Integrated fragmentomic profile and 5-Hydroxymethylcytosine 1 of capture-based low-pass sequencing data enables pan-cancer detection 2 3 via cfDNA 4 Zhidong Zhang^{1,†}, Xuenan Pi^{1,†}, Chang Gao¹, Jun Zhang², Lin Xia¹, Xiaogin 5 Yan², Xinlei Hu¹, Ziyue Yan², Shuxin Zhang³, Ailin Wei⁷, Yuer Guo¹, Jingfeng 6 Liu⁵, Ang Li⁴, Xiaolong Liu⁵, Wei Zhang⁶, Yanhui Liu³, Dan Xie1^{*} 7

¹National Frontier Center of Disease Molecular Network, State Key Laboratory 8

9 of Biotherapy, West China Hospital, Sichuan University, No. 37 Guoxue Alley,

10 Chengdu 610041, Sichuan Province, P. R. China.

²Tailai Inc., Shanghai 200233, P. R. China. 11

12 ³Department of Neurosurgery, West China Hospital, Sichuan University, No.

37 Guoxue Alley, Chengdu 610041, Sichuan Province, P. R. China. 13

⁴Department of Pancreatic Surgery, West China Hospital, Sichuan University, 14

15 No. 37 Guoxue Alley, Chengdu 610041, Sichuan Province, P. R. China.

⁵Mengchao Hepatobiliary Hospital of Fujian Medical University, Xihong Road 16

312, Fuzhou 350025, Fujian Province, P. R. China. 17

⁶Department of Respiratory and Critical Care Medicine, First Affiliated Hospital 18

of the Second Military Medical University, Shanghai 200433, Shanghai, P. R. 19

China. 20

⁷Guang'an People's Hospital, Guang'an, China 21

*Correspondence to: Dr Dan Xie, National Frontier Center of Disease 22 NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

Molecular Network, State Key Laboratory of Biotherapy, West China Hospital,
Sichuan University, No. 37 Guoxue Alley, Chengdu 610041, Sichuan Province,
P. R. China. Tel: +86-28-8558-2944; Fax: +86-28-8558-2944; E-mail:
danxie@scu.edu.cn (D. Xie).

²⁷ [†] These authors contributed equally to this work as joint first authors.

28

29 Abstract

Using epigenetic markers and fragmentomics of cell-free DNA for cancer 30 31 detection has been proven applicable. We further combine the two features and explore the diagnostic potential of the features on pan-cancer detection. 32 We extracted cfDNA fragmentomic features from 191 whole-genome 33 34 sequencing data and investigated them in 396 low-pass 5hmC sequencing data from four common cancer types and controls. We identified aberrant 35 ultra-long fragments (220-500bp) of cancer samples in 5hmC sequencing 36 data, both in size and coverage profile, and showed its dominant role in 37 cancer prediction. Since cfDNA hydroxymethylation and fragmentomic 38 markers can be detected simultaneously in low-pass 5hmC sequencing data, 39 we built an integrated model including 63 features of both fragmentomic 40 41 features and hydroxymethylation signatures for pan-cancer detection with high sensitivity and specificity (88.52% and 82.35%, respectively). We 42 showed that fragmentomic information in 5hmC sequencing data is an ideal 43 marker for cancer detection and that it shows high performance in low-pass 44

45 sequencing data.

46

47 Introduction

Displaying a trend of increase, more than 19.3 million individuals were 48 49 diagnosed with cancer, and 10.0 million individuals died of cancer worldwide in 2020 (1). Currently, tissue biopsy is the most widely used method for 50 cancer diagnosis and staging. However, owing to being an invasive 51 examination method, tissue biopsy has difficulties in tracking clonal evolution, 52 53 monitoring treatment response, and early detection of cancer and recurrence (2). Therefore, the need for the development of accurate liquid biopsy 54 methods is urgent, with the potential to detect cancer and identify resistance 55 56 mechanisms early, to quantify minimal residual disease, and to monitor treatment responses (2, 3). 57

Circulating cell-free tumor DNA (ctDNA) is a kind of tumor-derived material in 58 59 the blood of patients with cancer and is mainly derived from cancer cells via apoptosis or necrosis (3, 4). ctDNA carries cancer-specific genetic mutations, 60 epigenetic alterations, and fragmentomic aberrations, allowing for cancer 61 and tissue-of-origin prediction in a non-invasive 62 detection manner. Nevertheless, as cancer is a highly heterogeneous disease, biomarkers 63 based on genetic and epigenetic aberrations often have limitations in practical 64 applications. Recent studies have shown that cell-free DNA (cfDNA) 65 fragmentation is a non-random process (5, 6), and it is thus possible to 66

develop a generic approach for non-invasive cancer detection using cfDNA
 fragmentomics.

