secondary model (excluding MSI) within each of the four
211 external cohorts, based on mean prediction scores from the seven models.

212 Given the rarity of mutations and the limitations of metrics such as AUROC and the Area Under
213 the Precision-Recall Curve (AUPRC) in capturing imbalanced distributions³⁴, we employed diverse
214 metrics and analyzed individual prediction scores for the main prediction targets, stratified by MSI
215 status. Average prediction scores across seven folds were calculated per patient and categorized
216 into four subgroups per target based on microsatellite and mutational status: MSS/WT and

217 MSS/MUT represent MSS tumors with wild-type (WT) and mutated (MUT) targets, respectively,
218 while MSI/WT and MSI/MUT represent MSI tumors with wild-type and mutated targets (Fig. 1D).
219 Model mutation detection performance was evaluated through comparisons on training and testing
220 sets. After testing for normality with the Shapiro-Wilk test, the Mann-Whitney U test was used to
221 assess statistical significance between MSI and prediction target scores within each subgroup.
222 Additionally, four comparisons using the Wilcoxon test evaluated the model's ability to identify
223 prediction target mutations by comparing score distributions for MSS/WT vs MSS/MUT and
224 MSI/WT vs MSI/MUT. Finally, subgroup comparisons based on prediction target status were
225 performed: wild type (MSS/WT vs MSI/WT) or mutated (MSS/MUT vs MSI/MUT).

226 **Role of the funding source**

227 The funders of the study had no role in data collection, analysis, interpretation, writing of the
228 manuscript and the decision to submit.

229 **Results**

230 **Multi-Target Transformers Match and Partly Exceed Single-Target Models While** 231 **Enabling Simultaneous Prediction of Genetic Alterations**

232 Our approach aimed to reproduce prior DL-studies predicting genetic alterations in CRC from H&E
233 slides^{5-9,11-13,35,36}, expanding the scope to a broader set of alterations using a multi-target
234 transformer architecture (Fig. 1B). We compared the performance of our transformer for selected
235 targets including the main prediction targets (Fig. 2) with the external validation AUROCs reported
236 in the literature for single-target models, which used varying model architectures.

237

238 For MSI detection, single- and multi-target transformers achieved AUROCs of 0.91 (± 0.02) and
239 0.93 (± 0.01), respectively, on the primary test set ($p < 0.05$; refer to Tab. S4, Fig. 3A here and in
240 the following). The primary model achieved AUROCs from 0.87 (± 0.01 , TCGA) to 0.94 (± 0.01 ,
241 WHI) (Tab. S5-S8, Fig. S1; referred to here and in the following), consistent with the reported
242 literature range of 0.77–0.96.^{4,7,8,10,12,13,25,36} The following summarizes model performance for
243 external validation for selected targets, with AUROCs for single- and multi-target transformers
244 (Tab. S4, Fig. 3A) and comparisons between primary (multi-target with MSI) and secondary (multi-
245 target without MSI) models (Tab. S5-S12, Fig. S2). In *BRAF* MUT detection, the the AUROC was
246 0.72 (± 0.06 , single-target) and 0.78 (± 0.01 , multi-target, $p < 0.05$). The AUROCs for primary vs.
247 secondary models were lowest on WHI (0.77 (± 0.01) vs. 0.76 (± 0.01 , $p > 0.05$)) to highest on CRA
248 (0.83 (± 0.02) vs. 0.82 (± 0.03 , $p > 0.05$)). These results align with the literature range of 0.66 to
249 0.88.^{7,8,10-13,16,25,36} In *RNF43* MUT detection, the AUROC was 0.80 (± 0.05 , single-target) and 0.86
250 (± 0.01 , multi-target, $p < 0.05$). The AUROCs for primary vs. secondary models ranged from 0.80

251 (± 0.01) vs. 0.79 (± 0.01 , $p > 0.05$) on TCGA to 0.87 (± 0.02) vs. 0.85 (± 0.01 , $p > 0.05$) on CRA. These
252 results exceed the literature range of 0.63 to 0.72.^{11,16} For KRAS MUT detection, the AUROC was
253 0.65 (± 0.02) for both single- and multi-target models ($p > 0.05$). The primary vs. secondary models
254 ranged from 0.56 (± 0.02) vs. 0.55 (± 0.02 , TCGA, $p > 0.05$) to 0.69 (± 0.06) vs. 0.72 (± 0.03 , CPTAC,
255 $p > 0.05$). These results fall below and within the literature range of 0.60 to 0.80.^{8,11,12,16,25,36}
256 Interestingly, CPTAC had no KRAS MUT cases with MSI, while TCGA showed the highest number
257 of such cases among all cohorts, with up to 19 instances. For hypermutation and TP53 MUT
258 detection, the results fell within the literature ranges of 0.81–0.87^{8,10,13} and 0.60–0.75^{8,10,11,16}, with
259 no significant differences between models. For APC mutational status prediction, the results
260 ranged both above and below the single available literature AUROC of 0.67^{11,16}, with no statistical
261 differences between models, highlighting the limited number of comparable values. For BMPR2
262 and ZNRF43 there was no comparable data in the literature (AUROCs shown in Tab. S4, Fig. 3A,
263 Fig. S2). Additional prediction targets with AUROCs ≥ 0.75 are detailed in Tab. S4.

264 **Enhanced Predictive Performance for Genetic Alterations Associated with MSI**

265 Hierarchical clustering identified two primary genetic clusters: Cluster 2, comprising MSI-associ-
266 ated genes (e.g., BRAF, BMPR2, ZNRF3, RNF43) with strong co-occurrence with MSI and hy-
267 permutation, and Cluster 1, including MSS-associated genes such as TP53, KRAS, and APC (Fig.
268 2). Association rule mining reinforced robust MSI correlations in Cluster 2 (e.g., BMPR2, RNF43),
269 contrasting with inverse relationships in Cluster 1 MSS-linked genes (e.g., TP53, KRAS) (Fig. S3A,
270 Tab. S13).

271
272 MSI-associated alterations showed higher AUROC values for predicting genetic mutations com-
273 pared to MSS-associated alterations (Cluster 1) (Fig. 3B-C, Tab. S14, Fig. S1). Additional metrics
274 (Fig. 3B, Fig. S3B, Tab. S14) and patient-level prediction score distributions were analyzed for MSI
275 and genetic alterations within Clusters 1 and 2 (Fig. 1D, Fig. 4, Tab. S15-S16). Prediction scores
276 quantified model certainty, with high scores indicating mutation presence, low scores indicating
277 absence, and scores near 0.5 reflecting indecision. Comprehensive summaries supported internal
278 (Fig. S4) and external validation (Fig. S5). For MSS-associated Cluster 1 genes (TP53, APC,
279 KRAS) the external validation AUROCs ranged from 0.65 to 0.72, with MSI scores accurately
280 reflecting ground truth: low for MSS and high for MSI cases (Fig. 2, Fig. 4A, Tab. S4). High MSI
281 scores tended to align with WT target predictions, while low scores aligned with MUT predictions.
282 As a result, discrepancies were observed in MSS/WT and MSI/MUT subgroups, where target
283 predictions deviated from ground truth. Genetic Cluster 2, including MSI, hypermutation, and the
284 genes BMPR2, ZNRF3, RNF43, and BRAF, demonstrated higher AUROCs (0.75–0.88) during
285 external validation (Fig. 2, Fig. 4B, Tab. S4). In MSS patients, mutations in biomarkers like BMPR2
286 ($n=3$) and ZNRF3 ($n=9$) were rare, and MSS/MUT subgroups showed less distinct prediction
287 scores for RNF43 ($n=20$) and BRAF ($n=39$), reflecting higher prediction uncertainty. MSI scores

288 aligned with the ground truth (MSS/MSI), but unlike Cluster 1, alteration scores in all subgroups
289 correlated with MSI scores. This produced accurate prediction trends for MSS/WT and MSI/MUT
290 subgroups (low scores for WT, high for MUT), while MSS/MUT and MSI/WT subgroups showed
291 deviations from the target ground truth. These findings, partially reflected in AUROC results (Fig.
292 3B, Fig. S3B), suggest less effective differentiation between MUT and WT in MSS and MSI
293 subgroups compared to the combined group. Despite the strong influence of MSI-associated
294 morphology on predictions, aligning the prediction target scores with MSS (Cluster 1) or MSI-high
295 (Cluster 2), intrinsic phenotypes for some alterations cannot be fully confirmed or excluded.
296 Analyses of MSI-high and MSS subgroups (Tab. S17-S18, Fig. 3B, Fig. S3B-S3C) revealed
297 AUROCs of 0.60–0.70 for genes like *BRAF*, *RNF43*, and *TP53* in MSS patients, indicating modest
298 discrimination between MUT and WT. Statistical differences in score distributions for MSI and
299 these genetic alterations further supported these findings (Fig. 4A-B).

300 **MSI-Associated Morphological Patterns Drive Predictions Across Prediction Tar-** 301 **gets**

302 To investigate the morphological patterns underlying the prediction of genetic alterations in CRC
303 from H&E slides using DL, we analyzed WSI heatmaps (Fig. 5, Fig. S6-S10) and highly predictive
304 tiles ('top tiles', Fig. 6, Fig. S11-S23) for key prediction targets (Cluster 1: *TP53*, *APC*, *KRAS*;
305 Cluster 2: *RNF43*, *BRAF*, hypermutated) in 25 cases from CRA and WHI test cohorts. This was
306 assessed via manual histopathological review of WSI heatmaps and 20 top tiles per target (Tab.
307 S19). We found that our model predominantly focused on tumor regions, with minimal relevance
308 attributed to pen marks or non-tumor areas (e.g., Fig. 5A-D, Fig. S6B-D, Fig. S7A-D). Pen marks,
309 present on most slides, were rarely highlighted and appeared faintly in top tiles (e.g., Fig. S13, Fig.
310 S22), though slight highlighting in heatmaps occurred in rare cases (e.g., Fig. 5C, Fig. S9D).
311 Regions of high model attention primarily consisted of tumor tissue, not normal intestinal mucosa
312 or uninvolved connective/adipose tissue, despite the model not being explicitly trained for tumor
313 detection.

