

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 Although currently guided screening tests have each been shown to provide better treatment
76 outcomes and reduce cancer mortality, some of them are invasive, thus having low accessibility.
77 Importantly, most of them are single cancer screening tests, which may result in high false
78 positive rates when used sequentially (5). Multi-cancer early detection (MCED) tests can
79 potentially overcome these challenges by simultaneously detecting multiple cancer types from a
80 single test (6). Liquid biopsy, an emerging non-invasive approach for MCED, can capture a wide
81 range of tumor features, including cell free DNA (cfDNA), circulating tumor DNA (ctDNA),
82 exosomes, proteins, mRNA, and metabolites (7, 8). Among them, ctDNA has become a
83 promising biomarker for detecting early-stage cancers because it is a carrier of genetic and
84 epigenetic modifications from cancer-derived DNA (9). Indeed, ctDNA detection has

85 demonstrated several advantages in non-invasive diagnostic, prognostic, and monitoring of
86 cancer patients during and after treatment (10, 11). Furthermore, ctDNA carrying tumor-specific
87 alterations could be used to identify the corresponding unknown primary cancer and tumor
88 localization.

89 In recent years, there has been considerable interest in exploring the potential of ctDNA
90 alterations for early detection of cancer (11, 12). One such approach is the PanSeer test, which
91 uses 477 differentially methylated regions (DMRs) in ctDNA to detect five different types of
92 cancer up to four years prior to conventional diagnosis (13). The DELFI assay employs a
93 genome-wide analysis of ctDNA fragment profiles to increase sensitivity in early detection (14).
94 Recently, the Galleri test has emerged as a multi-cancer detection assay that analyses more than
95 100,000 methylation regions in the genome to detect over 50 cancer types and localize the tumor
96 site (15).

97 Despite their great potential, there remain several challenges that these assays must solve to
98 deliver accessible and reliable clinical adoption for the large population, including the low
99 fraction of ctDNA in the blood of early stage cancer patients, the heterogeneity of ctDNA
100 signatures from diverse cancer types, subtypes and stages (12), and the high sequencing depth
101 required. To address these challenges, recent studies have focused on multi-analyte approach -
102 combining genomic and nongenomic features such as methylomics and fragmentomics to
103 increase the detection of ctDNA and accuracy for TOO identification (12-14). Advances in
104 multimodal analysis approaches have led to the development of powerful screening tests that
105 enable high sensitivity and cost-effectiveness. For example, CancerSEEK uses a combined
106 approach of protein biomarkers and genetic alterations to detect and locate the presence of eight
107 types of cancers (15). In this assay, cancer-associated serum proteins play a complementary role

108 in tumor localization as cfDNA mutations are not tissue specific. However, detecting both
109 protein and genetic biomarkers are time-consuming and costly. Thus, the development of future
110 MCED tests should endeavor to deliver a screening approach with high sensitivity, specificity,
111 and TOO identification at cost-effective price to provide better clinical outcomes and treatment
112 opportunities for all cancer patients.

113 In an effort to address the challenges of early cancer detection, we have developed a multimodal
114 approach called SPOT-MAS (Screening for the Presence Of Tumor by DNA Methylation And
115 Size). This assay was previously applied to cohorts of colorectal (16) and breast cancer patients
116 (17) and demonstrated ability for early detection of these cancers at high sensitivity across
117 different cancer stages and patient age groups. In this study, we aimed to expand our multimodal
118 approach, SPOT-MAS, to comprehensively analyze methylomics, fragmentomics, DNA copy
119 number and end motifs of cfDNA and evaluate its utility to simultaneously detecting and
120 locating cancer from a single screening test. As proof of concept, we used 2,288 participants,
121 including 738 nonmetastatic patients and 1,550 healthy controls, to train and fully validate this
122 approach on five commonly diagnosed cancers, including breast, gastric, lung, colorectal, and
123 liver cancer. Our findings demonstrate that the multimodal approach of SPOT-MAS enables
124 profiling of multiple ctDNA signatures across the entire genome at low sequencing depth to
125 detect five different cancer types in their early stages. Beyond detecting the presence of cancer
126 signals, our assay was able to predict the tumor location, which is important for clinicians to fast-
127 track the follow-up diagnostic and guide necessary treatment. Thus, SPOT-MAS has the
128 potential to become a universal, simple, and cost-effective approach for early multi-cancer
129 detection in a large population.

130 **Methods**

131 **Patient enrollment**

132 This study recruited 738 cancer patients (223 breast cancer, 159 CRC, 122 liver cancer, 136 lung
133 cancer, 98 gastric cancer) and 1550 healthy subjects. All cancer patients were confirmed to have
134 one of the five cancers analyzed in this study. Cancer patients were confirmed to have cancer by
135 abnormal imaging examination and subsequent tissue biopsy confirmation of malignancy.
136 Cancer stages were determined by the TNM (Tumor, Node, Metastasis) system classification
137 according to the American Joint Committee on Cancer and the International Union for Cancer
138 Control. Our study only recruited cancer patients with non-systemic-metastatic stages (Stage I-
139 IIIA) in which cancer is localized to the primary sites and has not spread to other organs. We
140 excluded patients who were diagnosed with metastatic stage IIIB and IV cancer. All healthy
141 subjects were confirmed to have no history of cancer at the time of enrollment. They were
142 followed up at six months and one year after enrollment to ensure that they did not develop
143 cancer. Study subjects were recruited from the University of Medicine and Pharmacy, Thu Duc
144 City Hospital, Medic Medical Center, Medical Genetics Institute in Ho Chi Minh city, Vietnam,
145 National Cancer Hospital and Hanoi Medical University in Hanoi from May 2019 to December
146 2022.

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

151 **Isolation of cfDNA**

152 10 mL of blood was collected from each participant in a Cell-Free DNA BCT tube (Streck,
153 USA). Plasma was collected from blood samples after centrifugation with two rounds ($2,000 \times g$

154 for 10 min and then $16,000 \times g$ for 10 min). The plasma fraction was aliquoted for long-term
155 storage at -80°C . Cell free DNA (cfDNA) was extracted from 1 mL plasma aliquots using the
156 MagMAX Cell-Free DNA Isolation kit (ThermoFisher, USA), according to the manufacturer's
157 instructions. Extracted cfDNA was quantified by the QuantiFluor dsDNA system (Promega,
158 USA).

159 **Bisulfite conversion and library preparation**

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

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

167 DNA from library products were pooled equally, hybridized and captured using The XGen
168 hybridization and wash kit (Integrated DNA Technologies, 1072281, USA), together with our
169 customized panel of xGen Lockdown Probes including 450 regions across 18,000 CpG sites
170 (Integrated DNA Technologies, USA). The construction of panel was built as previously
171 described (16, 18, 19). After hybridization, the flow-through product was concentrated using
172 SpeedVac (N-Biotek, NB-503CIR, Korea) at 65°C . The samples were then added with the
173 hybridization master mixture (hybridization buffer, hybridization enhancer and H_2O) and
174 denatured. Biotinylated P5 and P7 probes (P5-biotin: $/5\text{Biosg}/\text{AATGATACGGCGACCACCGA}$,
175 P7-biotin: $/5\text{Biosg}/\text{CAAGCAGA AGACGGCATAACGAGAT}$) on streptavidin magnetic beads
176 (Invitrogen, CA, USA) were hybridized with the single-stranded DNA. The captured DNA

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

184 **Targeted methylation analysis (TM)**

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

$$190 \quad \text{Methylation ratio} = \frac{\text{methylated cytosine (C)}}{\text{methylated C} + \text{unmethylated C}}$$

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

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

195 The integrated bioinformatics pipeline Methy-pipe was used to analyse GWM. We carried out
196 the trimming step using Trimmomatic, removing adapter sequence and low-quality bases at
197 fragment ends (16). The methylation ratio for each bin was calculated as following equation.