69 The fragmentomics analysis of cfDNA encompasses the sizes (7), coverage (7), nucleosome footprints (8), end-points (9, 10) and end-motifs (11) of cfDNA. 70 71 For example, the size distribution of cfDNA fragments can indicate specific cellular biological changes (12, 13). Coverage in cfDNA sequencing can be 72 used to analyze nucleosome positions (8, 14) and identify expressed genes 73 (14). Fragmentation patterns of cfDNA can be characterized by deferentially 74 75 phased fragment end signals, which were preferentially found in tissuespecific open chromatin regions (15). Practically, Jiang et al. showed that 76 cfDNA fragments ending with particular genomic coordinates or motifs had 77 78 higher degrees of association with hepatocellular carcinoma (HCC) (10, 11). Recently, the epigenetic marker 5-hydroxymethylcytosine (5hmC) 79 has attracted attention in the field of cancer research due to its regulatory role in 80 81 tissue-specific gene expression (16). The cfDNA captured by 5hmC comprises a subset of all circulating cfDNA and, therefore, may contain fragmentomic 82 information. Integrative multi-omics analysis including the 5hmC profile of 83 cfDNA improves the sensitivity and specificity of cancer detection, as 84 demonstrated by recent studies (17, 18). 85

To highlight the potential of 5hmC sequencing data in cancer detection, we obtained a large cohort of cancer samples, including HCC, pancreatic ductal adenocarcinoma (PDAC), lung adenocarcinoma (LUAD), and glioblastoma

(GBM). We then established a pan-cancer preferred end map based on WGS data from these four cancer types and checked their characteristics in lowpass 5hmC sequencing data. We hypothesized that the fragmentomic information would be consistent between WGS and 5hmC sequencing data and tested it by comparison of each feature. We built classification models and investigated the effect of combining 5hmC signatures with fragmentomic information in pan-cancer screening.

96

97 Methods

98 Sample collection and study design

In total, low-pass (mean mapped reads: 38.32±8.02 million) 5hmC-seq data
from 85 healthy controls and 311 cancer patients were included in this study.
We obtained 5hmC sequencing data from 33 LUAD (18), 74 PDAC (37), 132
HCC (38), and 85 control (37) samples from three publications, and 72 GBM
samples from an unpublished paper.

Plasma samples from 15 healthy controls and 59 cancer patients were collected and subjected to WGS (mean: 14×, range: 8.6-23.7×) to establish the preferred end marker set, while the low-pass WGS data from 19 healthy controls and 98 cancer patients (mean: 3×, range: 2.3-6.2×) were obtained to cross-check with low-pass 5hmC sequencing data. 13 of the samples with low-pass WGS sequencing had above 10x WGS sequencing data. The demographics and clinical characteristics of the cohort are summarized in

111 Supplementary Tables 1 and 2. The design of the study is shown in Figure 1.

112

113 Plasma sample collection and cfDNA WGS sequencing

We collected plasma samples from 46 GBM patients, 29 PDAC patients and 114 34 controls from West China Hospital. Plasma samples of 37 HCC patients 115 and 32 LUAD patients were collected from Mengchao Hepatobiliary Hospital 116 of Fujian Medical University and the First Affiliated Hospital of the Second 117 Military Medical University, respectively. For every subject, the QIAseg cfDNA 118 119 All-in-one Kit (Qiagen) was applied on the 2 ml plasma sample for the cfDNA extraction and library construction. Pair-end 150 bp sequencing of the libraries 120 was performed on the Illumina Novaseg 6000 platform. 121

122

123 cfDNA 5hmC profiling and sequencing

cfDNA extraction and 5hmC library construction was performed as previously 124 125 described (28). Firstly, cfDNA was extracted from 2 ml plasma sample using QIAamp Circulating Nucleic Acid Kit (QIAGEN Inc., Valencia, CA, USA). Then, 126 cfDNA (5-10 ng) ligated with sequencing adaptors was incubated in a 25 µL 127 reaction solution containing HEPES buffer (50 mM, pH 8.0), 25 mM MgCl2, 60 128 μM N3-UDP-Glc (ActiveMotif, Carlsbad, CA, USA), and 12.5 U β-129 glucosyltransferase (NEB, Beverly, MA, USA) for 2 h at 37 °C. Next, 2.5 µL 130 DBCO-PEG4-biotin (Sigma, Carlsbad, CA, USA) was directly added and 131 incubated for 2 h at 37 °C. 10 µg sheared salmon sperm DNA (Life 132

Technologies, USA) was added, then the Micro Bio-Spin 30 Column (Bio-Rad, 133 Hercules, CA, USA) was used to purify the DNA following the instruction, and 134 the final volume was adjusted to 25µL. After that, the purified DNA was 135 incubated with 5 µL C1 streptavidin beads (Life Technologies, USA) in buffer 1 136 137 (5 mM Tris pH 7.5, 0.5 mM EDTA, 1 M NaCl and 0.2% Tween 20) for 30 min. The beads were subsequently undergone three 5-min washes each with 138 buffer 1, buffer 2 (buffer 1 without NaCl), buffer 3 (buffer 1 with pH 9) and 139 buffer 4 (buffer 3 without NaCl). Then, the beads were resuspended in water 140 141 and amplified with 11 cycles of PCR amplification (initial denaturing at 98 °C for 45 s, followed by 11 cycles of denaturing at 98 °C for 15 s, annealing at 142 60 °C for 30 s, extension at 72 °C for 30 s, and a final extension at 17 °C for 5 143 144 min). The amplified products were purified using 0.8 × AMPure XP beads (Beck-man Coulter, Fullerton, CA, USA). Pair-end 150 bp sequencing was 145 performed on the Illumina Novaseg 6000 platform. 146

147

148 Fragmentomic profiling of cfDNA and cancer prediction analysis

The details of the preferred end, 5hmC signatures, size profile and coverage profile calculation are described in the Supplementary Material. We used a random forest model to distinguish healthy people from cancer patients using fragmentomic features. To estimate the prediction error, we used five crossvalidations. Near-zero variance features were removed. The training, validation, and test sets account for 60%, 20%, and 20% of the data,

respectively. The samples were selected randomly in a balanced way to keep 155 the ratio of the number of cancer to non-cancer samples similar in the training. 156 157 validation, and test subsets. Feature importance was calculated using the training data in each cross-validation run, and we sorted the features 158 159 according to the mean value of feature importance. The number of features used in the final model was obtained by the highest AUC value in the 160 validation set. To build an integrative model, we fitted the features selected by 161 every feature-alone model into one random forest model, and performed 162 163 feature selection, training, testing as described above. Random forest machine learning implemented 164 was using the python package sklearn.ensemble.RandomForestClassifier with the following parameters: 165 166 n estimators = 100, criterion="gini".

