

1 **Multimodal analysis of methylomics and fragmentomics in plasma cell-free**
2 **DNA for multi-cancer early detection and localization**

3 Van Thien Chi Nguyen^{1,2#}, Trong Hieu Nguyen^{1,2#}, Nhu Nhat Tan Doan^{1,2}, Thi Mong Quynh
4 Pham^{1,2}, Giang Thi Huong Nguyen^{1,2}, Thanh Dat Nguyen^{1,2}, Thuy Thi Thu Tran^{1,2}, Duy Long Vo³,
5 Thanh Hai Phan⁴, Thanh Xuan Jasmine⁴, Van Chu Nguyen^{5,6}, Huu Thinh Nguyen³, Trieu Vu
6 Nguyen⁷, Thi Hue Hanh Nguyen^{1,2}, Le Anh Khoa Huynh^{1,8}, Trung Hieu Tran^{1,2}, Quang Thong
7 Dang³, Thuy Nguyen Doan³, Anh Minh Tran³, Viet Hai Nguyen³, Vu Tuan Anh Nguyen³, Le
8 Minh Quoc Ho³, Quang Dat Tran³, Thi Thu Thuy Pham⁴, Tan Dat Ho⁴, Bao Toan Nguyen⁴,
9 Thanh Nhan Vo Nguyen⁴, Thanh Dang Nguyen⁴, Dung Thai Bieu Phu⁴, Boi Hoan Huu Phan⁴, Thi
10 Loan Vo⁴, Thi Huong Thoang Nai⁴, Thuy Trang Tran⁴, My Hoang Truong⁴, Ngan Chau Tran⁴,
11 Trung Kien Le³, Thanh Huong Thi Tran^{5,6}, Minh Long Duong^{5,6}, Hoai Phuong Thi Bach^{5,6}, Van
12 Vu Kim^{5,6}, The Anh Pham^{5,6}, Duc Huy Tran³, Trinh Ngoc An Le³, Truong Vinh Ngoc Pham³,
13 Minh Triet Le³, Dac Ho Vo^{1,2}, Thi Minh Thu Tran^{1,2}, Minh Nguyen Nguyen^{1,2}, Thi Tuong Vi
14 Van^{1,2}, Anh Nhu Nguyen^{1,2}, Thi Trang Tran^{1,2}, Vu Uyen Tran^{1,2}, Minh Phong Le^{1,2}, Thi Thanh
15 Do^{1,2}, Thi Van Phan^{1,2}, Luu Hong Dang Nguyen^{1,2}, Duy Sinh Nguyen^{1,2}, Van Thinh Cao⁹, Thanh
16 Thuy Thi Do², Dinh Kiet Truong², Hung Sang Tang^{1,2}, Hoa Giang^{1,2}, Hoai Nghia Nguyen^{1,2}, Minh
17 Duy Phan^{1,2,*}, Le Son Tran^{1,2,*}

18

19 ¹Gene Solutions, Ho Chi Minh City, Vietnam

20 ²Medical Genetics Institute, Ho Chi Minh City, Vietnam

21 ³University Medical Center, Ho Chi Minh City, Vietnam

22 ⁴MEDIC Medical Center, Ho Chi Minh City, Vietnam

23 ⁵National Cancer Hospital, Hanoi, Vietnam

24 ⁶Hanoi Medical University, Hanoi, Vietnam

25 ⁷Thu Duc City Hospital, Ho Chi Minh City, Vietnam

26 ⁸Department of Biostatistics, Virginia Commonwealth University, School of Medicine,
27 Richmond, VA, USA

28 ⁹Pham Ngoc Thach University of Medicine, Ho Chi Minh City, Vietnam

29

30

31

32

33

34

35

36 # Van Thien Chi Nguyen and Trong-Hieu Nguyen contributed equally to this study.

37 *Correspondence: pmduy@yahoo.com; leson1808@gmail.com

38

39

40 **Key words:** liquid biopsy, multimodal analysis, methylation, fragment length, cell-free DNA,
41 circulating tumor DNA, multicancer early detection, tissue of origin, machine learning, graph
42 convolutional neural network.

43

44

45 **Abstract**

46 Despite their promise, circulating tumor DNA (ctDNA)-based assays for multi-cancer early
47 detection face challenges in test performance, due mostly to the limited abundance of ctDNA and
48 its inherent variability. To address these challenges, published assays to date demanded a very
49 high-depth sequencing, resulting in an elevated price of test. Herein, we developed a multimodal
50 assay called SPOT-MAS (Screening for the Presence Of Tumor by Methylation And Size) to
51 simultaneously profile methylomics, fragmentomics, copy number, and end motifs in a single
52 workflow using targeted and shallow genome-wide sequencing (~0.55X) of cell-free DNA. We
53 applied SPOT-MAS to 738 nonmetastatic patients with breast, colorectal, gastric, lung and liver
54 cancer, and 1,550 healthy controls. We then employed machine learning to extract multiple
55 cancer and tissue-specific signatures for detecting and locating cancer. SPOT-MAS successfully
56 detected the five cancer types with a sensitivity of 72.4% at 97.0% specificity. The sensitivities
57 for detecting early-stage cancers were 62.3% and 73.9% for stage I and II, respectively,
58 increasing to 88.3% for nonmetastatic stage IIIA. For tumor-of-origin, our assay achieved an
59 accuracy of 0.7. Our study demonstrates comparable performance to other ctDNA-based assays
60 while requiring significantly lower sequencing depth, making it economically feasible for
61 population-wide screening.

62 **Introduction**

63 The incidence of cancer-related morbidity and mortality is rapidly increasing globally, and
64 accounted for nearly one fifth of all deaths in 2020 (1). High-cost treatment is a significant
65 financial burden for cancer patients, with almost 286 billion dollars in 2021 and an increase of
66 8.2% to 581 billion dollars in 2030. In Vietnam, GLOBOCAN 2020 reported over 182,500
67 newly diagnosed cases and 122,690 cancer-related deaths (1). Among these, liver (14.5%), lung
68 (14.4%), breast (11.8%), gastric (9.8%), and colorectal cancer (9%) are the five most common
69 types. Up to 80% of cancer patients in Vietnam were diagnosed at stage III or stage IV, resulting
70 in a high rate of 1-year mortality (25%) and a low 5-year survival rate compared to other
71 countries (2). Diagnostic delays are associated with a lower chance of survival, greater
72 treatment-associated problems, and higher costs (3). Cancer detection at earlier stages can
73 improve the opportunity to control cancer progression, increase the patient survival rate, and
74 lower medical expenses (4).

75 Most current early cancer screening assays have limitations such as invasiveness, low
76 accessibility, and high false positive rates when used sequentially, resulting in overdiagnosis and
77 overtreatment. Multi-cancer early detection (MCED) tests can potentially overcome these
78 challenges by simultaneously detecting multiple cancer types from a single test (5). Liquid
79 biopsy, an emerging non-invasive approach for MCED, can capture a wide range of tumor
80 features, including cell free DNA (cfDNA), circulating tumor DNA (ctDNA), exosomes,
81 proteins, mRNA, and metabolites (6, 7). Among them, ctDNA has become a promising
82 biomarker for detecting early-stage cancers because it is a carrier of genetic and epigenetic
83 modifications from cancer-derived DNA (8). Indeed, ctDNA detection has demonstrated several
84 advantages in non-invasive diagnostic, prognostic, and monitoring of cancer patients during and

85 after treatment (9, 10). Furthermore, ctDNA carrying tumor-specific alterations could be used to
86 identify the corresponding unknown primary cancer and tumor localization.

87 In recent years, there has been considerable interest in exploring the potential of ctDNA
88 alterations for early detection of cancer and localization of the tissue of origin (TOO) (11, 12).

89 One such approach is the PanSeer test, which uses 477 differentially methylated regions (DMRs)
90 in ctDNA to detect five different types of cancer up to four years prior to conventional diagnosis
91 (13). However, this assay is limited in its ability to determine the TOO, as it only uses
92 methylation regions common to multiple cancers. The DELFI assay employs a genome-wide
93 analysis of ctDNA fragment profiles to increase sensitivity in early detection, but also lacks
94 accuracy in classifying the source of tumor-derived cfDNA (14). Recently, the Galleri test has
95 emerged as a multi-cancer detection assay that analyses more than 100,000 methylation regions
96 in the genome to detect over 50 cancer types and localize the tumor site (15). This approach
97 requires a large-scale target capture panel at very high-depth sequencing (with a depth coverage
98 of 30X), incurring high sequencing costs and limiting the accessibility of this test to the wider
99 population.

100 Despite their great potential, there remain several challenges that these assays must solve to
101 deliver accessible and reliable clinical adoption for the large population, including the low
102 fraction of ctDNA in the blood of early stage cancer patients, the heterogeneity of ctDNA
103 signatures from diverse cancer types, subtypes and stages (16), and the high sequencing depth
104 required. To address these challenges, recent studies have focused on multi-analyte approach -
105 combining genomic and nongenomic features such as methylomics and fragmentomics to
106 increase the detection of ctDNA and accuracy for TOO identification (16-18). Advances in
107 multimodal analysis approaches have led to the development of powerful screening tests that

108 enable high sensitivity and cost-effectiveness. For example, CancerSEEK uses a combined
109 approach of protein biomarkers and genetic alterations to detect and locate the presence of eight
110 types of cancers (19). In this assay, cancer-associated serum proteins play a complementary role
111 in tumor localization as cfDNA mutations are not tissue specific. However, detecting both
112 protein and genetic biomarkers are time-consuming and costly. Thus, the development of future
113 MCED tests should endeavor to deliver a screening approach with high sensitivity, specificity,
114 and TOO identification at cost-effective price to provide better clinical outcomes and treatment
115 opportunities for all cancer patients.

116 In an effort to address the challenges of early cancer detection, we have developed a multimodal
117 approach called SPOT-MAS (Screening for the Presence Of Tumor by DNA Methylation And
118 Size). This assay was previously applied to cohorts of colorectal (20) and breast cancer (under
119 review in *Frontiers in Oncology*) patients and demonstrated ability for early detection of these
120 cancers at high sensitivity across different cancer stages and patient age groups. In this study, we
121 expanded our multimodal approach, SPOT-MAS to comprehensively analyze methylomics,
122 fragmentomics, DNA copy number and end motifs of cfDNA for simultaneously detecting and
123 locating cancer from a single screening test. As proof of concept, we used 2,288 participants,
124 including 738 nonmetastatic patients and 1,550 healthy controls, to train and fully validate this
125 approach on five commonly diagnosed cancers, including breast, gastric, lung, colorectal, and
126 liver cancer. These cancer types accounted for more than 54% of new cancer cases and 57% of
127 cancer death worldwide as well as the most diagnosed cancer types in developing countries (1, 2).
128 By using targeted sequencing, shallow whole genome sequencing (with a depth coverage of
129 0.55X) and innovative machine learning algorithms, we could analyze a large multi-feature
130 datasets of cfDNA for multi-cancer early detection and tumor localization with high sensitivity

131 and cost-effectiveness. Our assay achieved sensitivities of early detection of 72.4% among the
132 five cancer types at specificity of 97.0%, with AUC of 0.95 . Moreover, we could identify the
133 tissue of origin with an accuracy of 0.7 in independent validation cohort using graph
134 convolutional neural network. Thus, SPOT-MAS has the potential to become a universal, simple,
135 and cost-effective approach for early multi-cancer detection in large populations.

136 **Methods**

137 **Patient enrollment**

138 This study recruited 738 cancer patients (223 breast cancer, 159 CRC, 122 liver cancer, 136 lung
139 cancer, 98 gastric cancer) and 1550 healthy subjects. All cancer patients were confirmed to have
140 one of the five cancers analyzed in this study. Cancer stages were determined following
141 guidelines from the American Joint Committee on Cancer and the International Union for Cancer
142 Control (21). Individuals were considered healthy if they had no history of cancer at the time of
143 enrollment and follow-up interviews were conducted by specialized physicians to confirm
144 noncancer status at 6 and 12 months after enrollment. Study subjects were recruited from the
145 University Medical Center, Thu Duc City Hospital, University of Medicine and Pharmacy,
146 Medic Medical Center and Medical Genetics Institute in Ho Chi Minh city, Vietnam, National
147 Cancer Hospital and Hanoi Medical University in Hanoi from May 2019 to December 2022.

148 Written informed consent was obtained from each participant in accordance with the Declaration
149 of Helsinki. This study was approved by the Ethics Committee of the Medic Medical Center,
150 University of Medicine and Pharmacy and Medical Genetics Institute, Ho Chi Minh city,
151 Vietnam. All cancer patients were treatment-naïve at the time of blood sample collection.

152 **Isolation of cfDNA**

153 10 mL of blood was collected from each participant in a Cell-Free DNA BCT tube (Streck,
154 USA). Plasma was collected from blood samples after centrifugation with two rounds ($2,000 \times g$
155 for 10 min and then $16,000 \times g$ for 10 min). The plasma fraction was aliquoted for long-term
156 storage at -80°C . Cell free DNA (cfDNA) was extracted from 1 mL plasma aliquots using the
157 MagMAX Cell-Free DNA Isolation kit (ThermoFisher, USA), according to the manufacturer's
158 instructions. Extracted cfDNA was quantified by the QuantiFluor dsDNA system (Promega,
159 USA).

160 **Bisulfite conversion and library preparation**

161 According to the manufacturer's instructions, bisulfite conversion and cfDNA purification were
162 prepared by EZ DNA Methylation-Gold Kit (Zymo research, D5006, USA). DNA library was
163 prepared from bisulfite-converted DNA samples using xGenTM Methyl-Seq DNA Library Prep
164 Kit (Integrated DNA Technologies, 10009824, USA) with AdaptaseTM technology, according to
165 the manufacture's instructions. The QuantiFlour dsDNA system (Promega, USA) was used to
166 analyse the concentration of DNA.

