

Deep Learning for Biomarker Discovery in Cancer Genomes

Michaela Unger (1), Chiara M. L. Loeffler (1, 2, 3), Laura Žigutyte (1), Srividhya Sainath (1), Tim Lenz (1), Julien Vibert (4), Andreas Mock (5, 6, 7), Stefan Fröhling (6, 7, 8, 9), Trevor A. Graham (10), Zunamys I. Carrero (1), Jakob Nikolas Kather (1, 2, 11, 12)

+ Correspondence to jakob-nikolas.kather@alumni.dkfz.de

1. Else Kroener Fresenius Center for Digital Health, University of Technology Dresden, Dresden, Germany
2. Medical Department 1, University Hospital and Faculty of Medicine Carl Gustav Carus, University of Technology Dresden, Dresden, Germany
3. National Center for Tumor Diseases Dresden (NCT/UCC), a partnership between DKFZ, Faculty of Medicine and University Hospital Carl Gustav Carus, TUD Dresden University of Technology, and Helmholtz-Zentrum Dresden - Rossendorf (HZDR), Dresden, Germany
4. Drug Development Department (DITEP), Gustave Roussy, Villejuif, France
5. Institute of Pathology, Ludwig-Maximilians-University München, Munich, Germany
6. Division of Translational Medical Oncology, German Cancer Research Center (DKFZ), Heidelberg, Germany
7. National Center for Tumor Diseases (NCT), NCT Heidelberg, a partnership between DKFZ and Heidelberg University Hospital, Heidelberg, Germany
8. German Cancer Consortium (DKTK), Core Center Heidelberg, Heidelberg, Germany
9. Division of Translational Precision Medicine, Institute of Human Genetics, Heidelberg University, Heidelberg, Germany
10. Centre for Evolution and Cancer, Institute of Cancer Research, London, UK
11. National Center for Tumor Diseases Dresden (NCT/UCC), a partnership between DKFZ, Faculty of Medicine and University Hospital Carl Gustav Carus, TUD Dresden University of Technology, and Helmholtz-Zentrum Dresden - Rossendorf (HZDR), Dresden, Germany
12. Medical Oncology, National Center for Tumor Diseases (NCT), University Hospital Heidelberg, Heidelberg, Germany

1 **Abstract**

2 **Background:** Genomic data is essential for clinical decision-making in precision oncology.
3 Bioinformatic algorithms are widely used to analyze next-generation sequencing (NGS) data,
4 but they face two major challenges. First, these pipelines are highly complex, involving multiple
5 steps and the integration of various tools. Second, they generate features that are human-
6 interpretable but often result in information loss by focusing only on predefined genetic
7 properties. This limitation restricts the full potential of NGS data in biomarker extraction and
8 slows the discovery of new biomarkers in precision oncology.

9 **Methods:** We propose an end-to-end deep learning (DL) approach for analyzing NGS data.
10 Specifically, we developed a multiple instance learning DL framework that integrates somatic
11 mutation sequences to predict two compound biomarkers: microsatellite instability (MSI) and
12 homologous recombination deficiency (HRD). To achieve this, we utilized data from 3,184
13 cancer patients obtained from two public databases: The Cancer Genome Atlas (TCGA) and
14 the Clinical Proteome Tumor Analysis Consortium (CPTAC).

15 **Results:** Our proposed deep learning method demonstrated high accuracy in identifying
16 clinically relevant biomarkers. For predicting MSI status, the model achieved an accuracy of
17 0.98, a sensitivity of 0.95, and a specificity of 1.00 on an external validation cohort. For
18 predicting HRD status, the model achieved an accuracy of 0.80, a sensitivity of 0.75, and a
19 specificity of 0.86. Furthermore, the deep learning approach significantly outperformed
20 traditional machine learning methods in both tasks (MSI accuracy, $p\text{-value} = 5.11 \times 10^{-18}$; HRD
21 accuracy, $p\text{-value} = 1.07 \times 10^{-10}$). Using explainability techniques, we demonstrated that the
22 model's predictions are based on biologically meaningful features, aligning with key DNA
23 damage repair mutation signatures.

24 **Conclusion:** We demonstrate that deep learning can identify patterns in unfiltered somatic
25 mutations without the need for manual feature extraction. This approach enhances the

- 1 detection of actionable targets and paves the way for developing NGS-based biomarkers
- 2 using minimally processed data.
- 3
- 4 **Keywords:** AI, Deep Learning, Biomarker, Microsatellite Instability, Homologous
- 5 Recombination Deficiency, Cancer Genomics

1 Background

2 Bulk and targeted sequencing of cancer are being progressively integrated into clinical routine
3 workflows. Today, knowledge about genetic variants or expression profiles are key pillars of
4 personalized oncology [1–5]. Several efforts and programs have been established to analyze
5 patients' cancer genomes, benefiting oncologists' decision-making [6,7]. In molecular tumor
6 boards, genomic biomarkers play a substantial role when determining diagnosis of a patient,
7 hereditary predispositions, treatment scheme and can monitor therapy response, or explain
8 resistances. Two important and highly clinically relevant biomarkers that are characterized by
9 sets of genomic alterations are microsatellite instability (MSI) [8] and homologous
10 recombination deficiency (HRD) [9]. In MSI, the mismatch repair (MMR) system of a cell is
11 impaired, leading to small-scale insertions and deletions (indels). MSI-high tumors are an ideal
12 target for immunotherapy [8,10]. In contrast, HRD introduces a different mutational scar in the
13 cancer genome [11,12]. This makes HRD tumors eligible for treatment with poly ADP ribose
14 polymerase (PARP) inhibitors in certain situations [13,14]. Both MSI and HRD leave
15 distinguishable and complex patterns in the cancer genome.

16 Since the start of the Human Genome Project [15], bioinformaticians have been involved in
17 designing computational algorithms to analyze genomes. Traditionally, extensive multi-step
18 pipelines process sequencing data for specific applications. For instance, gene expression is
19 determined by mapping transcriptome reads to the human genome and assigning counts to
20 the respective genes [16–20]. Similarly, small- and large-scale mutations are detected by
21 variant callers that compare the sequence of a cancer genome to a reference [21–24].
22 Information on gene expression or variants in cancer-related genes is then interpreted in the
23 clinical context. However, this approach is fundamentally limited by human expertise. By
24 constraining the analysis of next generation sequencing (NGS) data to pre-defined genomic
25 features, potentially relevant information is neglected. This limitation is particularly evident for
26 compound biomarkers such as MSI and HRD which are reflected by a large spectrum of

1 changes in the genome and can be defined through a range of molecular assays [25–27].
2 Therefore, having an objective, unbiased end-to-end tool that directly relates a disease-
3 specific genomic pattern to patient subgroups and outcomes could complement existing
4 bioinformatics methods in precision oncology.

5 One instrument that could serve such a purpose is deep learning (DL). With increasing
6 computational resources and decreasing sequencing costs [28], DL has been increasingly
7 applied in cancer bioinformatics [29,30]. For instance, DL has been applied to variant calling,
8 such as prognostication or drug response prediction [31–34]. Nevertheless, most of these
9 models still rely on heavily handcrafted features or make strong prior assumptions about
10 feature interactions. Here, we build upon previous attempts to use attention-based multiple
11 instance learning (attMIL) for genomic data analysis with minimal human intervention [35,36].
12 We extend and apply this approach on unfiltered somatic mutations to predict MSI status in
13 colon, rectum, gastric, and uterine cancers and HRD status in breast, ovarian, prostate, and
14 pancreatic cancers. To our knowledge, this is the first study to apply the attMIL DL method for
15 somatic mutations in multiple clinical cohorts with external validation. Finally, we propose a
16 set of explainability methods to gain insights into the potential biomedical features the model
17 uses for its predictions. In summary, our approach provides an interpretable, genomics-based
18 tool applicable for oncological tasks that could support clinical decision-making.

1 **Methods**

2 **Data acquisition**

3 Small-scale somatic mutations in form of single base substitutions (SBSs) and insertions and
4 deletions (indels) of 3,185 whole exome sequenced (WES) cancer patients from The Cancer
5 Genome Atlas Program (TCGA) (n=3,080 patients) and the Clinical Proteomic Tumor Analysis
6 Consortium (CPTAC) (n=105 patients) were utilized. (**Fig. 1a**) TCGA controlled-access and
7 CPTAC public mutation data were obtained in Mutation Annotation Format (.maf) files through
8 the Genomic Data Commons (GDC) (**Fig. 1b**). In total we acquired 8.3 million mutations for
9 TCGA and 70k mutations for CPTAC COAD.

10 To reconstruct the mutation sequence context, mutation entries were mapped back to their
11 genomic locations using the 'Chromosome', 'Start_Position', and 'End_Position' columns to
12 generate their sequence context. Human reference genome builds GRCh37 and GRCh38
13 were used for TCGA and CPTAC data respectively, with chromosome sequences obtained
14 from the Ensembl genome database (**Fig. 1c**).

15 For prediction of MSI status, four common cancer types were selected: colon adenocarcinoma
16 (COAD) (n=338), rectal adenocarcinoma (READ) (n=110), stomach adenocarcinoma (STAD)
17 (n=432), and uterine corpus endometrial carcinoma (UCEC) (n=450) (**Fig. 1a & S1**), resulting
18 in a dataset of 1,330 patients from TCGA. MSI status was obtained through cBioPortal
19 (<https://www.cbioportal.org/>) and previous studies conducted on TCGA data [37]. The MSI
20 status of all patients was determined using consensus calls from polymerase chain reaction
21 (PCR) -based assays, MANTIS [38], and MSIsensor [39], with previously described thresholds
22 set at binarized values of 0.4 and 3.5, respectively (**Additional File 1**). For validation, 104
23 COAD patients from CPTAC were used (**Supplementary Material 1**).