314
315 Frequent MSI-associated patterns in most MSI cases and Cluster 2 gene mutations included
316 medullary growth, high tumor-infiltrating lymphocytes (TILs), and mucinous differentiation, consis-
317 tent with prior studies.^{37–39} However, some MSI cases in our cohort and the literature display atypi-
318 cal morphologies not commonly linked to MSI (e.g., Fig. S8D, Fig. S9D).⁴⁰ In MSS-related Cluster
319 1, *KRAS* mutation predictions highlighted luminal tumor regions, especially villous adenomas with
320 high-grade dysplasia, as equally important as invasive adenocarcinoma areas (Fig. S16A-B, Fig.
321 S17A-B). This aligns with the known prevalence of *KRAS* mutations in villous adenomas⁴¹ and
322 carcinomas with close proximity to polyps⁴², reflecting the algorithm's ability to capture intratumoral
323 heterogeneity. Tiles with high MSI relevance displayed medullary carcinoma patterns, including
324 tumor cell sheets and high TILs, which corresponded to low *KRAS* MUT prediction scores (Fig. 6A-

325 B). Alterations such as *TP53* and *APC*, less common in MSI-high cases, relied on MSS-like
326 morphology (e.g., gland-forming conventional adenocarcinoma with dirty necrosis⁴⁰; Fig. S18-S20),
327 while MSI-associated alterations (e.g., *BRAF*, *RNF43*, hypermutation) depended on MSI-related
328 features, including medullary patterns, mucinous differentiation (with signet-ring cells; Fig. S23B),
329 and high TILs⁴⁰ (Fig. 6C-D, Fig. S12-S15, Fig. S21, Tab. S19). To explore target-specific patterns
330 beyond MSI, a pathological review of top tiles from selected slides (Fig. S22-S23) identified tumor
331 budding as a morphological feature associated with *BRAF* mutations in MSS cases (Fig. S23A-
332 C).⁴³

333 Discussion

334 Previous research in CRC has focused on high AUROCs for biomarker detection, showing MSI as
335 the most detectable, while often neglecting co-occurrence with other mutations and lacking in-
336 depth analyses of multiple genetic alterations and their links to MSI-associated morphologies. To
337 address these limitations, we developed a DL model to predict multiple genetic alterations in CRC
338 from H&E slides, using five GECCO cohorts with unified sequencing data with additional validation
339 on two public CRC-datasets, TCGA and CPTAC. The multi-target model achieved performance
340 within the literature range, outperforming single-target models for some alterations while enabling
341 the efficient and sustainable simultaneous prediction of numerous prediction targets. This in-depth
342 analysis was made possible by the comprehensive dataset and novel model architecture, both of
343 which are shared open-source with the scientific community.

344 Our cohort, comprising multiple sub-cohorts, is broadly representative of CRC patients⁴⁴, particu-
345 larly in Western populations, with data primarily from the US and Europe. It includes a higher pro-
346 portion of female patients (IWHS, WHI) and disproportionately includes White individuals, while
347 TCGA had the highest proportion of Black/African American patients, and CPTAC was limited to
348 fresh-frozen colon samples with the smallest cohort. MSI frequency ranged from 7–35%, aligning
349 with reported rates of 15–20%⁴⁵ (Tab. S1). Variability in target mutation frequencies across cohorts,
350 remained within plausible ranges, as did the experimental results. However, small sample sizes
351 for some prediction targets likely contributed to higher variability (Fig. S2), underscoring the need
352 for larger, more diverse datasets. The genetic clusters analyzed included key CRC biomarkers
353 (e.g., *TP53*, *APC*, *KRAS*, *BRAF*)^{7,8,11–13}, with mutational profiles and MSI co-occurrences aligning
354 with prior studies^{26,27,46}.

355 In line with the literature MSI was the most predictable alteration^{7,8,12,13,25,36}, though some MSI
356 cases scored low due to a lack of associated morphological features, which is plausible and
357 pathological assessment has shown. For the detection of multiple genetic alterations apart from
358 MSI, our multi-target model outperformed several of our own single-target models and models from
359 the literature, such as for detection of *RNF43*, using the AUROC as a common performance
360 measure.^{4–9,35,36} Notably, hypermutation prediction demonstrates a strong interdependence with

361 MSI due to their frequent co-occurrence and close biological linkage, with shared morphological
362 features leading to lower scores for hypermutated MSS cases. *BRAF* mutations, critical in CRC for
363 their prognostic and predictive significance⁴⁷, showed detectable phenotypic changes, including
364 mucinous differentiation and poorly differentiated clusters. The histopathological evaluation
365 revealed features such as mucinous differentiation and stroma-rich patterns dominating *BRAF*
366 MUT and *BRAF*WT tiles, respectively^{12,47}. Interestingly, we could identify mucinous/signet-ring dif-
367 ferentiation relevant for *BRAF* predictions.⁴⁷ However, these distinct features overlapped with MSI-
368 associated morphology, confounding specificity and occasionally leading to misclassification in
369 MSS cases. In contrast, some top tiles for *BRAF* MUT prediction showed poorly differentiated
370 clusters and tumor budding, which has already been linked to *BRAF* mutations⁴⁸, and is associated
371 with MSS.⁴³ While mutations like *BRAF* induce distinct phenotypic changes detectable by DL, MSI
372 morphology often overshadows fine-grained, target-specific patterns. Excluding MSI as a
373 prediction target did not affect the predictability of other prediction targets, and independent class
374 token behavior with no cross-token information flow in the decoder, a key feature of the model
375 architecture, was demonstrated (Fig. S24). In summary, we have shown that CRC mutation
376 predictability is primarily driven by MSI-associated morphology, the most discriminative feature for
377 mutation detection, though subgroup analyses of MSI and MSS suggest subtle mutation-specific
378 patterns warrant further study.

379 This study has several limitations. Despite its extent, the dataset was limiting the detection of rare
380 alterations and subtle mutation-specific morphologies, highlighting the need for larger and more
381 diverse datasets. Underrepresentation of Black/African-American individuals in the GECCO
382 dataset and missing or incomplete annotations further constrain generalizability and accuracy.
383 While monochrome tiles were excluded during preprocessing, some pen markings remained.
384 These markings, though occasionally attracting slight model attention and potentially introducing
385 bias, were largely transparent, limited to a few tiles, and may hold biologically relevant information,
386 particularly near tumor borders, warranting further investigation. Additionally, our study em-
387 phasizes the need for advanced explainability methods to validate DL-detected pathological pat-
388 terns and link them to specific prediction targets, fostering a bidirectional exchange between
389 pathology and DL. This could clarify whether MSI-associated morphology dominates model pre-
390 dictions due to the absence of other target-specific features or its overwhelming influence. Ap-
391 proaches such as modified loss functions could help mitigate MSI dominance and enhance de-
392 tection of subtle features⁴⁹ and integrating multimodal data could improve generalizability, accu-
393 racy, and understanding of mutation-specific morphologies. Our model expands on single-target
394 transformers and raises important questions regarding potential bias toward MSI-associated tar-
395 gets due to their larger subgroup size; however, as our model performs within the range reported
396 in the literature, further investigation is required to elucidate the influence of MSI- and non-MSI-
397 associated features, providing insights into the mechanisms of single- and multi-target transform-
398 ers. Similar to prior research^{34,50}, we show that AUROC can be a misleading metric for evaluating

399 biomarker detection; thus, we report it for comparability but include additional metrics. For future
400 studies, we recommend using extensive, diverse datasets, analyzing multiple prediction targets
401 and their interactions with various metrics (e.g. Accuracy, Precision, Sensitivity, Specificity, F1
402 Score), and addressing potential confounding factors by leveraging approaches such as our pro-
403 posed model architecture.

404 From a practical point of view, our multitarget DL-based biomarker prediction provides significant
405 value for research by identifying which prediction targets can be reliably predicted and the mor-
406 phological basis for these predictions, guiding future studies and highlighting prediction targets that
407 merit further investigation. For some alterations, such as *BRAF*, existing surrogate immuno-
408 histochemical assays could be applied⁵¹ to predicted regions in follow-up studies to deepen un-
409 derstanding of the biological mechanisms behind predictions. Clinically, this approach streamlines
410 colorectal cancer testing by eliminating the need for separate workflows and serving as a cost-
411 effective prescreening tool, particularly in resource-limited settings or earlier CRC stages.
412 Biologically, it sheds light on how genetic alterations influence tumor morphology and underscores
413 the dominant role of surrogate markers like MSI in CRC histology. As the comprehensive analysis
414 of co-occurrences in CRC is scarcely explored in relation to DL and is rare for other cancer types⁵²,
415 this efficient, scalable method can be extended to any cancer type, enabling simultaneous analysis
416 of multiple prediction targets and advancing both research and precision oncology.