167

168 Statistics

All statistical analyses were performed using R version 3.6.1 and python 3.7.0.
All tests were two-sided, and P values < 0.05 were considered statistically
significant. The ROC-AUC plot and AUC value were implemented using
sklearn.metrics.plot_roc_curve.

173

174 Ethics

This study was approved by the Ethics Committee of Sichuan Cancer Hospital
 (SCCHEC-02-2016-005). The written informed consent was obtained from all

- 177 participants.
- 178
- 179
- 180 Results
- 181
- 182

183 Constructing the cfDNA preferred end map

cfDNA fragments ending at specific genomic coordinates that were statistically 184 185 more abundant than the Poisson distribution was an essential fragmentation signature referred to as the preferred end (10). We utilized WGS data from 186 the four cancer types to build a preferred end map and checked the feature 187 188 using low-pass data. The basis of utilizing cfDNA preferred end coordinates to indicate the occurrence and location of cancer was its ability to locate 189 nucleosomes, the position of which is related to cell identity (19). To confirm 190 191 this, we defined upstream (U) ends and downstream (D) ends by their genomic coordinate order on the reference genome calculated via alignment. 192 We surveyed the distance distribution between adjacent U ends, D ends, and 193 nucleosome centers, and the result suggested the preferred ends were 194 enriched on the DNA linker (Supplementary Figure 1a-f). We further checked 195 the distribution of the predicted nucleosome location based on non-malignant 196 tissues at chr12: 34517269-34519122, and the preferred U and D ends of 197 WGS data from controls (n=15) and HCC patients (n=15). As expected, the 198

199 preferred U and D ends of the controls were enriched at the midpoint between two nucleosomes, which is the location of the DNA linker (Figure 2a), and U 200 201 ends were found downstream of their nearest D ends. However, in HCC patients, U ends were upstream of their nearest D ends, and only D ends 202 were found at one midpoint between two nucleosomes (Figure 2b), which 203 indicated inconsistent nucleosome positions between HCC patients with 204 controls. We observed the same abnormal preferred end distribution in GBM. 205 LUAD and PDAC patients (Supplementary Figure 1g-i). 206

²⁰⁷ By using WGS data, we constructed cfDNA fragment preferred end maps for

the four cancer types and control (Figure 2c-d). Preferred end enrichment was 208 observed in intron and open chromatin regions (Supplementary Figure 2). 209 210 Additionally, we noticed that for every cancer type, both U and D preferred ends were more enriched in health-specific open chromatin regions than in 211 cancer-specific open chromatin regions (Supplementary Figure 2). One 212 213 explanation was the nucleosome pattern in the health-specific open chromatin regions was unstable, as the enrichment of both U and D preferred ends in 214 215 healthy cfDNA were significantly lower than cancer patients (Mann–Whitney, P<0.05). Another was the heterogeneous nature of cfDNA, as cfDNA can be 216 derived from various tissue sources. To search for preferred ends specific to 217 each of the four cancer types, we determined genome-wide cancer-specific 218 preferred U/D end coordinates and health-specific preferred U/D end 219 coordinates based on the Youden index. In total, we got 81,083 D and 80,505 220

221 U preferred ends. Intriguingly, compared with all repeatable ends (preferred ends appeared in at least two samples in each cancer type), cancer-specific 222 223 preferred U/D ends were no longer enriched in health-specific open chromatin regions, but they showed enrichment in promoter regions and cancer-specific 224 225 open chromatin regions (Figure 2e-f, Supplementary Figure 3-4). Furthermore, health-specific preferred U/D ends showed enrichment in intergenic regions 226 and were not found in promoters, as cancer-specific preferred ends were 227 (Supplementary Figure 4). To cross check the cancer specificity of the 228 229 preferred ends in low-pass WGS data, we calculated a ratio, which was the ratio of detected cancer-specific preferred end numbers to health-specific 230 preferred end numbers for each sample. The results showed significantly 231 higher ratios of cancer-specific to health-specific preferred U/D ends in the 232 three cancer types (except GBM) than in the control (Figure 2g, 233 Mann–Whitney, p value<0.001). In addition to the difference between cancer 234 and health, statistically significant differences in the ratios between cancer 235 types were observed (Supplementary Figure 5). 236

237

238 Fragmentomic features in 5hmC data

We checked the batch effect of the 5hmC sequencing data based on size profile and coverage profile, and found no clustering of samples according to cancer types (Supplementary Figure 6).