167 **Target region capture, whole genome hybridization & sequencing**

168 DNA from library products were pooled equally, hybridized and captured using The XGen
169 hybridization and wash kit (Integrated DNA Technologies, 1072281, USA), together with our
170 customized panel of xGen Lockdown Probes including 450 regions across 18,000 CpG sites
171 (Integrated DNA Technologies, USA). The construction of panel was built as previously
172 described (13, 20, 22). After hybridization, the flow-through product was concentrated using
173 SpeedVac (N-Biotek, NB-503CIR, Korea) at 65°C . The samples were then added with the
174 hybridization master mixture (hybridization buffer, hybridization enhancer and H_2O) and
175 denatured. Biotinylated P5 and P7 probes (P5-biotin: /5Biosg/AATGATACGGCGACCACCGA,

176 P7-biotin: /5Biosg/CAAGCAGA AGACGGCATACGAGAT) on streptavidin magnetic beads
177 (Invitrogen, CA, USA) were hybridized with the single-stranded DNA. The captured DNA
178 products were amplified by a PCR reaction with free P5 and P7 primers (P5 primer:
179 AATGATACGGCGACCACCGA, P7 primer: CAAGCAGAAGACGGCATACGA). The
180 concentrations of DNA libraries were determined using the QuantiFluor dsDNA system
181 (Promega, USA). Both target and flow-through fraction were sequenced on the DNBSEQ-G400
182 DNA system (MGI Tech, Shenzhen, China) with 100-bp paired-end reads at a sequencing depth
183 of 20 million reads per fraction. Data was demultiplexed by bcl2fastq (Illumina, CA, USA).
184 FASTQ files were then examined using FastQC v. 0.11.9 and MultiQC v. 1.12.

185 **Targeted methylation analysis (TM)**

186 All paired-end reads were processed by Trimmomatic v 0.32 with the option HEADCROP. The
187 trimmed reads were then aligned by Bismark v. 0.22.3. Deduplication and sorting of BAM files
188 were conducted using Samtools v. 1.15. Reads falling into our 450 target regions were filtered
189 using Bedtools v. 2.28. Methylation calling was performed using Bismark methylation extractor
190 (20). Briefly, methylation ratio was measured for each target region:

$$191 \text{Methylation ratio} = \frac{\text{methylated cytosine (C)}}{\text{methylated C} + \text{unmethylated C}}$$

192 Methylation fold change from cancer to control was calculated for each target region. For
193 analyzing differential methylated regions, significance level was set at $p \leq 0.05$, corresponding to
194 a $-\log_{10}$ adjusted p-value ≥ 1.301 (Benjamini-Hochberg correction).

195 **Genome-wide methylation analysis (GWM)**

196 The integrated bioinformatics pipeline Methy-pipen was used to analyse GWM. We carried out
197 the trimming step using Trimmomatic, removing adapter sequence and low-quality bases at
198 fragment ends (20). The methylation ratio for each bin was calculated as following equation.

199
$$\text{Methylation ratio} = \frac{\text{methylated cytosine (C)}}{\text{methylated C} + \text{unmethylated C}}$$

200 Mean methylation ratio was calculated for each bin and subsequently used to plot GWM density
201 curves. To identify bins with significant methylation changes between cancer and control group,
202 methylation ratio in each bin of cancer samples were compared with corresponding values in
203 control samples using Wilcoxon rank sum test. Bins with adjusted p-value (Benjamini-Hochberg
204 correction) ≤ 0.05 were considered significant. Those with \log_2 fold change (cancer vs control) $>$
205 0 were categorized as hypermethylated bins. Those with \log_2 fold change (cancer vs control) < 0
206 were categorized as hypomethylated bins.

207 **Copy number aberration analysis (CNA)**

208 CNA analysis was performed using the R-package QDNaseq (23). We also used 1-Mb
209 segmentation strategy to analyse CNA. We excluded bins that fell into the low mappability and
210 Duke blacklist regions(24). The number of reads mapped to each bin was measured by the
211 function “binReadCounts”, and GC-content correction was conducted by the functions
212 “estimateCorrection” and “correctBins”. The final CNA feature was derived by bin-wise
213 normalizing and outlier smoothing with the functions “normalizeBins” and “smoothOutlierBins”.
214 This process resulted in a feature vector of a length of 2691 bins.

215 To identify significant DNA gain or loss between cancer and control group, CNA values in each
216 bin of cancer samples were compared with corresponding values in control samples using
217 Wilcoxon rank sum test. Bins with adjusted p-value (Benjamini-Hochberg correction) ≤ 0.05

218 were considered significant. Those with \log_2 fold change (cancer vs control) > 0 were
219 categorized as significant increase. Those with \log_2 fold change (cancer vs control) < 0 were
220 categorized as significant decrease.

221 **Fragment length analysis (FLEN, SHORT, LONG, TOTAL, RATIO)**

222 We used an in-house python script to convert the.bsaligned files into BAM files and collected the
223 fragment length from 100 to 250 bp, resulting in 151 possible fragment lengths for further
224 analysis. The fragment frequency in each length (%) was measured by getting the proportion of
225 reads with that length to the total read count in the range of 100 to 250 bp. Fragment length (bp)
226 against fragment frequency (%) was plotted to obtain a FLEN distribution curve.

227 We divided the whole genome into 588 non-overlapping bins of 5Mb (5 million bases) long and
228 then extracted the read counts regarding these bins. Short fragments have lengths from 100 to
229 150 bp and long fragments have lengths from 151 to 250 bp. The ratio of short and long
230 fragments was calculated by dividing the number of each fragment. All the short, long and total
231 read counts for each sample in 588 bins were normalized using z-score normalization. The short,
232 long and total normalized read counts and short/long ratios were chosen as features analyzed
233 (SHORT, LONG, TOTAL, RATIO).

234 **End motif analysis (EM)**

235 AdaptaseTM technology (Integrated DNA Technologies, USA) was used during library
236 preparation to ligate adapters to ssDNA fragments in a template-independent reaction (25). This
237 step involved adding a random tail to the 5' end of reverse reads. Although median length of the
238 tail was 8 bp and thus allowed trimming to obtain information for other analysis, the random-
239 length tails did not allow exact determination of the 5' end of the reverse reads. Therefore, EM

240 features were determined based on the genomic coordinate of the 5' end of the forward reads.
241 We determined the first 4-mer sequence based on the human reference genome hg19. In 256
242 possible 4-mer motifs, the frequency of each motif was calculated by dividing the number of
243 reads carrying that motif by the total number of reads, generating an EM feature vector of a
244 length of 256 for each sample.

245 **Construction of machine learning models**

246 All samples in the discovery cohort were used for model training to classify if a sample is
247 cancerous or not. For every feature type (TM, GWM, CNA, FLEN, SHORT, LONG, TOTAL,
248 RATIO and EM), three machine learning algorithms, including Logistic regression (LR),
249 Random Forest (RF) and Extreme Gradient Boosting (XGB), were applied. By using the
250 “GridSearchCV” function in the scikit-learn (v.1.0.2), model hyperparameters with the best
251 performance were chosen with ‘CV’ parameter (cross-validation) set to 5. The best
252 hyperparameters for each algorithm were found using function ‘best_params_’ implemented in
253 GridSearchCV. Subsequently, feature selection was performed for each algorithm as follows: (1)
254 for LR, the “penalty” parameter with ‘l1’ (LASSO regression), ‘l2’ (Ridge regression) and
255 ‘none’ (no penalty) were examined to select the setting with the best performance; (2) for RF and
256 XGB, a “SelectFromModel” function with the ‘threshold’ was set at 0.0001 to get all features.
257 Then, the three algorithms (LR, RF, XGB) trained with the best hyperparameters and selected
258 features were validated using k-fold cross validation approach on the dataset of training cohort
259 with k-fold set to 20-fold, and ‘scoring’ parameter set to ‘roc_auc’. This split the data into 20
260 groups, in which 19 groups were model-fitted and the remaining group was tested, which
261 resulted in 20 ‘roc_auc’ scores. The average of these scores was used to obtain the prediction
262 performance of each model. The model with the highest ‘roc_auc’ average score was chosen

263 (either LR, RF or XGB). Ensembled models were constructed by combining probability scores
264 of nine single-feature base models (TM, GWM, CNA, FLEN, SHORT, LONG, TOTAL,
265 RATIO, EM) with different combination using LR, resulting in one probability score for every
266 sample. An extensive search was performed to evaluate the performance of all possible
267 combinations (n= 511) and the combination with highest AUC was selected as the final model.
268 The model cut-off was set at the threshold specificity of >95%. This combination model
269 performance was evaluated on an independent validation dataset to examine the model
270 classification power.

271 In addition the stacking ensemble, another combinatory strategy was examined. Instead of
272 combining nine base models , we generated a single dataframe consisting of raw data of all nine
273 features. The model hyperparameters tuning and features selecting were followed the same
274 strategy as described above. After choosing the best algorithm, the model performance was also
275 evaluated using the same external validation dataset.

276 **Construction of models for TOO**

277 **Strategy 1: Random Forest (RF) model**

278 A single data frame of nine features in discovery cohort was used to train the Random Forest
279 (RF) to classify 5 cancer types. By using the “GridSearchCV” function in the scikit-learn
280 (v.1.0.2), model hyperparameters with the best performance were chosen with ‘CV’ parameter
281 (cross-validation) set to 3 and “class_weight” parameter set to “balanced”. The best
282 hyperparameters were found by function ‘best_params_’. Then, the model was validated using k-
283 fold cross validation approach on the training cohort with k-fold set to 10-fold and its
284 performance was evaluated on the validation cohort.

285 **Strategy 2: Deep neural network (DNN) model**

286 Backpropagation trained the H₂O deep neural network (DNN) (multi-layer feedforward artificial
287 neural network) (H₂O package, version 3.36.1.2) with stochastic gradient descent. The random
288 grid search was selected as previously described (20).

289 **Strategy 3: Graph Convolutional Neural Network (GCNN) model**

290 The model training utilized an input graph formed from a discovery dataset and a validation
291 dataset as transudative setting (13) comprising patients diagnosed with five types of cancer:
292 breast, colorectal (CRC), gastric, liver, and lung. The discovery dataset contains a set of sample-
293 label pairs $\mathcal{J} = \{(X_i, Y_i) | i = 1, \dots, N\}$ where X_i represents the i th sample and Y_i represents i th
294 label, and N is the number of sample-label pairs. For each X_i in the discovery dataset, a node's
295 feature vector $f = \{F_0, \dots, F_d\} \in \mathbb{R}^d$ is constructed by combining groups of features, where F_i is
296 the i th feature, d is the number of features. The same procedure was applied for the independent
297 validation dataset. To construct an interaction graph between cancer nodes, we employed the k -
298 nearest neighbors' algorithm. An interaction graph defined as $G = (V, E)$ where $V = \{X_i | i =$
299 $1, \dots, N\}$ is a node set formed by the discovery samples, and $E = \{e_{ij}\}$ is an edge set, where e_{ij}
300 denotes an edge. Given N nodes in the node set, i.e. $|V| = N$, a graph topology $A \in \mathbb{R}^{N \times N}$ is
301 defined by:

$$A_{ij} = \begin{cases} 1, & e_{ij} \in E \text{ and } d_{ij} < \delta \\ 0, & \text{otherwise} \end{cases}$$

302 where d_{ij} is the Euclidean distance of node i and j , and δ is set to 0.8.

303 In accordance with (22), a Graph Convolutional Neural Network (GCNN) was constructed for
304 the purpose of tissue of origin classification. The network comprised three message-passing

305 layers, each with a hidden size of 44 and a head number of 4. Tissue of origin classification was
306 approached as a node classification problem, wherein the model assigned each node to one of
307 five cancer types: breast, colorectal, gastric, liver or lung cancer. Focal loss was employed for
308 multi-class classification optimization and the Adam optimizer was utilized for gradient-based
309 optimization. A 10-fold cross-validation approach was implemented on the discovery dataset;
310 nine groups were used for model training and one group for evaluation. The optimal model was
311 selected based on its ability to achieve the highest accuracy on the validation set during 10-fold
312 cross-validation. This model was subsequently applied to an independent validation dataset
313 consisting of 239 cancer patients across five cancer types to obtain the performance of tissue of
314 origin classification.

315 Given the predictions of trained model and the graph topology, we estimated the feature
316 importance score by the GNN Explainer [4]. The feature was considered important if it satisfied:

$$F_i > \delta_f$$

317 where F_i is the important score of i th feature estimated by the GNN Explainer, δ_f is the chosen
318 cut-off and was set to 0.9.

319 **Statistical analysis**

320 This study used either the Wilcoxon Rank Sum test or t-test to find statistically significant
321 differences between cancer and control. The Kolmogorov-Smirnov test was used to decide
322 whether two cohorts have the same statistical distribution. The Benjamini-Hochberg correction
323 was used to correct p-value for multiple comparisons (with a corrected p-value cutoff $\alpha \leq 0.05$).
324 DeLong's test was used to compare the differences between AUCs. All statistical analyses were

325 performed using R (4.1.0) packages, including ggplot2, pROC, and caret. 95% confident interval
326 (95% CI) was presented in a bracket next to a value accordingly.

327 **Results**

328 **Clinical characteristics of cancer and healthy participants.**

329 This study recruited 738 patients with five common cancer types, including breast cancer
330 (n=223), CRC (n=159), gastric cancer (n=98), liver cancer (n=122), lung cancer (n=136) and
331 1,550 healthy participants (Table S1). Cancer patients were diagnosed by either imaging and/or
332 histology analysis, depending on cancer type. All cancer patients were treatment-naïve at the
333 time of blood collection. Healthy participants had no history of cancer at the time of sample
334 collection and remained cancer-free at the 6- and 12-month follow-ups. Cancer patients and
335 healthy participants were randomly assigned to the discovery and validation cohorts (Table 1 and
336 Table S2). The discovery cohort was used to profile multiple cancer- and tissue-specific
337 signatures and to construct machine learning algorithm while the validation cohort was used
338 solely to external evaluation of the performance of machine learning models.