198
$$\text{Methylation ratio} = \frac{\text{methylated cytosine (C)}}{\text{methylated C} + \text{unmethylated C}}$$

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

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

207 CNA analysis was performed using the R-package QDNaseq (20). We also used 1-Mb
208 segmentation strategy to analyse CNA. We excluded bins that fell into the low mappability and
209 Duke blacklist regions(21). The number of reads mapped to each bin was measured by the
210 function “binReadCounts”, and GC-content correction was conducted by the functions
211 “estimateCorrection” and “correctBins”. The final CNA feature was derived by bin-wise
212 normalizing and outlier smoothing with the functions “normalizeBins” and “smoothOutlierBins”.
213 This process resulted in a feature vector of a length of 2691 bins.

214 To identify significant DNA gain or loss between cancer and control group, CNA values in each
215 bin of cancer samples were compared with corresponding values in control samples using
216 Wilcoxon rank sum test. Bins with adjusted p-value (Benjamini-Hochberg correction) ≤ 0.05
217 were considered significant. Those with \log_2 fold change (cancer vs control) > 0 were
218 categorized as significant increase. Those with \log_2 fold change (cancer vs control) < 0 were
219 categorized as significant decrease.

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

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

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

233 **End motif analysis (EM)**

234 AdaptaseTM technology (Integrated DNA Technologies, USA) was used during library
235 preparation to ligate adapters to ssDNA fragments in a template-independent reaction (22). This
236 step involved adding a random tail to the 5' end of reverse reads. Although median length of the
237 tail was 8 bp and thus allowed trimming to obtain information for other analysis, the random-
238 length tails did not allow exact determination of the 5' end of the reverse reads. Therefore, EM
239 features were determined based on the genomic coordinate of the 5' end of the forward reads.
240 We determined the first 4-mer sequence based on the human reference genome hg19. In 256
241 possible 4-mer motifs, the frequency of each motif was calculated by dividing the number of

242 reads carrying that motif by the total number of reads, generating an EM feature vector of a
243 length of 256 for each sample.

244 **Construction of machine learning models**

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

265 sample. An extensive search was performed to evaluate the performance of all possible
266 combinations (n= 511) and the combination with highest AUC was selected as the final model.
267 The model cut-off was set at the threshold specificity of >95%. This combination model
268 performance was evaluated on an independent validation dataset to examine the model
269 classification power.

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

275 **Construction of models for TOO**

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

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

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

285 Backpropagation trained the H₂O deep neural network (DNN) (multi-layer feedforward artificial
286 neural network) (H₂O package, version 3.36.1.2) with stochastic gradient descent. The random
287 grid search was selected as previously described (16).

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

289 The model training utilized an input graph formed from a discovery dataset and a validation
290 dataset as transudative setting (19) comprising patients diagnosed with five types of cancer:
291 breast, colorectal (CRC), gastric, liver, and lung. The discovery dataset contains a set of sample-
292 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
293 label, and N is the number of sample-label pairs. For each X_i in the discovery dataset, a node's
294 feature vector $f = \{F_0, \dots, F_d\} \in \mathbb{R}^d$ is constructed by combining groups of features, where F_i is
295 the i th feature, d is the number of features. The same procedure was applied for the independent
296 validation dataset. To construct an interaction graph between cancer nodes, we employed the k -
297 nearest neighbors' algorithm. An interaction graph defined as $G = (V, E)$ where $V = \{X_i | i =$
298 $1, \dots, N\}$ is a node set formed by the discovery samples, and $E = \{e_{ij}\}$ is an edge set, where e_{ij}
299 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
300 defined by:

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

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

302 In accordance with (18), a Graph Convolutional Neural Network (GCNN) was constructed for
303 the purpose of tissue of origin classification. The network comprised three message-passing
304 layers, each with a hidden size of 44 and a head number of 4. Tissue of origin classification was

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

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

$$F_i > \delta_f$$

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

318 **Statistical analysis**

319 This study used either the Wilcoxon Rank Sum test or t-test to find statistically significant
320 differences between cancer and control. The Kolmogorov-Smirnov test was used to decide
321 whether two cohorts have the same statistical distribution. The Benjamini-Hochberg correction
322 was used to correct p-value for multiple comparisons (with a corrected p-value cutoff $\alpha \leq 0.05$).
323 DeLong's test was used to compare the differences between AUCs. All statistical analyses were
324 performed using R (4.1.0) packages, including ggplot2, pROC, and caret. 95% confident interval
325 (95% CI) was presented in a bracket next to a value accordingly.

326 **Results**

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

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

338 The discovery cohort comprised of 499 cancer patients (156 breast, 106 CRC, 67 gastric, 77 liver
339 and 93 lung, Table S1) and 1,076 healthy participants. The cancer group had a median age of 58
340 (range 25 to 97, Table 1) and consisted of 279 females and 220 males. The discovery healthy
341 group consisted of 599 females and 477 males, with a median age of 47 (range 18 to 84, Table
342 1). In the discovery cohort, gender ratios were similar between cancer and healthy control
343 groups, whereas cancer patients were older than controls ($p < 0.0001$, Mann-Whitney test, Table
344 1). Of the cancer patients, 10.4% were at stage I, 33.9% were at stage II, and 30.1% were at non-
345 metastatic stage IIIA. Staging information was not available for 25.7% of cancer patients, who
346 were confirmed by specialized clinicians to have non-metastatic tumors (Table 1).

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

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

356 In our recent study of SPOT-MAS, we have demonstrated that the integration of ctDNA
357 methylation and fragmentomic features can significantly improve the early detection of
358 colorectal cancer (16) and breast cancer (17). Here, we expanded the breadth of ctDNA analyses
359 by adding two sets of features including DNA copy number and end motif into SPOT-MAS to
360 maximize cancer detection rate and identify TOO. Briefly, a novel and cost-effective workflow
361 of SPOT-MAS was developed involving three main steps (Figure 1). In step 1, cfDNA was
362 isolated from peripheral blood and subjected to bisulfite conversion and adapter ligation to create
363 a single whole-genome bisulfite library of cfDNA. From this library, in step 2, a hybridization
364 reaction was performed to collect the target capture fraction (450 cancer specific regions), then
365 the whole-genome fraction was retrieved by collecting the ‘flow-through’ and hybridizing with
366 probes specific for adapter sequences of DNA library. Both the target capture fraction and whole-
367 genome fraction were sequenced to the depth of ~52X and 0.55X, respectively (Table S3). Data
368 pre-processing was performed to generate five different sets of cfDNA features, including
369 methylation changes at target regions (TM), genome-wide methylation (GWM), fragment length

370 patterns (Flen), copy number aberrations (CNA) and end motif (EM). In step 3, these features
371 were used as inputs for a two-stage model to obtain prediction outcomes. Stage 1 of our model
372 comprised of a stacking ensemble machine learning model for binary classification of cancer
373 versus healthy. Then the samples predicted as cancer were passed to stage 2 where graph
374 convolution neural network (GCNN) was adopted to predict TOO (Figure 1).