242 The cfDNA captured by 5hmC marks comprised a subset of DNA in whole

plasma. In theory, one sequenced fragment should contain at least one 5hmC. 243 Hence, we investigated the consistency of the size profile between WGS data 244 245 and 5hmC sequencing data. According to the results, fragment sizes of 5hmC cfDNA in the four cancer samples were significantly shorter than those in the 246 control (Mann–Whitney test, p value<0.001, pan-cancer Median N50 = 174 bp, 247 control Median N50 = 210 bp), though the difference between cancer and 248 control groups in WGS data was only observed for HCC and LUAD (Figure 249 3a). Next, we confirmed that the cfDNA size profile of the 5hmC sequencing 250 251 data was consistent with that of the WGS data in the size range of 0-220 bp (Figure 3b), with both data types displaying a dominant peak at 167 bp 252 (Supplementary Figure 7a). A difference between the cancer and control 253 254 samples in both short (100-150 bp) and long (151-220 bp) cfDNA fragments (7) was found in the 5hmC data, as in the WGS data (Supplementary Figure 255 7b-c), though the 5hmC data showed a unique secondary peak at ~320 bp 256 (Figure 3c, Supplementary Figure 7a). In the WGS data, cfDNA fragments 257 ranging from 0-220 bp accounted for nearly 100% of fragments ranging from 258 0-800 bp, but the percent value of the same size interval in the 5hmC 259 sequencing data ranged from 31.8% to 99.6% (Figure 3f). Thus, we confirmed 260 the presence of ultra-long fragments (221-500 bp) in the 5hmC cfDNA. We 261 suspected that the existence of ultra-long fragments may have resulted from 262 the capture-based technique applied in 5hmC sequencing, 263 and we consequently assessed the relationship between fragment length and 5hmC 264

peak number. In ultra-long fragments, the percentage of fragments with more 265 than one peak was significantly higher than that in non-ultra-long fragments 266 (0-220 bp) (Mann-Whitney test, p value<0.001) (Figure 3d). Indeed, 267 enrichment analysis of the ultra-long fragments with more than one peak 268 revealed primary enrichment at CpG islands (Supplementary Figure 7d). This 269 indicated that the capture-based sequencing technique is more likely than 270 WGS to capture ultra-long cfDNA fragments. Another finding was that in all 271 cancer types, the percentage of ultra-long fragments was significantly lower 272 273 than that in the control (Supplementary Figure 7b). We examined the percentage of ultra-long fragments with more than one 5hmC peak in the four 274 cancer types and the control and found that the value was significantly lower 275 in the former (Figure 3e). As reported by a recent study, the stability of 276 circulating DNA derives mostly from the nucleosome structure (21). We hence 277 evaluated the proportion of 5hmC cfDNA fragments at 146 bp, 166 bp, 312 bp, 278 279 and 332 bp, which indicates the length of one mono-nucleosome and one dinucleosome (plus the linker size). The result showed a higher proportion of 280 mono-nucleosome-sized cfDNA fragments and a lower proportion of di-281 nucleosome-sized cfDNA fragments in the cancer group than in the control 282 (Supplementary Figure 7e). 283

We also explored the cancer specificity of the pan-cancer preferred end set built upon WGS data in the 5hmC sequencing data by calculating the ratio of cancer-specific to health-specific preferred U/D ends. The ratios were

significantly increased in cfDNA samples in three cancer types (except PDAC) 287 compared with healthy individuals (Figure 3g, Supplementary Figure 5c-d, 288 Mann-Whitney, p value<0.005). Coverage profiles in 5hmC sequencing data 289 also showed a higher variance than WGS data in samples from control and 290 cancer (Supplementary Figure 8). We explored the coverage patterns of short, 291 long, and ultra-long cfDNA fragments at the 5 Mb bins of the genome. Similar 292 to WGS data (Supplementary Material), cancer patients exhibited multiple 293 unstable genomic regions: coverage of short/long/ultra-long fragments was 294 295 inconsistent in cancer patients but consistent in controls (Figure 4a, Supplementary Figure 9a). We analyzed the genome-wide coverage profile 296 correlation of cancer patients to controls in the 5hmC data and found that the 297 298 coverage profiles were consistent among controls and that the correlation value in cancer patients was significantly lower (Wilcoxon rank-sum test, p 299 value<0.005). Importantly, the largest difference between cancer samples and 300 controls was found with respect to the coverage profile of ultra-long fragments 301 (Wilcoxon rank-sum test, p value<0.005) (Figure 4b). 302

303

304 Genomic distribution of 5hmC in cancer and control cohorts

To gain an understanding of the genomic regions associated with hydroxymethylation in cfDNA, we first determined 5hmC-enriched loci in the control and four types of cancer, which were detected as peaks via MACS2 (22). 5hmC-modified regions among samples were compared in 1 kb bins on

the reference genome, and thus a consensus list of the absence and 309 presence of 5hmC-modified peaks among the samples was obtained. 5hmC 310 311 signatures were enriched over genic features, most significantly in the promoter, 5'UTRs, 3'UTRs, exons, and transcription end sites (TESs) (Figure 312 313 5a). Comparison of 5hmC signature enrichment in the four types of cancers with the control revealed significant differences, whereby 5hmC peak 314 enrichment was lower in all cancer types than in the control over promoters, 315 5'UTRs, exons, and CpG islands (Wilcoxon rank-sum test, p value<0.001) 316 317 (Figure 5a).