1 Furthermore, to simulate how our models perform on mutations from targeted sequencing
2 panels, we filtered the WES-derived variant data using the gene lists from FoundationOne
3 CDx (324 genes) [40,41] and TruSight Oncology 500 (523 genes) [42,43]. Variants were
4 filtered by matching their HUGO nomenclature entries to the gene lists. Since the gene panels
5 both include only coding exonic regions, we kept mutations with the following variant
6 classification: Missense_Mutation, Silent, Nonsense_Mutation, Frame_Shift_Del,
7 Frame_Shift_Ins, In_Frame_Del, In_Frame_Ins. The full CPTAC dataset contained 70,861
8 variants from 16,478 genes. After filtering by Trusight Oncology 2,718 variants from 432 genes
9 remained after the FoundationOne Dx panel 1,758 variants from 280 genes remained
10 (**Supplementary Material 2**). To match the mutation load of simulated CPTAC panel
11 sequencing with WES, we imputed the mutation distribution by repeating the variants ten
12 times.

13 For prediction of HRD status, three, cancer types where HRD is prevalent were selected [40]:
14 breast (BRCA) (n=967), ovarian (OV) (n=182), pancreatic (PAAD) (n=139) and prostate
15 cancer (PRAD) (n=462) from TCGA (**Fig. 1a & S2**). The genomic HRD scar score [41] consists
16 of three numerical properties: loss of heterozygosity (LOH) [42], telomeric allelic imbalance
17 (TAI) [43], and large-scale transitions (LST) [44] from which a sum HRD status was
18 determined. Using tissue-specific cutoffs [45], we generated binarized prediction targets with
19 cutoffs set at 45 for BRCA, 54 for OV, 37 for PAAD, and 21 for PRAD [45] (**Additional File 2**).
20 Due to the limited number of HR-deficient samples in CPTAC for these cancer types, external
21 validation was not feasible for HRD predictions.

22 **Data Preparation**

23 In our experiments we compared our DL model to a state-of-the-art (SOTA) machine learning
24 (ML) tool to have a baseline reference regarding performance. Data preparation for both
25 models differed in fundamental parts.

1 The input data for the DL model consisted of the sequence context surrounding each mutation.
2 Four 20-nucleotide sequences were extracted: the mutation itself, the upstream sequence (5'),
3 the downstream sequence (3'), and the reference sequence at the mutation site. For mutations
4 shorter than 20 nucleotides - such as single-nucleotide substitutions - the remaining positions
5 were padded with gap characters ('-'). Similarly, shorter reference sequences were also
6 padded to ensure consistent length. For all experiments, a window size of 20 nucleotides was
7 defined to capture the specific indel signatures associated with MSI and HRD mechanisms
8 (**Fig. 1d**).

9 To account for the double stranded nature of DNA, the reverse complement of each mutation
10 sequence was provided as well and processed as blocks as described above. Additionally, a
11 binary indicator was included to specify the DNA strand on which the mutation occurred. All
12 sequences were numerically encoded by mapping nucleotide bases (A, C, G, T) and gap
13 characters ('-').

14 For the ML model, mutations were represented in their SBS and indel forms, comparable to
15 Catalogue Of Somatic Mutations In Cancer (COSMIC) v3.4 [46] features. Therefore, we used
16 178 input features for the ML model, 96 for SBS mutations and 83 for indels. Feature values
17 were normalized by mutation count of the respective patient (**Supplementary Methods**).

18 **Data Splitting**

19 To enhance generalizability and mitigate potential site-specific biases, a site-specific training
20 split based on the Tissue Source Site (TSS) codes within the TCGA patient identifiers was
21 conducted [47]. TSS codes were mapped to their respective institutions (**Supplementary**
22 **Material 3, Additional File 3**). The test set for each model comprised patients from specific
23 institutions excluded from the training set to ensure comparable results. The training data was
24 divided for a 5-fold cross-validation in an approximate 70:15:15 split (**Fig. S1 & S2**).
25 Additionally, to this source-specific split, for MSI experiments we validated our findings on

1 CPTAC. For the training data positive samples for MSI ranged from 20 to 22 % between folds,
2 and from 15 to 17 % for HRD (**Table S1 & S2**).

3 **Baseline Gradient Boosting**

4 To establish a SOTA ML classifier on somatic mutations as baseline for performance, Extreme
5 Gradient Boosting (XGB) models were employed using the scikit-learn v1.2.0 and xgboost
6 v2.1.1 libraries in Python. For MSI and HRD classification tasks, default parameters of the
7 `xgb.XGBClassifier()` and `fit()` function were used.

8 **Deep Learning Model Setup**

9 The DL model employed in this study was directly adopted from the implementation provided
10 by Anaya et al. [35], in TensorFlow v2.12.0. The model comprises two main components: a
11 trainable mutation encoder and an attention multiple instance learning (attMIL) module
12 (**Supplementary Methods**) (**Fig. 1d & 1e**).

13 The mutation encoder consists of a 2D convolutional layer followed by a dense layer,
14 transforming each mutation's sequence context into a feature vector. Convolutions are applied
15 separately to the forward and reverse complement sequences, and the resulting features are
16 concatenated and passed through the dense layer. The output dimensionality is set to 128
17 (**Fig. 1d**).

18 In the attMIL module [48], mutation-level feature vectors are aggregated into a patient-level
19 representation using an attention mechanism. Weighted sum pooling is employed, where
20 feature vectors are summed and weighted by their attention scores. The aggregated patient
21 vector is then passed through two dense layers, with the final layer using a sigmoid activation
22 function for binary classification tasks. For MSI and HRD, binary cross-entropy was used as a
23 loss function (**Fig. 1e**).

1 Models were trained with the Adam optimizer and a batch size of 128. For MSI models, training
2 was conducted for up to 300 epochs with a learning rate of 0.001. HRD models were trained
3 for up to 500 epochs with a learning rate of 0.01. The training and deployment of the models
4 were carried out on an NVIDIA RTX A6000 with 46 GB of RAM.

5 Model performance was evaluated using accuracy, F1 score, receiver operating characteristic
6 (ROC) area under the curve (AUC), and precision-recall (PR) AUC, specificity, sensitivity and
7 absolute values of true positives (TP), false positives (FP), true negatives (TN) and false
8 negatives (FN) for classification tasks. To estimate significance between model ROC AUCs
9 we utilized the DeLong's test, for accuracy McNemar, for count data (TP, FP, TN, FN) the
10 Wilcoxon Signed-Rank test and for all other metrics a paired t-test was applied. Furthermore,
11 for patient level significance tests, such as DeLong's test and McNemar we aggregated the p-
12 value for all five folds with Fisher's method. To generate a binary classification from continuous
13 prediction values we used the best cutoff of the ROC AUC on the training set.

14 **Explainability Techniques**

15 Several explainability techniques were applied to interpret the model's predictions. This
16 approach allowed us to interpret how the attMIL model encodes and evaluates mutations, as
17 well as how these encodings contribute to patient-level predictions. Dimensionality reduction
18 was performed using Uniform Manifold Approximation and Projection (UMAP) (python library
19 umap-learn v0.5.5), which facilitated the visualization of mutation- and patient-level feature
20 vectors by projecting them into lower-dimensional space to identify sample clusters. Mutation-
21 level features were extracted from a layer that encodes individual mutations into vector
22 representations, while patient-level features were obtained from the aggregation layer, which
23 integrates information across all mutations for a given patient. Mutation encodings were then
24 stratified based on attention value and grouped using k-means clustering (k=7) with the
25 sklearn.cluster library, allowing the analysis of the importance of similar mutations to the
26 model. The choice of k=7 balanced sufficient granularity to capture distinct mutation patterns

1 with avoiding excessive fragmentation. To further understand the relationship between these
2 clusters and biological processes, mutation catalogs from the clusters were generated and
3 compared to known COSMIC mutational signatures associated with defective mismatch repair
4 (dMMR) and HRD.

5 Finally, the relationship between the model's prediction scores and various tumor properties
6 was examined for potential trends. For MSI, these included associations with cancer type,
7 tumor mutational burden (TMB), indel mutation density, driver mutations in dMMR genes (e.g.,
8 *MLH1*, *MSH2*, *MSH3*, *MSH6*, *PMS2*) and *POLE*, as well as the methylation status of *MLH1*
9 promoter. For this, TCGA patients included in the model's validation and test sets were
10 selected, and the binarization threshold was determined based on the optimal cutoff for the
11 ROC AUC derived from the training set. For HRD, correlations were evaluated in relation to
12 cancer type, *BRCA1/2* status, continuous genomic scarHRD scores and their components
13 (i.e., LOH, TAI, and LST), TMB, and mutations in additional HR pathway genes (e.g., *ATM*,
14 *BRIP1*, *CHEK2*, *NBN*, *PALB2*, *RAD51C*).

15 For XGB models, Shapley additive explanations (SHAP) values were calculated with the
16 python library shap (v0.45.1) for each input feature. By default, 20 most important features are
17 ranked by SHAP value and displayed in increasing order.