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

377 DNA methylation is an important epigenetic signature responsible for major changes in
378 regulating expression of cancer associated genes by impacting the binding of transcription
379 factors to regulatory sites and the structure of chromatin (23, 24). Of the 450 target regions
380 associated with cancer that were selected from public data (18, 19, 25), 402 regions were
381 identified as differentially methylated regions (DMRs) in cancer patients when compared to
382 healthy participants from the discovery cohort (Wilcoxon rank-sum test, p-values < 0.05, Figure
383 2A and Table S4). Of those, 339 (84.3%) regions were identified as hypermethylated ($\log_{2}FC >$
384 0), and 63 (15.7%) regions as hypomethylated in cancer samples ($\log_{2}FC < 0$, Figure 2A). We
385 next examined the genomic location of the 402 DMRs and found 100, 108, 107 and 87 DMRs
386 that were mapped to promoter, exon, intron and intergenic regions, respectively (Figure 2B). To
387 understand the relationship between the differences in methylation regions and biological
388 pathways, we performed pathway enrichment analysis using g:Profiler on hypermethylated
389 DMRs. We detected 36 enriched pathways, including 14 from Kyoto Encyclopedia of Genes and
390 Genomes (KEGG) and 22 from WikiPathway (WP) (Figure 2C and Table S5). These significant
391 pathways were known to regulate tumorigenesis of breast, gastric, hepatocellular, and colorectal
392 cancer. Therefore, the methylation changes in the targeted regions, particularly the

393 hypermethylated DMRs, mostly occur early in tumorigenesis and are crucial for distinguishing
394 early-stage cancer patients from healthy individuals.

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

396 In addition to site-specific hypermethylation, hypomethylation is a significant genome-wide
397 change that has been identified in many types of cancers (21, 26, 27). To investigate the
398 methylation changes at genome-wide level, bisulfite sequencing reads from the whole-genome
399 fraction were mapped to the human genome, split into bins of 1Mb (2,734 bins across the
400 genome), and the reads from each bin were used to calculate methylation ratio. As expected, we
401 observed a left-ward shift in the distribution of methylation ratio in cancer samples compared to
402 healthy controls, indicating global hypomethylation in the cancer genome ($p < 0.0001$, two-
403 sample Kolmogorov-Smirnov test, Figure 3A). Of these bins, we identified 1,715 (62.7%) bins
404 as significantly hypomethylated in cancer, located across 22 autosomes of the genome (Figure
405 3B, Wilcoxon rank sum test with Benjamini-Hochberg adjusting p-value < 0.05). In contrast,
406 there were only 10 bins identified as hypermethylated and mapped to chromosome 1, 2, 3, 5, 6, 7
407 and 12 in the cancer genome (Figure 3B). Therefore, our data confirmed the widespread
408 hypomethylation across the genome and this would potentially serve to distinguish cancer
409 patients from healthy controls.

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

411 Somatic copy number aberrations (CNAs) in the cancer genome are associated with the initiation
412 and progression of numerous cancers by altering transcriptional levels of both oncogenes and
413 tumor suppressor genes (28). Recent studies have shown that CNAs detection could identify and
414 quantify the fraction of ctDNA in plasma cfDNA (29-31). To examine CNAs at genome-wide

415 scale, we used 1Mb bin to determine the percentage of bins that showed significant copy number
416 gains or losses between cancer and control group. We identified 729 bins (27.1%) with a
417 significant gain and 976 bins (36.3%) with a significant loss in copy number across 22
418 chromosomes of the cancer genome (Benjamini-Hochberg adjusting p-value <0.05, Wilcoxon
419 rank sum test, Figure 4A). We noted that chromosome 8 had the highest proportion of bins with
420 CNA gains, while chromosome 22 showed the highest proportion of bins with CNA losses
421 (Figure 4B).

422 It is thought that the abnormal hypomethylation at genome-wide level is linked with somatic
423 copy number aberration (CNA), resulting in genome instability, which is an important
424 tumorigenic event (32-34). Indeed, our data showed a significant increase in levels of CNA in
425 hypomethylated bins compared to bins with unchanged methylation (p=0.024, Figure S1A).
426 Consistently, bins with CNA gains showed significant decreases in methylation as compared to
427 those with CNA losses or unchanged CNA (p<0.01, Figure S1B). In summary, SPOT-MAS
428 enables comprehensive profiling of both global differences in methylation and somatic CNA as
429 individual feature types, as well as exploring their functional links during cancer initiation and
430 development, rendering them ideal biomarkers for cancer detection.

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

432 Several studies have shown that the fragmentation pattern of cfDNA is a non-random event
433 mediated by apoptotic-dependent caspases and ctDNA fragments tend to be shorter than non-
434 cancer cfDNA (20, 35-38). One novel technical aspect of SPOT-MAS is the use of bisulfite
435 sequencing data not only for methylation but also for fragment length analysis. Certain studies
436 showed evidence of DNA degradation followed bisulfite treatment, possibly due to high
437 temperature and low pH conditions of the bisulfite conversion procedure, while other showed

438 that bisulfite sequencing affects large genomic DNA but not small size cfDNA (39-42).
439 Therefore, to demonstrate the use of bisulfite treated cfDNA for fragment length analysis, we
440 randomly selected 3 healthy controls and 9 cancer samples to perform pair-wise comparison
441 between bisulfite and non-bisulfite sequencing results. We observed a strong correlation between
442 fragment length profile of non-bisulfite and bisulfite sequencing (Pearson correlation, $R^2 > 0.9$,
443 $p < 0.0001$, Figure S2A) for all 12 tested samples, indicating the feasibility of using bisulfite
444 sequencing data for cfDNA fragment length analysis. Indeed, the fragment size distributions of
445 bisulfite-treated cfDNA in both cancer patients and control subjects showed a peak at 167 bp
446 (Figure 5A), corresponding to the length of DNA wrapped around histone (~147 bp) plus
447 linker regions (~2x10 bp), which was in good agreement with previous studies using non-
448 bisulfite cfDNA (22, 37). Importantly, our results showed that cfDNA of cancer patients was
449 more fragmented than that of healthy participants, with a higher frequency of fragments ≤ 150 bp
450 and a lower frequency of fragment > 150 bp (Figure 5A).

451 To examine whether the fragment length variation in cancer-derived cfDNA and non-cancer
452 cfDNA could be position-dependent (37), we calculated the ratios of short (≤ 150 bp) to long
453 fragments (> 150 bp) across the genome in cancer patients and healthy controls. The mean ratio
454 of short to long fragments in cancer patients was 0.29 (range 0.28 to 0.42), which was higher
455 than the mean ratio of 0.27 (range 0.26 to 0.39) for healthy controls (Figure 5B). The changes of
456 mean ratio were across 22 autosomes of the genome. Our results indicate that the SPOT-MAS
457 technology can effectively capture differences in fragmentation patterns between cancer and
458 healthy participants across the entire genome, making them potential biomarkers for the
459 detection of circulating tumor DNA in plasma.

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

461 Associated with differences in fragment length is the differences in the DNA motifs at the end of
462 each fragment as the consequences of differential cleavage between DNA in cancer cells and
463 normal cells during apoptosis (22, 43). Here, we calculated the frequencies of 256 4-mer end
464 motifs (EMs) of cfDNA fragments and compared them between cancer patients and healthy
465 participants. Consistent with the fragment length features, we also confirmed the correlation of
466 EM frequency between bisulfite and non-bisulfite sequencing results of 12 randomly selected
467 samples, suggesting that EM profiles were reserved in bisulfite treated cfDNA (Figure S2B). Of
468 the 256 4-mer EMs, we detected 78 motifs with increased frequencies and 106 motifs with
469 decreased frequencies between cancer and healthy controls (Figure 6A and Table S6).