Differential analysis of 5hmC-modified peaks between cancer patients and 318 controls by Fisher's exact test detected 1,010, 6,395 and 773 differentially 319 320 modified 5hmC peaks in GBM, HCC and PDAC, respectively. Moreover, this pan-cancer 5hmC peak set was able to separate most cancer samples from 321 the control (Figure 5b). KEGG enrichment analysis of genes located within the 322 323 differentially modified 5hmC peaks revealed associated oncogenesis pathways for each cancer type (Supplementary Figure 10), such as the 324 neurotrophin signaling pathway and platelet activation for GBM (23, 24), the 325 MAPK signaling pathway and focal adhesion for HCC (25, 26), and the FoxO 326 signaling pathway and insulin signaling pathway for PDAC (27). 327

328

Cancer detection by combining fragmentomic features and 5hmC signatures using
 5hmC sequencing data

331	As 5hmC sequencing data retain both 5hmC signatures and fragmentomic
332	information, it is expected to theoretically show high sensitivity and specificity
333	in cancer detection. A receiver operator characteristic (ROC) curve was used
334	to evaluate the performance of the classifier. In total, 53 5hmC signatures
335	were selected, and the AUC of the classifier was 0.876 in the validation set
336	(sensitivity = 82.26%, specificity = 82.35%) and 0.872 in the test set
337	(sensitivity = 81.97%, specificity = 82.35%) (Figure 5d-e, Supplementary
338	Figure 11 c).

Overall, the fragmentomic information in the 5hmC sequencing data 339 performed well. Using the size profile as features, 40 were selected to build a 340 pan-cancer prediction model, among which 32 features were of the ultra-long 341 342 and longer size range (>= 220 bp) (Figure 5c). The AUC value was 0.981 in the validation set (sensitivity = 93.55%, specificity = 94.12%) and 0.882 in the 343 test set (sensitivity = 71.43%, specificity = 88.24%) (Figure 5d-e). Using 344 preferred ends as features, 37 were selected to build a pan-cancer prediction 345 model. ROC analysis showed an AUC value of 0.940 in the validation set 346 (sensitivity = 83.87%, specificity = 94.12%) and 0.899 in the test set 347 (sensitivity = 73.02%, specificity = 94.12%) (Figure 5d-e, Supplementary 348 Figure 11a). Regarding the coverage profile of the 5hmC sequencing data, we 349 explored the genomic distribution of short, long, and ultra-long cfDNA 350 fragments on a 100 kb window of the genome. This model achieved an AUC 351 value of 0.946 in the validation set (sensitivity = 80.65%, specificity = 94.12%) 352

and 0.882 in the test set (sensitivity = 71.43%, specificity = 88.24%) (Figure
5d-e).

355 Finally, we constructed an integrated cancer screening model by combining fragmentomic features and 5hmC signatures in 5hmC data. Sixty-three 356 features, including 10 5hmC signatures, 24 size profile features, 21 coverage 357 profile features, and 8 preferred end features, were selected to build the pan-358 cancer prediction model. The AUC value was 0.927 in the validation set 359 (sensitivity = 93.44%, specificity = 88.24%) and 0.920 in the test set 360 361 (sensitivity = 88.52%, specificity = 82.35%) (Figure 5d-e). The performance of single-cancer detection is depicted in Supplementary Figure 11b. 362

With data from the four cancer types, we further explored tissue-of-origin prediction. The performance of the random forest model with 5hmC fragmentomic features and the integrated model are shown in Supplementary Tables 3-5.

367

368 **Discussion**

cfDNA 5hmC signatures are reported to have potential for cancer detection, such as in lung cancer (28), hepatocellular carcinoma (28), colon cancer (29), gastric cancer (29), and pancreatic cancer (17, 30). The large-scale cohort of 5hmC sequencing data utilized in this study additionally shows the potential of cfDNA 5hmC signatures for cancer detection in glioblastoma and pan-cancer. Besides, we further examined the classification effect using fragmentomic

features in tissue-of-origin prediction of cancer. Most importantly, we explored 375 cfDNA fragmentomic information in 5hmC sequencing data and found that it is 376 377 possible to detect cfDNA hydroxymethylation and fragmentomic markers simultaneously. The size profile, preferred end, and coverage profile in 5hmC 378 sequencing data showed a large difference between cancer patients and 379 healthy individuals. The integrated model covering the 5hmC signature, size 380 profile, preferred end, and coverage profile contained more information than 381 the model with the 5hmC signature alone, as revealed by higher sensitivity 382 383 and specificity in pan-cancer detection.

384

Previous work reported a difference in the size profile between specific cancer 385 386 patients or pan-cancer patients and healthy controls (12, 13, 31), though the analyses only focused on lengths less than 320 bp. Here, we characterized 387 ultra-long fragments in 5hmC sequencing data and identified their size and 388 coverage profile aberrations in cancer samples. Although it is commonly 389 acknowledged that ctDNA tends to be more fragmented than cfDNA in normal 390 391 tissue (12, 31), we found an even larger size difference in cfDNA between cancer samples and controls in 5hmC data owing to ultra-long fragments. 392 Hence, we further expanded the size range to 800 bp by setting 80 adjacent, 393 non-overlapping bins to capture 10-bp interval signals (32). The coverage 394 profile is another major improvement we made with regard to previously 395 reported models based on fragmentation information (9). Indeed, only the 396

ratio of short to long cfDNA fragments has been analyzed in previously 397 reported models, with information on cfDNA ultra-long fragments being lost. 398 399 Our results suggest that the coverage distribution of size-selected cfDNA fragments in ultra-long fragments has the strongest inconsistency between 400 401 control and cancer samples (Figure 4a). Of note, the classification models selected the majority of features from ultra-long and longer size range 402 fragments in the coverage profile model (34/42), size profile model (32/40) 403 (Figure 5c, Supplementary Figure 12c) and integrated model (41/63), 404 405 suggesting that ultra-long features are a major contribution to the classification model. Further research on the mechanism of the aberration of 406 the captured ultra-long fragments in cfDNA from cancer samples may provide 407 408 insights into the fundamental biological properties of plasma cfDNA from patients with cancer. 409