1 Results

2 DL detects microsatellite instability from sequencing panels

3 MSI arises because of defects in the mismatch repair (MMR) pathway [49,50], and leads to
4 characteristic small-scale indels in repetitive sequences as well as distinctive patterns of SBSs
5 (**Fig. 2a**), which in principle are recognizable by DL-based models [11,51–53]. We trained an
6 attMIL model to detect MSI status in 1,330 patients from TCGA, including patients from COAD,
7 READ, STAD and UCEC. After training, we externally validated this model on the COAD
8 cohort from CPTAC (n=105 patients). The attMIL model achieved an accuracy of 0.98 ± 0.01 ,
9 a ROC AUC of 1.00 ± 0.00 , sensitivity of 0.95 ± 0.02 and a specificity of 1.00 ± 0.00 confirming its
10 ability to detect MSI-related patterns (**Table S3**). When analyzing the absolute number of
11 misclassified patients across 5-fold cross-validation (ensemble of 5 models), we observed no
12 false positive patients and only 1.20 ± 0.45 false negatives in the total population of 105
13 patients. These results were based on an optimal threshold determined in the training set. For
14 real-world application, the threshold could be adapted to yield a high positive predictive value.
15 For example, in the external test set, a threshold corresponding to a sensitivity of 0.99 would
16 yield a specificity of 0.99 ± 0.01 , which is much higher than for many clinically used tests.

17 Next, we compared the DL performance to that of a SOTA ML model to investigate whether
18 DL is indeed more powerful than standard ML for this application. As expected, XGB also
19 achieved excellent performance on the CPTAC datasets, with an accuracy of 0.89 ± 0.07 , a
20 ROC AUC of 0.99 ± 0.01 , sensitivity of 0.99 ± 0.02 and a specificity of 0.80 ± 0.16 (**Fig. 2b, Table**
21 **S3**). However, a significant difference in performance between DL and ML models was seen
22 when comparing F1 scores (DL 0.97 ± 0.01 ; ML 0.76 ± 0.15) with a p-value of 0.04 (**Table S3**).
23 Furthermore, when investigating misclassified patients, we found that the XGB model
24 produces a large amount of false positive patients (16.40 ± 12.93 of 105 patients). These data

1 show that MSI prediction is a relatively easy task which even classical ML solves well,
2 however, our new DL approach was still substantially, and clinically meaningfully, better.

3 In the real world, most patients do not undergo whole exome sequencing (WES), but just panel
4 sequencing of a few hundred genes. We evaluated the performance of DL and ML models on
5 pseudo panel sequencing data. To this end, the number of somatic mutations was filtered to
6 include only genes present in two targeted sequencing panels: FoundationOne Dx [54,55] and
7 TruSight Oncology [56,57]. Reducing the mutation set to the TruSight Oncology 523 gene
8 panel revealed a significant divergence in performance between DL and ML models, with the
9 attMIL model achieving an accuracy of 0.96 ± 0.04 outperforming XGB's 0.61 ± 0.05 (p-value
10 1.67×10^{-38}) and a F1 score of attMIL of 0.96 ± 0.04 compared to 0.50 ± 0.02 of XGB (p-value
11 7.56×10^{-5}) (**Fig. 2b**). As expected, the task became more challenging with fewer mutations
12 available, indicated by lower overall accuracy - however, the DL model maintained a
13 reasonably good performance. In addition, we found that especially the XGB model had
14 reduced performance regarding specificity, from 0.80 ± 0.16 to 0.54 ± 0.28 , when lowering the
15 number of variants for model input, due to high false positive number (24.60 ± 15.26 of 105
16 patients) (**Table S4**). This suggests that the classical ML model is considerably biased towards
17 positive predictions. When the genes were further narrowed to 324 from the FoundationOne
18 Dx panel, performance differences persisted despite increased variance (accuracy attMIL
19 0.83 ± 0.10 vs. XGB 0.59 ± 0.06 , p-value 7.94×10^{-21} ; F1 scores attMIL 0.78 ± 0.17 vs. XGB
20 0.47 ± 0.11 , p-value 0.06) (**Fig. 2b, Table S5**). These results indicate that the DL model
21 outperforms classical ML in clinically relevant genomic subsets present in commercial
22 sequencing panels.

23 **DL detects biomedically relevant MSI patterns**

24 Subsequently, we investigated if the feature representations obtained by the DL model result
25 in a clinically relevant clustering of patients, which would provide further proof that the DL
26 model learned clinically relevant patterns. We investigated a possible clustering of at both the

1 mutation and patient levels (**Fig. 2c**). While no substantial differences were apparent at the
2 mutation level, a clear separation emerged at the patient level. MSI and MSS cancers
3 segregated distinctively, indicating that the DL model was able to encode their respective
4 patients differently (**Fig. 2c**).

5 We then evaluated which mutations were most influential for the model's predictions by
6 extracting mutation encodings from the embedding layer. These encodings were clustered
7 using k-means (k=7) and ranked by the mean attention values of their mutations to identify
8 highly and minimally influential groups (**Fig. 2c**). When examining the mutation catalogs of the
9 clusters with the highest attention scores (clusters II+III), they were enriched in C and G indels
10 (**Fig. 2c**). Comparing this to the indel mutation signature ID7 [46] (**Fig. 2a**) associated with
11 MSI, which depicts deletions of mononucleotide stretches but also of dinucleotide repeats, this
12 suggests that the model is capturing relevant MSI-associated mutational patterns.

13 When analyzing a mutation cluster with average attention scores (cluster I), we observed
14 several similarities to the SBS6 and SBS15 mutation signatures associated with defective
15 MMR (dMMR) [46] (**Fig. S3a**). Similarities encompassed the features of the C>T mutation
16 class, with comparable peaks in the sequence context pairs of an up- and downstream GG,
17 AG and CG. This could further indicate that the model may have learned to group and assign
18 weights to MSI-specific mutations. Other SBS clusters mostly contain similar mutations with
19 C>A mutations in a TT context (**Fig. S3b**). Comparing this to the Shapley additive explanation
20 (SHAP) analysis of the handcrafted features of XGB, XGB also utilizes deletions of a single
21 C/G similar to the DL model (**Fig. S4**). However, other features that strongly influenced the
22 XGB model's predictions, such as indels at repetitive sites or certain C>G/C>A mutations, did
23 not resemble classical dMMR signatures, suggesting either the identification of novel patterns
24 or a potential bias towards irrelevant features.

1 **DL identifies MSI-associated patterns beyond mutation counts**

2 Finally, we investigated whether the model's predictions were driven by subtle patterns within
3 the mutations or if it was simply confounded by factors such as mutation counts. To this end,
4 we investigated if the model's predictions are associated with other relevant properties of the
5 patients' cancers related to MSI (Fig. 2d). These properties included the ground truth label,
6 features of genomic instability, driver mutations, and the methylation status of MMR-related
7 genes. The model's normalized prediction scores were highly indicative of MSI and MSS
8 status across all cancer types, effectively reflecting the underlying distinction between the two
9 groups. Next, as expected, a correlation trend between the prediction score and TMB and
10 indel density was observed. However, patients with high TMB/indel count in an MSS context
11 were not predicted as MSI by the DL model. This suggests that the model learns subtle,
12 clinically relevant patterns and does not simply count mutations as indicators for MSI. We also
13 observed that driver mutations in key MMR pathway proteins are not necessarily indicative of
14 a positive model prediction, nor were they always correlated with consensus labels from
15 MANTIS, MSIsensor or PCR (Fig. 2d). A close relationship was found between the model
16 scores and the MLH1 promoter methylation, as silencing of this gene has a more significant
17 impact on MMR function than mutations alone [58].

18 These results emphasize DL's potential as an alternative tool for analyzing NGS data,
19 specifically for panel sequencing data. Our results demonstrate high performance
20 generalizability, and explainability of the DL-based attMIL model.

21 **DL predicts HRD and highlights mutational patterns in alternative** 22 **repair pathways**

23 We next investigated HRD, which arises from defects in key homologous recombination (HR)
24 pathway genes, as well as from other associated genomic defects. When HR fails to repair
25 double-strand breaks (DSBs), cells rely on alternative pathways, such as microhomology-

1 mediated end joining (MMEJ), which in turn induce specific mutational patterns we detected
2 with our DL model. (**Fig. 3a**) [59,60]. We trained an attMIL model to detect HRD using data
3 from 1,750 patients across BRCA, OV, PAAD, and PRAD cohorts from TCGA, with internal
4 validation performed in a source-site-specific manner. After 5-fold cross-validation, the attMIL
5 model performed slightly better than XGB (attMIL accuracy: 0.80 ± 0.01 , XGB accuracy:
6 0.75 ± 0.02 (p-value 1.07×10^{-10}); attMIL ROC AUC: 0.88 ± 0.01 , XGB ROC AUC: 0.86 ± 0.01 (p-
7 value 0.72); DL PR AUC: 0.67 ± 0.01 , XGB PR AUC: 0.75 ± 0.02) (**Fig. 3b, Fig. S6a**). Regarding
8 specificity, both models performed well with (attMIL 0.86 ± 0.05 , XGB 0.98 ± 0.00 (p-value
9 0.049)), however sensitivity was comparably low with 0.74 ± 0.06 for attMIL and 0.52 ± 0.03 for
10 XGB (p-value 0.003), indicating that predicting HRD is a more challenging task than predicting
11 MSI (**Table S6**).

12 We then investigated the explainability of DL-based HRD predictions. The patient level UMAP
13 shows a gradient between HRD and HR- proficient (HRP) patients, but no clear separation as
14 in the MSI case (Fig. 3c). Furthermore, the mutation level features for HRD and HRP also
15 cluster very much together. In the k-mean clustering heatmap of mutation features, we again
16 identified groups of mutations with high attention scores. The second highest attention group
17 (cluster II) consists of indel mutations, primarily monomer stretch insertions and deletions of
18 C and G, as well as various other indels in repetitive regions. Notably, the highest attention
19 scores (cluster III) were assigned to microhomologies and deletions at 5-nucleotide repeats,
20 consistent with indel signatures ID6 [46] (Fig. 3a) and ID8 [46] (Fig. S4a). Both signatures are
21 associated with alternative repair pathways to HRD, namely MMEJ and non-homologous end
22 joining (NHEJ). A group of mutations with mean attention attribution (cluster I) displays mostly
23 SBSs of various classes related to SBS3 [46], reflecting the general high mutation load of HRD
24 cancers (Fig. S4a). These findings are in line with the highly ranked indel features of the XGB
25 model, which mostly consist of microhomologies or deletions of single C/G nucleotides (Fig.
26 S4b). Additionally, the XGB model identified several C>G mutations as important predictors,
27 which, to our knowledge, previously have not been brought into context with HRD.