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

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

480 In order to increase the sensitivity of early cancer detection while avoiding the high cost of deep
481 sequencing, a screening test should survey a wide range of ctDNA signatures (12). Therefore, we
482 utilized multiple ctDNA signatures to construct classification models for distinguishing cancer
483 patients from healthy individuals. To expand the feature space, we generated four additional

484 features based on fragment length, including short, long, total fragment count, and short-to-long
485 ratio, resulting in nine input feature groups (Figure 7A). For each feature group, we tested three
486 different algorithms, including random forest (RF), logistic regression (LR), or extreme gradient
487 boosting (XGB), to tune hyperparameters and select the optimal algorithms (Figure 7A). To
488 evaluate the performance of these single-feature models, we performed 20-fold cross-validation
489 on the discovery dataset and calculated “Area Under the Curve” (AUC) of the “Receiver
490 Operating Characteristic” (ROC) curve. Among the nine features, EM-based model showed the
491 highest AUC of 0.90 (95% CI: 0.89-0.92, Figure 7B) while the SHORT-based model had the
492 lowest AUC of 0.71 (95% CI: 0.69-0.74, Figure 7B).

493 To assess whether combining features could improve classification, we used two strategies to
494 construct multi-feature models. In the first strategy, all nine feature groups were concatenated
495 into a single data frame before being fed into the RF, LR, or XGB algorithms. Of the three
496 algorithms, the XGB model exhibited the best performance with an AUC of 0.88 (95% CI: 0.87-
497 0.90, Figure 7B). However, this AUC is still lower than that of the EM-based model (0.88 versus
498 0.90, Figure 7B). In the second strategy, we constructed an ensemble stacking model using
499 logistic regression to combine the prediction results of the single-feature models. We conducted
500 an exhaustive search approach to evaluate the performance of 511 possible combinations. The
501 stacking ensemble model based on combining eight features, including TM, GW, CNA, FLEN,
502 LONG, TOTAL, RATIO and EM, exhibited the best performance and outperformed the single-
503 feature models (Table S7), with an AUC of 0.93 (95% CI: 0.92-0.95, Figure 7B and Figure S3).
504 In the independent validation cohort, we obtained similar results, where the ensemble model also
505 outperformed single-feature models, with an AUC of 0.95 (95% CI: 0.93-0.96, Figure 7C).

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

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

521 Upon stratifying our dataset by gender, we found that there was no significant difference in the
522 prediction of healthy status between males and females (Figure S4A and S4C). However, in the
523 case of cancer prediction, our model demonstrated higher accuracy in males than females in both
524 the discovery and validation cohorts (Figure S4A and S4C). Notably, when breast cancer
525 samples were removed from our analysis, there was no difference in the detection rates between
526 male and female patients (Figure S4B and S4D), suggesting that the observed gender bias may
527 be attributed to the high proportion of breast cancer patients (all females) in our cohort, who
528 exhibited the lowest detection rate among the five cancer groups.

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

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

547 The ability to predict the tissue origin of ctDNA is critical for early cancer detection as this can
548 guide subsequent diagnostic tests and treatment. Previous studies have attempted to use either
549 fragment length or methylation landscapes to achieve this goal (6, 37, 44). In this study, we
550 demonstrated the ability of SPOT-MAS to identify the TOO using low-depth bisulfite
551 sequencing to generate multiple sets of cfDNA features. We first concatenated the nine sets of

552 cfDNA features into a single data frame and focused our analysis on 499 cancer patients with
553 five cancer types in the discovery cohort. We then constructed a Random Forest (RF) and two
554 neural network models (convolutional neural network and graph convolutional neural network)
555 to predict the TOO and used 10-fold cross-validation to estimate and compare the performance
556 of these models (Figure 8A and Figure S5A). The Graph Convolutional Neural Network
557 (GCNN) was chosen due to its superior performance and stability (Figure S5B and S5C and
558 Table S9).

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

569 The median accuracy for TOO identification among the five cancer types by the GCNN-based
570 multi-feature model was 0.73 (range 0.54 to 0.87) in the discovery cohort (Figure 8E). The
571 accuracy in the discovery cohort was highest for breast (0.87) and liver cancer (0.82) and lowest
572 for gastric cancer (0.54). In the validation cohort, we obtained a slightly lower accuracy with a
573 median of 0.70 (range 0.55 to 0.78). The accuracies for individual cancer types were 0.78 for
574 breast, 0.76 for liver, 0.66 for colorectal, 0.63 for lung and 0.55 for gastric cancer (Figure 8F).

575 Among the 5 cancer types, breast cancer showed the highest TOO accuracy, possibly due to the
576 highest number of important features detected by the model. In contrast, CRC and gastric cancer
577 exhibited the lowest TOO accuracy with high misprediction rates between these two cancer types
578 (0.11 and 0.19 for CRC versus gastric and gastric versus CRC, respectively). Together, our study
579 highlights the benefits of integrating multimodal analysis with the GCNN model to capture the
580 broad landscape of tissue-specific markers in different cancer types.

581 **Discussion**

582 In an era marked by a global rise in cancer-related morbidity and mortality, the development of
583 liquid biopsy screening tests that can detect and localize cancer at an early stage holds
584 tremendous potential to revolutionize cancer diagnosis and therapy. However, the low amount of
585 ctDNA fragments in plasma samples of patients with early-stage cancer as well as the molecular
586 heterogeneity of different cancer types are known as the major challenges for liquid biopsy based
587 multi-cancer detection assays. Thus, sequencing at high depth coverages is required to capture
588 enough informative cancer DNA fragments in the finite plasma sample to achieve early cancer
589 detection. In support to this notion, many groups (6, 15, 37, 45) have developed assays that
590 exploited high depth coverage of sequencing to detect ctDNA fragments in plasma of early stage
591 cancer patients. However, this strategy might not be cost effective and feasible for population
592 wide screening in developing countries. Alternatively, we argued that increasing breadth of
593 ctDNA analysis could maximize the ability to detect ctDNA fragments with heterogeneous
594 genetic and epigenetic changes at shallow sequencing depth, thus improving the sensitivity for
595 multicancer detection. To demonstrate the feasibility of this approach, we built a stacking
596 ensemble model to combine nine different ctDNA signatures and demonstrated its superior
597 performance on cancer detection in comparison to single-feature models (Figure 7B and 7C).

598 SPOT-MAS achieved a sensitivity of 72.4 % at a specificity of 97.0 % for detecting five
599 common cancer types using shallow depth sequencing. Furthermore, it can predict the tissue of
600 origin with an accuracy of 70%.