417 In conclusion, our study highlights the utility of multi-target transformers for detecting biomarker-
418 specific patterns in CRC H&E images, with MSI phenotype as the dominant factor influencing
419 predictability. The model's reliance on MSI- and MSS-associated morphology underscores the
420 importance of morphology and co-occurrence patterns over AUROC metrics. By enabling the effi-
421 cient investigation of numerous prediction targets simultaneously, this model facilitates a more
422 comprehensive understanding of biomarker interactions. These findings emphasize the need to
423 consider individual prediction scores and tumor-specific confounders in biomarker studies, laying
424 the foundation for future research to address these effects across other cancer types.

425 **List of abbreviations**

426 AUROC: Area under the Receiver Operating Characteristic Curve; AUPRC: Area Under the Pre-
427 cision-Recall Curve; CIN: Chromosomal instability; CPTAC: Clinical Proteomic Tumor Analysis
428 Consortium; CRC: Colorectal cancer; DL: Deep Learning; H&E: Hematoxylin and eosin; Mb:
429 Megabases; MSI: Microsatellite instability; MSS: Microsatellite stable; MUT: Mutated; NOS: Not
430 otherwise specified; px: Pixel; ROC: Receiver Operating Characteristic Curve; TCGA: The Cancer
431 Genome Atlas; TILs: Tumor infiltrating lymphocytes; ViT: Vision Transformer; WSI: Whole Slide
432 Image; WT: Wild type

433 **Declarations**

434 **Ethics approval and consent to participate**

435 This study was performed in accordance with the Declaration of Helsinki. This study is a retro-
436 spective analysis of scanned images of anonymized tissue samples of various cohorts of cancer
437 patients. Data were collected and anonymized and ethical approval was obtained. The overall
438 analysis was approved by the Ethics board of the Medical Faculty of Technical University Dresden
439 under the ID BO-EK-444102022.

440 **Contributions**

441 MG and JNK conceptualized the study. CS, PH, LAB, AJF, ELG, AG, MJG, PL, NM, ST and UP
442 provided clinical and digital histopathological data. MG curated the source data and conducted the
443 literature research. MvT implemented the source code for the deep learning pipeline and the
444 heatmaps. MG developed the code for data analysis and visualization. MG planned and conducted
445 the experiments, visualizations and analysis of data and results. ARM assisted with the statistical
446 evaluation. MG, CMLL, ZIC and NGR interpreted the data. NGR and BM did the pathological
447 examination of the samples and analyzed the heatmaps and top/bottom tiles. MG wrote the first
448 draft of the manuscript. All authors revised the manuscript draft, jointly interpreted the data and
449 agreed to the submission of this article. All authors had access to all the data, and MG and JNK
450 have verified the data.

451 **Declaration of interest**

452 JNK declares consulting services for Biopimus, France; Owkin, France; DoMore Diagnostics,
453 Norway; Panakeia, UK; AstraZeneca, UK; Mindpeak, Germany; and MultiplexDx, Slovakia. Fur-
454 thermore, he holds shares in StratifAI GmbH, Germany, Synagen GmbH, Germany; has received
455 a research grant by GSK; and has received honoraria by AstraZeneca, Bayer, Daiichi Sankyo,
456 Eisai, Janssen, Merck, MSD, BMS, Roche, Pfizer, and Fresenius. MG has received honoraria for
457 lectures sponsored by Techniker Krankenkasse (TK) and AstraZeneca. SF has received honoraria
458 for lectures by BMS and MSD. UP declares consulting services for AbbVie and her husband is
459 holding individual stocks for the following companies: BioNTech SE – ADR, Amazon, CureVac BV,
460 NanoString Technologies, Google/Alphabet Inc Class C, NVIDIA Corp, Microsoft Corp.. No other
461 potential conflicts of interest are reported by any of the authors.

462 **Data sharing**

463 All source code we used to conduct this study is publicly available with publication. The code for
464 image preprocessing is publicly available at <https://github.com/KatherLab>. The tessellation script is
465 available at <https://github.com/KatherLab/preprocessing-ng>. Extraction of CTransPath features
466 was conducted with scripts from <https://github.com/KatherLab/marugoto>. The scripts for the multi-

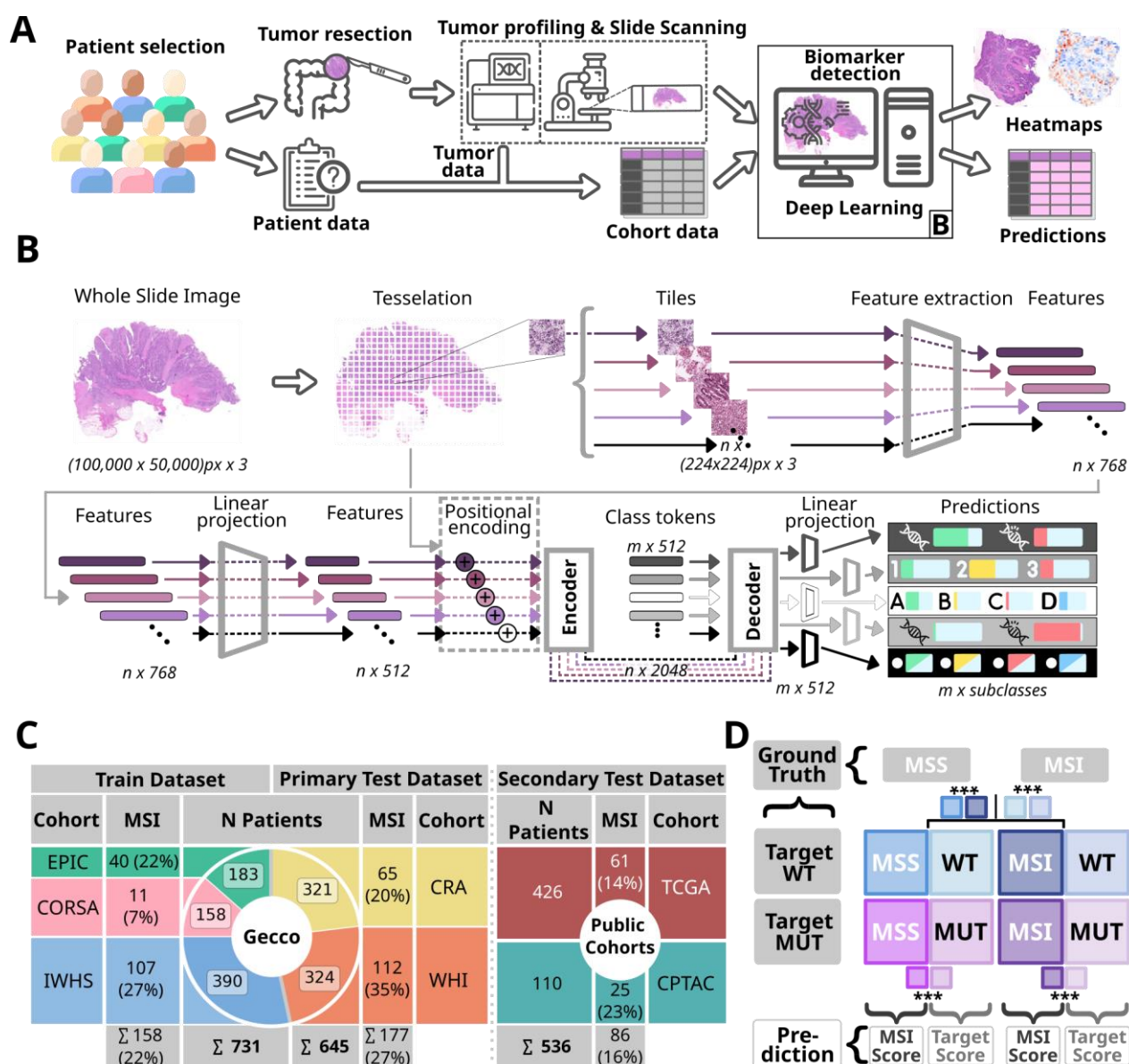
467 target transformer and the heatmaps used for explainability can be accessed at <https://github.com/LocalToasty/barspoon-transformer/>. The saved checkpoints for the trained models with the
468 lowest validation loss for all seven folds can be obtained from
469 <https://github.com/gustavmarco/barspoon-transformer/releases/tag/gustav2024>. Links to the exact repository versions used can be found in Tab. S3. The digitized whole slide images (WSIs) for
470 the TCGA²² cohort are publicly accessible at <https://portal.gdc.cancer.gov/> and
471 <https://www.cbioportal.org/>. The digitized WSIs for the CPTAC cohort are publicly accessible at
472 The Cancer Imaging Archive (TCIA)⁵³. The molecular and clinical data for TCGA and CPTAC is
473 publicly accessible at <https://portal.gdc.cancer.gov/> and <https://www.cbioportal.org/>. The datasets
474 from the GECCO consortium are available from the corresponding author on reasonable request.
475 All data generated or analyzed during this study are included in this published article and its supplementary information files.
476
477
478