1 **DL predicts HRD with biologically meaningful features**

2 Following the same approach as before with a MSI, we again compared the model's
3 predictions to target-specific tumor properties. (**Fig. 3d**). We observed that higher prediction
4 scores correlated with a greater likelihood of patients being HRD-positive, suggesting that the
5 model's prediction scores reflect a ranking of HRD probability among samples. Upon
6 examining the model's prediction across cancer types, we noted that the proportion of patients
7 predicted as HRD varied: approximately 1:4 BRCA patients and 3:5 OV patients were
8 predicted as HRD positive, whereas there were no positive cases in PAAD and only 1:20 in
9 PRAD, aligning with the distribution of the HRD status of patients in these cancer types (**Table**
10 **S2**). It remains unclear whether this is a result of the model learning to differentiate tissue
11 types or solely predicting HRD status. The model predictions were not correlated to LST, LOH,
12 TAI individually, but only to the HRD scar signature, which is composed of all three
13 components [45]

14 In summary, we show that deep learning models can effectively predict key biomarkers like
15 MSI and HRD across cancer types while using biologically meaningful features and
16 generalizing well to external cohorts, outperforming traditional machine learning approaches,
17 especially in panel sequencing data.

18

1 **Discussion**

2 In this study, we demonstrated that DL models can learn to predict clinically relevant
3 biomarkers solely from sequences of small-scale somatic mutation of cancer genomes,
4 without using any prior assumptions or expert knowledge. By comparing an attMIL DL
5 framework with a XGB ML model, we found that for both targets, MSI and HRD, DL
6 outperforms ML. The DL model was explainable and specific, biologically plausible, genomic
7 features were associated with known mutational patterns. Although this model still relies on
8 variants called by classical bioinformatic tools as input, it represents a first step towards end-
9 to-end DL workflows for cancer genomes.

10 Our findings indicate that DL models offer advantages over traditional ML approaches. For
11 MSI prediction, the DL model outperformed an ML model, especially when applied to limited
12 data during inference. While both models performed similarly on full datasets, the DL model
13 maintained excellent performance even with targeted sequencing panels, such as the TruSight
14 Oncology 500 gene panel. This suggests that DL models are more robust to data sparsity and
15 could streamline clinical workflows by reducing the need for additional MSI testing. In addition,
16 as sequencing panels are widely used in clinical routine, DL models performing well not only
17 on WES/WGS data, could help bring genomic advancements to a broader population.

18 In predicting HRD, the DL model was slightly better than that of the ML model with a
19 performance comparable to other tools [61–63]. HRD is primarily characterized by large-scale
20 genomic alterations rather than small indels alone. This suggests that small-scale somatic
21 mutations may not provide sufficient information for accurate HRD prediction. Incorporating
22 further information about the sequencing data, such as variant allele frequency, genomic
23 region or affected gene as well as data on structural variations and large-scale
24 rearrangements could enhance model performance.

1 DL models excel in adapting to diverse datasets without predefined features, making them
2 ideal for the heterogeneity of tumor genomes. Unlike classical ML models, they process
3 sequencing data end-to-end, uncovering both known and novel patterns.

4 This study is not without limitations. External validation for HRD predictions was not feasible
5 due to limited data availability, restricting the generalizability of our findings. Furthermore, in
6 the case of MSI, DL predictions based on NGS and PCR labels should be compared to the
7 histological approach of diagnosis - immunohistochemistry. Generally, small cohort sizes and
8 class imbalances may have introduced biases to our predictions. Unlike most studies, we
9 meticulously attempted to mitigate biases through site-aware data splitting [47], confounding
10 factors cannot be entirely excluded. Future studies should aim to replicate and validate our
11 findings in larger and more diverse patient cohorts to establish robustness.

12 As sequencing technologies advance, end-to-end DL approaches could become integral to
13 clinical genomics, particularly when WES/WGS data is unavailable and analysis relies on
14 routinely used targeted sequencing panels. These models hold the potential to streamline
15 workflows by bypassing lengthy preprocessing pipelines, possibly accelerating time to
16 diagnosis and treatment, and with this might reduce costs. For this purpose, future
17 advancements could focus on incorporating raw sequencing reads directly into DL models,
18 making them fully end-to-end. Finally, the incorporation of multi-omics data - whole
19 genome/exome, RNA and bisulfite sequencing - could enable DL to fully capture the
20 complexity and heterogeneity inherent in cancer genomes [64].

1 **Conclusion**

2 In summary, our study displays that DL models can predict clinically relevant biomarkers from
3 genomic data and while capturing complex mutational patterns associated with MSI and HRD.
4 By learning to identify and prioritize mutations influencing a given phenotype, DL models
5 reduce the need for human intervention and complement human expertise in precision
6 oncology. We provide an open-source toolkit to enable reproducibility and broader application
7 of these methods across cancer types and other diseases. Future efforts should focus on
8 expanding datasets, integrating raw multi-omics data, and refining model architectures to
9 maximize the impact of DL in advancing personalized medicine.

1 **Abbreviations**

- 2 AI – Artificial Intelligence
- 3 AUC – Area under the curve
- 4 AUROC – Area under the receiver operating characteristic
- 5 BRCA – Breast invasive carcinoma
- 6 COAD – Colon adenocarcinoma
- 7 CRC – Colorectal carcinoma
- 8 DL – Deep learning
- 9 HRD – Homologous Recombination Deficiency
- 10 HRP – Homologous Recombination Proficiency
- 11 LUAD – Lung adenocarcinoma
- 12 LUSC – Lung squamous cell carcinoma
- 13 MLP – Multilayer perceptron
- 14 MMEJ – Microhomology-mediated end joining
- 15 (d)MMR – (Defective) mismatch repair
- 16 MIL – Multiple instance learning
- 17 MSI – Microsatellite instability
- 18 MSS – Microsatellite stability
- 19 MUT – Mutated
- 20 NHEJ – Non-homologous end joining
- 21 NGS – Next-generation sequencing
- 22 PAAD – Pancreatic adenocarcinoma
- 23 PRAD – Prostate adenocarcinoma
- 24 OV – Ovarian serous cystadenocarcinoma
- 25 TCGA – The Cancer Genome Atlas
- 26 READ – Rectal adenocarcinoma
- 27 ROC – Receiver operating characteristic

- 1 UCEC – Uterine corpus endometrial carcinoma
- 2 WGS – Whole genome sequencing
- 3 WT – Wildtype
- 4 WXS – Whole exome sequencing
- 5 XGB – Extreme gradient boosting

1 **References**

- 2 1. Andre F, Filleron T, Kamal M, Mosele F, Arnedos M, Dalenc F, et al. Genomics to select
3 treatment for patients with metastatic breast cancer. *Nature*. 2022;610:343–8.
- 4 2. Mateo J, Steuten L, Aftimos P, André F, Davies M, Garralda E, et al. Delivering precision
5 oncology to patients with cancer. *Nat Med*. 2022;28:658–65.
- 6 3. Hodder A, Leiter SM, Kennedy J, Addy D, Ahmed M, Ajithkumar T, et al. Benefits for
7 children with suspected cancer from routine whole-genome sequencing. *Nat Med*.
8 2024;30:1905–12.
- 9 4. Pleasance ED, Cheetham RK, Stephens PJ, McBride DJ, Humphray SJ, Greenman CD,
10 et al. A comprehensive catalogue of somatic mutations from a human cancer genome.
11 *Nature*. 2010;463:191–6.
- 12 5. Horak P, Heining C, Kreutzfeldt S, Hutter B, Mock A, Hüllein J, et al. Comprehensive
13 genomic and transcriptomic analysis for guiding therapeutic decisions in patients with rare
14 cancers. *Cancer Discov*. 2021;11:2780–95.
- 15 6. Sosinsky A, Ambrose J, Cross W, Turnbull C, Henderson S, Jones L, et al. Insights for
16 precision oncology from the integration of genomic and clinical data of 13,880 tumors from
17 the 100,000 Genomes Cancer Programme. *Nat Med*. 2024;30:279–89.
- 18 7. ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analysis
19 of whole genomes. *Nature*. 2020;578:82–93.
- 20 8. Le DT, Durham JN, Smith KN, Wang H, Bartlett BR, Aulakh LK, et al. Mismatch repair
21 deficiency predicts response of solid tumors to PD-1 blockade. *Science*. 2017;357:409–13.
- 22 9. Chopra N, Tovey H, Pearson A, Cutts R, Toms C, Proszek P, et al. Homologous
23 recombination DNA repair deficiency and PARP inhibition activity in primary triple negative