601 Previous studies have reported that methylation changes at target regions could be exploited for
602 detecting ctDNA in plasma of patients with early-stage cancer (46, 47). Consistently, in TM
603 analysis, out of 450 TM regions chosen from previous publications (18, 19), we identified 402
604 regions as significant differentially methylated regions (DMRs) in cancer patients (Figure 2A).
605 These DMRs were enriched for regulatory regions of well-known cancer-related gene families
606 such as PAX family genes, TBX family genes, FOX family genes and HOX family genes, and
607 some have previously been reported as biomarkers for noninvasive cancer diagnosis, such as
608 *SEPT9* and *SHOX2* (48, 49). In addition to the targeted hypermethylation regions, our study also
609 showed widespread hypomethylation patterns across 22 autosomes of cancer patients (Figure 3),
610 a hallmark of cancer (50). In addition to methylation alterations, recent studies have revealed that
611 the DNA copy number , fragmentomics profile (37) and end motif profile (22) at genome wide
612 scales have been shown as useful features for healthy-cancer classification. Therefore, we
613 propose that the combination of these markers might provide added value to increase the
614 performance of liquid biopsy assays. We demonstrated that the same bisulfite sequencing data
615 could be used to identify somatic CNA (Figure 4), cancer-associated fragment length (Figure 5)
616 and end motifs (Figure 6), highlighting the advantage of SPOT-MAS in capturing the broad
617 landscape of ctDNA signatures without high cost deep sequencing. For cancer-associated
618 fragment length, we pre-processed this data into five different feature tables to better reflect the
619 information embedded within the data. Overall, we integrated multiple features of ctDNA
620 including methylation, fragment length, end motif and copy number changes into a multi-cancer

621 detection model and demonstrated that this approach could distinguish healthy individuals with
622 patients from five popular cancer types. This strategy enables increased breadth of ctDNA
623 analysis at shallow sequencing depth to overcome the limitation of low amount of ctDNA
624 fragments in plasma samples as well as molecular heterogeneity of cancers.

625 The involvement and orthogonal links of the above features in the transcriptional regulation of
626 cancer-associated genes during carcinogenesis prompted us to examine whether the combination
627 of multiple cancer-specific signatures in cfDNA could improve the efficiency of cancer detection
628 (51, 52). We first determined the performance of models constructed using individual type of
629 cfDNA features. Next, by performing exhaustive searches for all possible combinations of
630 single-feature models, we identified that the stacking ensemble of seven features could achieve
631 the AUC of 0.95 (95% CI: 0.93-0.96, Figure 8C and Figure S3), which is superior to all single-
632 feature models. Moreover, this study showed that the feature of EM achieved the highest
633 performance among the five examined ctDNA signatures in discriminating cancer from healthy
634 controls (Figure S6). Importantly, we found that combining EM with other ctDNA signatures in
635 a stack model could further improve the sensitivity for detecting cancer samples, with significant
636 improvement for lung cancer patients (Figure S6A and S6B). These findings highlighted that the
637 multimodal analysis of multiple ctDNA signatures by SPOT-MAS could increase the breadth of
638 ctDNA feature analysis, thus enhancing the detection sensitivity while maintaining the low cost
639 of sample preparation and sequencing. Among the five cancer types, breast cancer showed the
640 lowest detection rate of 58.3% and 49.3% in the discovery and validation cohort, respectively.
641 Variations in detection rates among different cancer types have been previously reported (6, 15,
642 44). Consistently, it has been reported that the detection of breast cancer, particularly in early
643 stages, is challenging due to the low levels of ctDNA shedding and heterogeneity of molecular

644 subtypes of breast tumors (6). In contrast, we obtained the highest detection rate for liver cancer
645 patients with the sensitivity of 89.6% and 91.1% in the discovery and validation cohort,
646 respectively. Our finding is in good agreement with the literature showing that liver tumors shed
647 high amounts of ctDNA (53). Despite a slightly higher AUC value in the validation cohort
648 compared to the discovery cohort, no significant differences in AUC values were observed
649 between the two cohorts at CV of 10 or 50 ($p=0.1277$, DeLong's test). This result demonstrated
650 the advantage of a multimodal approach to enhance ctDNA detection in plasma. We also
651 conducted a survey of liquid biopsy assays to put our SPOT-MAS into the context of current
652 state-of-the art in the field. Table S11 showed that SPOT-MAS is using the lowest sequencing
653 depth approach (with a depth coverage of $\sim 0.55X$) and making up for this by integrating the
654 greatest number of cfDNA features to achieve comparable performance to other assays.

655 For TOO identification, our results showed that the graph convolutional neural network (GCNN)
656 performed the best among the models tested (Figure S5 and Figure 8). GCNN has the ability to
657 explore the similarity and mutual representation among samples, therefore achieving great
658 success in multi-class classification tasks (54, 55). Unlike the reference-based deconvolution
659 approaches (56, 57), our GCNN approach is independent of a reference methylation atlas, which
660 was developed from tissue or cell type specific methylation markers and thus may introduce bias
661 due to discordance between the methylomes of tissue gDNA and plasma cfDNA (12, 58).
662 Although the methylation changes were reported as most predictive for TOO in previous studies
663 (56, 57), our results showed the contribution of each of the 9 features for TOO identification
664 (Figure 8C). In addition to GWM, fragment ratio (RATIO) and CNA are the major contributors
665 to the discrimination of different tissue types. This finding provided additional evidence that the
666 multimodal approach capturing the breadth of tissue-specific signatures could improve the

667 accuracy of TOO identification (6). Our GCNN model achieved an accuracy of 0.70 for TOO
668 prediction in validation cohort. This was comparable to the performance of CancerLocator,
669 which was based a probabilistic distribution model of tissue specific methylation markers (59).
670 Recently, Liu et al. (6) developed a methylation atlas based method, which achieved a higher
671 accuracy of 93% for locating 50 types of cancer. However, this approach is based on deep
672 genome-wide sequencing with high depth coverage of 30X (Table S11), thus might not be a cost
673 effective approach for cancer screening in large populations, especially in low-income countries.

674 For an effective screening test, careful consideration of disease prevalence, cancer in this
675 context, is imperative. Given the low prevalence of cancers, even a small proportion of false-
676 positive test results arising from reduced assay specificity, if extrapolated to a national
677 population, could significantly escalate the need for confirmatory imaging and biopsy procedures
678 for benign abnormalities detected during screening. Thus, false positives can have substantial
679 implications for both healthcare resources and patient well-being. Conversely, a screening test
680 with high sensitivity ensures that most cancer cases are detected and minimizes delays in
681 diagnosis. To address potential limitations posed by low sensitivity in cancer screening tests, we
682 suggest that current liquid biopsy tests should be employed as a complementary approach to
683 existing diagnostic methods to enhance cancer detection rates. To be used a stand-alone test,
684 further work is required to improve its performance, with a particular emphasis on improving
685 sensitivity while preserving high specificity.