339 The discovery cohort comprised of 499 cancer patients (156 breast, 106 CRC, 67 gastric, 77 liver
340 and 93 lung, Table S1) and 1,076 healthy participants. The cancer group had a median age of 58
341 (range 25 to 97, Table 1) and consisted of 279 females and 220 males. The discovery healthy
342 group consisted of 599 females and 477 males, with a median age of 47 (range 18 to 84, Table
343 1). In the discovery cohort, gender ratios were similar between cancer and healthy control
344 groups, whereas cancer patients were older than controls ($p < 0.0001$, Mann-Whitney test, Table
345 1). Of the cancer patients, 10.4% were at stage I, 33.9% were at stage II, and 30.1% were at non-

346 metastatic stage IIIA. Staging information was not available for 25.7% of cancer patients, who
347 were confirmed by specialized clinicians to have non-metastatic tumors (Table 1).

348 The validation cohort consisted of 239 cancer patients (67 breast, 53 CRC, 31 gastric, 45 liver
349 and 43 lung, Table S1) and 474 healthy participants (Table 1). Consistent with the discovery
350 cohort, the gender distribution was comparable between the cancer and healthy control groups,
351 and the cancer group was older than the control group, with a median age of 59 and 48 years old,
352 respectively ($p < 0.0001$, Mann-Whitney test, Table 1). The percentage of cancer patients with
353 each stage was similar to that of the discovery cohort, with 9.6% at stage I, 28.9% at stage II and
354 32.2% at stage IIIA. Staging information was unavailable for 29.3% of non-metastatic cancer
355 patients (Table 1).

356 **The multimodal SPOT-MAS assay for multi-cancer and tissue of origin detection**

357 In our recent study of SPOT-MAS, we have demonstrated that the integration of ctDNA
358 methylation and fragmentomic features can significantly improve the early detection of
359 colorectal cancer (20). Here, we expanded the breadth of ctDNA analyses by adding two sets of
360 features including DNA copy number and end motif into SPOT-MAS to maximize cancer
361 detection rate and identify TOO. Briefly, a novel and cost-effective workflow of SPOT-MAS
362 was developed involving three main steps (Figure 1). In step 1, cfDNA was isolated from
363 peripheral blood and subjected to bisulfite conversion and adapter ligation to create a single
364 whole-genome bisulfite library of cfDNA. From this library, in step 2, a hybridization reaction
365 was performed to collect the target capture fraction (450 cancer specific regions), then the whole-
366 genome fraction was retrieved by collecting the ‘flow-through’ and hybridizing with probes
367 specific for adapter sequences of DNA library. Both the target capture fraction and whole-
368 genome fraction were sequenced to the depth of ~52X and 0.55X, respectively (Table S3). Data

369 pre-processing was performed to generate five different sets of cfDNA features, including
370 methylation changes at target regions (TM), genome-wide methylation (GWM), fragment length
371 patterns (Flen), copy number aberrations (CNA) and end motif (EM). In step 3, these features
372 were used as inputs for a two-stage model to obtain prediction outcomes. Stage 1 of our model
373 comprised of a stacked ensemble machine learning model for binary classification of cancer
374 versus healthy. Then the samples predicted as cancer were passed to stage 2 where graph
375 convolution neural network (GCNN) was adopted to predict TOO (Figure 1).

376 **Identification of differentially methylated regions (DMRs) in cancer patients from target** 377 **capture fraction**

378 DNA methylation is an important epigenetic signature responsible for major changes in
379 regulating expression of cancer associated genes by impacting the binding of transcription
380 factors to regulatory sites and the structure of chromatin (26, 27). Of the 450 target regions
381 associated with cancer that were selected from public data (13, 22), 402 regions were identified
382 as differentially methylated regions (DMRs) in cancer patients when compared to healthy
383 participants from the discovery cohort (Wilcoxon rank-sum test, p-values < 0.05, Figure 2A and
384 Table S4). Of those, 339 (84.3%) regions were identified as hypermethylated ($\log_{2}FC > 0$), and
385 63 (15.7%) regions as hypomethylated in cancer samples ($\log_{2}FC < 0$, Figure 2A). We next
386 examined the genomic location of the 402 DMRs and found 100, 108, 107 and 87 DMRs that
387 were mapped to promoter, exon, intron and intergenic regions, respectively (Figure 2B). To
388 understand the relationship between the differences in methylation regions and biological
389 pathways, we performed pathway enrichment analysis using g:Profiler on hypermethylated
390 DMRs. We detected 36 enriched pathways, including 14 from Kyoto Encyclopedia of Genes and
391 Genomes (KEGG) and 22 from WikiPathway (WP) (Figure 2C and Table S5). These significant

392 pathways were known to regulate tumorigenesis of breast, gastric, hepatocellular, and colorectal
393 cancer. Therefore, the methylation changes in the targeted regions, particularly the
394 hypermethylated DMRs, mostly occur early in tumorigenesis and are crucial for distinguishing
395 early-stage cancer patients from healthy individuals.

396 **Genome-wide methylation changes in cfDNA of cancer patients**

397 In addition to site-specific hypermethylation, hypomethylation is a significant genome-wide
398 change that has been identified in many types of cancers (24, 28, 29). To investigate the
399 methylation changes at genome-wide level, bisulfite sequencing reads from the whole-genome
400 fraction were mapped to the human genome, split into bins of 1Mb (2,734 bins across the
401 genome), and the reads from each bin were used to calculate methylation ratio. As expected, we
402 observed a left-ward shift in the distribution of methylation ratio in cancer samples compared to
403 healthy controls, indicating global hypomethylation in the cancer genome ($p < 0.0001$, two-
404 sample Kolmogorov-Smirnov test, Figure 3A). Of these bins, we identified 1,715 (62.7%) bins
405 as significantly hypomethylated in cancer, located across 22 autosomes of the genome (Figure
406 3B, Wilcoxon rank sum test with Benjamini-Hochberg adjusting p-value < 0.05). In contrast,
407 there were only 10 bins identified as hypermethylated and mapped to chromosome 1, 2, 3, 5, 6, 7
408 and 12 in the cancer genome (Figure 3B). Therefore, our data confirmed the widespread
409 hypomethylation across the genome and this would potentially serve to distinguish cancer
410 patients from healthy controls.

411 **Increase DNA copy number aberrations (CNAs) in cfDNA of cancer patients**

412 Somatic copy number aberrations (CNAs) in the cancer genome are associated with the initiation
413 and progression of numerous cancers by altering transcriptional levels of both oncogenes and

414 tumor suppressor genes (30). Recent studies have shown that CNAs detection could identify and
415 quantify the fraction of ctDNA in plasma cfDNA (31-33). To examine CNAs at genome-wide
416 scale, we used 1Mb bin to determine the percentage of bins that showed significant copy number
417 gains or losses between cancer and control group. We identified 729 bins (27.1%) with a
418 significant gain and 976 bins (36.3%) with a significant loss in copy number across 22
419 chromosomes of the cancer genome (Benjamini-Hochberg adjusting p-value <0.05, Wilcoxon
420 rank sum test, Figure 4A). We noted that chromosome 8 had the highest proportion of bins with
421 CNA gains, while chromosome 22 showed the highest proportion of bins with CNA losses
422 (Figure 4B).

423 It is thought that the abnormal hypomethylation at genome-wide level is linked with somatic
424 copy number aberration (CNA), resulting in genome instability, which is an important
425 tumorigenic event (34-36). Indeed, our data showed a significant increase in levels of CNA in
426 hypomethylated bins compared to bins with unchanged methylation (p=0.024, Figure S1A).
427 Consistently, bins with CNA gains showed significant decreases in methylation as compared to
428 those with CNA losses or unchanged CNA (p<0.01, Figure S1B). In summary, SPOT-MAS
429 enables comprehensive profiling of both global differences in methylation and somatic CNA as
430 individual feature types, as well as exploring their functional links during cancer initiation and
431 development, rendering them ideal biomarkers for cancer detection.

432 **Fragment length analysis captured patterns of ctDNA in plasma**

433 Several studies have shown that the fragmentation pattern of cfDNA is a non-random event
434 mediated by apoptotic-dependent caspases and ctDNA fragments tend to be shorter than non-
435 cancer cfDNA (14, 23, 37-39). One novel technical aspect of SPOT-MAS is the use of bisulfite
436 sequencing data not only for methylation but also for fragment length analysis. Certain studies

437 showed evidence of DNA degradation followed bisulfite treatment, possibly due to high
438 temperature and low pH conditions of the bisulfite conversion procedure, while other showed
439 that bisulfite sequencing affects large genomic DNA but not small size cfDNA (40-43).
440 Therefore, to demonstrate the use of bisulfite treated cfDNA for fragment length analysis, we
441 randomly selected 3 healthy controls and 9 cancer samples to perform pair-wise comparison
442 between bisulfite and non-bisulfite sequencing results. We observed a strong correlation between
443 fragment length profile of non-bisulfite and bisulfite sequencing (Pearson correlation, $R^2 > 0.9$,
444 $p < 0.0001$, Figure S2A) for all 12 tested samples, indicating the feasibility of using bisulfite
445 sequencing data for cfDNA fragment length analysis. Indeed, the fragment size distributions of
446 bisulfite-treated cfDNA in both cancer patients and control subjects showed a peak at 167 bp
447 (Figure 5A), corresponding to the length of DNA wrapped around histone (~ 147 bp) plus
448 linker regions ($\sim 2 \times 10$ bp), which was in good agreement with previous studies using non-
449 bisulfite cfDNA (14, 25). Importantly, our results showed that cfDNA of cancer patients was
450 more fragmented than that of healthy participants, with a higher frequency of fragments ≤ 150 bp
451 and a lower frequency of fragment > 150 bp (Figure 5A).

452 To examine whether the fragment length variation in cancer-derived cfDNA and non-cancer
453 cfDNA could be position-dependent (14), we calculated the ratios of short (≤ 150 bp) to long
454 fragments (> 150 bp) across the genome in cancer patients and healthy controls. The mean ratio
455 of short to long fragments in cancer patients was 0.29 (range 0.28 to 0.42), which was higher
456 than the mean ratio of 0.27 (range 0.26 to 0.39) for healthy controls (Figure 5B). The changes of
457 mean ratio were across 22 autosomes of the genome. Our results indicate that the SPOT-MAS
458 technology can effectively capture differences in fragmentation patterns between cancer and

459 healthy participants across the entire genome, making them potential biomarkers for the
460 detection of circulating tumor DNA in plasma.

461 **Profile of 4-mer end motifs reflecting differences between cancer and healthy cfDNA**

462 Associated with differences in fragment length is the differences in the DNA motifs at the end of
463 each fragment as the consequences of differential cleavage between DNA in cancer cells and
464 normal cells during apoptosis (25, 44). Here, we calculated the frequencies of 256 4-mer end
465 motifs (EMs) of cfDNA fragments and compared them between cancer patients and healthy
466 participants. Consistent with the fragment length features, we also confirmed the correlation of
467 EM frequency between bisulfite and non-bisulfite sequencing results of 12 randomly selected
468 samples, suggesting that EM profiles were reserved in bisulfite treated cfDNA (Figure S2B). Of
469 the 256 4-mer EMs, we detected 78 motifs with increased frequencies and 106 motifs with
470 decreased frequencies between cancer and healthy controls (Figure 6A and Table S6).

471 Interestingly, EMs beginning with cytosine (C) exhibited the highest number of EMs with
472 significant changes of frequency in cancer samples (Figure 6A). Figure 6B shows the top ten
473 EMs exhibiting significant differences. Specifically, the frequencies of five motifs (CAAA,
474 TAGA, CAGA, CAAG, and CAAT) were found to be significantly increased, while the
475 frequencies of another five motifs (CGCT, CGCC, CGCA, GCCT, and CGTT) were
476 significantly decreased in cancer patients (Figure 6B). Therefore, the differences in end motif
477 frequency identified by SPOT-MAS between cancer patients and healthy participants may serve
478 as a promising target for the identification of ctDNA.

479 **SPOT-MAS assay combining different features of cfDNA to enhance the accuracy of** 480 **cancer detection**

481 In order to increase the sensitivity of early cancer detection while avoiding the high cost of deep
482 sequencing, a screening test should survey a wide range of ctDNA signatures (16). Therefore, we
483 utilized multiple ctDNA signatures to construct classification models for distinguishing cancer
484 patients from healthy individuals. To expand the feature space, we generated four additional
485 features based on fragment length, including short, long, total fragment count, and short-to-long
486 ratio, resulting in nine input feature groups (Figure 7A). For each feature group, we tested three
487 different algorithms, including random forest (RF), logistic regression (LR), or extreme gradient
488 boosting (XGB), to tune hyperparameters and select the optimal algorithms (Figure 7A). To
489 evaluate the performance of these single-feature models, we performed 20-fold cross-validation
490 on the discovery dataset and calculated “Area Under the Curve” (AUC) of the “Receiver
491 Operating Characteristic” (ROC) curve. Among the nine features, EM-based model showed the
492 highest AUC of 0.90 (95% CI: 0.89-0.92, Figure 7B) while the SHORT-based model had the
493 lowest AUC of 0.71 (95% CI: 0.69-0.74, Figure 7B).