479 Acknowledgements

480 JNK is supported by the German Federal Ministry of Health (DEEP LIVER, ZMVI1-2520DAT111),
481 the Max-Eder-Programme of German Cancer Aid (grant #70113864), the German Federal Ministry
482 of Education and Research (PEARL, 01KD2104C; CAMINO, 01EO2101; SWAG, 01KD2215A;
483 TRANSFORM LIVER, 031L0312A; TANGERINE, 01KT2302 through ERA-NET Transcan), the
484 German Academic Exchange Service (SECAI, 57616814), the German Federal Joint Committee
485 (Transplant.KI, 01VSF21048) the European Union's Horizon Europe and innovation programme
486 (ODELIA, 101057091; GENIAL, 101096312) and the National Institute for Health and Care
487 Research (NIHR, NIHR213331) Leeds Biomedical Research Centre. The views expressed are
488 those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health
489 and Social Care. SF is supported by the German Federal Ministry of Education and Research
490 (SWAG, 01KD2215C), the German Cancer Aid (DECADE, 70115166 and TargHet, 70115995) and
491 the German Research Foundation (504101714). The Genetics and Epidemiology of Colorectal
492 Cancer Consortium (GECCO) is funded by: National Cancer Institute, National Institutes of Health,
493 U.S. Department of Health and Human Services (U01 CA137088, R01 CA488857, P20
494 CA252733). Genotyping/Sequencing services were provided by the Center for Inherited Disease
495 Research (CIDR) contract number HHSN268201700006I. This research was funded in part
496 through the NIH/NCI Cancer Center Support Grant P30 CA015704. Scientific Computing
497 Infrastructure at Fred Hutch funded by ORIP grant S10OD028685. The CORSA study was funded
498 by Austrian Research Funding Agency (FFG) BRIDGE (grant 829675, to Andrea Gsur), the
499 "Herzfelder'sche Familienstiftung" (grant to Andrea Gsur) and was supported by COST Action
500 BM1206. CRA was supported by the National Institutes of Health grant R01 CA068535. The
501 coordination of EPIC is financially supported by the International Agency for Research on Cancer
502 (IARC) and also by the Department of Epidemiology and Biostatistics, School of Public Health,
503 Imperial College London which has additional infrastructure support provided by the NIHR Imperial
504 Biomedical Research Centre (BRC). The national cohorts are supported by: Danish Cancer
505 Society (Denmark); Ligue Contre le Cancer, Institut Gustave Roussy, Mutuelle Générale de
506 l'Éducation Nationale, Institut National de la Santé et de la Recherche Médicale (INSERM)
507 (France); German Cancer Aid, German Cancer Research Center (DKFZ), German Institute of
508 Human Nutrition Potsdam- Rehbruecke (DIfE), Federal Ministry of Education and Research
509 (BMBF) (Germany); Associazione Italiana per la Ricerca sul Cancro-AIRC-Italy, Compagnia di
510 SanPaolo and National Research Council (Italy); Dutch Ministry of Public Health, Welfare and
511 Sports (VWS), Netherlands Cancer Registry (NKR), LK Research Funds, Dutch Prevention Funds,
512 Dutch ZON (Zorg Onderzoek Nederland), World Cancer Research Fund (WCRF), Statistics
513 Netherlands (The Netherlands); Health Research Fund (FIS) - Instituto de Salud Carlos III (ISCIII),
514 Regional Governments of Andalucía, Asturias, Basque Country, Murcia and Navarra, and the
515 Catalan Institute of Oncology - ICO (Spain); Swedish Cancer Society, Swedish Research Council
516 and and Region Skåne and Region Västerbotten (Sweden); Cancer Research UK (14136 to EPIC-
517 Norfolk; C8221/A29017 to EPIC-Oxford), Medical Research Council (1000143 to EPIC-Norfolk;
518 MR/M012190/1 to EPIC-Oxford) (United Kingdom). The IWHS study was supported by NIH grants
519 CA107333 (R01 grant awarded to P.J. Limburg) and HHSN261201000032C (N01 contract
520 awarded to the University of Iowa). The WHI program is funded by the National Heart, Lung, and
521 Blood Institute, National Institutes of Health, U.S. Department of Health and Human Services

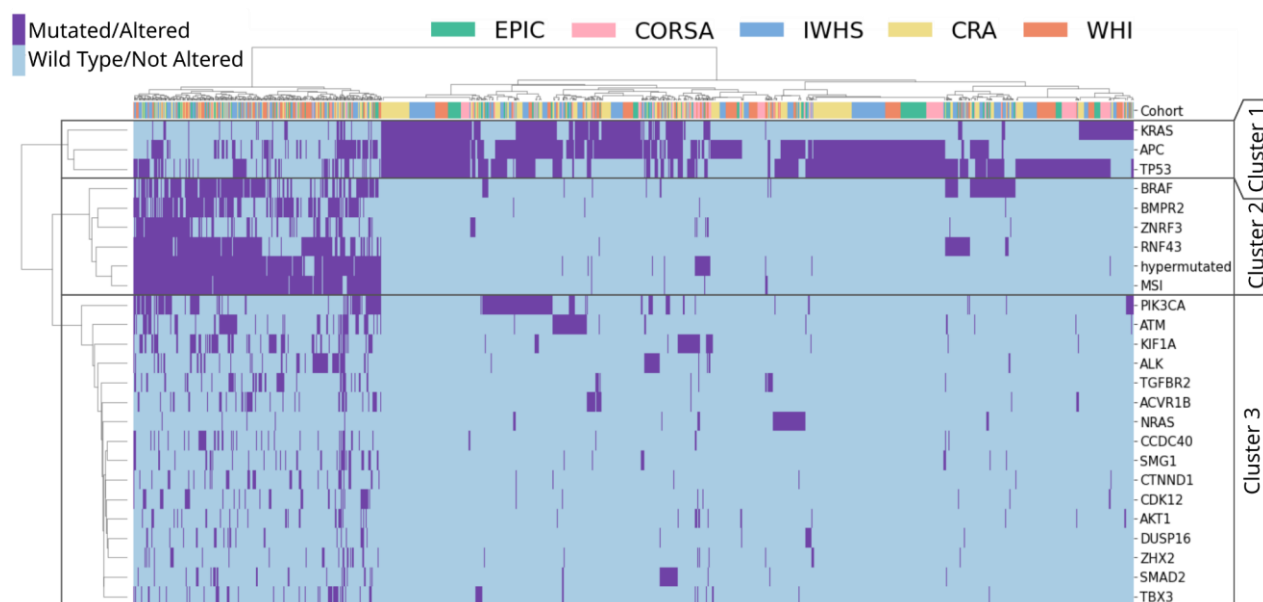
522 through contracts 75N92021D00001, 75N92021D00002, 75N92021D00003, 75N92021D00004,
523 75N92021D00005. We kindly thank all individuals who agreed to participate in the CORSA study.
524 Furthermore, we thank all cooperating physicians and students and the Biobank Graz of the
525 Medical University of Graz. We also acknowledge the TCGA Research Network and the Clinical
526 Proteomic Tumor Analysis Consortium (CPTAC), which generated the data on which some of the
527 results shown in this study are based. Where authors are identified as personnel of the International
528 Agency for Research on Cancer/World Health Organization, the authors alone are responsible for
529 the views expressed in this article and they do not necessarily represent the decisions, policy or
530 views of the International Agency for Research on Cancer/World Health Organization. The authors
531 thank the WHI investigators and staff for their dedication, and the study participants for making the
532 program possible. A full listing of WHI investigators can be found at:
533 [http://www.whi.org/researchers/Docu-](http://www.whi.org/researchers/Documents%20%20Write%20a%20Paper/WHI%20Investigator%20Short%20List.pdf)
534 [ments%20%20Write%20a%20Paper/WHI%20Investigator%20Short%20List.pdf](http://www.whi.org/researchers/Documents%20%20Write%20a%20Paper/WHI%20Investigator%20Short%20List.pdf). The readability
535 and language of the work were improved using ChatGPT-4o and DeepL.

536 **Figures**



537 Fig. 1: **Experimental design, cohort characterization, and schematic for predictive analysis.**

538 **A.** Tissue samples from colorectal cancer (CRC) patients across five independent cohorts were obtained via
 539 surgical resection, with associated demographic, clinical, and sequencing data collected. Upon Hematoxylin
 540 and Eosin (H&E) staining, tumor tissues are digitized into Whole Slide Images (WSIs) for profiling genetic
 541 alterations. The WSIs are then used to train and test a deep learning (DL) algorithm for biomarker detection,
 542 to simultaneously predict multiple mutational statuses and provide heatmap explanations. **B.** The DL pipeline
 543 tessellates the WSIs into smaller tiles while rejecting background and blurry areas, extracting n feature
 544 vectors from n tiles. Feature vectors are compressed and processed in a multi-target transformer, employing
 545 an attention mechanism in an encoder-decoder structure for class token learning. The transformer generates
 546 individual scores for the respective amount of classes per target. The code is able to comprise positional tile
 547 embedding (dashed lines), which did not result in improved performance and were therefore excluded from
 548 our study. **C.** Overview of the five GECCO and two public cohorts, including patient numbers, slides,
 549 extracted features, and MSI case proportions. The cohorts are divided into train datasets and test datasets.
 550 **D.** Schematic for interpreting result plots and statistics, delineating dataset partitioning based on
 551 microsatellite (MSS: microsatellite stability, MSI: microsatellite instability) and gene mutational status (MUT:
 552 mutated, WT: wild type). The diagram illustrates distinct groups by color, with the left side representing MSI
 553 prediction scores and the right side for prediction target scores. True ground truth labels of samples guide
 554 the group organization, with model-generated scores depicted in corresponding colors.