686 There are several limitations in our study. First, despite using a large dataset of 738 cancer
687 samples, there was an unequal distribution of samples among cancer types, with breast cancer
688 accounting for 30.2% (223/738, Table S2) of the total samples and gastric cancer having a much
689 smaller representation (13.3%, Table S2). As a result, our models may have been influenced by

690 this imbalance, potentially introducing bias in the training and evaluation process. Therefore,
691 future studies should consider incorporating more samples to better estimate the overall
692 performance of the SPOT-MAS test. Second, tumor staging information was not available for
693 26.8% of cancer patients (198/738) in our study. For patients with unavailable staging
694 information, their initial imaging examinations were conducted at the study hospitals. However,
695 subsequent tests and surgical procedures were performed at a different hospital, as per the
696 patients' preferences. Consequently, the original study hospitals lacked access to comprehensive
697 tumor staging data. To address this limitation, the metastasis status of these patients was obtained
698 via communication channels between the clinicians at the study hospitals and those at the surgery
699 hospitals. This enabled the retrieval of limited information, adhering to an established data-
700 sharing agreement between the two institutions. To maintain the robustness of our analysis,
701 patients diagnosed with metastatic cancer or those with indeterminate metastatic status were
702 subsequently excluded from the study. Therefore, all cancer patients recruited in this study were
703 confirmed to have non-metastatic tumors. Third, the cancer patients in both the discovery and
704 validation cohort were older than the healthy participants. Age differences could be a
705 confounding variable of methylation and could affect the model performance (60, 61). However,
706 we observed no significant association between the participants' age and model prediction scores
707 (Figure S4). Fourth, the ability of SPOT-MAS to differentiate cancer patients from those with
708 benign lesions has not been examined in this study. Fifth, this study only focused on the top 5
709 common cancer types, thus the current version of SPOT-MAS might misidentify cancer patients
710 of other types, resulting in lower sensitivity to real world application. At each research sites,
711 blood samples from both cancer patients and healthy subjects were collected in Streck Cell-
712 Free DNA BCT tubes and subsequently transported to a central laboratory located in Medical

713 Genetics Institute for cfDNA isolation, library preparation and sequencing. In a recent
714 publication (62), we have investigated the impact of logistic time and hemolysis rates of blood
715 samples collected from different clinical sites on cfDNA concentration and sequencing quality.
716 We did not observe any noticeable impact of such variations on cfDNA concentrations or
717 sequencing library yields. However, future analytical validation studies using a larger sample
718 size are required to evaluate the impact of variation in sampling technique across different
719 clinical sites on the robustness or accuracy of assay results. Lastly, this was a retrospective
720 cohort study and may be biased by the nature of this study design. In an interim 6-month report
721 of a prospective study named K-DETEK, we were encouraged by the preliminary data
722 demonstrated the ability of SPOT-MAS to detect cancer patients who exhibited no symptoms at
723 the time of testing (62). Despite these promising results, the performance of SPOT-MAS as an
724 early cancer screening test remains to be fully validated in a large, multi-center prospective study
725 with 1 to 2 years of follow up.

726 **Conclusions**

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

734

735 **References**

- 736 1. Sung H, Ferlay J, Siegel RL, Laversanne M, Soerjomataram I, Jemal A, et al. Global
737 Cancer Statistics 2020: GLOBOCAN Estimates of Incidence and Mortality Worldwide for 36
738 Cancers in 185 Countries. *CA: A Cancer Journal for Clinicians*. 2021;71(3):209-49.
- 739 2. Pham T, Bui L, Kim G, Hoang D, Tran T, Hoang M. Cancers in Vietnam-Burden and
740 Control Efforts: A Narrative Scoping Review. *Cancer Control*. 2019;26(1):1073274819863802.
- 741 3. Hawkes N. Cancer survival data emphasise importance of early diagnosis. *BMJ*.
742 2019;364:l408.
- 743 4. Kakushadze Z, Raghubanshi R, Yu W. Estimating Cost Savings from Early Cancer
744 Diagnosis. *Data [Internet]*. 2017; 2(3).
- 745 5. Sasieni P, Smittenaar R, Hubbell E, Broggio J, Neal RD, Swanton C. Modelled mortality
746 benefits of multi-cancer early detection screening in England. *Br J Cancer*. 2023;129(1):72-80.
- 747 6. Liu MC, Oxnard GR, Klein EA, Swanton C, Seiden MV. Sensitive and specific multi-
748 cancer detection and localization using methylation signatures in cell-free DNA. *Ann Oncol*.
749 2020;31(6):745-59.
- 750 7. Li J, Han X, Yu X, Xu Z, Yang G, Liu B, et al. Clinical applications of liquid biopsy as
751 prognostic and predictive biomarkers in hepatocellular carcinoma: circulating tumor cells and
752 circulating tumor DNA. *J Exp Clin Cancer Res*. 2018;37(1):213.
- 753 8. Nguyen HT, Luong BA, Tran DH, Nguyen TH, Ngo QD, Le LGH, et al. Ultra-Deep
754 Sequencing of Plasma-Circulating DNA for the Detection of Tumor-Derived Mutations in
755 Patients with Nonmetastatic Colorectal Cancer. *Cancer Invest*. 2022;40(4):354-65.
- 756 9. Gao Q, Zeng Q, Wang Z, Li C, Xu Y, Cui P, et al. Circulating cell-free DNA for cancer
757 early detection. *The Innovation*. 2022;3(4):100259.
- 758 10. Pascual J, Attard G, Bidard FC, Curigliano G, De Mattos-Arruda L, Diehn M, et al.
759 ESMO recommendations on the use of circulating tumour DNA assays for patients with cancer:
760 a report from the ESMO Precision Medicine Working Group. *Ann Oncol*. 2022;33(8):750-68.
- 761 11. Nguyen HT, Tran DH, Ngo QD, Pham HT, Tran TT, Tran VU, et al. Evaluation of a
762 Liquid Biopsy Protocol using Ultra-Deep Massive Parallel Sequencing for Detecting and
763 Quantifying Circulation Tumor DNA in Colorectal Cancer Patients. *Cancer Invest*.
764 2020;38(2):85-93.
- 765 12. Moser T, Kühberger S, Lazzeri I, Vlachos G, Heitzer E. Bridging biological cfDNA
766 features and machine learning approaches. *Trends Genet*. 2023;39(4):285-307.
- 767 13. Im YR, Tsui DWY, Diaz LA, Jr., Wan JCM. Next-Generation Liquid Biopsies:
768 Embracing Data Science in Oncology. *Trends Cancer*. 2021;7(4):283-92.
- 769 14. Zhou Q, Kang G, Jiang P, Qiao R, Lam WKJ, Yu SCY, et al. Epigenetic analysis of cell-
770 free DNA by fragmentomic profiling. *Proceedings of the National Academy of Sciences*.
771 2022;119(44):e2209852119.
- 772 15. Cohen JD, Li L, Wang Y, Thoburn C, Afsari B, Danilova L, et al. Detection and
773 localization of surgically resectable cancers with a multi-analyte blood test. *Science*.
774 2018;359(6378):926-30.
- 775 16. Nguyen HT, Khoa Huynh LA, Nguyen TV, Tran DH, Thu Tran TT, Khang Le ND, et al.
776 Multimodal analysis of ctDNA methylation and fragmentomic profiles enhances detection of
777 nonmetastatic colorectal cancer. *Future Oncol*. 2022;18(35):3895-912.