494 To assess whether combining features could improve classification, we used two strategies to
495 construct multi-feature models. In the first strategy, all nine feature groups were concatenated
496 into a single data frame before being fed into the RF, LR, or XGB algorithms. Of the three
497 algorithms, the XGB model exhibited the best performance with an AUC of 0.88 (95% CI: 0.87-
498 0.90, Figure 7B). However, this AUC is still lower than that of the EM-based model (0.88 versus
499 0.90, Figure 7B). In the second strategy, we constructed an ensemble stacking model using
500 logistic regression to combine the prediction results of the single-feature models. We conducted
501 an exhaustive search approach to evaluate the performance of 511 possible combinations. The
502 stacking ensemble model based on combining eight features, including TM, GW, CNA, FLEN,
503 LONG, TOTAL, RATIO and EM, exhibited the best performance and outperformed the single-

504 feature models (Table S7), with an AUC of 0.93 (95% CI: 0.92-0.95, Figure 7B and Figure S3).
505 In the independent validation cohort, we obtained similar results, where the ensemble model also
506 outperformed single-feature models, with an AUC of 0.95 (95% CI: 0.93-0.96, Figure 7C).
507 In order to ensure cost-effectiveness and minimize psychological impact of cancer screening
508 tests in a large population, high specificity is a crucial requirement. Accordingly, we established
509 the cutoff value for each constructed model based on a minimum specificity threshold of 95%.
510 Of the nine single-feature models, EM and GWM models exhibited the highest sensitivities, at
511 59.5% and 60.9%, respectively. The stacking ensemble model achieved a sensitivity of 73.8%
512 and a specificity of 95.1% with a cutoff value of 0.546 in the discovery cohort (Figure 7D), and a
513 mean sensitivity of 72.4% and a specificity of 97.0% in the validation cohort (Figure 7E).
514 Stratification of samples by cancer types revealed that the ensemble model performed most
515 accurately in predicting liver cancer (89.6% sensitivity), followed by CRC (82.1% sensitivity),
516 lung cancer (78.5% sensitivity) and gastric cancer (71.6% sensitivity) (Figure 7F, Table S8).
517 Breast cancer had the lowest detection rate of 58.3% (91/156 patients). Importantly, the
518 performance of our ensemble model remained consistent in the validation cohort, with liver
519 cancer again showing the highest sensitivity (91.1%), followed by lung cancer (83.7%), CRC
520 (83.0%), gastric cancer (61.3%), and breast cancer (49.3%) (Figure 7G, Table S8).

521 **Influence of clinical features on model prediction**

522 Upon stratifying our dataset by gender, we found that there was no significant difference in the
523 prediction of healthy status between males and females (Figure S4A and S4C). However, in the
524 case of cancer prediction, our model demonstrated higher accuracy in males than females in both
525 the discovery and validation cohorts (Figure S4A and S4C). Notably, when breast cancer
526 samples were removed from our analysis, there was no difference in the detection rates between

527 male and female patients (Figure S4B and S4D), suggesting that the observed gender bias may
528 be attributed to the high proportion of breast cancer patients (all females) in our cohort, who
529 exhibited the lowest detection rate among the five cancer groups.

530 We next evaluated the potential confounding effect of age on our prediction model by examining
531 the correlation between the model prediction scores and the participants' ages. The results
532 revealed no significant correlation, suggesting that age differences are unlikely to affect the
533 accuracy of our model (Figure S4E and S4F). With regards to cancer burden (ie. tumor size), our
534 model performed better for cancers with higher burden, as reflected by the higher cancer scores
535 assigned to these cases (Figure S4G and S4H). Specifically, patients with tumor diameter ≥ 3.5
536 cm were more likely to be detected than those with a diameter < 3.5 cm (Figure S4G and S4H).
537 Similarly, cancer stages also influence the performance of our stacking ensemble model,
538 showing increasing detection accuracy as the stages get more advanced. In the discovery cohort,
539 the model's accuracy was highest for stage IIIA cancers, with an AUC of 0.95 (95% CI 0.93-
540 0.97), and lowest for stage I cancer, with an AUC of 0.90 (95% CI 0.86-0.95) (Figure S4I and
541 S4J). Consistently, our model performance was lower with an AUC of 0.94 (95% CI 0.89-0.98)
542 and 0.93 (95% CI 0.90-0.96) for stage I and II cancer, respectively, increasing to 0.98 (95% CI
543 0.97-0.99) for stage IIIA in the validation cohort (Figure S4K and S4L). These results
544 demonstrated that our ensemble model can detect cancers at all stages found in our cohorts,
545 despite a slightly lower performance in early stages (stage I and II) compared to non-metastatic
546 stage (IIIA).

547 **SPOT-MAS enables prediction of cancer types**

548 The ability to predict the tissue origin of ctDNA is critical for early cancer detection as this can
549 guide subsequent diagnostic tests and treatment. Previous studies have attempted to use either

550 fragment length or methylation landscapes to achieve this goal (5, 14, 45). In this study, we
551 demonstrated the ability of SPOT-MAS to identify the TOO using low-depth bisulfite
552 sequencing to generate multiple sets of cfDNA features. We first concatenated the nine sets of
553 cfDNA features into a single data frame and focused our analysis on 499 cancer patients with
554 five cancer types in the discovery cohort. We then constructed a Random Forest (RF) and two
555 neural network models (convolutional neural network and graph convolutional neural network)
556 to predict the TOO and used 10-fold cross-validation to estimate and compare the performance
557 of these models (Figure 8A and Figure S5A). The Graph Convolutional Neural Network
558 (GCNN) was chosen due to its superior performance and stability (Figure S5B and S5C and
559 Table S9).

560 We then used the GNNExplainer tool to measure the importance of different cfDNA features.
561 Our results showed that breast cancer had the highest number of features with an important score
562 >0.9 (497 features), while lung cancer had the lowest number of important features (126
563 features) (Figure 8B). Colorectal, gastric, and liver cancers had 363, 309, and 204 important
564 features, respectively (Figure 8B and Table S10). Genome-wide methylation and copy number
565 aberration were the most important features for differentiating breast, colorectal, CRC, gastric
566 and liver cancer from other cancer types, while the end motif had the lowest contribution to
567 distinguish cancer types (Figure 8C). Visualization of the 3D GCNN showed that this set of
568 discriminative features could segregate the five different cancer types (Figure 8D), highlighting
569 the benefits of a multimodal approach for predicting TOO.

570 The median accuracy for TOO identification among the five cancer types by the GCNN-based
571 multi-feature model was 0.73 (range 0.54 to 0.87) in the discovery cohort (Figure 8E). The
572 accuracy in the discovery cohort was highest for breast (0.87) and liver cancer (0.82) and lowest

573 for gastric cancer (0.54). In the validation cohort, we obtained a slightly lower accuracy with a
574 median of 0.70 (range 0.55 to 0.78). The accuracies for individual cancer types were 0.78 for
575 breast, 0.76 for liver, 0.66 for colorectal, 0.63 for lung and 0.55 for gastric cancer (Figure 8F).
576 Among the 5 cancer types, breast cancer showed the highest TOO accuracy, possibly due to the
577 highest number of important features detected by the model. In contrast, CRC and gastric cancer
578 exhibited the lowest TOO accuracy with high misprediction rates between these two cancer types
579 (0.11 and 0.19 for CRC versus gastric and gastric versus CRC, respectively). Together, our study
580 highlights the benefits of integrating multimodal analysis with the GCNN model to capture the
581 broad landscape of tissue-specific markers in different cancer types.

582 **Discussion**

583 In an era marked by a global rise in cancer-related morbidity and mortality, the development of
584 liquid biopsy screening tests that can detect and localize cancer at an early stage holds
585 tremendous potential to revolutionize cancer diagnosis and therapy. Despite this, challenges in
586 test performance and cost must still be overcome, due mostly to the limited abundance of ctDNA
587 and its inherent variability. To address these, published liquid biopsy assays to date demanded a
588 very high-depth sequencing (15), or a combination of protein and genetic biomarkers (19),
589 resulting in an elevated price of test. In the current study, we present the SPOT-MAS assay as a
590 single workflow with comparable performance to current tests while requiring a much lower
591 sequencing depth (Table S11). SPOT-MAS achieved a sensitivity of 72.4 % at a specificity of
592 97.0 % for detecting five common cancer types using shallow depth sequencing. Furthermore, it
593 can predict the tissue of origin with an accuracy of 70%.

594 The SPOT-MAS assay allows comprehensive investigation of multiple biomarkers in cfDNA,
595 including targeted methylation (TM), genome-wide methylation (GWM), copy number

596 aberration (CNA), end motif (EM) and fragment length profiles (Flen). In TM analysis, out of
597 450 TM regions chosen from previous publications (13, 22), we identified 402 regions as
598 significant differentially methylated regions (DMRs) in cancer patients (Figure 2A). These
599 DMRs were enriched for regulatory regions of well-known cancer-related gene families such as
600 PAX family genes, TBX family genes, FOX family genes and HOX family genes, and some
601 have previously been reported as biomarkers for noninvasive cancer diagnosis, such as *SEPT9*
602 and *SHOX2* (46, 47). In addition to the targeted hypermethylation regions, our study also showed
603 widespread hypomethylation patterns across 22 autosomes of cancer patients (Figure 3), a
604 hallmark of cancer (48). Importantly, we demonstrated that the same bisulfite sequencing data
605 could be used to identify somatic CNA (Figure 4), cancer-associated fragment length (Figure 5)
606 and end motifs (Figure 6), highlighting the advantage of SPOT-MAS in capturing the broad
607 landscape of ctDNA signatures without high cost deep sequencing. For cancer-associated
608 fragment length, we pre-processed this data into five different feature tables to better reflect the
609 information embedded within the data. Overall, nine feature tables are available for model
610 construction.

611 The involvement and orthogonal links of the above features in the transcriptional regulation of
612 cancer-associated genes during carcinogenesis prompted us to examine whether the combination
613 of multiple cancer-specific signatures in cfDNA could improve the efficiency of cancer detection
614 (49, 50). We first determined the performance of models constructed using individual type of
615 cfDNA features. Next, by performing exhaustive searches for all possible combinations of
616 single-feature models, we identified that the stacking ensemble of seven features could achieve
617 the AUC of 0.95 (95% CI: 0.93-0.96, Figure 8C and Figure S3), which is superior to all single-
618 feature models. Among the five cancer types, breast cancer showed the lowest detection rate of

619 58.3% and 49.3% in the discovery and validation cohort, respectively. Variations in detection
620 rates among different cancer types have been previously reported (5, 19, 45). Consistently, it has
621 been reported that the detection of breast cancer, particularly in early stages, is challenging due
622 to the low levels of ctDNA shedding and heterogeneity of molecular subtypes of breast tumors
623 (5). In contrast, we obtained the highest detection rate for liver cancer patients with the
624 sensitivity of 89.6% and 91.1% in the discovery and validation cohort, respectively. Our finding
625 is in good agreement with the literature showing that liver tumors shed high amounts of ctDNA
626 (51). This result demonstrated the advantage of a multimodal approach to enhance ctDNA
627 detection in plasma. We also conducted a survey of liquid biopsy assays to put our SPOT-MAS
628 into the context of current state-of-the art in the field. Table S11 showed that SPOT-MAS is
629 using the lowest sequencing depth approach (with a depth coverage of $\sim 0.55X$) and making up
630 for this by integrating the greatest number of cfDNA features to achieve comparable
631 performance to other assays.

632 For TOO identification, our results showed that the graph convolutional neural network (GCNN)
633 performed the best among the models tested (Figure S5 and Figure 8). GCNN has the ability to
634 explore the similarity and mutual representation among samples, therefore achieving great
635 success in multi-class classification tasks (52, 53). Unlike the reference-based deconvolution
636 approaches (54, 55), our GCNN approach is independent of a reference methylation atlas, which
637 was developed from tissue or cell type specific methylation markers and thus may introduce bias
638 due to discordance between the methylomes of tissue gDNA and plasma cfDNA (16, 56).
639 Although the methylation changes were reported as most predictive for TOO in previous studies
640 (54, 55), our results showed the contribution of each of the 9 features for TOO identification
641 (Figure 8C). In addition to GWM, fragment ratio (RATIO) and CNA are the major contributors

642 to the discrimination of different tissue types. This finding provided additional evidence that the
643 multimodal approach capturing the breadth of tissue-specific signatures could improve the
644 accuracy of TOO identification (5). Our GCNN model achieved an accuracy of 0.70 for TOO
645 prediction in validation cohort. This was comparable to the performance of CancerLocator,
646 which was based a probabilistic distribution model of tissue specific methylation markers (57).
647 Recently, Liu et al. (5) developed a methylation atlas based method, which achieved a higher
648 accuracy of 93% for locating 50 types of cancer. However, this approach is based on deep
649 genome-wide sequencing with high depth coverage of 30X (Table S11), thus might not be a cost
650 effective approach for cancer screening in large populations, especially in low-income countries.

651 There are several limitations in our study. First, despite using a large dataset of 738 cancer
652 samples, there was an unequal distribution of samples among cancer types, with breast cancer
653 accounting for 30.2% (223/738, Table S2) of the total samples and gastric cancer having a much
654 smaller representation (13.3%, Table S2). As a result, our models may have been influenced by
655 this imbalance, potentially introducing bias in the training and evaluation process. Therefore,
656 future studies should consider incorporating more samples to better estimate the overall
657 performance of the SPOT-MAS test. Second, tumor staging information was not available for
658 26.8% of cancer patients (198/738) in our study. This is due to the patients' decision to select
659 different hospitals for diagnostics and treatment, leading to missing histopathological
660 information at the hospitals where they were originally recruited. However, all cancer patients
661 recruited in this study were confirmed to have non-metastatic tumors. Third, the cancer patients
662 in both the discovery and validation cohort were older than the healthy participants. Age
663 differences could be a confounding variable of methylation and could affect the model
664 performance (58, 59). However, we observed no significant association between the participants'

665 age and model prediction scores (Figure S4). Fourth, the ability of SPOT-MAS to differentiate
666 cancer patients from those with benign lesions has not been examined in this study. Fifth, this
667 study only focused on the top 5 common cancer types, thus the current version of SPOT-MAS
668 might misidentify cancer patients of other types, resulting in lower sensitivity to real world
669 application. Lastly, this was a retrospective cohort study and may be biased by the nature of this
670 study design. In an interim 6-month report of a prospective study named K-DETEK, we were
671 encouraged by the preliminary data demonstrated the ability of SPOT-MAS to detect cancer
672 patients who exhibited no symptoms at the time of testing (60). Despite these promising results,
673 the performance of SPOT-MAS as an early cancer screening test remains to be fully validated in
674 a large, multi-center prospective study with 1 to 2 years of follow up.