410

411 In our work, we constructed a single-base resolution preferred end set based on a set of WGS data. Our results indicated that the preferred end set had 412 good generalization ability because the differences between patients with 413 cancer and controls were also shown in other independent data sets, e.g., 414 low-pass WGS data and 5hmC data. It has been reported that referred end 415 coordinates differ between HCC patients and healthy people at a whole-416 genome scale (10). Apart from HCC, our results suggest the cfDNA preferred 417 end can also be used to distinguish patients with GBM, LUAD, and PDAC. We 418

419 herein discussed the distribution of cfDNA preferred ends in genomic features, open chromatin regions, and nucleosome structure. The distribution of cfDNA 420 421 preferred ends can reflect nucleosome positioning, though there is no apparent preference in genomic features and open chromatin regions. 422 423 consistent with previous work (15). However, informative cancer-specific preferred ends were found to be significantly enriched in promoter and far 424 health-specific chromatin 425 away from open regions (Figure 2e-f. Supplementary Figure 4c-d, Supplementary Figure 5c-d). 426

427

In summary, our results indicate that cfDNA fragmentomic analysis with a 428 nonlinear classification algorithm using low-pass 5hmC sequencing data may 429 430 provide a simple, low-cost, and highly effective method for cancer detection. The success of utilizing ultra-long fragments in 5hmC sequencing data as 431 biomarkers indicates the potential of other capture-based sequencing 432 433 approaches, such as cell-free methylated DNA immunoprecipitation and highthroughput sequencing (cfMeDIP-seg) (33), in the diagnosis of cancer. 434 435 Notably, third-generation sequencing methods are expected to detect long cfDNA molecules, which have been utilized in non-invasive prenatal testing 436 (34). As third-generation sequencing methods can detect sequence context 437 and epigenetic modification at the same time, they can provide simultaneous 438 analysis of genome-wide genetic, fragmentomic and epigenetic detection (35, 439 36). Moreover, standardized and effective workflows for analysis of cfDNA 440

441 fragments need to be developed.

442

443 Data availability

444 The 5hmC sequencing data for controls, LUAD, HCC and PDAC were publicly

445 available as described in the Method section. The 5hmC sequencing data for

446 GBM can be accessed at https://ngdc.cncb.ac.cn/gsa-human/s/15mfS980.

447

448 Acknowledgments

449 This study was supported by the 1.3.5 project for disciplines of excellence, West China Hospital, Sichuan University (ZYYC20006) to D. Xie; Thousand 450 Talents Program of the West China Hospital (0040205401F58) to D. Xie; 451 452 Sichuan Provincial Foundation of Science and Technology (2020YFS0051) to D. Xie, (2017SZ0006) to Y. Liu; Clinical Research Innovation Project, West 453 China Hospital, Sichuan University (19HXCX009) to Y. Liu; the San Hang 454 455 Program of the Second Military Medical University, Medical basic research project of the First Affiliated Hospital, the Second Military Medical University 456 (2021JCMS16) to W. Zhang: the Science and technology project of Sichuan 457 Province (2021YFS0109) to A. Li; Post-Doctor Research Project, West China 458 459 Hospital, Sichuan University (20HXBH035) to S. Zhang; the high quality development of Guang 'an People's Hospital (21FZ003) to A. Wei. 460

461 **Author Contributions**

462 Conception and design: Dan Xie, Zhidong Zhang, Xuenan Pi and Lin Xia.

463	Colle	ction and assembly of data: Jun Zhang, Xiaoqin Yan, Shuxing Zhang,
464	Ailin	Wei, Jingfeng Liu, Ang Li, Xiaolong Liu, Wei Zhang and Yanhui Liu. Data
465	analy	sis and interpretation: Zhidong Zhang, Xuenan Pi, Chang Gao, Xinlei Hu,
466	Xiyue	Yan and Yuer Guo. Manuscript writing: Zhidong Zhang and Xuenan Pi.
467	All au	thors read and approved the final version of the manuscript.
468	Com	peting interests
469	The a	authors declare no competing interests.
470		
471	Refe	rences
472	1.	Sung, H., et al. Global cancer statistics 2020: GLOBOCAN estimates of
473		incidence and mortality worldwide for 36 cancers in 185 countries. CA
474		<i>Cancer. J. Clin.</i> 0 , 1–41 (2021).
475	2.	Heitzer, E., Haque, I. S., Roberts, C. E. S. & Speicher, M. R. Current
476		and future perspectives of liquid biopsies in genomics-driven oncology.
477		<i>Nat. Rev. Genet.</i> 20 , 71–88 (2019).
478	3.	Siravegna, G., Marsoni, S., Siena, S. & Bardelli, A. Integrating liquid
479		biopsies into the management of cancer. Nat. Rev. Clin. Oncol. 14,
480		531–48 (2017).
481	4.	Jahr, S., et al. DNA fragments in the blood plasma of cancer patients:
482		quantitations and evidence for their origin from apoptotic and necrotic
483		cells. <i>Cancer. Res.</i> 61 , 1659-65 (2001).
484	5.	Serpas, L., et al. Dnase113 deletion causes aberrations in length and

485	end-motif frequencies i	n nlasma DNA	Proc Natl	Acad Scill S A
405	chu-moui nequencies i	$1 \mu a \sin a \nu \sin \pi$.	. <i>i i 00. ivali.</i>	

486 **116**, 641-49 (2019).