- 778 17. Pham TMQ, Phan TH, Jasmine TX, Tran TTT, Huynh LAK, Vo TL, et al. Multimodal
779 analysis of genome-wide methylation, copy number aberrations, and end motif signatures
780 enhances detection of early-stage breast cancer. *Front Oncol.* 2023;13:1127086.
- 781 18. Nguyen H-N, Cao N-PT, Van Nguyen T-C, Le KND, Nguyen DT, Nguyen Q-TT, et al.
782 Liquid biopsy uncovers distinct patterns of DNA methylation and copy number changes in
783 NSCLC patients with different EGFR-TKI resistant mutations. *Scientific Reports.*
784 2021;11(1):16436.
- 785 19. Chen X, Gole J, Gore A, He Q, Lu M, Min J, et al. Non-invasive early detection of
786 cancer four years before conventional diagnosis using a blood test. *Nature Communications.*
787 2020;11(1):3475.
- 788 20. Mouliere F, Chandrananda D, Piskorz AM, Moore EK, Morris J, Ahlborn LB, et al.
789 Enhanced detection of circulating tumor DNA by fragment size analysis. *Sci Transl Med.*
790 2018;10(466).
- 791 21. Ehrlich M. DNA methylation in cancer: too much, but also too little. *Oncogene.*
792 2002;21(35):5400-13.
- 793 22. Jiang P, Sun K, Peng W, Cheng SH, Ni M, Yeung PC, et al. Plasma DNA End-Motif
794 Profiling as a Fragmentomic Marker in Cancer, Pregnancy, and Transplantation. *Cancer*
795 *Discovery.* 2020;10(5):664-73.
- 796 23. Yin Y, Morgunova E, Jolma A, Kaasinen E, Sahu B, Khund-Sayeed S, et al. Impact of
797 cytosine methylation on DNA binding specificities of human transcription factors. *Science.*
798 2017;356(6337).
- 799 24. Buitrago D, Labrador M, Arcon JP, Lema R, Flores O, Esteve-Codina A, et al. Impact of
800 DNA methylation on 3D genome structure. *Nature Communications.* 2021;12(1):3243.
- 801 25. Phan TH, Chi Nguyen VT, Thi Pham TT, Nguyen VC, Ho TD, Quynh Pham TM, et al.
802 Circulating DNA methylation profile improves the accuracy of serum biomarkers for the
803 detection of nonmetastatic hepatocellular carcinoma. *Future Oncol.* 2022;18(39):4399-413.
- 804 26. Das PM, Singal R. DNA methylation and cancer. *J Clin Oncol.* 2004;22(22):4632-42.
- 805 27. Hoffmann MJ, Schulz WA. Causes and consequences of DNA hypomethylation in
806 human cancer. *Biochem Cell Biol.* 2005;83(3):296-321.
- 807 28. Shao X, Lv N, Liao J, Long J, Xue R, Ai N, et al. Copy number variation is highly
808 correlated with differential gene expression: a pan-cancer study. *BMC Med Genet.*
809 2019;20(1):175.
- 810 29. Baldacchino S, Grech G. Somatic copy number aberrations in metastatic patients: The
811 promise of liquid biopsies. *Semin Cancer Biol.* 2020;60:302-10.
- 812 30. Knuutila S, Aalto Y, Autio K, Björkqvist AM, El-Rifai W, Hemmer S, et al. DNA copy
813 number losses in human neoplasms. *Am J Pathol.* 1999;155(3):683-94.
- 814 31. Dereli-Öz A, Versini G, Halazonetis TD. Studies of genomic copy number changes in
815 human cancers reveal signatures of DNA replication stress. *Mol Oncol.* 2011;5(4):308-14.
- 816 32. Brennan K, Flanagan JM. Is there a link between genome-wide hypomethylation in blood
817 and cancer risk? *Cancer Prev Res (Phila).* 2012;5(12):1345-57.
- 818 33. Zhang W, Klinkebiel D, Barger CJ, Pandey S, Guda C, Miller A, et al. Global DNA
819 Hypomethylation in Epithelial Ovarian Cancer: Passive Demethylation and Association with
820 Genomic Instability. *Cancers (Basel).* 2020;12(3).
- 821 34. Chan KCA, Jiang P, Chan CWM, Sun K, Wong J, Hui EP, et al. Noninvasive detection
822 of cancer-associated genome-wide hypomethylation and copy number aberrations by plasma

- 823 DNA bisulfite sequencing. *Proceedings of the National Academy of Sciences*.
824 2013;110(47):18761-8.
- 825 35. Underhill HR, Kitzman JO, Hellwig S, Welker NC, Daza R, Baker DN, et al. Fragment
826 Length of Circulating Tumor DNA. *PLOS Genetics*. 2016;12(7):e1006162.
- 827 36. Lo YMD, Han DSC, Jiang P, Chiu RWK. Epigenetics, fragmentomics, and topology of
828 cell-free DNA in liquid biopsies. *Science*. 2021;372(6538):eaaw3616.
- 829 37. Cristiano S, Leal A, Phallen J, Fiksel J, Adleff V, Bruhm DC, et al. Genome-wide cell-
830 free DNA fragmentation in patients with cancer. *Nature*. 2019;570(7761):385-9.
- 831 38. Nguyen VC, Nguyen TH, Phan TH, Tran TT, Pham TTT, Ho TD, et al. Fragment length
832 profiles of cancer mutations enhance detection of circulating tumor DNA in patients with early-
833 stage hepatocellular carcinoma. *BMC Cancer*. 2023;23(1):233.
- 834 39. Raizis AM, Schmitt F, Jost JP. A bisulfite method of 5-methylcytosine mapping that
835 minimizes template degradation. *Anal Biochem*. 1995;226(1):161-6.
- 836 40. Tanaka K, Okamoto A. Degradation of DNA by bisulfite treatment. *Bioorg Med Chem*
837 *Lett*. 2007;17(7):1912-5.
- 838 41. Kint S, De Spiegelaere W, De Kesel J, Vandekerckhove L, Van Crieckinge W. Evaluation
839 of bisulfite kits for DNA methylation profiling in terms of DNA fragmentation and DNA
840 recovery using digital PCR. *PLoS One*. 2018;13(6):e0199091.
- 841 42. Ehrich M, Zoll S, Sur S, van den Boom D. A new method for accurate assessment of
842 DNA quality after bisulfite treatment. *Nucleic Acids Res*. 2007;35(5):e29.
- 843 43. Jin C, Liu X, Zheng W, Su L, Liu Y, Guo X, et al. Characterization of fragment sizes,
844 copy number aberrations and 4-mer end motifs in cell-free DNA of hepatocellular carcinoma for
845 enhanced liquid biopsy-based cancer detection. *Mol Oncol*. 2021;15(9):2377-89.
- 846 44. Klein EA, Richards D, Cohn A, Tummala M, Lapham R, Cosgrove D, et al. Clinical
847 validation of a targeted methylation-based multi-cancer early detection test using an independent
848 validation set. *Ann Oncol*. 2021;32(9):1167-77.
- 849 45. Stackpole ML, Zeng W, Li S, Liu C-C, Zhou Y, He S, et al. Cost-effective methylome
850 sequencing of cell-free DNA for accurately detecting and locating cancer. *Nature*
851 *Communications*. 2022;13(1):5566.
- 852 46. Xu RH, Wei W, Krawczyk M, Wang W, Luo H, Flagg K, et al. Circulating tumour DNA
853 methylation markers for diagnosis and prognosis of hepatocellular carcinoma. *Nat Mater*.
854 2017;16(11):1155-61.
- 855 47. Luo H, Zhao Q, Wei W, Zheng L, Yi S, Li G, et al. Circulating tumor DNA methylation
856 profiles enable early diagnosis, prognosis prediction, and screening for colorectal cancer. *Sci*
857 *Transl Med*. 2020;12(524).
- 858 48. Ilse P, Biesterfeld S, Pomjanski N, Wrobel C, Schramm M. Analysis of
859 SHOX2 Methylation as an Aid to Cytology in Lung Cancer Diagnosis.
860 *Cancer Genomics - Proteomics*. 2014;11(5):251.
- 861 49. Warren JD, Xiong W, Bunker AM, Vaughn CP, Furtado LV, Roberts WL, et al. Septin 9
862 methylated DNA is a sensitive and specific blood test for colorectal cancer. *BMC Med*.
863 2011;9:133.
- 864 50. Jones PA, Ohtani H, Chakravarthy A, De Carvalho DD. Epigenetic therapy in immune-
865 oncology. *Nat Rev Cancer*. 2019;19(3):151-61.
- 866 51. Ulz P, Perakis S, Zhou Q, Moser T, Belic J, Lazzeri I, et al. Inference of transcription
867 factor binding from cell-free DNA enables tumor subtype prediction and early detection. *Nat*
868 *Commun*. 2019;10(1):4666.