675 **Conclusions**

676 In conclusion, we have developed the SPOT-MAS assay to comprehensively profile methylomic,
677 fragmentomic, copy number aberrations, and motif end signatures of plasma cfDNA. Our large-
678 scale case-control study demonstrated that SPOT-MAS, with its unique combination of
679 multimodal analysis of cfDNA signatures and innovative machine-learning algorithms, can
680 successfully detect and localize multiple types of cancer at a low-cost sequencing. These findings
681 provided important supporting evidence for the incorporation of SPOT-MAS into clinical
682 settings as a complementary cancer screening method for at-risk populations.

683

684 **References**

- 685 1. Sung H, Ferlay J, Siegel RL, Laversanne M, Soerjomataram I, Jemal A, et al. Global
686 Cancer Statistics 2020: GLOBOCAN Estimates of Incidence and Mortality Worldwide for 36
687 Cancers in 185 Countries. *CA: A Cancer Journal for Clinicians*. 2021;71(3):209-49.
- 688 2. Pham T, Bui L, Kim G, Hoang D, Tran T, Hoang M. Cancers in Vietnam-Burden and
689 Control Efforts: A Narrative Scoping Review. *Cancer Control*. 2019;26(1):1073274819863802.
- 690 3. Hawkes N. Cancer survival data emphasise importance of early diagnosis. *BMJ*.
691 2019;364:l408.
- 692 4. Kakushadze Z, Raghubanshi R, Yu W. Estimating Cost Savings from Early Cancer
693 Diagnosis. *Data [Internet]*. 2017; 2(3).
- 694 5. Liu MC, Oxnard GR, Klein EA, Swanton C, Seiden MV. Sensitive and specific multi-
695 cancer detection and localization using methylation signatures in cell-free DNA. *Ann Oncol*.
696 2020;31(6):745-59.
- 697 6. Li J, Han X, Yu X, Xu Z, Yang G, Liu B, et al. Clinical applications of liquid biopsy as
698 prognostic and predictive biomarkers in hepatocellular carcinoma: circulating tumor cells and
699 circulating tumor DNA. *J Exp Clin Cancer Res*. 2018;37(1):213.
- 700 7. Nguyen HT, Luong BA, Tran DH, Nguyen TH, Ngo QD, Le LGH, et al. Ultra-Deep
701 Sequencing of Plasma-Circulating DNA for the Detection of Tumor- Derived Mutations in
702 Patients with Nonmetastatic Colorectal Cancer. *Cancer Invest*. 2022;40(4):354-65.
- 703 8. Gao Q, Zeng Q, Wang Z, Li C, Xu Y, Cui P, et al. Circulating cell-free DNA for cancer
704 early detection. *The Innovation*. 2022;3(4):100259.

- 705 9. Pascual J, Attard G, Bidard FC, Curigliano G, De Mattos-Arruda L, Diehn M, et al.
706 ESMO recommendations on the use of circulating tumour DNA assays for patients with cancer:
707 a report from the ESMO Precision Medicine Working Group. *Ann Oncol.* 2022;33(8):750-68.
- 708 10. Nguyen HT, Tran DH, Ngo QD, Pham HT, Tran TT, Tran VU, et al. Evaluation of a
709 Liquid Biopsy Protocol using Ultra-Deep Massive Parallel Sequencing for Detecting and
710 Quantifying Circulation Tumor DNA in Colorectal Cancer Patients. *Cancer Invest.*
711 2020;38(2):85-93.
- 712 11. Constantin N, Sina AA, Korbie D, Trau M. Opportunities for Early Cancer Detection:
713 The Rise of ctDNA Methylation-Based Pan-Cancer Screening Technologies. *Epigenomes.*
714 2022;6(1).
- 715 12. Phan TH, Chi Nguyen VT, Thi Pham TT, Nguyen VC, Ho TD, Quynh Pham TM, et al.
716 Circulating DNA methylation profile improves the accuracy of serum biomarkers for the
717 detection of nonmetastatic hepatocellular carcinoma. *Future Oncol.* 2022;18(39):4399-413.
- 718 13. Chen X, Gole J, Gore A, He Q, Lu M, Min J, et al. Non-invasive early detection of
719 cancer four years before conventional diagnosis using a blood test. *Nature Communications.*
720 2020;11(1):3475.
- 721 14. Cristiano S, Leal A, Phallen J, Fiksel J, Adleff V, Bruhm DC, et al. Genome-wide cell-
722 free DNA fragmentation in patients with cancer. *Nature.* 2019;570(7761):385-9.
- 723 15. Jamshidi A, Liu MC, Klein EA, Venn O, Hubbell E, Beausang JF, et al. Evaluation of
724 cell-free DNA approaches for multi-cancer early detection. *Cancer Cell.* 2022;40(12):1537-
725 49.e12.
- 726 16. Moser T, Kühberger S, Lazzeri I, Vlachos G, Heitzer E. Bridging biological cfDNA
727 features and machine learning approaches. *Trends Genet.* 2023;39(4):285-307.

- 728 17. Im YR, Tsui DWY, Diaz LA, Jr., Wan JCM. Next-Generation Liquid Biopsies:
729 Embracing Data Science in Oncology. *Trends Cancer*. 2021;7(4):283-92.
- 730 18. Zhou Q, Kang G, Jiang P, Qiao R, Lam WKJ, Yu SCY, et al. Epigenetic analysis of cell-
731 free DNA by fragmentomic profiling. *Proceedings of the National Academy of Sciences*.
732 2022;119(44):e2209852119.
- 733 19. Cohen JD, Li L, Wang Y, Thoburn C, Afsari B, Danilova L, et al. Detection and
734 localization of surgically resectable cancers with a multi-analyte blood test. *Science*.
735 2018;359(6378):926-30.
- 736 20. Nguyen HT, Khoa Huynh LA, Nguyen TV, Tran DH, Thu Tran TT, Khang Le ND, et al.
737 Multimodal analysis of ctDNA methylation and fragmentomic profiles enhances detection of
738 nonmetastatic colorectal cancer. *Future Oncol*. 2022;18(35):3895-912.
- 739 21. Edge SB, Compton CC. The American Joint Committee on Cancer: the 7th Edition of the
740 AJCC Cancer Staging Manual and the Future of TNM. *Annals of Surgical Oncology*.
741 2010;17(6):1471-4.
- 742 22. Nguyen H-N, Cao N-PT, Van Nguyen T-C, Le KND, Nguyen DT, Nguyen Q-TT, et al.
743 Liquid biopsy uncovers distinct patterns of DNA methylation and copy number changes in
744 NSCLC patients with different EGFR-TKI resistant mutations. *Scientific Reports*.
745 2021;11(1):16436.
- 746 23. Mouliere F, Chandrananda D, Piskorz AM, Moore EK, Morris J, Ahlborn LB, et al.
747 Enhanced detection of circulating tumor DNA by fragment size analysis. *Sci Transl Med*.
748 2018;10(466).
- 749 24. Ehrlich M. DNA methylation in cancer: too much, but also too little. *Oncogene*.
750 2002;21(35):5400-13.

- 751 25. Jiang P, Sun K, Peng W, Cheng SH, Ni M, Yeung PC, et al. Plasma DNA End-Motif
752 Profiling as a Fragmentomic Marker in Cancer, Pregnancy, and Transplantation. *Cancer*
753 *Discovery*. 2020;10(5):664-73.
- 754 26. Yin Y, Morgunova E, Jolma A, Kaasinen E, Sahu B, Khund-Sayeed S, et al. Impact of
755 cytosine methylation on DNA binding specificities of human transcription factors. *Science*.
756 2017;356(6337).
- 757 27. Buitrago D, Labrador M, Arcon JP, Lema R, Flores O, Esteve-Codina A, et al. Impact of
758 DNA methylation on 3D genome structure. *Nature Communications*. 2021;12(1):3243.
- 759 28. Das PM, Singal R. DNA methylation and cancer. *J Clin Oncol*. 2004;22(22):4632-42.
- 760 29. Hoffmann MJ, Schulz WA. Causes and consequences of DNA hypomethylation in
761 human cancer. *Biochem Cell Biol*. 2005;83(3):296-321.
- 762 30. Shao X, Lv N, Liao J, Long J, Xue R, Ai N, et al. Copy number variation is highly
763 correlated with differential gene expression: a pan-cancer study. *BMC Med Genet*.
764 2019;20(1):175.
- 765 31. Baldacchino S, Grech G. Somatic copy number aberrations in metastatic patients: The
766 promise of liquid biopsies. *Semin Cancer Biol*. 2020;60:302-10.
- 767 32. Knuutila S, Aalto Y, Autio K, Björkqvist AM, El-Rifai W, Hemmer S, et al. DNA copy
768 number losses in human neoplasms. *Am J Pathol*. 1999;155(3):683-94.
- 769 33. Dereli-Öz A, Versini G, Halazonetis TD. Studies of genomic copy number changes in
770 human cancers reveal signatures of DNA replication stress. *Mol Oncol*. 2011;5(4):308-14.
- 771 34. Brennan K, Flanagan JM. Is there a link between genome-wide hypomethylation in blood
772 and cancer risk? *Cancer Prev Res (Phila)*. 2012;5(12):1345-57.

- 773 35. Zhang W, Klinkebiel D, Barger CJ, Pandey S, Guda C, Miller A, et al. Global DNA
774 Hypomethylation in Epithelial Ovarian Cancer: Passive Demethylation and Association with
775 Genomic Instability. *Cancers (Basel)*. 2020;12(3).
- 776 36. Chan KCA, Jiang P, Chan CWM, Sun K, Wong J, Hui EP, et al. Noninvasive detection
777 of cancer-associated genome-wide hypomethylation and copy number aberrations by plasma
778 DNA bisulfite sequencing. *Proceedings of the National Academy of Sciences*.
779 2013;110(47):18761-8.
- 780 37. Underhill HR, Kitzman JO, Hellwig S, Welker NC, Daza R, Baker DN, et al. Fragment
781 Length of Circulating Tumor DNA. *PLOS Genetics*. 2016;12(7):e1006162.
- 782 38. Lo YMD, Han DSC, Jiang P, Chiu RWK. Epigenetics, fragmentomics, and topology of
783 cell-free DNA in liquid biopsies. *Science*. 2021;372(6538):eaaw3616.
- 784 39. Nguyen VC, Nguyen TH, Phan TH, Tran TT, Pham TTT, Ho TD, et al. Fragment length
785 profiles of cancer mutations enhance detection of circulating tumor DNA in patients with early-
786 stage hepatocellular carcinoma. *BMC Cancer*. 2023;23(1):233.
- 787 40. Raizis AM, Schmitt F, Jost JP. A bisulfite method of 5-methylcytosine mapping that
788 minimizes template degradation. *Anal Biochem*. 1995;226(1):161-6.
- 789 41. Tanaka K, Okamoto A. Degradation of DNA by bisulfite treatment. *Bioorg Med Chem*
790 *Lett*. 2007;17(7):1912-5.
- 791 42. Kint S, De Spiegelaere W, De Kesel J, Vandekerckhove L, Van Criekinge W. Evaluation
792 of bisulfite kits for DNA methylation profiling in terms of DNA fragmentation and DNA
793 recovery using digital PCR. *PLoS One*. 2018;13(6):e0199091.
- 794 43. Ehrich M, Zoll S, Sur S, van den Boom D. A new method for accurate assessment of
795 DNA quality after bisulfite treatment. *Nucleic Acids Res*. 2007;35(5):e29.

- 796 44. Jin C, Liu X, Zheng W, Su L, Liu Y, Guo X, et al. Characterization of fragment sizes,
797 copy number aberrations and 4-mer end motifs in cell-free DNA of hepatocellular carcinoma for
798 enhanced liquid biopsy-based cancer detection. *Mol Oncol.* 2021;15(9):2377-89.
- 799 45. Klein EA, Richards D, Cohn A, Tummala M, Lapham R, Cosgrove D, et al. Clinical
800 validation of a targeted methylation-based multi-cancer early detection test using an independent
801 validation set. *Ann Oncol.* 2021;32(9):1167-77.
- 802 46. Ilse P, Biesterfeld S, Pomjanski N, Wrobel C, Schramm M. Analysis of
803 SHOX2 Methylation as an Aid to Cytology in Lung Cancer Diagnosis.
804 *Cancer Genomics - Proteomics.* 2014;11(5):251.
- 805 47. Warren JD, Xiong W, Bunker AM, Vaughn CP, Furtado LV, Roberts WL, et al. Septin 9
806 methylated DNA is a sensitive and specific blood test for colorectal cancer. *BMC Med.*
807 2011;9:133.
- 808 48. Jones PA, Ohtani H, Chakravarthy A, De Carvalho DD. Epigenetic therapy in immune-
809 oncology. *Nat Rev Cancer.* 2019;19(3):151-61.
- 810 49. Ulz P, Perakis S, Zhou Q, Moser T, Belic J, Lazzeri I, et al. Inference of transcription
811 factor binding from cell-free DNA enables tumor subtype prediction and early detection. *Nat*
812 *Commun.* 2019;10(1):4666.
- 813 50. Ivanov M, Baranova A, Butler T, Spellman P, Mileyko V. Non-random fragmentation
814 patterns in circulating cell-free DNA reflect epigenetic regulation. *BMC Genomics.* 2015;16
815 Suppl 13(Suppl 13):S1.
- 816 51. Caggiano C, Celona B, Garton F, Mefford J, Black BL, Henderson R, et al.
817 Comprehensive cell type decomposition of circulating cell-free DNA with CelFiE. *Nature*
818 *Communications.* 2021;12(1):2717.