- 1 breast cancer. *Nat Commun.* 2020;11:2662.
- 2 10. Germano G, Lamba S, Rospo G, Barault L, Magri A, Maione F, et al. Inactivation of DNA
3 repair triggers neoantigen generation and impairs tumour growth. *Nature.* 2017;552:116–20.
- 4 11. Alexandrov LB, Kim J, Haradhvala NJ, Huang MN, Tian Ng AW, Wu Y, et al. The
5 repertoire of mutational signatures in human cancer. *Nature.* 2020;578:94–101.
- 6 12. Polak P, Kim J, Braunstein LZ, Karlic R, Haradhavala NJ, Tiao G, et al. A mutational
7 signature reveals alterations underlying deficient homologous recombination repair in breast
8 cancer. *Nat Genet.* 2017;49:1476–86.
- 9 13. Moore Kathleen, Colombo Nicoletta, Scambia Giovanni, Kim Byoung-Gie, Oaknin Ana,
10 Friedlander Michael, et al. Maintenance Olaparib in Patients with Newly Diagnosed
11 Advanced Ovarian Cancer. *N Engl J Med.* 2018;379:2495–505.
- 12 14. Tutt ANJ, Garber JE, Kaufman B, Viale G, Fumagalli D, Rastogi P, et al. Adjuvant
13 Olaparib for Patients with BRCA1- or BRCA2-Mutated Breast Cancer. *N Engl J Med.*
14 2021;384:2394–405.
- 15 15. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial
16 sequencing and analysis of the human genome. *Nature.* 2001;409:860–921.
- 17 16. Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat*
18 *Rev Genet.* 2009;10:57–63.
- 19 17. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
20 universal RNA-seq aligner. *Bioinformatics.* 2013;29:15–21.
- 21 18. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. Mapping and quantifying
22 mammalian transcriptomes by RNA-Seq. *Nat Methods.* 2008;5:621–8.
- 23 19. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq

- 1 quantification. *Nat Biotechnol.* 2016;34:525–7.
- 2 20. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-
3 aware quantification of transcript expression. *Nat Methods.* 2017;14:417–9.
- 4 21. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A framework
5 for variation discovery and genotyping using next-generation DNA sequencing data. *Nat*
6 *Genet.* 2011;43:491–8.
- 7 22. Cibulskis K, Lawrence MS, Carter SL, Sivachenko A, Jaffe D, Sougnez C, et al.
8 Sensitive detection of somatic point mutations in impure and heterogeneous cancer
9 samples. *Nat Biotechnol.* 2013;31:213–9.
- 10 23. Ge H, Liu K, Juan T, Fang F, Newman M, Hoek W. FusionMap: detecting fusion genes
11 from next-generation sequencing data at base-pair resolution. *Bioinformatics.*
12 2011;27:1922–8.
- 13 24. Rausch T, Zichner T, Schlattl A, Stütz AM, Benes V, Korbel JO. DELLY: structural
14 variant discovery by integrated paired-end and split-read analysis. *Bioinformatics.*
15 2012;28:i333–9.
- 16 25. Li K, Luo H, Huang L, Luo H, Zhu X. Microsatellite instability: a review of what the
17 oncologist should know. *Cancer Cell Int.* 2020;20:16.
- 18 26. Dedeurwaerdere F, Claes KB, Van Dorpe J, Rottiers I, Van der Meulen J, Breyne J, et al.
19 Comparison of microsatellite instability detection by immunohistochemistry and molecular
20 techniques in colorectal and endometrial cancer. *Sci Rep.* 2021;11:12880.
- 21 27. Stewart MD, Merino Vega D, Arend RC, Baden JF, Barbash O, Beaubier N, et al.
22 Homologous Recombination Deficiency: Concepts, Definitions, and Assays. *Oncologist.*
23 2022;27:167–74.
- 24 28. Kris A. Wetterstrand MS. DNA sequencing costs: Data [Internet]. *Genome.gov. NHGRI;*

- 1 2019 [cited 2024 Sep 9]. Available from: [https://www.genome.gov/about-genomics/fact-](https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data)
- 2 [sheets/DNA-Sequencing-Costs-Data](https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data)
- 3 29. Unger M, Kather JN. Deep learning in cancer genomics and histopathology. *Genome*
- 4 *Med.* 2024;16:44.
- 5 30. Unger M, Kather JN. A systematic analysis of deep learning in genomics and
- 6 histopathology for precision oncology. *BMC Med Genomics.* 2024;17:48.
- 7 31. Poplin R, Chang P-C, Alexander D, Schwartz S, Colthurst T, Ku A, et al. A universal
- 8 SNP and small-indel variant caller using deep neural networks. *Nat Biotechnol.*
- 9 2018;36:983–7.
- 10 32. Chiu Y-C, Chen H-IH, Zhang T, Zhang S, Gorthi A, Wang L-J, et al. Predicting drug
- 11 response of tumors from integrated genomic profiles by deep neural networks. *BMC Med*
- 12 *Genomics.* 2019;12:18.
- 13 33. Elmarakeby HA, Hwang J, Arafah R, Crowdis J, Gang S, Liu D, et al. Biologically
- 14 informed deep neural network for prostate cancer discovery. *Nature.* 2021;598:348–52.
- 15 34. Friedman S, Gauthier L, Farjoun Y, Banks E. Lean and deep models for more accurate
- 16 filtering of SNP and INDEL variant calls. *Bioinformatics.* 2020;36:2060–7.
- 17 35. Anaya J, Sidhom J-W, Mahmood F, Baras AS. Multiple-instance learning of somatic
- 18 mutations for the classification of tumour type and the prediction of microsatellite status. *Nat*
- 19 *Biomed Eng.* 2024;8:57–67.
- 20 36. Sanjaya P, Maljanen K, Katainen R, Waszak SM, Genomics England Research
- 21 Consortium, Aaltonen LA, et al. Mutation-Attention (MuAt): deep representation learning of
- 22 somatic mutations for tumour typing and subtyping. *Genome Med.* 2023;15:47.
- 23 37. Cancer Genome Atlas Network. Comprehensive molecular characterization of human
- 24 colon and rectal cancer. *Nature.* 2012;487:330–7.

- 1 38. Kautto EA, Bonneville R, Miya J, Yu L, Krook MA, Reeser JW, et al. Performance
2 evaluation for rapid detection of pan-cancer microsatellite instability with MANTIS.
3 *Oncotarget*. 2017;8:7452–63.
- 4 39. Niu B, Ye K, Zhang Q, Lu C, Xie M, McLellan MD, et al. MSIsensor: microsatellite
5 instability detection using paired tumor-normal sequence data. *Bioinformatics*.
6 2014;30:1015–6.
- 7 40. Nguyen L, W M Martens J, Van Hoeck A, Cuppen E. Pan-cancer landscape of
8 homologous recombination deficiency. *Nat Commun*. 2020;11:5584.
- 9 41. Sztupinski Z, Diossy M, Krzystanek M, Reiniger L, Csabai I, Favero F, et al. Migrating
10 the SNP array-based homologous recombination deficiency measures to next generation
11 sequencing data of breast cancer. *NPJ Breast Cancer*. 2018;4:16.
- 12 42. Abkevich V, Timms KM, Hennessy BT, Potter J, Carey MS, Meyer LA, et al. Patterns of
13 genomic loss of heterozygosity predict homologous recombination repair defects in epithelial
14 ovarian cancer. *Br J Cancer*. 2012;107:1776–82.
- 15 43. Birkbak NJ, Wang ZC, Kim J-Y, Eklund AC, Li Q, Tian R, et al. Telomeric allelic
16 imbalance indicates defective DNA repair and sensitivity to DNA-damaging agents. *Cancer*
17 *Discov*. 2012;2:366–75.
- 18 44. Popova T, Manié E, Rieunier G, Caux-Moncoutier V, Tirapo C, Dubois T, et al. Ploidy
19 and large-scale genomic instability consistently identify basal-like breast carcinomas with
20 BRCA1/2 inactivation. *Cancer Res*. 2012;72:5454–62.
- 21 45. Rempel E, Kluck K, Beck S, Ourailidis I, Kazdal D, Neumann O, et al. Pan-cancer
22 analysis of genomic scar patterns caused by homologous repair deficiency (HRD). *NPJ*
23 *Precis Oncol*. 2022;6:36.
- 24 46. Cosmic. COSMIC [Internet]. 2020 [cited 2024 Jul 30]. Available from:

- 1 <https://cancer.sanger.ac.uk/signatures/id/>
- 2 47. Howard FM, Dolezal J, Kochanny S, Schulte J, Chen H, Heij L, et al. The impact of site-
3 specific digital histology signatures on deep learning model accuracy and bias. *Nat*
4 *Commun.* 2021;12:4423.
- 5 48. Ilse M, Tomczak J, Welling M. Attention-based Deep Multiple Instance Learning. In: Dy
6 J, Krause A, editors. *Proceedings of the 35th International Conference on Machine Learning*.
7 PMLR; 10--15 Jul 2018. p. 2127–36.
- 8 49. Cortes-Ciriano I, Lee S, Park W-Y, Kim T-M, Park PJ. A molecular portrait of
9 microsatellite instability across multiple cancers. *Nat Commun.* 2017;8:15180.
- 10 50. Dámaso E, Castillejo A, Arias MDM, Canet-Hermida J, Navarro M, Del Valle J, et al.
11 Primary constitutional MLH1 epimutations: a focal epigenetic event. *Br J Cancer.*
12 2018;119:978–87.
- 13 51. Aaltonen LA, Peltomäki P, Leach FS, Sistonen P, Pylkkänen L, Mecklin JP, et al. Clues
14 to the pathogenesis of familial colorectal cancer. *Science.* 1993;260:812–6.
- 15 52. Guo Q, Househam J, Lakatos E, Nowinski S, Al Bakir I, Grant H, et al. Long deletion
16 signatures in repetitive genomic regions track somatic evolution and enable sensitive
17 detection of microsatellite instability [Internet]. *Bioinformatics.* bioRxiv; 2024. Available from:
18 <https://www.biorxiv.org/content/10.1101/2024.10.03.616572v1>
- 19 53. Meier B, Volkova NV, Hong Y, Schofield P, Campbell PJ, Gerstung M, et al. Mutational
20 signatures of DNA mismatch repair deficiency in *C. elegans* and human cancers. *Genome*
21 *Res.* 2018;28:666–75.
- 22 54. TUMOR TYPES BIOMARKERS FDA-APPROVED THERAPY. Table 1: Companion
23 diagnostic indications [Internet]. [cited 2024 Nov 15]. Available from:
24 [27](https://www.foundationmedicine.com/sites/default/files/media/documents/2024-</div><div data-bbox=)