556

557

558

559

560

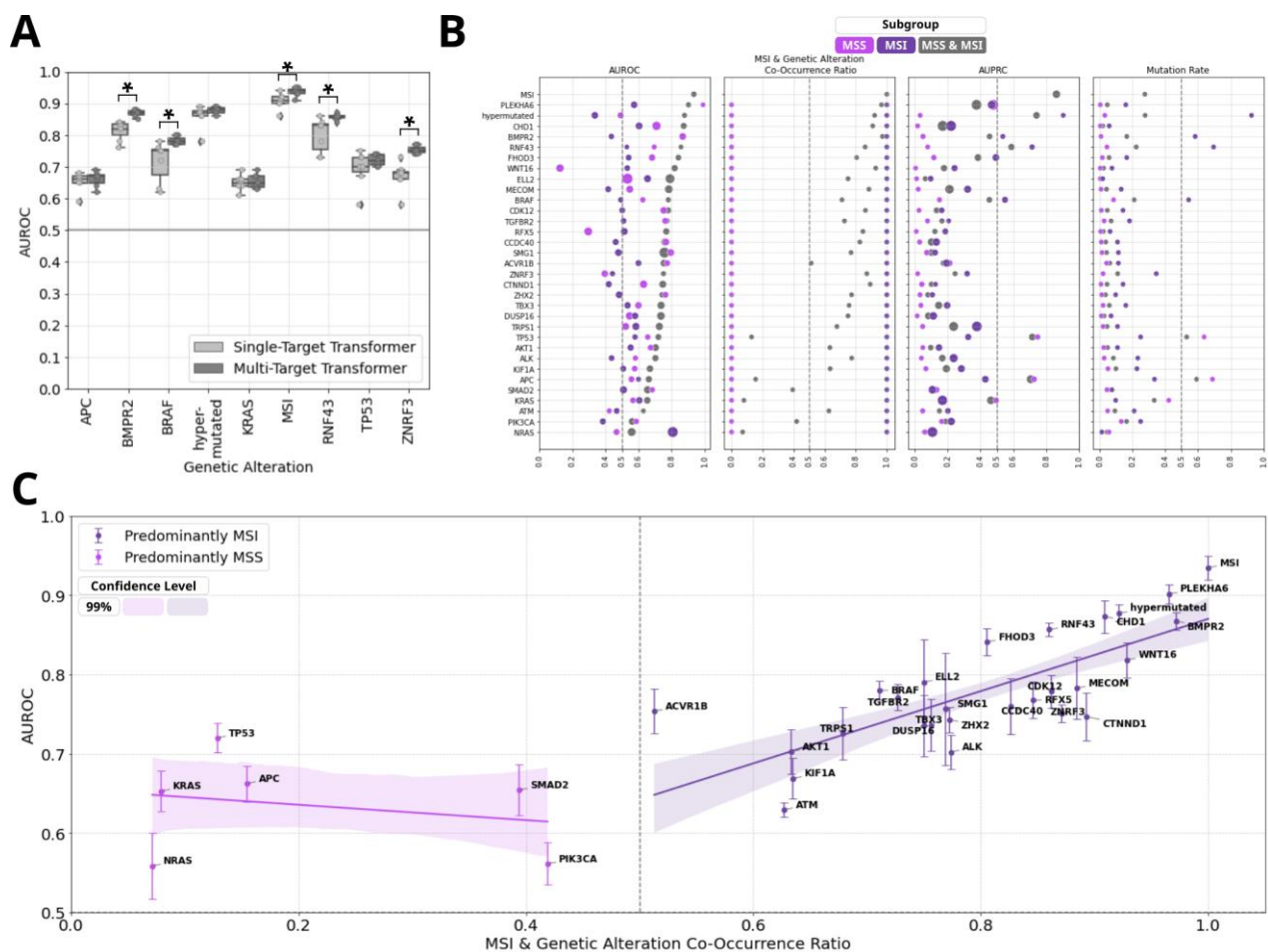
561

562

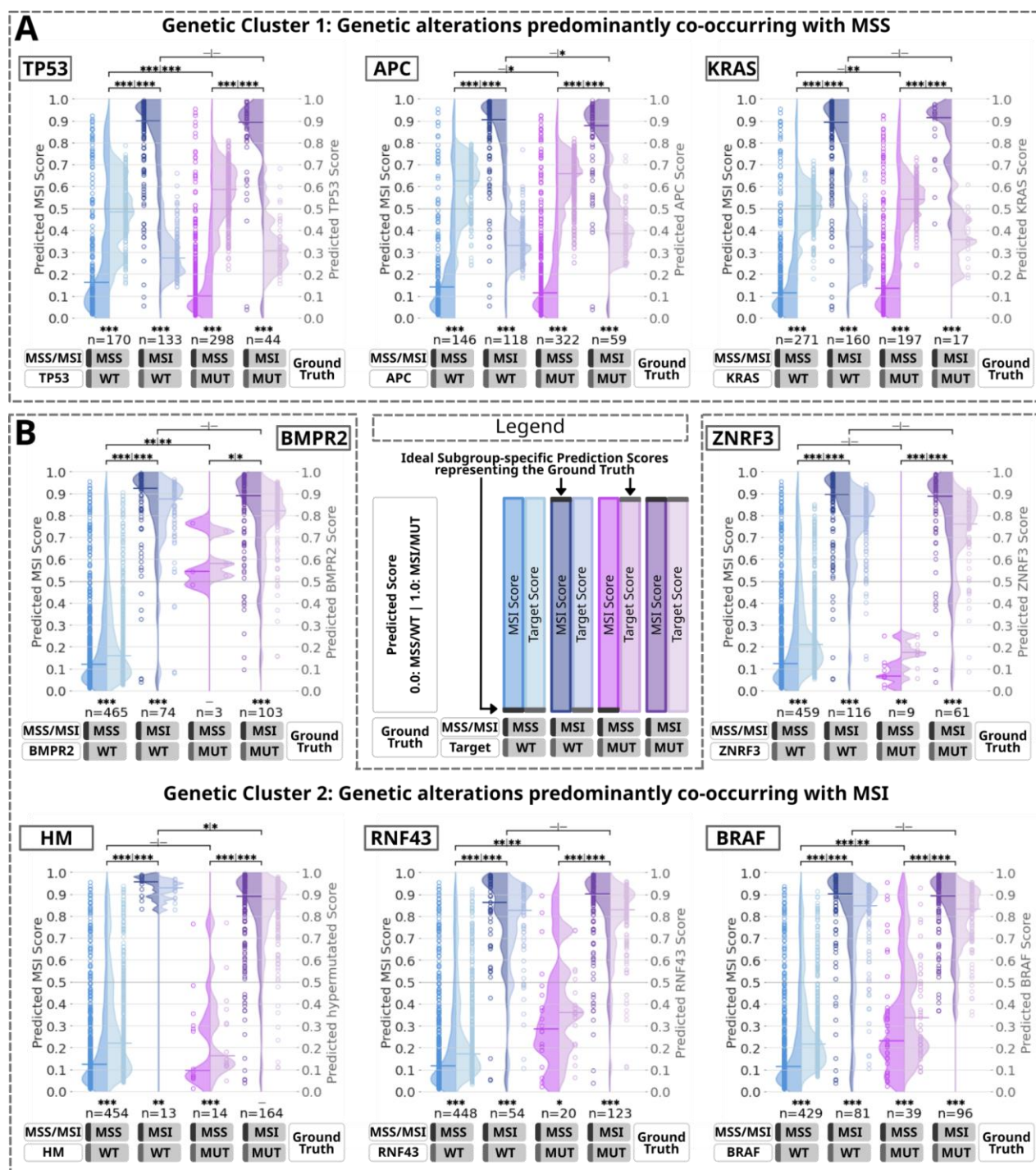
563

564

Fig. 2: Analysis of genetic alterations co-occurrence in CRC for GECCO cohorts. Hierarchical clustering analysis was conducted on the ground truth of genetic alterations with fully available mutational information. Each row corresponds to a genetic alteration, and each column represents a patient from the dataset. The top row indicates the distribution of patients from various cohorts within genetic clusters. The distance calculation was performed using the 'Euclidean' metric, and the 'Ward' method was applied to clustering. Three unique genetic clusters were created and marked. The patient clustering shows a diverse distribution of samples across all five cohorts and genetic clusters (top row).