- 819 52. Yin C, Cao Y, Sun P, Zhang H, Li Z, Xu Y, et al. Molecular Subtyping of Cancer Based
820 on Robust Graph Neural Network and Multi-Omics Data Integration. *Front Genet.*
821 2022;13:884028.
- 822 53. Huang Y, Chung ACS. Disease prediction with edge-variational graph convolutional
823 networks. *Med Image Anal.* 2022;77:102375.
- 824 54. Moss J, Magenheimer J, Neiman D, Zemmour H, Loyfer N, Korach A, et al.
825 Comprehensive human cell-type methylation atlas reveals origins of circulating cell-free DNA in
826 health and disease. *Nature Communications.* 2018;9(1):5068.
- 827 55. Loyfer N, Magenheimer J, Peretz A, Cann G, Bredno J, Klochendler A, et al. A DNA
828 methylation atlas of normal human cell types. *Nature.* 2023;613(7943):355-64.
- 829 56. Zhou X, Cheng Z, Dong M, Liu Q, Yang W, Liu M, et al. Tumor fractions deciphered
830 from circulating cell-free DNA methylation for cancer early diagnosis. *Nature Communications.*
831 2022;13(1):7694.
- 832 57. Kang S, Li Q, Chen Q, Zhou Y, Park S, Lee G, et al. CancerLocator: non-invasive cancer
833 diagnosis and tissue-of-origin prediction using methylation profiles of cell-free DNA. *Genome*
834 *Biol.* 2017;18(1):53.
- 835 58. Yusipov I, Bacalini MG, Kalyakulina A, Krivonosov M, Pirazzini C, Gensous N, et al.
836 Age-related DNA methylation changes are sex-specific: a comprehensive assessment. *Aging*
837 *(Albany NY).* 2020;12(23):24057-80.
- 838 59. Field AE, Robertson NA, Wang T, Havas A, Ideker T, Adams PD. DNA Methylation
839 Clocks in Aging: Categories, Causes, and Consequences. *Mol Cell.* 2018;71(6):882-95.
- 840 60. Nguyen THH, Lu YT, Le VH, Bui VQ, Nguyen LH, Pham NH, et al. Clinical validation
841 of a ctDNA-Based Assay for Multi-Cancer Detection: An Interim Report from a Vietnamese

842 Longitudinal Prospective Cohort Study of 2795 Participants. Cancer Investigation.

843 2023;41(3):232-48.

844

845

846 **List of abbreviations**

MCED	Multi-cancer early detection
TOO	tissue of origin
cfDNA	Circulating cell-free DNA
ctDNA	Circulating tumor DNA
SPOT-MAS	Screening for the Presence Of Tumor by DNA Methylation And Size
AUC	Area Under the Curve
TM	Targeted methylation
GWM	Genome-wide methylation
CNA	Copy number aberration
FLEN	Fragment length
EM	End motif
LB	Liquid Biopsy
LR	Logistic regression
RF	Random Forest
XGB	Extreme Gradient Boosting
DNN	Deep neural network
GCNN	Graph Convolutional Neural Network

CRC Colorectal cancer

ROC Receiver Operating Characteristic

847 **Declarations**

848 **Ethics approval and consent to participate:**

849 This study was approved by the Ethics Committee of the Medic Medical Center, University of
850 Medicine and Pharmacy and Medical Genetics Institute, Ho Chi Minh city, Vietnam. Written
851 informed consent was obtained from each participant in accordance with the Declaration of
852 Helsinki.

853 **Consent for publication:**

854 Not applicable.

855 **Availability of data and materials:**

856 Sequencing data will be deposited in a public portal database (NCBI SRA) upon acceptance and
857 are available on request from the corresponding author, LST. The data are not publicly available
858 due to ethical restrictions.

859 **Competing interests:**

860 The authors declare no conflict of interest.

861 **Funding:**

862 The study was funded by Gene Solutions

863 **Disclosure statement:**

864 The authors including LST, HNN, HG, MDP, HHN and DSN hold equity in Gene Solutions. The
865 funder Gene Solutions provided support in the form of salaries for authors are inventors on the
866 patent application (USPTO 17930705). We also confirm that this does not alter our adherence to
867 Cancer Investigation policies on sharing data and materials.

868 **Author contribution:**

869 Conceptualization: DLV, THP, TXJ, VCN, HTN, TVN, HG, HNN, MDP, LST

870 Patient consultancy and screening: DLV, THP, TXJ, VCN, HTN, TVN, QTD, TND, AMT,
871 VHN, VTAN, LMQH, QDT, TTTP, TDH, BTN, TNVN, TDN, DTBP, BHHP, TLV, THTN,
872 TTT, MHT, NCT, TKL, THTT, MLD, HPTB, VVK, TAP, DHT, TNAL, TVNP, MTL, DSN,
873 VTC, TTTD, HST

874 Formal analysis: VTCN, THN, NNTD, TMQP, TDN, THHN, LAKH, THT, DHV, TMTT,
875 MNN, TTVV, ANN, TTT, VUT, MPL, TTD, TVP, LHDN

876 Supervision: DKT

877 Writing-original draft: VTCN, GTHN, TTTT, LST

878 Writing-review and editing: VTCN, GTHN, TTTT, MDP, LST

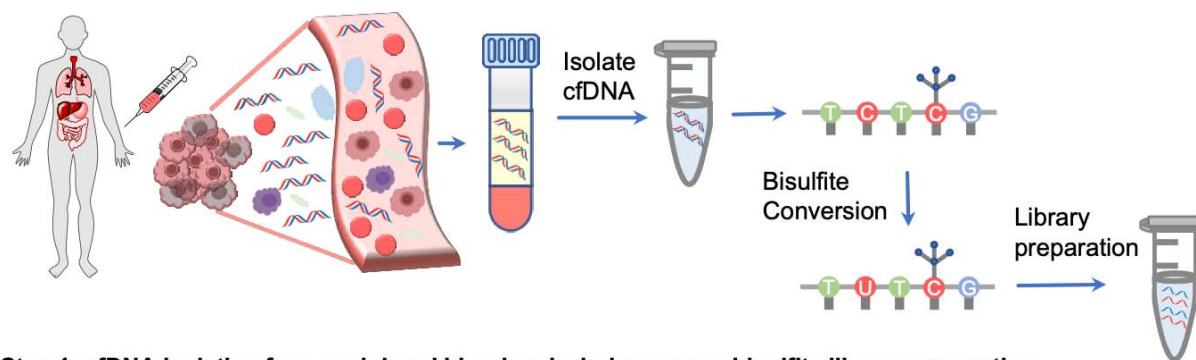
879 **Acknowledgments:**

880 We thank all participants who agreed to take part in this study, and all the clinics and hospitals
881 who assisted in patient consultation and sample collection.

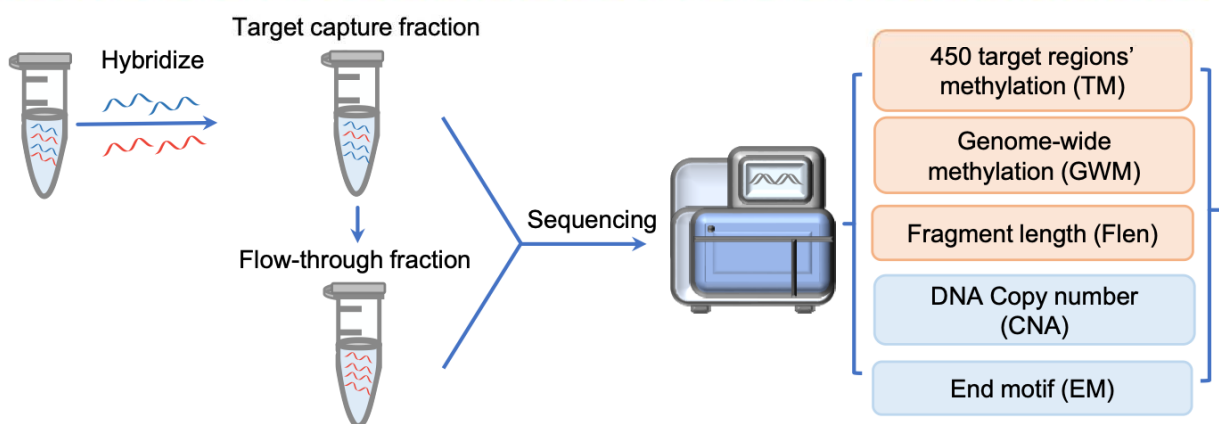
882

883

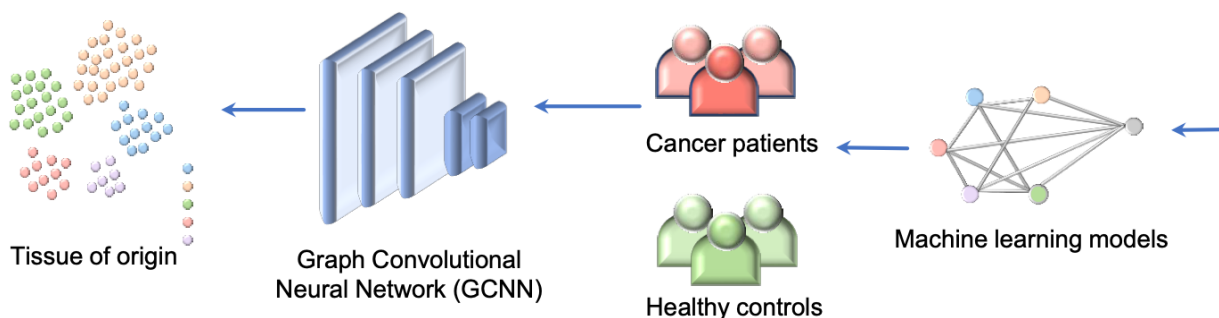
884 **Figures and Tables**



Step 1: cfDNA isolation from peripheral blood and whole-genome bisulfite library preparation



Step 2: Target and whole-genome fraction separation and sequencing



Step 3: Analysis of cfDNA signatures and model construction

885

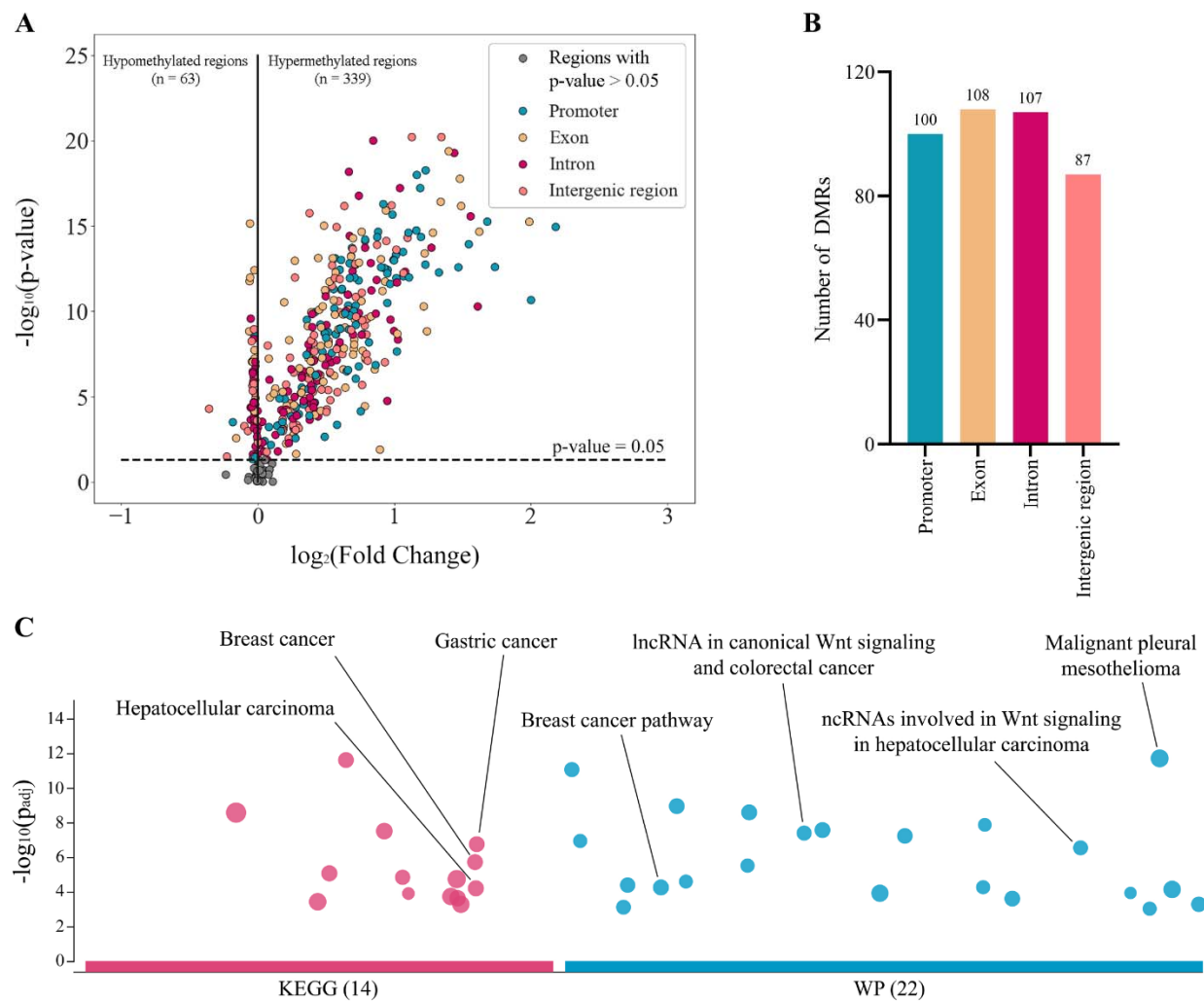
886 **Figure 1. Workflow of SPOT-MAS assay for multi-cancer detection and localization.** There

887 are three main steps in the SPOT-MAS assay. Firstly, cfDNA is isolated from peripheral blood,

888 then treated with bisulfite conversion and adapter ligation to make whole-genome bisulfite