- 487 6. Han, D. S. C., et al. The Biology of Cell-free DNA Fragmentation and
- the Roles of DNASE1, DNASE1L3, and DFFB. Am. J. Hum. Genet. 106,
- 489 202-14 (2020).
- 490 7. Cristiano, S., et al. Genome-wide cell-free DNA fragmentation in
- 491 patients with cancer. *Nature*. **570**, 385-89 (2019).
- 492 8. Snyder, M. W., Kircher M., Hill A. J., Daza R. M. & Shendure J. Cell-free
- 493 DNA Comprises an In Vivo Nucleosome Footprint that Informs Its
- 494 Tissues-Of-Origin. *Cell.* **164**, 57-68 (2016).
- 495 9. Chan, K. C., et al. Second generation noninvasive fetal genome
- 496 analysis reveals de novo mutations, single-base parental inheritance,
- 497 and preferred DNA ends. *Proc. Natl. Acad. Sci. U. S. A.* **113**, E8159-68
- 498 **(2016)**.
- 499 10. Jiang, P., et al. Preferred end coordinates and somatic variants as
- signatures of circulating tumor DNA associated with hepatocellular
- 501 carcinoma. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E10925-33 (2018).
- 502 11. Jiang, P., et al. Plasma DNA End-Motif Profiling as a Fragmentomic
- Marker in Cancer, Pregnancy, and Transplantation. *Cancer. Discov.* 10,
 664-73 (2020).
- Mouliere, F., et al. Enhanced detection of circulating tumor DNA by
 fragment size analysis. *Sci. Transl. Med.* **10**, eaat4921 (2018).

507	13.	Jiang, P., et al. Lengthening and shortening of plasma DNA in
508		hepatocellular carcinoma patients. Proc. Natl. Acad. Sci. U. S. A. 112,
509		E1317-25 (2015).
510	14.	Ulz, P., et al. Inferring expressed genes by whole-genome sequencing
511		of plasma DNA. <i>Nat. Genet.</i> 48 , 1273-8 (2016).
512	15.	Sun, K., et al. Orientation-aware plasma cell-free DNA fragmentation
513		analysis in open chromatin regions informs tissue of origin. Genome
514		Res. 29 , 418-27 (2019).
515	16.	He, B., et al. Tissue-specific 5-hydroxymethylcytosine landscape of the
516		human genome. <i>Nat. Commun.</i> 12 , 4249 (2021).
517	17.	Chen, L., et al. Genome-scale profiling of circulating cell-free DNA
518		signatures for early detection of hepatocellular carcinoma in cirrhotic
519		patients. <i>Cell. Res.</i> 31 , 589-92 (2021).
520	18.	[dataset]* Hu, X., et al. Integrated 5-hydroxymethylcytosine and
521		fragmentation signatures as enhanced biomarkers in lung cancer. Clin.
522		<i>Epigenetics.</i> 14 , 15 (2022).
523	19.	Roadmap Epigenomics Consortium, et al. Integrative analysis of 111
524		reference human epigenomes. Nature. 518, 317-30 (2015).
525	20.	Gaffney, D. J., et al. Controls of nucleosome positioning in the human
526		genome. <i>PLoS. Genet.</i> 8, e1003036 (2012).
527	21.	Sanchez, C., et al. Circulating nuclear DNA structural features, origins,
528		and complete size profile revealed by fragmentomics. JCI. Insight. 6,

- e144561 (2021). 529
- Zhang, Y., et al. Model-based analysis of ChIP-Seq (MACS). Genome 22. 530 531 Biol. 9, R137 (2008).
- Lawn, S., et al. Neurotrophin signaling via TrkB and TrkC receptors 23. 532
- 533 promotes the growth of brain tumor-initiating cells. J. Biol. Chem. 290,
- 3814-24 (2015). 534

549

- Nolte, I., Przibylla, H., Bostel, T., Groden, C. & Brockmann, M. A. 535 24.
- 536 Tumor-platelet interactions: glioblastoma growth is accompanied by
- 537 increasing platelet counts. Clin. Neurol. Neurosurg. 110, 339-42 (2008).
- 25. Moon, H. & Ro, S. W. MAPK/ERK Signaling Pathway in Hepatocellular 538

Carcinoma. Cancers (Basel). 13, 3026 (2021). 539

- 540 26. Gnani, D., et al. Focal adhesion kinase depletion reduces human
- hepatocellular carcinoma growth by repressing enhancer of zeste 541

homolog 2. Cell. Death. Differ. 24, 889-902 (2017). 542

543 27. Li, J., et al. Knockdown of FOXO3a induces epithelial-mesenchymal

transition and promotes metastasis of pancreatic ductal 544

- 545 adenocarcinoma by activation of the β -catenin/TCF4 pathway through
- SPRY2. J. Exp. Clin. Cancer. Res. 38, 38 (2019). 546
- 547 28. Song, C. X., et al. 5-Hydroxymethylcytosine signatures in cell-free DNA
- provide information about tumor types and stages. Cell. Res. 27, 1231-548 42 (2017).
- 29. Li, W., et al. 5-Hydroxymethylcytosine signatures in circulating cell-free 550