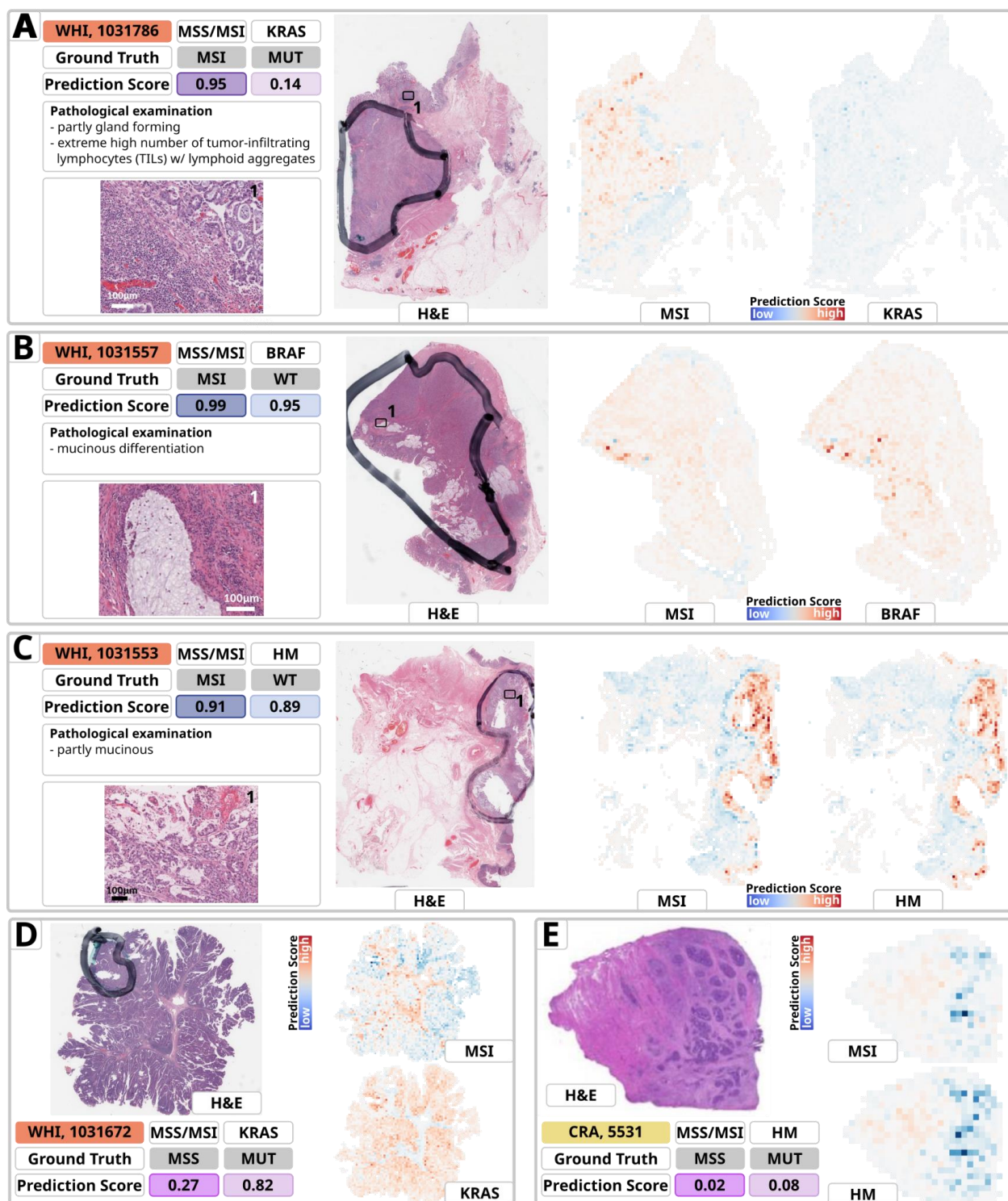


565
 566
 567 **Fig. 3: Evaluation of the performance of the Multi-Target Transformer on selected prediction**
 568 **targets for the external cohorts from GECCO. A.** The comparison of Single-Target Transformer
 569 versus Multi-Target Transformers shows the Area Under the Receiver Operating Characteristic
 570 curve (AUROC) from each of the 7 folds of external cross-validation, with the median value
 571 highlighted with a horizontal line in each box. The figure includes selected representative potential
 572 biomarkers of genetic alterations associated with MSS (Fig. 2, genetic Cluster 1) and MSI (Fig. 2,
 573 genetic Cluster 2). The test set cohorts consist of CRA and WHI (Fig. 1C). The horizontal line
 574 positioned at an AUROC of 0.50 represents a random guess of the model. Significance was
 575 determined through a two-sided DeLong test with a p-value threshold of less than 0.05. **B.**
 576 Performance metrics of Multi-Target Transformers for external validation. The mean (center of dot)
 577 and standard deviation (diameter of dot) for relevant selected prediction targets for the whole
 578 external set, as well as the MSI and MSS subgroups, are displayed based on the 7 folds of cross-
 579 validation. The threshold for binary classification is pre-defined as 0.50. The evaluation metrics
 580 include the Area Under the Receiver Operating Characteristic Curve (AUROC), and the Area
 581 Under the Precision-Recall Curve (AUPRC), along with the corresponding mutation rates in
 582 external cohorts. The Mutation Rate refers to the fraction of instances with a specific mutation in
 583 the subgroup. The MSI & Genetic Alteration Co-Occurrence Ratio is the fraction of cases harboring
 584 MSI among all cases with a particular genetic mutation. The data is sorted for AUROC and shown
 585 in Tab. S14 and Tab. S17-S18. An extended version of this panel with more metrics is shown in
 586 Fig. S3B. **C.** Distribution of Areas under the Receiver Operating Characteristic Curve (AUROCs,
 587 mean \pm standard deviation) for selected prediction targets and their co-occurrence with MSI with
 588 corresponding values and further metrics shown in Tab. S14. An extended version of this panel
 with MSS/MSI-subgroup specific AUROCs is shown in Fig. S3C.



589
590
591 **Fig. 4: Evaluation of prediction scores based on the multi-target transformer in external**
592 **validation on the GECCO test set subgrouped by the co-occurrence of the prediction targets**
593 **with MSI. A-B.** Violin plots representing individual patient scores from the test set cohorts for MSI
594 and representative genetic alterations in four subgroups based on microsatellite and alteration
595 mutational status. The left y-axis represents the MSI score scale (left violin halves) and the right y-
596 axis corresponds to the prediction target scores (right violin halves). The legend displays gray
597 horizontal lines in the concept violins that represent the optimal position of the prediction scores
598 based on ground truth. The selection of prediction targets includes *TP53*, *APC*, and *KRAS* from
599 genetic Cluster 1 (A.), and *BMPR2*, *ZNRF3*, Hypermethylation (HM), *RNF43*, and *BRAF* from genetic
600 Cluster 2 (B.) (Fig. 2). The data encompasses both external cohorts CRA and WHI (Fig. 1C). Each
601 dot represents the mean value of individual patient prediction scores calculated from 7 folds, with
602 the horizontal line on each side of the violin indicating the median of all individual mean patient
603 scores. A horizontal line at 0.50 denotes the line of model uncertainty. The sample count for each
604 subgroup is indicated below the violins. Statistical significance is denoted in the figures as follows:
* for $p < 0.05$, ** for $p < 0.01$, *** for $p < 0.001$, with more details provided in Fig. 1D. After testing

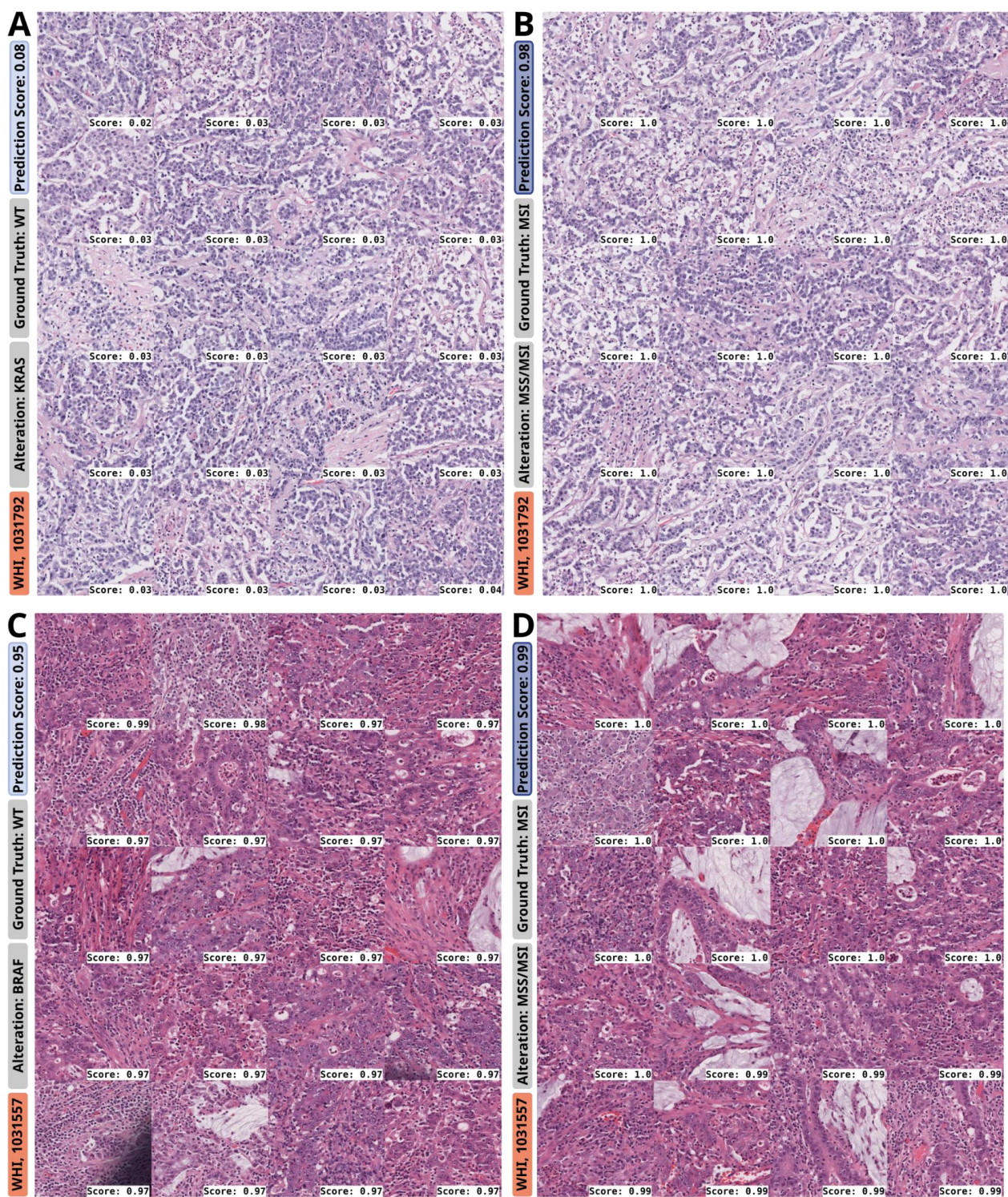
605 for normal distribution (Tab. S20), the Mann-Whitney U test was used for within-group
606 comparisons, and the Wilcoxon test was used for between-group comparisons. Abbreviations: HM:
607 Hypermethylation; MSI: Microsatellite instability; MSS: Microsatellite stability; MUT: Mutated; WT:
608 Wild type.