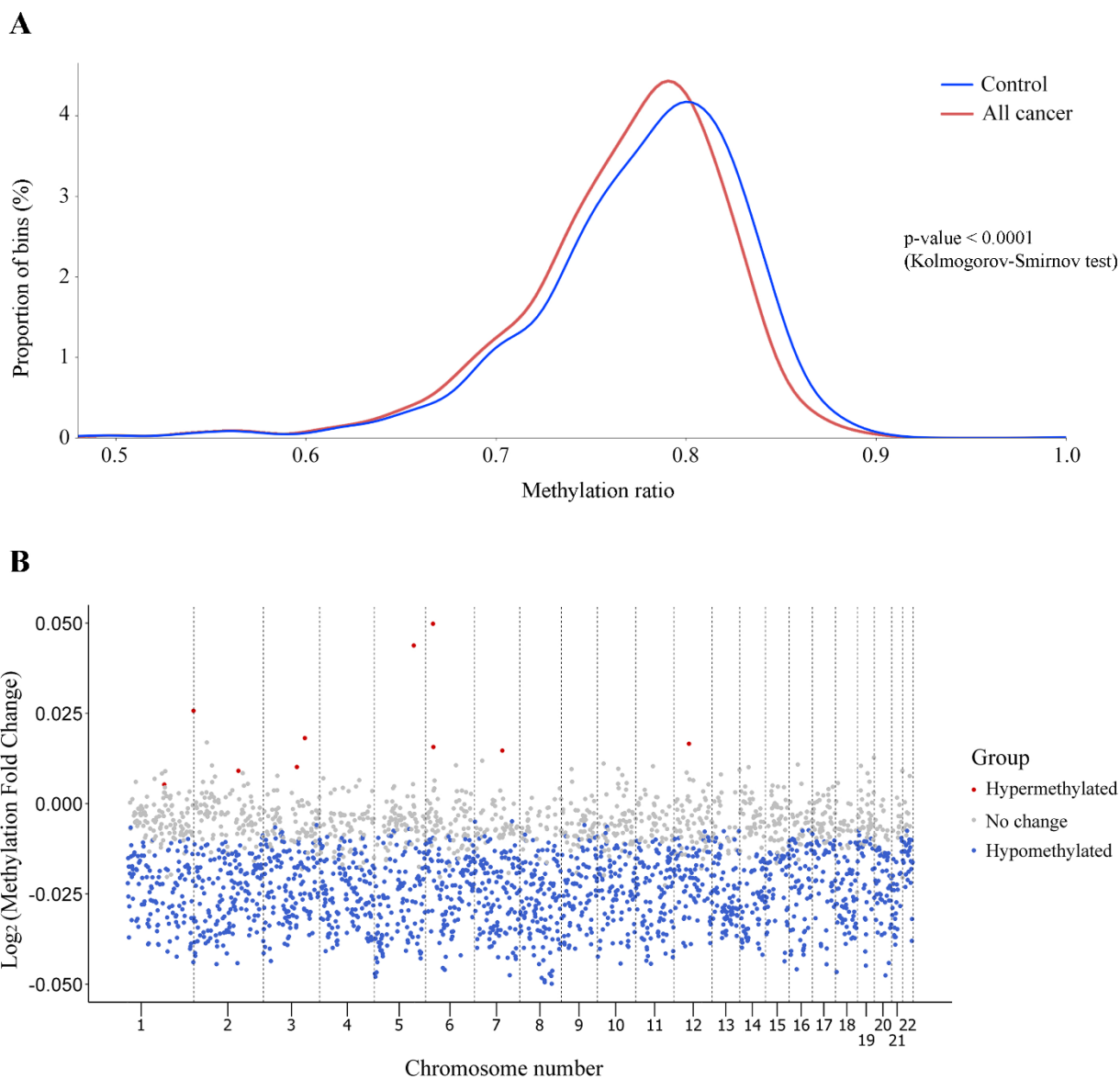
889 cfDNA library. Secondly, whole-genome bisulfite cfDNA library is subjected to hybridization by
890 probes specific for 450 target regions to collect the target capture fraction. The whole-genome
891 fraction was retrieved by collecting the ‘flow-through’ and hybridized with probes specific for
892 adapter sequences of DNA library. Both the target capture and whole-genome fractions were
893 subjected to massive parallel sequencing and the resulting data were pre-processed into five
894 different features of cfDNA: Target methylation (TM), genome-wide methylation (GWM),
895 fragment length profile (Flen), DNA copy number (CNA) and end motif (EM). Finally, machine
896 learning models and graph convolutional neural networks are adopted for classification of cancer
897 status and identification tissue of origin.

898



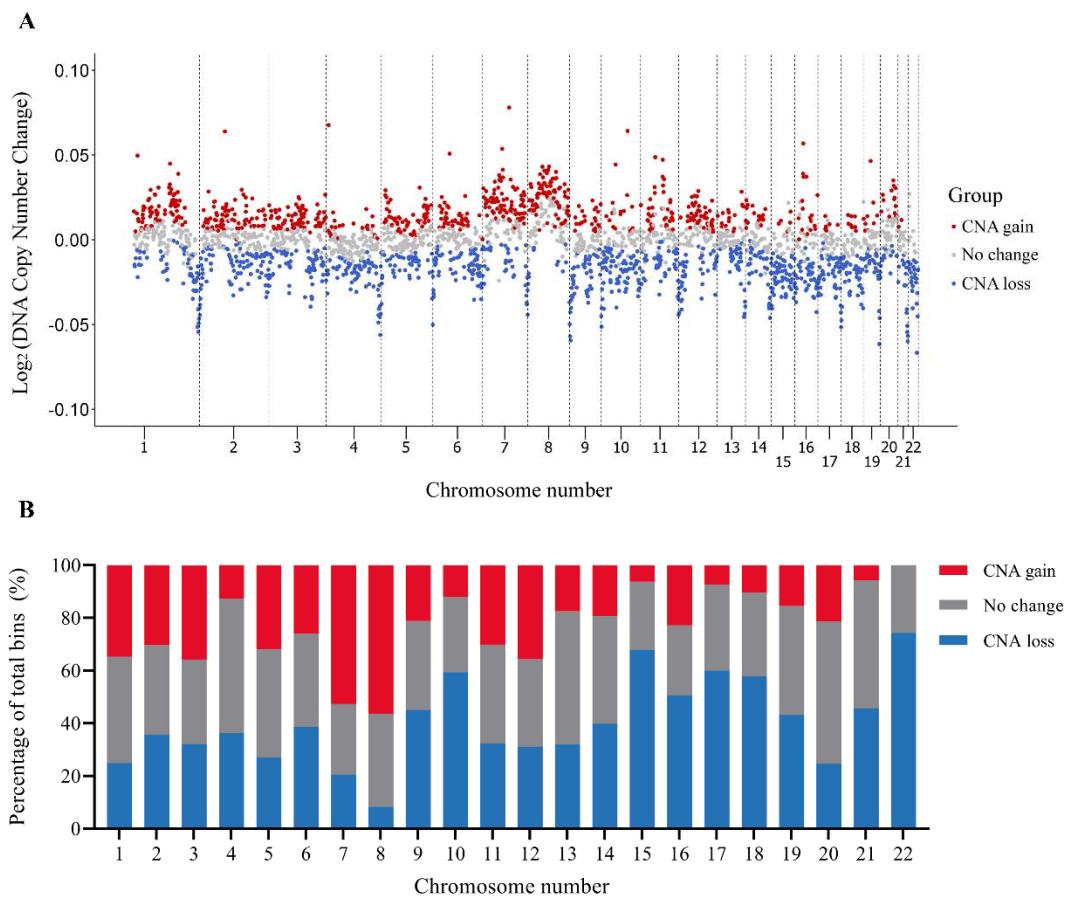
899
 900 **Figure 2. Analysis of targeted methylation in cfDNA.** (A) Volcano plot shows log₂ fold
 901 change (logFC) and significance (-log₁₀ Benjamini-Hochberg adjusted p-value from Wilcoxon
 902 rank-sum test) of 450 target regions when comparing 499 cancer patients and 1,076 healthy
 903 controls in the discovery cohort. There are 402 DMRs (p-value < 0.05), color-coded by genomic
 904 locations. (B) Number of DMRs in the four genomic locations. (C) KEGG and WP pathway
 905 enrichment analysis using g:Profiler for genes associated with the DMRs. A total of 36 pathways
 906 are enriched, suggesting a link between differences in methylation regions and tumorigenesis.

907



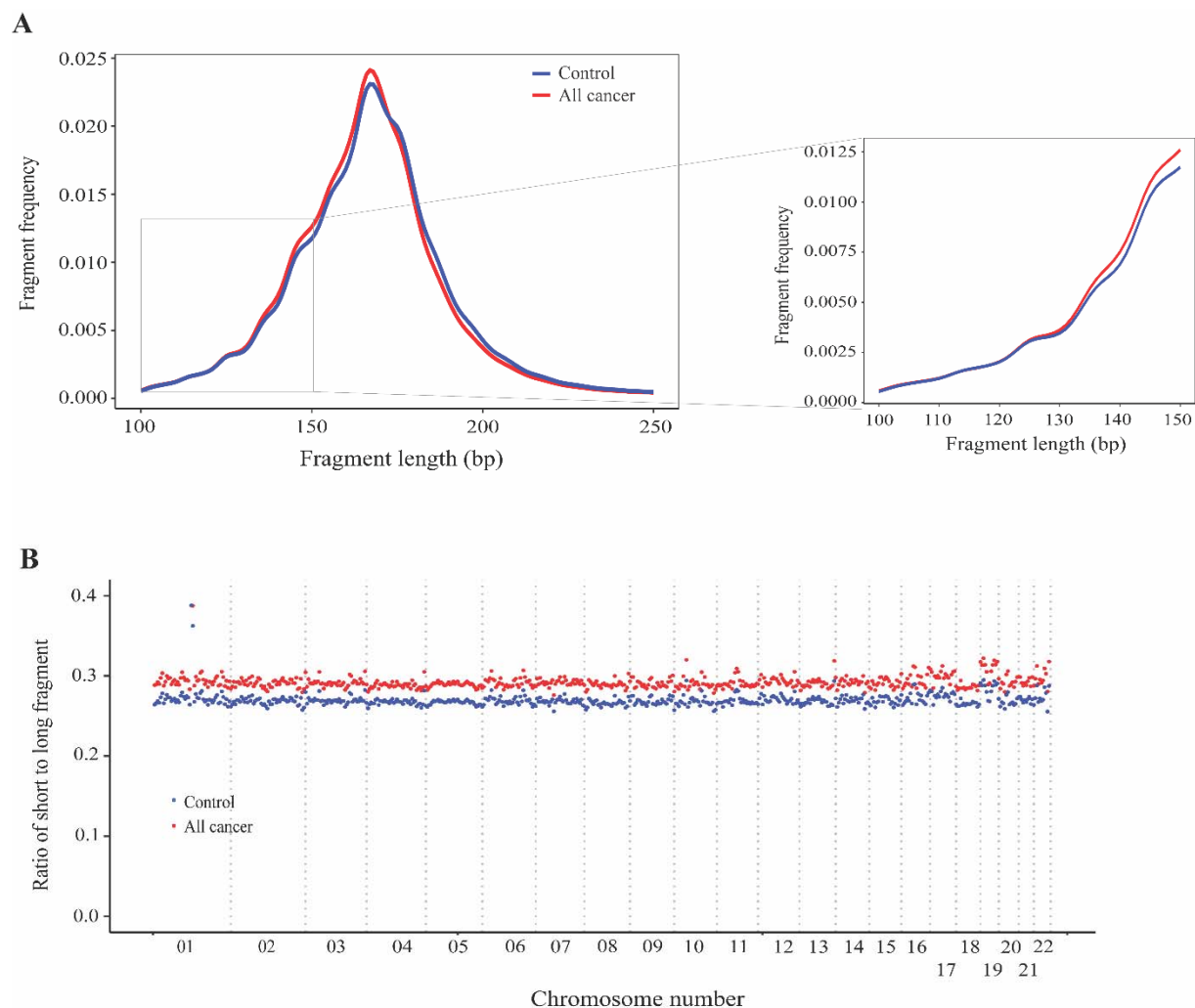
908

909 **Figure 3. Genome-wide methylation changes in cfDNA of cancer patients.** (A) Density plot
910 showing the distribution of genome-wide methylation ratio for all cancer patients (red curve, n=
911 499) and healthy participants (blue curve, n= 1,076). The left-ward shift in cancer samples
912 indicates global hypomethylation in the cancer genome (p < 0.0001, two-sample Kolmogorov-
913 Smirnov test). (B) Log₂ fold change of methylation ratio between cancer patients and healthy
914 participants in each bin across 22 chromosomes. Each dot indicates a bin, identified as
915 hypermethylated (red), hypomethylated (blue), or no significant change in methylation (grey).