551		DNA as diagnostic biomarkers for human cancers. Cell. Res. 27, 1243-
552		57 (2017).
553	30.	Guler, G. D., et al. Detection of early stage pancreatic cancer using 5-
554		hydroxymethylcytosine signatures in circulating cell free DNA. Nat.
555		Commun. 11, 5270 (2020).
556	31.	Lam, W. K. J., et al. Sequencing-based counting and size profiling of
557		plasma Epstein-Barr virus DNA enhance population screening of
558		nasopharyngeal carcinoma. Proc. Natl. Acad. Sci. U. S. A. 115, E5115-
559		24 (2018).
560	32.	Lo, Y. M., et al. Maternal plasma DNA sequencing reveals the genome-
561		wide genetic and mutational profile of the fetus. Sci. Transl. Med. 2,
562		61ra91 (2010).
563	33.	Shen, S. Y., et al. Sensitive tumour detection and classification using
564		plasma cell-free DNA methylomes. Nature. 563, 579-83 (2018).
565	34.	Yu, S. C. Y., et al. Single-molecule sequencing reveals a large
566		population of long cell-free DNA molecules in maternal plasma. Proc.
567		Natl. Acad. Sci. U. S. A. 118 , e2114937118 (2021).
568	35.	Tse, O. Y. O., et al. Genome-wide detection of cytosine methylation by
569		single molecule real-time sequencing. Proc. Natl. Acad. Sci. U. S. A.
570		118 , e2019768118 (2021).
571	36.	Liu, Q., et al. Detection of DNA base modifications by deep recurrent
572		neural network on Oxford Nanopore sequencing data. Nat. Commun. 10,

- **2449 (2019)**.
- 574 37. [dataset]* Cao, F., et al. Integrated epigenetic biomarkers in circulating
- 575 cell-free DNA as a robust classifier for pancreatic cancer. *Clin.*
- *Epigenetics.* **12**, 112 (2020).
- 577 38. [dataset]* Cai, Z., et al. Liquid biopsy by combining 5-
- 578 hydroxymethylcytosine signatures of plasma cell-free DNA and protein
- 579 biomarkers for diagnosis and prognosis of hepatocellular carcinoma.
- 580 ESMO Open. 6, 100021 (2021).

- . . .



Figure 1 Overview of the study design.

cfDNA from cancer patients and controls were sequenced with WGS.

601	Available 5hmC sequencing data from the same four cancer types and
602	controls were combined with WGS data to get the accurate fragmentomic
603	information and upon it a machine learning model was built to distinguish
604	healthy controls and cancer patients.





Figure 2 Genome-wide cfDNA fragment preferred end map construction. 606

cfDNA preferred end signals in a nucleosome array region (chr12: 34517269-607 34519122) in healthy (a) and HCC (b) subjects. Brown dots at the bottom 608

609	represent the predicted nucleosome center loci reported in a previous study (8,
610	22). (c) Venn diagram showing the intersection of the preferred U and D ends
611	between the cancer cohort and the control cohort. (d) Circle diagram showing
612	the density of preferred ends in 1 Mb windows in four types of cancer and in
613	healthy controls (the order of sample sources from outer ring to inner ring are
614	control, GBM, HCC, LUAD and PDAC). (e) Enrichment of preferred U/D ends
615	in promoter regions. (f) Enrichment of preferred U/D ends in health-specific
616	open chromatin regions. (g) Boxplots showing the ratios of the cancer-specific
617	to health-specific preferred U/D ends in low-pass WGS data for the four
618	cancer types (Wilcoxon rank-sum test).



619

Figure 3 cfDNA fragmentomic information in 5hmC sequencing data. 620

(a) Violin plots showing fragment size profile comparison between cancer 621 622 samples and controls in terms of N50 (Mann-Whitney test). (b) tSNE plot showing 5hmC sequencing data and WGS data by fragments below 220 bp. 623 (c) The cfDNA size profile of 5hmC sequencing data. (d) Comparison of the 624 percentage of fragments with more than one 5hmC peak between ultra-long 625 and non-ultra-long fragments in every sample. The dashed line indicates the 626 diagonal line. (e) Violin plots showing the percentage of ultra-long fragments 627 with more than one 5hmC peak in cancer samples and controls 628 (Mann-Whitney test). (f) Boxplots showing the percentage of fragments 629

- ranging from 0-220 bp and 221-500 bp among fragments ranging from 0-800
- 631 bp in WGS data and 5hmC sequencing data (Wilcoxon rank-sum test). (g)
- 632 Boxplots showing the ratios of the cancer-specific to health-specific preferred
- 633 U/D ends in each of the four cancer types (Wilcoxon rank-sum test).







of cfDNA fragments in 5hmC sequencing data. 636

(a) Genome-wide coverage profile of short, long, and ultra-long cfDNA 637 fragments; color indicates sample-wise correlation to median health. (b) 638

- Coverage profile correlation to median healthy in control and four types of 639
- cancer in short, long, and ultra-long cfNDA fragments (Wilcoxon rank-sum test, 640

p value<0.005). 641



Figure 5 Performance of cancer detection combining fragmentomic 643 features and 5hmC signatures using 5hmC sequencing data. 644

(a) Enrichment of 5hmC in genomic features in control and four types of 645 cancer (Wilcoxon rank-sum test, p value<0.001). (b) Heatmap of differentially 646 modified 5hmC peaks. (c) Selected size profile features in the size profile 647 model. (d) ROC curves for the validation set and test set in pan-cancer 648 diagnosis. (e) Performance evaluation of validation set and test set in pan-649 cancer diagnosis. 650

651