- 1 07/F1CDxTechnical_Specifications_Commercial_SPEC-01197_V3.0.pdf
- 2 55. Takeda M, Takahama T, Sakai K, Shimizu S, Watanabe S, Kawakami H, et al. Clinical
3 application of the FoundationOne CDx assay to therapeutic decision-making for patients with
4 advanced solid tumors. *Oncologist*. 2021;26:e588–96.
- 5 56. TruSight Oncology 500 [Internet]. [cited 2024 Nov 15]. Available from:
6 [https://emea.illumina.com/products/by-type/clinical-research-products/trusight-oncology-](https://emea.illumina.com/products/by-type/clinical-research-products/trusight-oncology-500.html#tabs-48377e6865-item-d4e16e33c1-documentation)
7 [500.html#tabs-48377e6865-item-d4e16e33c1-documentation](https://emea.illumina.com/products/by-type/clinical-research-products/trusight-oncology-500.html#tabs-48377e6865-item-d4e16e33c1-documentation)
- 8 57. Zhao C, Jiang T, Hyun Ju J, Zhang S, Tao J, Fu Y, et al. TruSight Oncology 500:
9 Enabling Comprehensive Genomic Profiling and Biomarker Reporting with Targeted
10 Sequencing [Internet]. *Bioinformatics*. bioRxiv; 2020. Available from:
11 <https://www.biorxiv.org/content/10.1101/2020.10.21.349100v1>
- 12 58. Sinicrope FA, Sargent DJ. Molecular pathways: microsatellite instability in colorectal
13 cancer: prognostic, predictive, and therapeutic implications. *Clin Cancer Res*. 2012;18:1506–
14 12.
- 15 59. Black SJ, Ozdemir AY, Kashkina E, Kent T, Rusanov T, Ristic D, et al. Publisher
16 Correction: Molecular basis of microhomology-mediated end-joining by purified full-length
17 Polθ. *Nat Commun*. 2020;11:1831.
- 18 60. Bennardo N, Cheng A, Huang N, Stark JM. Alternative-NHEJ is a mechanistically distinct
19 pathway of mammalian chromosome break repair. *PLoS Genet*. 2008;4:e1000110.
- 20 61. Davies H, Glodzik D, Morganella S, Yates LR, Staaf J, Zou X, et al. HRDetect is a
21 predictor of BRCA1 and BRCA2 deficiency based on mutational signatures. *Nat Med*.
22 2017;23:517–25.
- 23 62. Lee JJ, Kang HJ, Kim D, Lim SO, Kim SS, Kim G, et al. expHRD: an individualized,
24 transcriptome-based prediction model for homologous recombination deficiency assessment

- 1 in cancer. BMC Bioinformatics [Internet]. 2024;25. Available from:
- 2 <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-024-05854-y#Sec2>
- 3 63. Abbasi A, Steele CD, Bergstrom EN, Khandekar A, Farswan A, McKay RR, et al.
- 4 Detecting HRD in whole-genome and whole-exome sequenced breast and ovarian cancers.
- 5 medRxiv [Internet]. 2024; Available from:
- 6 <https://pmc.ncbi.nlm.nih.gov/articles/PMC11261949/>
- 7 64. Lipkova J, Chen RJ, Chen B, Lu MY, Barbieri M, Shao D, et al. Artificial intelligence for
- 8 multimodal data integration in oncology. Cancer Cell. 2022;40:1095–110.

1 **Declarations**

2 **Ethics approval and consent to participate**

3 This study was carried out in accordance with the Declaration of Helsinki. Datasets from
4 CPTAC and TCGA do not require formal ethics approval for a retrospective study since
5 samples are already anonymised. Furthermore, the analysis was approved by the Ethics
6 commission of the Medical Faculty of the Technical University Dresden (BO-EK-444102022).

7 **Consent for publication**

8 Not applicable.

9 **Availability of data and materials**

10 The mutation data used in this study can be accessed through the following sources: TCGA
11 data is available at <https://gdc.cancer.gov/about-data/publications/mc3-2017>, and CPTAC
12 data can be found at <https://portal.gdc.cancer.gov/>. TCGA access can be granted through eRA
13 Commons (https://public.era.nih.gov/commonsplus/home.era?menu_itemPath=600) and
14 dbGaP (<https://www.ncbi.nlm.nih.gov/gap/>). Chromosome information is accessible at
15 <https://ftp.ensembl.org/pub/>. Information regarding MSI and HRD status, as well as mutation
16 statuses of driver genes, is available through cBioPortal at <https://www.cbioportal.org/>.

17 Scripts for data preparation and model implementation are hosted on GitHub and can be found
18 at <https://github.com/mxsunc/Biomarker-ATGC>.

19 **Competing Interests**

20 **JNK** declares consulting services for Bioptimus, France; Owkin, France; DoMore Diagnostics,
21 Norway; Panakeia, UK; AstraZeneca, UK; Mindpeak, Germany; and MultiplexDx, Slovakia.
22 Furthermore, he holds shares in StratifAI GmbH, Germany, Synagen GmbH, Germany; has

1 received a research grant by GSK; and has received honoraria by AstraZeneca, Bayer, Daiichi
2 Sankyo, Eisai, Janssen, Merck, MSD, BMS, Roche, Pfizer, and Fresenius. No other competing
3 financial interests are declared by any of the remaining authors. **TG** is named as a coinventor
4 on patent applications that describe a method for TCR sequencing (GB2305655.9), and a
5 method to measure evolutionary dynamics in cancers using DNA methylation (GB2317139.0).
6 TG has received honorarium from Genentech and consultancy fees from DAiNA therapeutics.
7 The remaining authors have no competing interests to declare.

8 **Funding**

9 JNK is supported by the German Cancer Aid (DECADE, 70115166), the German Federal
10 Ministry of Education and Research (PEARL, 01KD2104C; CAMINO, 01EO2101; SWAG,
11 01KD2215A; TRANSFORM LIVER, 031L0312A; TANGERINE, 01KT2302 through ERA-NET
12 Transcan; Come2Data, 16DKZ2044A; DEEP-HCC, 031L0315A), the German Academic
13 Exchange Service (SECAI, 57616814), the German Federal Joint Committee (TransplantKI,
14 01VSF21048) the European Union's Horizon Europe and innovation programme (ODELIA,
15 101057091; GENIAL, 101096312), the European Research Council (ERC; NADIR,
16 101114631), the National Institutes of Health (EPICO, R01 CA263318) and the National
17 Institute for Health and Care Research (NIHR, NIHR203331) Leeds Biomedical Research
18 Centre. The views expressed are those of the author(s) and not necessarily those of the NHS,
19 the NIHR or the Department of Health and Social Care. This work was funded by the European
20 Union. Views and opinions expressed are however those of the author(s) only and do not
21 necessarily reflect those of the European Union. Neither the European Union nor the granting
22 authority can be held responsible for them.

23 **Author contributions**

24 MU and JNK conceptualized the study. MU performed the DL experiments and explainability
25 analysis. MU prepared the original draft and created the figures. JNK acquired funding and

1 supervised the study. MU, CMLL, LZ, SS, TL, JV, AM, SF, TG, ZIC, JNK provided scientific
2 input, reviewed and edited the manuscript.

3 **Acknowledgements**

4 We acknowledge the TCGA (<https://www.cancer.gov/tcga>) and the CPTAC
5 (<https://proteomics.cancer.gov/programs/cptac>) research networks, which generated the data
6 on which the results shown in this manuscript are based on.

7 In accordance with the COPE (Committee on Publication Ethics) position statement of 13
8 February 2023 (<https://publicationethics.org/cope-position-statements/ai-author>), the authors
9 hereby disclose the use of the following artificial intelligence models during the writing of this
10 article: GPT-4 (OpenAI) for checking spelling and grammar.