917

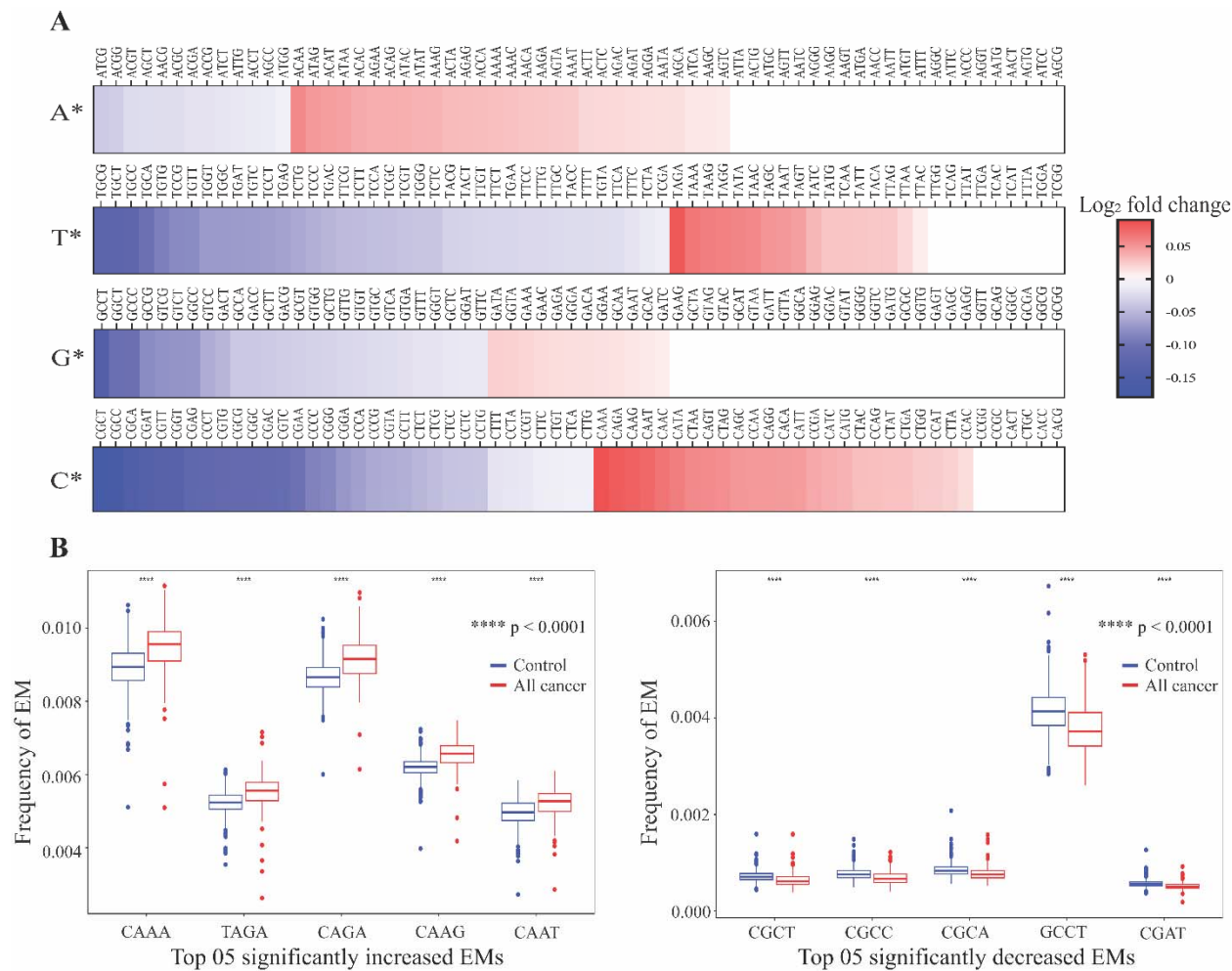
918 **Figure 4. Analysis of copy number aberration (CNA) in cfDNA.** (A) Log_2 fold change of
919 DNA copy number in each bin across 22 autosomes between 499 cancer patients and 1,076
920 healthy participants in the discovery cohort. Each dot represents a bin identified as gain (red),
921 loss (blue) or no change (grey) in copy number. (B) Proportions of different CNA bins in each
922 autosomes.



923

924 **Figure 5. Analysis of fragment length patterns of ctDNA in plasma.** (A) Density plot of
925 fragment length between cancer patients (red, n=499) and healthy participants (blue, n=1,076) in
926 the discovery cohort. Inset corresponds to an x-axis expansion of short fragment (<150 bp). (B)
927 Ratio of short to long fragments across 22 autosomes. Each dot indicates a mean ratio for each
928 bin in cancer patients (red) and healthy participants (blue).

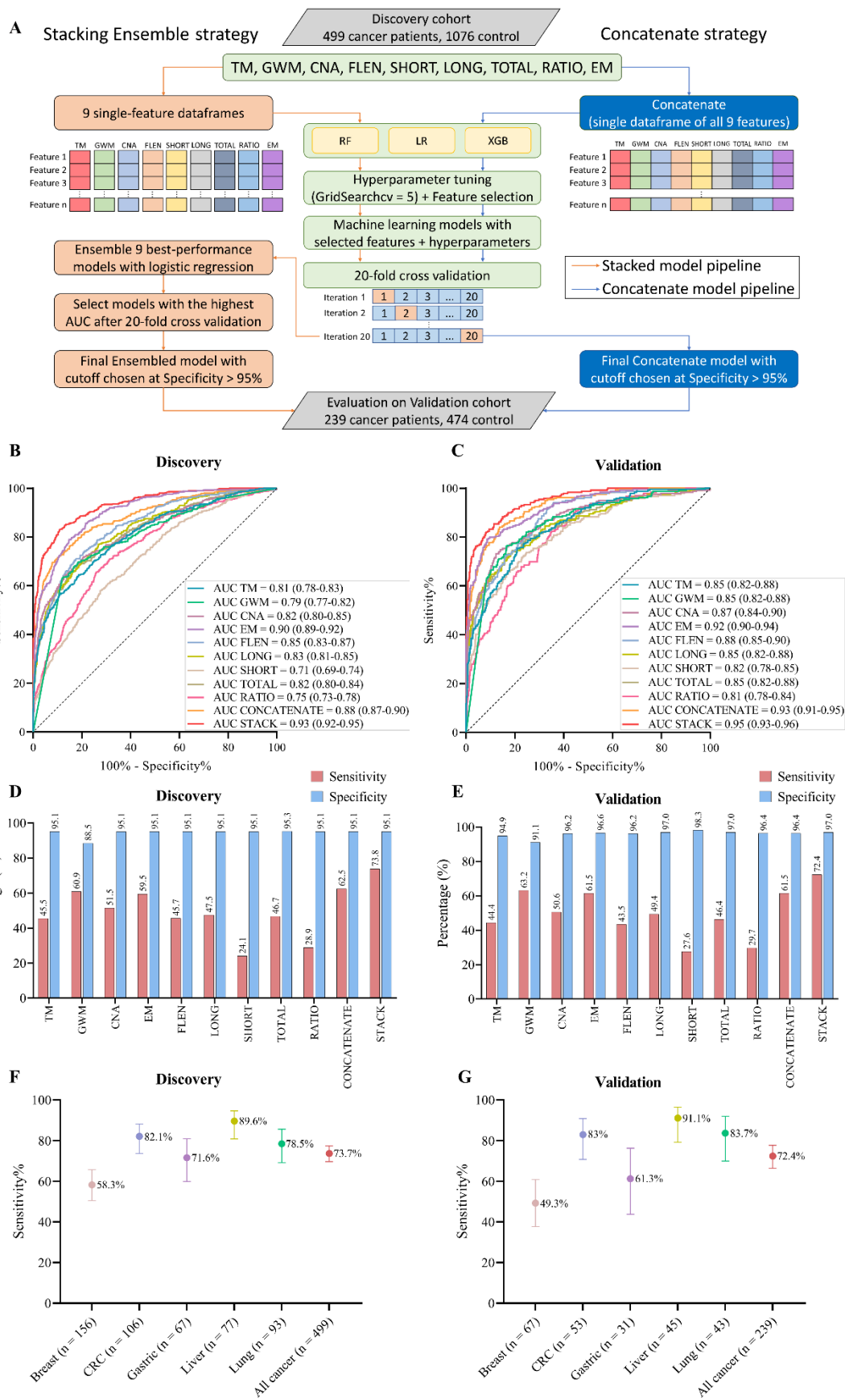
929



930

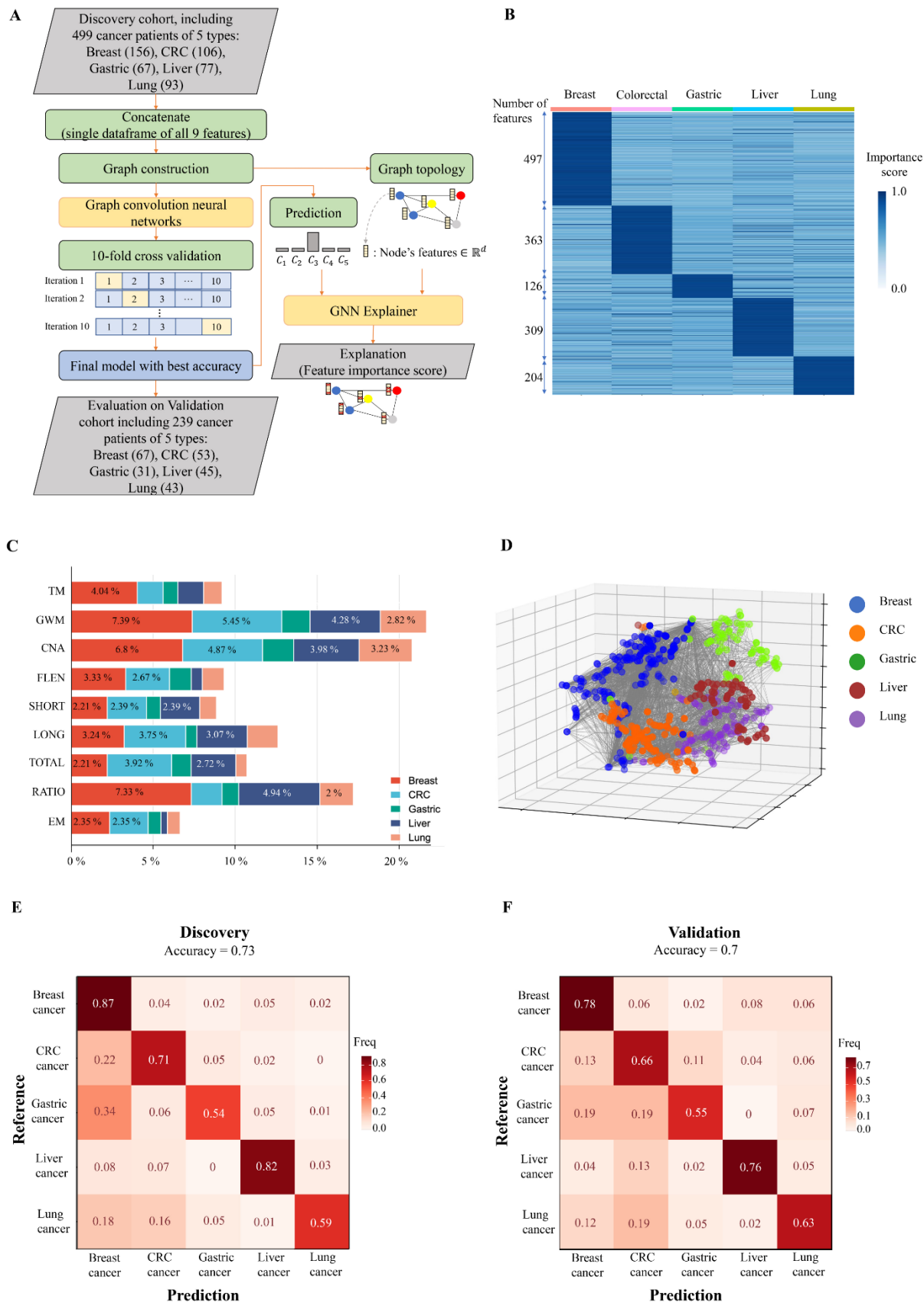
931 **Figure 6. Differences in 4-mer end motif between cancer and healthy cfDNA.** (A) Heatmap
 932 shows log₂ fold change of 256 4-mer end motifs in cancer patients (n=499) compared to healthy
 933 controls (n=1,076). (B) Box plots showing the top ten motifs with significant differences in
 934 frequency between cancer patients (red) and healthy controls (blue) using Wilcoxon rank-sum
 935 test with Bonferroni-adjusted p-value < 0.0001.

936



938 **Figure 7. Model construction and performance validation for SPOT-MAS.** (A) Two model
939 construction strategies for cancer detection. (B) and (C) ROC curves comparing the performance
940 of single-feature models, and two combination models (concatenate and ensemble stacking) in
941 the discovery (B) and validation cohorts (C). (D) and (E) Bar charts showing the specificity and
942 sensitivity of single-feature models and two combination models (concatenate and ensemble
943 stacking) in the discovery (D) and validation cohorts (E). (F) and (G) Dot plots showing the
944 sensitivity of SPOT-MAS assay in detection of 5 different cancer types in the discovery (F) and
945 validation cohorts (G). The points and error bars represent the average sensitivity over 20 runs
946 and 95% confidence intervals. Feature abbreviations as follows: TM – target methylation
947 density, GWM – genome-wide methylation density, CNA – copy number aberration, EM – 4-
948 mer end motif, FLEN – fragment length distribution, LONG – long fragment count, SHORT –
949 short fragment count, TOTAL – all fragment count, RATIO – ratio of short/long fragment.

950



952 **Figure 8. The performance of SPOT-MAS assay in prediction of the tissue of origin.** (A)
953 Model construction strategy to predict tissue of origin by combining nine sets of cfDNA features
954 using graph convolutional neural networks. (B) Heatmap shows feature important scores of five
955 cancer types. (C) Bar chart indicates the contribution of important features for classifying five
956 different cancers. (D) Three dimensions graph represents the classification of five cancer types.
957 (E) and (F) Cross-tables show agreement between the prediction (x-axis) and the reference (y-
958 axis) to predict tissue of origin in the discovery cohort (E) and validation cohort (F).

959

960 **Table 1.** Summary of clinical features of 738 cancer patients and 1,550 healthy controls in
 961 discovery and validation cohorts.

Clinical features		Discovery cohort (N=1,575)					Validation cohort (N=713)				
		Cancer (N = 499)		Healthy (N = 1,076)		p-value (Cancer vs Healthy)	Cancer (N = 239)		Healthy (N = 474)		p-value (Cancer vs Healthy)
		N	Percentage	N	Percentage		N	Percentage	N	Percentage	
Gender	Female	279	55.9%	599	55.7%	0.9281 [#]	126	52.72%	270	56.1%	0.2818 [#]
	Male	220	44.1%	477	44.3%		113	47.28%	204	43.9%	
Age	Median	58		47		< 0.0001 ##	59		48		< 0.0001 ##
	Min	25		18			28		19		
	Max	97		84			92		85		
Stage	I	52	10.4%			0.4947 [#]	23	9.6%			
	II	169	33.9%				69	28.9%			
	IIIA	150	30.1%				77	32.2%			
	Non- metastasis with unknown staging information	128	25.7%				70	29.3%			

962 [#] P-values from Chi-square test; ^{##} P-values from Mann-Whitney test

963