609
610
611
612
613
614
615
616
617
618
619
620
621
622

Fig. 5: Heatmaps of representative samples for prediction of MSI, KRAS, BRAF and Hypermethylation (HM) from the external GECCO validation dataset. The heatmaps are derived from the model with the median AUROC for MSI detection and the majority of prediction targets evaluated by sevenfold cross-validation. The cohort, Sample-ID, ground truth and prediction scores for MSI, along with the individual mutational status of the target, a brief pathological evaluation and magnified views of specific areas are provided for in-depth analysis. The heatmaps indicate relevant areas for the various predictions. The red areas are of high importance and indicate a mutant type (MUT), while the blue areas are of low importance and indicate a wild type (WT). The color intensity showcases the model's attention to that distinct area. **A**. The tumor exhibits both gland-forming and more solid components and extremely high numbers of tumor-infiltrating lymphocytes (TILs) with dense lymphoid aggregates. The pathological examination confirms the plausibility of a high MSI score indicating MSI which is also the ground truth. A low KRAS score indicates KRASWT but the ground truth is KRAS MUT. The heatmap highlights similar tumor areas

623 but with diverging scores: where MSI map is red indicating high score, *KRAS* map is blue indicating
624 low score. **B.** The presence of mucinous differentiation in a MSI, *BRAF* WT case results in high
625 MSI and *BRAF* scores. The MSI score is pathologically plausible whereas the *BRAF* score
626 indicates a contrary prediction tendency than the ground truth holds. For both predictions, the
627 model focuses on similar tumor areas with similar scores indicating MSI/*BRAF* MUT. **C.** Partly
628 mucinous morphology indicates the possibility of MSI, with a high score predicting MSI. HM is also
629 predicted MUT with a high score, even though HM is WT for this sample. Both heatmaps primarily
630 label the tumor and the same region with comparable significance. **D.** Villous adenoma with high
631 grade dysplasia is a common precursor lesion associated with high frequency of *KRAS* mutations
632 ^{41,42} The heatmaps highlight similar large scale tumor areas but with converging scores: where the
633 MSI map is red indicating a high score, the *KRAS* map is blue indicating low score. **E.** The tumor
634 area appears to be mainly MSS, and the heatmap predicts a low score, indicating this. Although it
635 is being mutated in the ground truth, it is still predicted as non-hypermuted. This is a rare MSS
636 case with HM ⁵⁴ Both heatmaps predominantly mark the tumor area and the same region with
637 comparable relevance. Abbreviations: HM: Hypermuation; MSI: Microsatellite instability; MSS:
638 Microsatellite stability; MUT: Mutated; WT: Wild type; w/: with



639

640

641

642

643

644

645

646

647

648

649

650

651

652

Fig. 6: Top tiles for prediction of genetic alterations (left column) and MSI (right column) for two selected slides from the GECCO test set. A.-B. WHI, 1031792: Medullary carcinoma with sheets of tumor cells, low stroma content, high number of tumor-infiltrating lymphocytes in a *KRAS* WT and MSI case, leading to high MSI prediction scores (**B.**) and low *KRAS* MUT prediction scores (**A.**). *KRAS* mutations show a lower frequency in MSI CRCs. Medullary carcinoma is a key morphological feature of MSI CRCs. **C.-D.** WHI 1031557: Top tiles for *BRAF* MUT as well as MSI prediction, both predictions with high prediction scores, both displaying a mixed morphology with partly medullary, partly mucinous, partly gland-forming histology, and high number of tumor-infiltrating/associated lymphocytes. Medullary growth pattern with lymphocytic infiltration and mucinous differentiation are typical features of a MSI-like morphology. Accordingly, the case was correctly predicted as MSI with a really high prediction score (**D.**). As *BRAF* MUT and MSI often co-occur and share morphologic overlap, the case was misclassified with regards to *BRAF*-status, resulting in high prediction scores for *BRAF* MUT (**C.**), even though the ground truth was *BRAF*

653 WT. Abbreviations: CRC: Colorectal cancer; MSI: Microsatellite instability; MSS: Microsatellite
654 stability; MUT: Mutated; WT: Wild type.

655 References

- 656 1. Tsimberidou AM, Fountzilias E, Nikanjam M, Kurzrock R. Review of precision cancer
657 medicine: Evolution of the treatment paradigm. *Cancer Treat Rev.* 2020 Jun;86:102019.
- 658 2. Xiao W, Ren L, Chen Z, Fang LT, Zhao Y, Lack J, et al. Toward best practice in cancer
659 mutation detection with whole-genome and whole-exome sequencing. *Nat Biotechnol.* 2021
660 Sep;39(9):1141–50.
- 661 3. Phillips KA, Douglas MP, Wordsworth S, Buchanan J, Marshall DA. Availability and funding
662 of clinical genomic sequencing globally. *BMJ Glob Health [Internet].* 2021 Feb;6(2). Available
663 from: <http://dx.doi.org/10.1136/bmjgh-2020-004415>
- 664 4. Kather JN, Pearson AT, Halama N, Jäger D, Krause J, Loosen SH, et al. Deep learning can
665 predict microsatellite instability directly from histology in gastrointestinal cancer. *Nat Med.*
666 2019 Jul;25(7):1054–6.
- 667 5. Echle A, Grabsch HI, Quirke P, van den Brandt PA, West NP, Hutchins GGA, et al. Clinical-
668 Grade Detection of Microsatellite Instability in Colorectal Tumors by Deep Learning.
669 *Gastroenterology.* 2020 Oct;159(4):1406–16.e11.
- 670 6. Lee SH, Song IH, Jang HJ. Feasibility of deep learning-based fully automated classification
671 of microsatellite instability in tissue slides of colorectal cancer. *Int J Cancer.* 2021 Aug
672 1;149(3):728–40.
- 673 7. Saldanha OL, Quirke P, West NP, James JA, Loughrey MB, Grabsch HI, et al. Swarm
674 learning for decentralized artificial intelligence in cancer histopathology. *Nat Med.* 2022
675 Jun;28(6):1232–9.
- 676 8. Bilal M, Raza SEA, Azam A, Graham S, Ilyas M, Cree IA, et al. Development and validation
677 of a weakly supervised deep learning framework to predict the status of molecular pathways
678 and key mutations in colorectal cancer from routine histology images: a retrospective study.
679 *Lancet Digit Health.* 2021 Dec;3(12):e763–72.
- 680 9. Saillard C, Dubois R, Tchita O, Loiseau N, Garcia T, Adriansen A, et al. Validation of
681 MSIntuit as an AI-based pre-screening tool for MSI detection from colorectal cancer
682 histology slides. *Nat Commun.* 2023 Nov 6;14(1):6695.
- 683 10. Zamanitajeddin N, Jahanifar M, Bilal M, Eastwood M, Rajpoot N. Social network analysis of
684 cell networks improves deep learning for prediction of molecular pathways and key
685 mutations in colorectal cancer. *Med Image Anal.* 2024 Apr;93(103071):103071.
- 686 11. Kather JN, Heij LR, Grabsch HI, Loeffler C, Echle A, Muti HS, et al. Pan-cancer image-
687 based detection of clinically actionable genetic alterations. *Nat Cancer.* 2020 Aug;1(8):789–
688 99.
- 689 12. Niehues JM, Quirke P, West NP, Grabsch HI, van Treeck M, Schirris Y, et al. Generalizable
690 biomarker prediction from cancer pathology slides with self-supervised deep learning: A
691 retrospective multi-centric study. *Cell Rep Med.* 2023 Apr 18;4(4):100980.
- 692 13. Guo B, Li X, Yang M, Jonnagaddala J, Zhang H, Xu XS. Predicting microsatellite instability
693 and key biomarkers in colorectal cancer from H&E-stained images: achieving state-of-the-art
694 predictive performance with fewer data using Swin Transformer. *Hip Int.* 2023 May;9(3):223–
695 35.
- 696 14. Bilal M, Tsang YW, Ali M, Graham S, Hero E, Wahab N, et al. Development and validation of
697 artificial intelligence-based prescreening of large-bowel biopsies taken in the UK and
698 Portugal: a retrospective cohort study. *Lancet Digit Health.* 2023 Nov 1;5(11):e786–97.