Figure 1

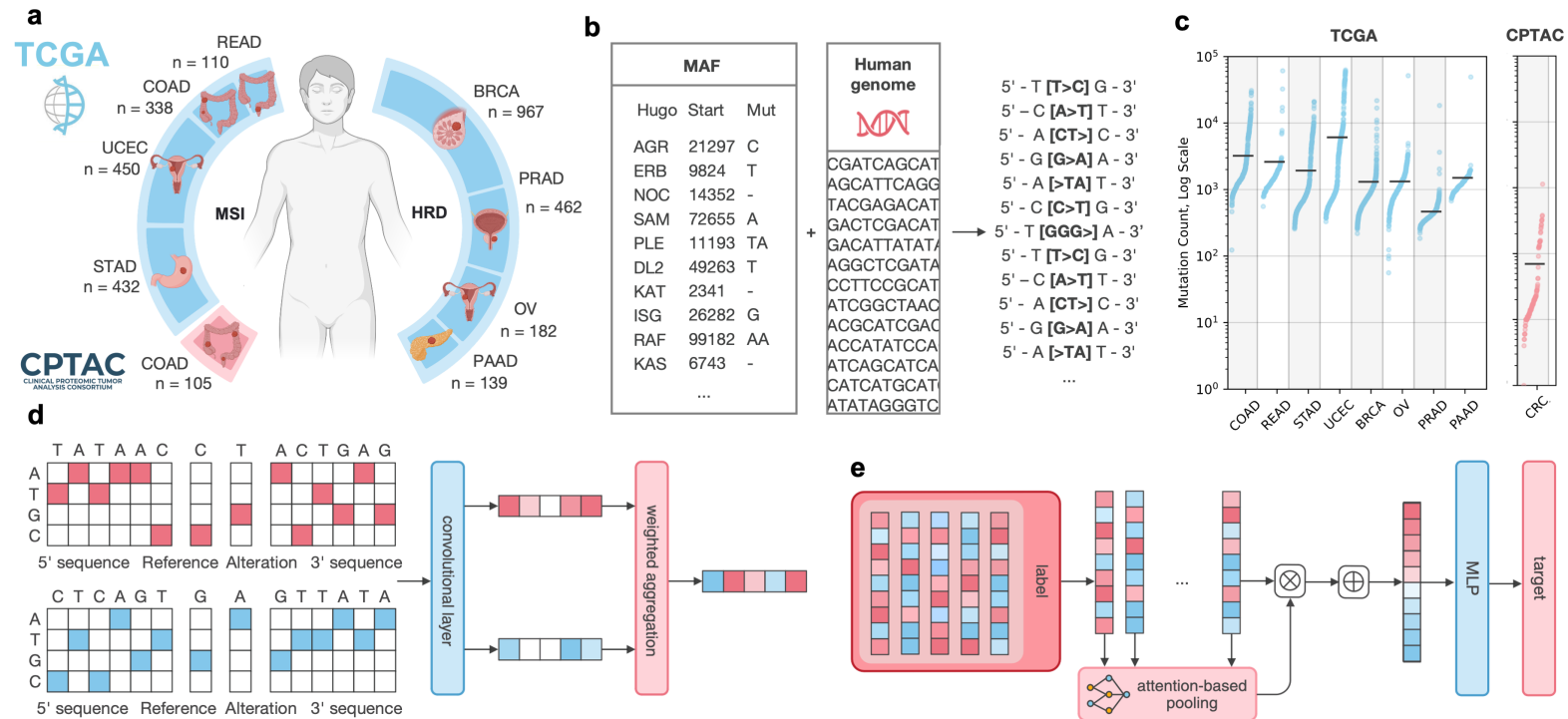
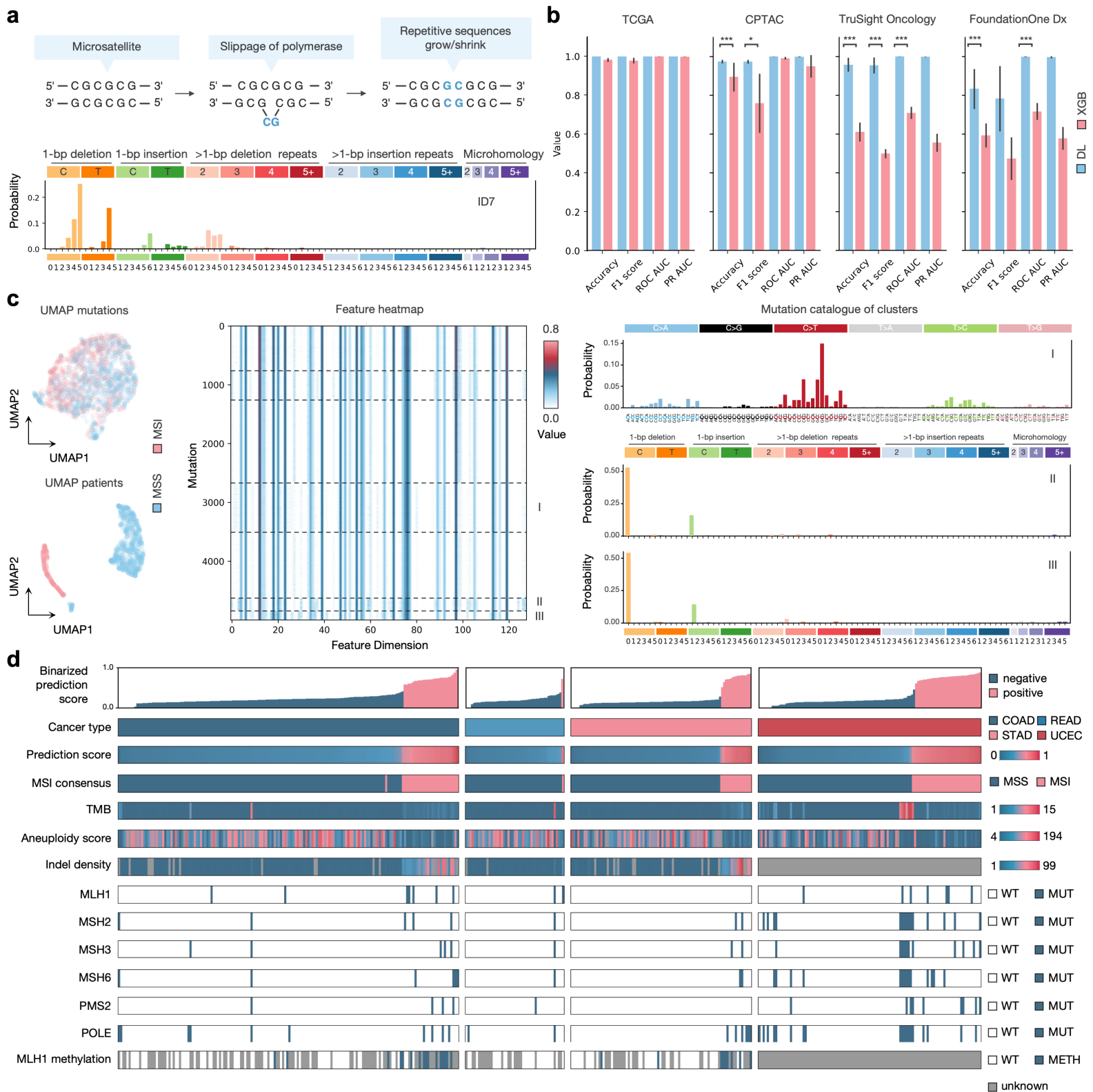


Figure 2

1 **Figure Captions**

2 **Fig. 1 | Overview of the study setup. a** Datasets (TCGA and CPTAC) used in this study with
3 cancer types and patient counts stratified by biomarkers. Cancer types from TCGA are
4 accentuated in blue, red indicates CPTAC. **b** Preprocessing of .maf files and reference
5 genome to mutation information. **c** Mutation counts by cohort and cancer type used in this
6 study. The black line indicates the mean mutation count. **d** Encoder part of the attMIL model.
7 Mutations from (c) are separated into four blocks: upstream, reference, alteration and
8 downstream element. Each block contains 20 nucleotides or gaps. Reverse strand is modeled
9 in the same manner by the reverse complement of the mutation. Both strand encodings are
10 passed through a 2D-convolution and a dense layer producing a mutation feature vector. **e**
11 Mutation vectors are gathered at the patient level and are aggregated by the attMIL
12 mechanism to a patient feature vector. Patient features are used for the classification task.
13 (COAD - colon adenocarcinoma, READ - rectal adenocarcinoma, STAD - stomach
14 adenocarcinoma, UCEC - uterine corpus endometrial carcinoma, OV - Ovarian serous
15 cystadenocarcinoma, BRCA - Breast invasive carcinoma, PRAD - Prostate adenocarcinoma,
16 PAAD - Pancreatic adenocarcinoma, MLP - multilayer perceptron)

17 **Fig. 2 | DL predicts MSI with high accuracy in panel sequencing. a** Mutation mechanism
18 and COSMIC Indel signature of MSI. **b** Bar chart of performance metrics to compare attMIL
19 and XGB models in the internal test dataset of TCGA, the external full CPTAC dataset and
20 CPTAC data filtered by targeted sequencing panels of FoundationOne and Trusight Oncology.
21 Significance is indicated by * (p-value < 0.05), ** (p-value < 0.005) and *** (p-value < 0.0005).
22 **c** Explainability of the attMIL model. UMAP of mutation and patient level features. Heatmap of
23 5,000 random mutation features clustered by k-means (k=7) and ranked by mean attention
24 score. Indel/SBS mutation catalogues of three mutation clusters. **d** Patient features sorted by
25 increasing prediction score of the attMIL model and separated by cancer type. (COAD - colon
26 adenocarcinoma, READ - rectal adenocarcinoma, STAD - stomach adenocarcinoma, UCEC

1 - uterine corpus endometrial carcinoma, MSI - Microsatellite instability, MSS - microsatellite
2 stability, WT - wildtype, MUT - mutated, METH - methylated)

3 **Fig. 3 | DL predicts HRD. a** Mutation mechanism of MMEJ and COSMIC Indel signature of
4 HRD. **b** ROC AUC and PR AUC plots for attMIL HRD predictions. **c** Explainability of the attMIL
5 model. UMAP of mutation and patient level features. Heatmap of 5,000 random mutation
6 features clustered by k-means (k=7) and ranked by mean attention score. Indel/SBS mutation
7 catalogues of three mutation clusters. **d** Patient features sorted by increasing prediction score
8 of the attMIL model and separated by cancer type. (OV - Ovarian serous cystadenocarcinoma,
9 BRCA - Breast invasive carcinoma, PRAD - Prostate adenocarcinoma, PAAD - Pancreatic
10 adenocarcinoma, HRD - Homologous Recombination Deficiency, HRP - Homologous
11 Recombination Proficiency, WT - wildtype, MUT - mutated)

1 **Supplementary Material**

2 **Supplementary Figures**

3 **Supplementary Fig. 1 | MSI data preparation.** **a** Bar chart of site-specific data split for MSI.
4 Dark blue indicates training data, light blue validation data and red the test set. The test set
5 stays the same over all folds. **b** Doughnut plot of the data splits regarding class distribution of
6 MSI vs MSS and number of patients. **c** Bar chart of number of patients per cancer type.

