

1 **File S1 Supplementary methods**

2 *Additional fragmentomics feature profiling*

3 To generate the Griffin feature, a profiling framework was adopted from Doebley and
 4 colleagues [1] to evaluate nucleosome occupancy and the corresponding protection of cell-free
 5 DNA (cfDNA). The GC-corrected coverage profile was quantified by employing three observable
 6 characteristics: the central coverage value measured at 30 base pairs from the specified location,
 7 the average coverage calculated within a 1000 base pair radius of the location, and the amplitude
 8 ascertained through the application of a Fast-Fourier Transform analysis. A total of 854 Griffin
 9 features representing transcription factor binding sites were generated from the low-pass whole-
 10 genome sequencing (WGS) data.

11 The neomer features were defined as short DNA sequences, which recur in tumor genomes
 12 but are absent from the human reference genome, by Georgakopoulos-Soares and colleagues [2].
 13 A total of 977 recurrent single-nucleotide polymorphisms (SNPs) were identified from 2577
 14 cancer patient samples using the PCAWG database (<https://dcc.icgc.org/releases/PCAWG/>). In
 15 total, 4,616 neomers of 16bp length were extracted from these SNPs, which were then filtered
 16 against common population variants compiled in the Genome Aggregation Database (gnomAD v2)
 17 [3], resulting in a final total of 1,758 neomer feature. The neomer features were profiled as the
 18 ratio of neomer-detecting reads over the total reads and the read count of each of the 1,758 neomers.

19 The motif breakpoint (MBP) feature [4] examined the frequencies of the 6bp motif at the 5'
 20 breakpoints on the human reference genome hg19, which extended 3bp to each direction. A total
 21 of 4,096 (4^6) MBP features were generated from the low-pass WGS data.

22 *cfDNA fragmentomics (cfFrag) score construction*

23 An automated machine-learning (autoML) process that utilizes five different algorithms,
 24 including generalized linear model (GLM), gradient boosting machine (GBM), random forest (RF),
 25 deep learning (DL), and eXtreme gradient boosting (XGBoost) [5], was employed to generate
 26 optimal base learners. The autoML utilized a randomized search for automatic algorithm selection,
 27 as well as for hyperparameter tuning. For each cfDNA fragmentomics feature type, a total of 200
 28 base learners were constructed using an autoML procedure, which performs hyperparameter
 29 tuning via random grid search.

30 The area under the curves was calculated using the training dataset via a 5-fold cross-
 31 validation approach for all base learners. For each feature type, the base learners were ranked by
 32 their AUCs, and the top 8 performing were then selected for constructing the cfFrag score. The
 33 final cfFrag score for each sample was then generated by calculating the mean predict score of the
 34 total 24 (3 × 8) optimal base learners.

$$\begin{aligned}
 & \text{cfFrag Score} \\
 & = \left(\frac{1}{8}\right) * \sum(\text{Score}(\text{FSR}_i)) + \left(\frac{1}{8}\right) * \sum(\text{Score}(\text{FSD}_i)) + \left(\frac{1}{8}\right) \\
 & * \sum(\text{Score}(\text{CNV}_i))
 \end{aligned}$$

38 The base learner predicts the score, which ranges from 0 to 1, representing the probability of
 39 a sample being breast cancer (0 = perfect benign nodule, 1 = perfect breast cancer; CNV, copy
 40 number variation; FSD, fragment size distribution; FSR, fragment size ratio.).

41 *Fragmentomics features evaluation*

42 We further evaluated the contribution of individual fragmentomics features by ranking them
 43 according to their importance in the final cfFrag model. For each base learner, we calculated the
 44 relative importance of individual features and sorted them from highest to lowest importance
 45 (using the maximum rank method for any tied values). For each of the three feature types, including
 46 copy number variation (CNV), fragment size distribution (FSD), and fragment size ratio (FSR),
 47 we determined the final importance of individual features by ranking the summed importance
 48 ranks of the eight selected top-base learners. A lower summed importance rank indicates a higher
 49 feature importance in the final cfFrag model.

50 **References**

- 51 [1] Doebley A-L, Ko M, Liao H, Cruikshank AE, Santos K, Kikawa C, et al. A framework for
 52 clinical cancer subtyping from nucleosome profiling of cell-free DNA. Nature Communications
 53 2022;13.
- 54 [2] Georgakopoulos-Soares I, Yizhar-Barnea O, Mouratidis I, Hemberg M, Ahituv N. Absent from
 55 DNA and protein: genomic characterization of nullomers and nullpeptides across functional
 56 categories and evolution. Genome Biol 2021;22:245.
- 57 [3] Gudmundsson S, Singer-Berk M, Watts NA, Phu W, Goodrich JK, Solomonson M, et al. Variant
 58 interpretation using population databases: Lessons from gnomAD. Hum Mutat 2022;43:1012-30.

- 59 [4] Ma X, Chen Y, Tang W, Bao H, Mo S, Liu R, et al. Multi-dimensional fragmentomic assay for
60 ultrasensitive early detection of colorectal advanced adenoma and adenocarcinoma. *J Hematol*
61 *Oncol* 2021;14:175.
- 62 [5] Bao H, Wang Z, Ma X, Guo W, Zhang X, Tang W, et al. Letter to the Editor: An ultra-sensitive
63 assay using cell-free DNA fragmentomics for multi-cancer early detection. *Mol Cancer*
64 2022;21:129.