- 699 15. Kacew AJ, Strohbehn GW, Saulsberry L, Laiteerapong N, Cipriani NA, Kather JN, et al.
700 Artificial Intelligence Can Cut Costs While Maintaining Accuracy in Colorectal Cancer
701 Genotyping [Internet]. Vol. 11, *Frontiers in Oncology*. 2021. Available from:
702 <http://dx.doi.org/10.3389/fonc.2021.630953>
- 703 16. Saldanha OL, Loeffler CML, Niehues JM, van Treeck M, Seraphin TP, Hewitt KJ, et al. Self-
704 supervised attention-based deep learning for pan-cancer mutation prediction from
705 histopathology. *NPJ Precis Oncol*. 2023 Mar 28;7(1):35.
- 706 17. Common Variant GWAS, Genetics & Epidemiology of Colorectal Cancer Consortium
707 (GECCO) [Internet]. [cited 2023 Sep 28]. Available from:
708 https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001078.v1.p1
- 709 18. Riboli E, Kaaks R. The EPIC Project: rationale and study design. *European Prospective*
710 *Investigation into Cancer and Nutrition*. *Int J Epidemiol*. 1997;26 Suppl 1:S6–14.
- 711 19. Gsur A, Baierl A, Brezina S. Colorectal Cancer Study of Austria (CORSA): A Population-
712 Based Multicenter Study. *Biology* [Internet]. 2021 Jul 28;10(8). Available from:
713 <http://dx.doi.org/10.3390/biology10080722>
- 714 20. Folsom AR, Kushi LH, Anderson KE, Mink PJ, Olson JE, Hong CP, et al. Associations of
715 general and abdominal obesity with multiple health outcomes in older women: the Iowa
716 Women’s Health Study. *Arch Intern Med*. 2000 Jul 24;160(14):2117–28.
- 717 21. Hays J, Hunt JR, Hubbell FA, Anderson GL, Limacher M, Allen C, et al. The Women’s
718 Health Initiative recruitment methods and results. *Ann Epidemiol*. 2003 Oct;13(9
719 Suppl):S18–77.
- 720 22. Cancer Genome Atlas Network. Comprehensive molecular characterization of human colon
721 and rectal cancer. *Nature*. 2012 Jul 18;487(7407):330–7.
- 722 23. Edwards NJ, Oberti M, Thangudu RR, Cai S, McGarvey PB, Jacob S, et al. The CPTAC
723 Data Portal: A resource for cancer proteomics research. *J Proteome Res*. 2015 Jun
724 5;14(6):2707–13.
- 725 24. Liu Y, Sethi NS, Hinoue T, Schneider BG, Cherniack AD, Sanchez-Vega F, et al.
726 Comparative Molecular Analysis of Gastrointestinal Adenocarcinomas. *Cancer Cell*. 2018
727 Apr 9;33(4):721–35.e8.
- 728 25. Wagner SJ, Reisenbüchler D, West NP, Niehues JM, Zhu J, Foersch S, et al. Transformer-
729 based biomarker prediction from colorectal cancer histology: A large-scale multicentric
730 study. *Cancer Cell*. 2023 Sep 11;41(9):1650–61.e4.
- 731 26. Zhuang Y, Wang H, Jiang D, Li Y, Feng L, Tian C, et al. Multi gene mutation signatures in
732 colorectal cancer patients: predict for the diagnosis, pathological classification, staging and
733 prognosis. *BMC Cancer*. 2021 Apr 9;21(1):380.
- 734 27. Bao X, Zhang H, Wu W, Cheng S, Dai X, Zhu X, et al. Analysis of the molecular nature
735 associated with microsatellite status in colon cancer identifies clinical implications for
736 immunotherapy. *J Immunother Cancer* [Internet]. 2020 Oct;8(2). Available from:
737 <http://dx.doi.org/10.1136/jitc-2020-001437>
- 738 28. Murtagh F, Legendre P. Ward’s Hierarchical Clustering Method: Clustering Criterion and
739 Agglomerative Algorithm [Internet]. *arXiv [stat.ML]*. 2011. Available from:
740 <http://arxiv.org/abs/1111.6285>
- 741 29. Han J, Kamber M, Pei J. *Data Mining: Concepts and Techniques*. 3rd ed. Oxford, England:
742 Morgan Kaufmann; 2011. 744 p. (The Morgan Kaufmann Series in Data Management
743 Systems).

- 744 30. Canny J. A computational approach to edge detection. *IEEE Trans Pattern Anal Mach Intell.*
745 1986 Jun;8(6):679–98.
- 746 31. Wang X, Yang S, Zhang J, Wang M, Zhang J, Yang W, et al. Transformer-based
747 unsupervised contrastive learning for histopathological image classification. *Med Image*
748 *Anal.* 2022 Oct;81:102559.
- 749 32. Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, et al. An Image
750 is Worth 16x16 Words: Transformers for Image Recognition at Scale [Internet]. *arXiv.* 2020.
751 Available from: <http://arxiv.org/abs/2010.11929>
- 752 33. Selvaraju RR, Cogswell M, Das A, Vedantam R, Parikh D, Batra D. Grad-CAM: Visual
753 Explanations from Deep Networks via Gradient-based Localization [Internet]. *arXiv [cs.CV].*
754 2016. Available from: <http://arxiv.org/abs/1610.02391>
- 755 34. McDermott MBA, Hansen LH, Zhang H, Angelotti G, Gallifant J. A Closer Look at AUROC
756 and AUPRC under Class Imbalance [Internet]. *arXiv [cs.LG].* 2024. Available from:
757 <http://arxiv.org/abs/2401.06091>
- 758 35. Yamashita R, Long J, Longacre T, Peng L, Berry G, Martin B, et al. Deep learning model for
759 the prediction of microsatellite instability in colorectal cancer: a diagnostic study. *Lancet*
760 *Oncol.* 2021 Jan;22(1):132–41.
- 761 36. Schrammen PL, Ghaffari Laleh N, Echle A, Truhn D, Schulz V, Brinker TJ, et al. Weakly
762 supervised annotation-free cancer detection and prediction of genotype in routine
763 histopathology. *J Pathol.* 2022 Jan;256(1):50–60.
- 764 37. Tuppurainen K, Mäkinen JM, Junttila O, Liakka A, Kyllönen AP, Tuominen H, et al.
765 Morphology and microsatellite instability in sporadic serrated and non-serrated colorectal
766 cancer. *J Pathol.* 2005 Nov;207(3):285–94.
- 767 38. Greenson JK, Huang SC, Herron C, Moreno V, Bonner JD, Tomsho LP, et al. Pathologic
768 predictors of microsatellite instability in colorectal cancer. *Am J Surg Pathol.* 2009
769 Jan;33(1):126–33.
- 770 39. Malik A, Bhatia JK, Sahai K, Boruah D, Sharma A. Evaluating morphological features for
771 predicting microsatellite instability status in colorectal cancer. *Armed Forces Med J India.*
772 2022 Sep;78(Suppl 1):S96–104.
- 773 40. Shia J, Schultz N, Kuk D, Vakiani E, Middha S, Segal NH, et al. Morphological
774 characterization of colorectal cancers in The Cancer Genome Atlas reveals distinct
775 morphology-molecular associations: clinical and biological implications. *Mod Pathol.* 2017
776 Apr;30(4):599–609.
- 777 41. Zauber P, Marotta S, Sabbath-Solitare M. KRAS gene mutations are more common in
778 colorectal villous adenomas and in situ carcinomas than in carcinomas. *Int J Mol Epidemiol*
779 *Genet.* 2013 Mar 18;4(1):1–10.
- 780 42. Rosty C, Young JP, Walsh MD, Clendenning M, Walters RJ, Pearson S, et al. Colorectal
781 carcinomas with KRAS mutation are associated with distinctive morphological and molecular
782 features. *Mod Pathol.* 2013 Jun;26(6):825–34.
- 783 43. Hatthakarnkul P, Quinn JA, Matly AAM, Ammar A, van Wyk HC, McMillan DC, et al.
784 Systematic review of tumour budding and association with common mutations in patients
785 with colorectal cancer. *Crit Rev Oncol Hematol.* 2021 Nov;167:103490.
- 786 44. Siegel RL, Wagle NS, Cercek A, Smith RA, Jemal A. Colorectal cancer statistics, 2023. *CA*
787 *Cancer J Clin.* 2023 May 1;73(3):233–54.

- 788 45. Nojadeh JN, Behrouz Sharif S, Sakhinia E. Microsatellite instability in colorectal cancer.
789 EXCLI J. 2018 Jan 22;17:159–68.
- 790 46. Giannakis M, Hodis E, Jasmine Mu X, Yamauchi M, Rosenbluh J, Cibulskis K, et al. RNF43
791 is frequently mutated in colorectal and endometrial cancers. Nat Genet. 2014
792 Dec;46(12):1264–6.
- 793 47. Pai RK, Jayachandran P, Koong AC, Chang DT, Kwok S, Ma L, et al. BRAF-mutated,
794 microsatellite-stable adenocarcinoma of the proximal colon: an aggressive adenocarcinoma
795 with poor survival, mucinous differentiation, and adverse morphologic features. Am J Surg
796 Pathol. 2012 May;36(5):744–52.
- 797 48. Trinh A, Ladrach C, Dawson HE, Ten Hoorn S, Kuppen PJK, Reimers MS, et al. Tumour
798 budding is associated with the mesenchymal colon cancer subtype and RAS/RAF mutations:
799 a study of 1320 colorectal cancers with Consensus Molecular Subgroup (CMS) data. Br J
800 Cancer. 2018 Nov;119(10):1244–51.
- 801 49. Nguyen TT, Huynh TT, Le Nguyen P, Liew AWC, Yin H, Nguyen QVH. A Survey of Machine
802 Unlearning [Internet]. arXiv [cs.LG]. 2022. Available from: <http://arxiv.org/abs/2209.02299>
- 803 50. Kleppe A. Area under the curve may hide poor generalisation to external datasets. ESMO
804 Open. 2022 Apr;7(2):100429.
- 805 51. Bledsoe JR, Kamionek M, Mino-Kenudson M. BRAF V600E immunohistochemistry is
806 reliable in primary and metastatic colorectal carcinoma regardless of treatment status and
807 shows high intratumoral homogeneity. Am J Surg Pathol. 2014 Oct;38(10):1418–28.
- 808 52. Dawood M, Eastwood M, Jahanifar M, Young L, Ben-Hur A, Branson K, et al. Cross-linking
809 breast tumor transcriptomic states and tissue histology. Cell Rep Med. 2023 Dec
810 19;4(12):101313.
- 811 53. Clark K, Vendt B, Smith K, Freymann J, Kirby J, Koppel P, et al. The Cancer Imaging
812 Archive (TCIA): maintaining and operating a public information repository. J Digit Imaging.
813 2013 Dec 25;26(6):1045–57.
- 814 54. Zaidi SH, Harrison TA, Phipps AI, Steinfeld R, Trinh QM, Qu C, et al. Landscape of
815 somatic single nucleotide variants and indels in colorectal cancer and impact on survival. Nat
816 Commun. 2020 Jul 20;11(1):3644.