7 **Supplementary Fig. 2 | HRD data preparation.** **a** Bar chart of site-specific data split for HRD.
8 Dark blue indicates training data, light blue validation data and red the test set. The test set
9 stays the same over all folds. **b** Doughnut plot of the data splits regarding class distribution of
10 HRD vs HRP and number of patients. **c** Bar chart of number of patients per cancer type.

11 **Supplementary Fig. 3 | MSI features.** **a** COSMIC SBS mutation signatures 6 and 15
12 associated with MSI. **b** SBS mutation catalogues of mutation feature clusters not displayed in
13 Fig. 2c. Only mutation catalogues of the majority mutation class are displayed.

14 **Supplementary Fig. 4 | XGB features used in MSI predictions.** 20 most important features
15 of XGB to classify MSI. Features are sorted by contribution to prediction (high to low SHAP
16 values). SBS mutations are displayed as, for example, a C>T mutation with an upstream G
17 and a downstream C (GC>TC). Indel mutations can be read as first the indel pattern length
18 and then the number of repetitions. For example, an insertion pattern of at least five
19 nucleotides without a repetition (5InsRep0) or a mononucleotide deletion of a C with the length
20 one (1DelC1).

21 **Supplementary Fig. 5 | HRD features.** **a** COSMIC SBS mutation signature three and indel
22 signature eight associated with HRD and NHEJ. **b** SBS and indel mutation catalogues of
23 mutation feature clusters not displayed in Fig. 3c. Only mutation catalogues of the majority
24 mutation class are displayed.

1 **Supplementary Fig. 6 | XGB features used in HRD predictions. a** ROC AUC and PR AUC
2 of XGB predicting HRD in a five-fold cross validation. **b** 20 most important features of XGB to
3 classify HRD. Features are sorted by contribution to prediction (high to low SHAP values).
4 Microhomology mutations can be read as first the deletion length and then the length of the
5 microhomology. For example, a deletion of at least five nucleotides containing a
6 microhomology of two nucleotides to at least one of the flanking sites (5+DelMH2).

1 Supplementary Tables

2 **Supplementary Table. 1 | MSI patient counts.** MSI and MSS cases in the TCGA data split
 3 by cancer type and in total

	COAD	READ	STAD	UCEC	total
MSI cases	59	3	82	136	280
MSS cases	279	107	350	314	1050

4
 5 **Supplementary Table. 2 | HRD patient counts.** HRD and HRP cases in the TCGA data split
 6 by cancer type and in total

	BRCA	OV	PAAD	PRAD	total
HRD cases	149	74	6	79	308
HRP cases	818	108	133	383	1442

7
 8 **Supplementary Table. 3 | MSI performance metrics.** Performance values on the full
 9 CPTAC external test dataset of attMIL and XGB predicting MSI. The p-values for
 10 corresponding model comparison are stated below the performances. (accuracy - acc, sens -
 11 sensitivity, spec - specificity)

	acc	F1 score	ROC AUC	PR AUC	sens	spec	TP	FP	TN	FN
attMIL	0.98± 0.01	0.97± 0.01	1.00± 0.00	1.00± 0.00	0.95± 0.02	1.00± 0.00	22.80 ± 0.45	0.00± 0.00	82.00 ± 0.00	1.20± 0.45
XGB	0.89± 0.07	0.76± 0.15	0.99± 0.01	0.95± 0.06	0.99± 0.02	0.80± 0.16	22.80 ±0.45	16.40 ± 12.93	64.60 ± 12.93	0.20± 0.45
p-value	5.11x 10 ⁻¹⁸	0.04	0.32	-	0.04	0.047	1.00	0.06	0.06	0.06

12
 13 **Supplementary Table. 4 | MSI performance metrics on TruSight Oncology 500 CPTAC**
 14 **variants.** Performance values on the TruSight Oncology 500 filtered CPTAC external test

1 dataset of attMIL and XGB predicting MSI. The p-values for corresponding model comparison
 2 are stated below the performances. (accuracy - acc, sens - sensitivity, spec - specificity)

	acc	F1 score	ROC AUC	PR AUC	sens	spec	TP	FP	TN	FN
attMIL	0.96± 0.04	0.96± 0.04	1.00± 0.00	1.00± 0.00	0.92± 0.07	1.00± 0.00	22.00 ± 1.73	0.00± 0.00	82.00 ± 0.00	2.00± 1.73
XGB	0.61± 0.05	0.50± 0.02	0.71± 0.03	0.56± 0.05	0.68± 0.19	0.54± 0.28	15.60 ±4.28	24.60 ± 15.26	29.40 ± 15.26	7.40± 4.28
p-value	1.67x 10 ⁻³⁸	7.56x 10 ⁻⁵	4.66x 10 ⁻¹⁶	-	0.04	0.02	0.06	0.06	0.06	0.59

3
 4 **Supplementary Table. 5 | MSI performance metrics on FoundationOne Dx CPTAC**
 5 **variants.** Performance values on the FoundationOne Dx filtered CPTAC external test dataset
 6 of attMIL and XGB predicting MSI. The p-values for corresponding model comparison are
 7 stated below the performances. (accuracy - acc, sens - sensitivity, spec - specificity)

	acc	F1 score	ROC AUC	PR AUC	sens	spec	TP	FP	TN	FN
attMIL	0.83± 0.10	0.78± 0.17	1.00± 0.00	1.00± 0.00	0.67± 0.20	1.00± 0.00	16.00 ± 4.95	0.00± 0.00	82.00 ± 0.00	8.00± 4.95
XGB	0.59± 0.06	0.47± 0.11	0.72± 0.04	0.58± 0.06	0.63± 0.26	0.55± 0.28	14.60 ±6.02	22.40 ± 14.21	27.60 ± 14.21	8.40± 6.02
p-value	7.94x 10 ⁻²¹	0.06	2.88x 10 ⁻¹⁷	-	0.85	0.02	0.6	0.06	0.06	0.06

8
 9 **Supplementary Table. 6 | HRD performance metrics.** Performance values on the TCGA
 10 test dataset of attMIL and XGB predicting HRD. The p-values for corresponding model
 11 comparison are stated below the performances. (accuracy - acc, sens - sensitivity, spec -
 12 specificity)

	acc	F1 score	ROC AUC	PR AUC	sens	spec	TP	FP	TN	FN
--	-----	----------	---------	--------	------	------	----	----	----	----

attMIL	0.80± 0.01	0.65± 0.02	0.88± 0.01	0.67± 0.01	0.74± 0.06	0.86± 0.05	33.40 ± 2.88	25.00 ± 9.38	159.0 ± 9.38	11.60 ± 2.88
XGB	0.75± 0.02	0.66± 0.03	0.86± 0.01	0.75± 0.02	0.52± 0.03	0.98± 0.00	25.40 ± 1.52	2.80± 0.84	183.2 0± 0.84	23.60 ± 1.52
p-value	1.07x 10 ⁻¹⁰	0.51	0.72	-	0.003	0.049	0.06	0.06	0.06	0.06

1 **Supplementary Methods**

2 **Feature Catalogues for ML model**

3 To generate the SBS mutation catalogues, patient mutations were grouped by 96 SBS
4 features, encompassing six mutation classes: C>A, C>G, C>T, T>A, T>C, T>G (including the
5 reverse complement). Each class was further stratified by the adjacent 3' and 5' nucleotides.
6 Since there are four possible options for the 3' and 5' positions, this results in $6 * 4 * 4 = 96$
7 features.

8 For indel feature generation, mutations were grouped according to insertion or deletion length,
9 repetition pattern, and a possible homology to flanking sites. For mononucleotide deletions of
10 C or T (including complementary base), two classes were created, along with two classes for
11 corresponding mononucleotide insertions. Additionally, four classes were defined based on
12 repetition patterns for deletions and insertions. Each of these twelve initial subclasses was
13 then subdivided into six subclasses based on the number of repeats deleted or inserted,
14 resulting in 72 subclasses. Finally, four microhomology classes were assigned, describing
15 mutations in which either the 3' or 5' flanking site has partial homology to the deletion
16 sequence. These classes were further subdivided based on homology length, adding 11
17 additional classes to the 72 remaining classes, resulting in a total of 83 indel classes.

18 **Multiple Instance Learning**

19 Multiple Instance Learning (MIL) is a deep learning framework designed to handle scenarios
20 where labels are associated with sets of instances, referred to as "bags," rather than individual
21 instances. In MIL, only a subset of instances within a bag contribute to the overall class
22 prediction, while others may be irrelevant or unrelated. This makes MIL particularly suitable
23 for contexts where identifying which instances are informative is inherently challenging or
24 unknown [48]

1 Here, we applied MIL to model patient-level predictions, where each patient is represented as
2 a "bag" containing multiple somatic mutations ("instances") [35,36]. Not all mutations
3 contribute to the phenotype of interest (e.g., MSI or HRD), as some mutations may not be
4 associated with the underlying mutational mechanisms driving these biomarkers. MIL allows
5 the model to autonomously learn which mutations are most relevant by assigning attention
6 weights to individual instances, rather than relying on predefined rules or human-curated
7 knowledge.

1 **Additional Files**

2 **Additional file 1 | MSI status.** Excel file in which the patient identifiers of TCGA and CPTAC
3 patients with corresponding tissue type and MSI status is stored.

4 **Additional file 2 | HRD status.** Excel file in which the patient identifiers of TCGA patients with
5 corresponding information is stored. Information includes: tissue type, scar HRD scores with
6 subscores (LST, LOH, TAI) and HRD status with cutoffs at 42 and tissue specific cutoffs.

7 **Additional file 3 | Tissue source sites.** Excel file in which the tissue source site codes within
8 the TCGA patient identifiers and their mapping to the actual source site is stored.