- 869 52. Ivanov M, Baranova A, Butler T, Spellman P, Mileyko V. Non-random fragmentation
870 patterns in circulating cell-free DNA reflect epigenetic regulation. *BMC Genomics*. 2015;16
871 Suppl 13(Suppl 13):S1.
- 872 53. Caggiano C, Celona B, Garton F, Mefford J, Black BL, Henderson R, et al.
873 Comprehensive cell type decomposition of circulating cell-free DNA with CelFiE. *Nature*
874 *Communications*. 2021;12(1):2717.
- 875 54. Yin C, Cao Y, Sun P, Zhang H, Li Z, Xu Y, et al. Molecular Subtyping of Cancer Based
876 on Robust Graph Neural Network and Multi-Omics Data Integration. *Front Genet*.
877 2022;13:884028.
- 878 55. Huang Y, Chung ACS. Disease prediction with edge-variational graph convolutional
879 networks. *Med Image Anal*. 2022;77:102375.
- 880 56. Moss J, Magenheimer J, Neiman D, Zemmour H, Loyfer N, Korach A, et al.
881 Comprehensive human cell-type methylation atlas reveals origins of circulating cell-free DNA in
882 health and disease. *Nature Communications*. 2018;9(1):5068.
- 883 57. Loyfer N, Magenheimer J, Peretz A, Cann G, Bredno J, Klochendler A, et al. A DNA
884 methylation atlas of normal human cell types. *Nature*. 2023;613(7943):355-64.
- 885 58. Zhou X, Cheng Z, Dong M, Liu Q, Yang W, Liu M, et al. Tumor fractions deciphered
886 from circulating cell-free DNA methylation for cancer early diagnosis. *Nature Communications*.
887 2022;13(1):7694.
- 888 59. Kang S, Li Q, Chen Q, Zhou Y, Park S, Lee G, et al. CancerLocator: non-invasive cancer
889 diagnosis and tissue-of-origin prediction using methylation profiles of cell-free DNA. *Genome*
890 *Biol*. 2017;18(1):53.
- 891 60. Yusipov I, Bacalini MG, Kalyakulina A, Krivonosov M, Pirazzini C, Gensous N, et al.
892 Age-related DNA methylation changes are sex-specific: a comprehensive assessment. *Aging*
893 (Albany NY). 2020;12(23):24057-80.
- 894 61. Field AE, Robertson NA, Wang T, Havas A, Ideker T, Adams PD. DNA Methylation
895 Clocks in Aging: Categories, Causes, and Consequences. *Mol Cell*. 2018;71(6):882-95.
- 896 62. Nguyen THH, Lu YT, Le VH, Bui VQ, Nguyen LH, Pham NH, et al. Clinical validation
897 of a ctDNA-Based Assay for Multi-Cancer Detection: An Interim Report from a Vietnamese
898 Longitudinal Prospective Cohort Study of 2795 Participants. *Cancer Investigation*.
899 2023;41(3):232-48.

900

901

902 **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

903 **Declarations**

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

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

909 **Consent for publication:**

910 Not applicable.

911 **Availability of data and materials:**

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

915 **Competing interests:**

916 The authors declare no conflict of interest.

917 **Funding:**

918 The study was funded by Gene Solutions

919 **Disclosure statement:**

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

924 **Author contribution:**

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

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

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

932 Supervision: DKT

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

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

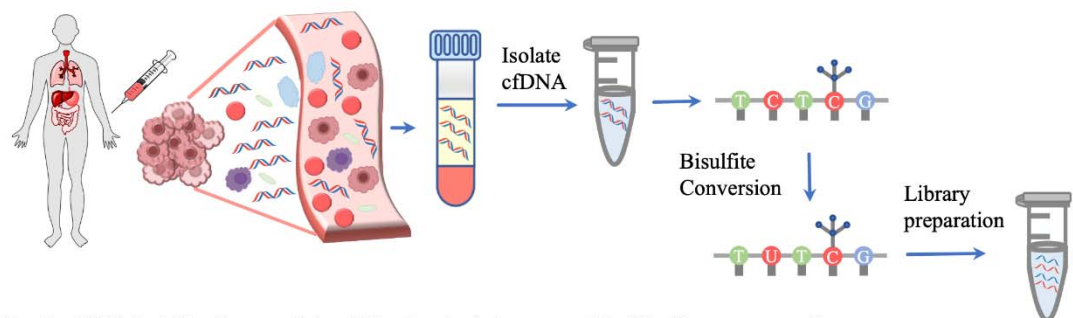
935 **Acknowledgments:**

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

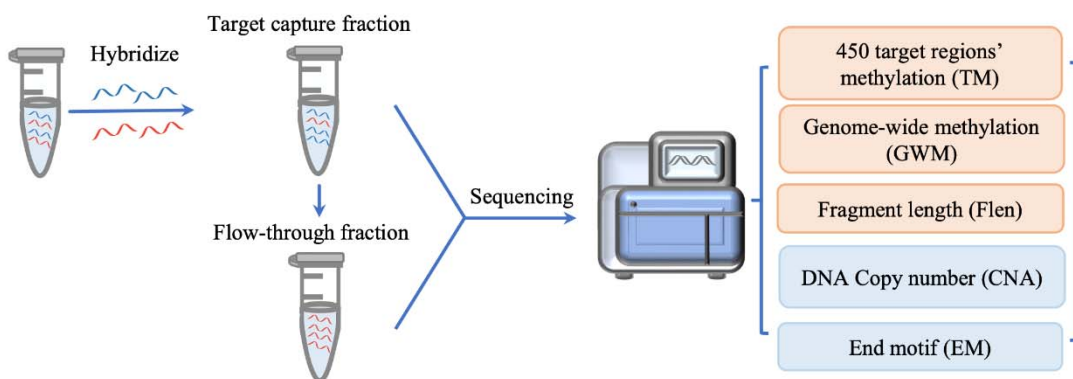
938

939

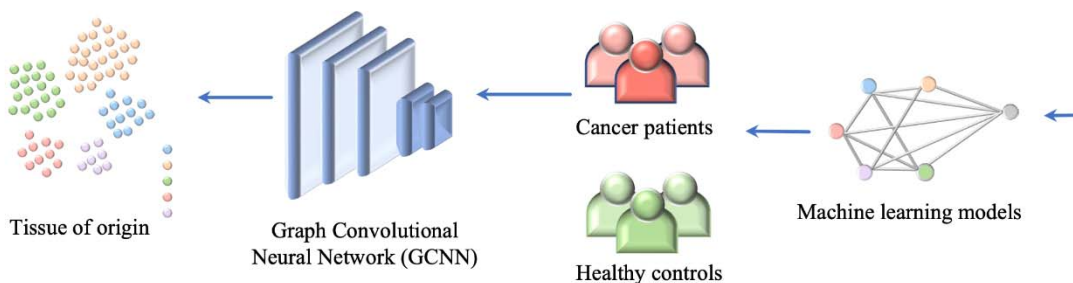
940 **Figures and Tables**



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



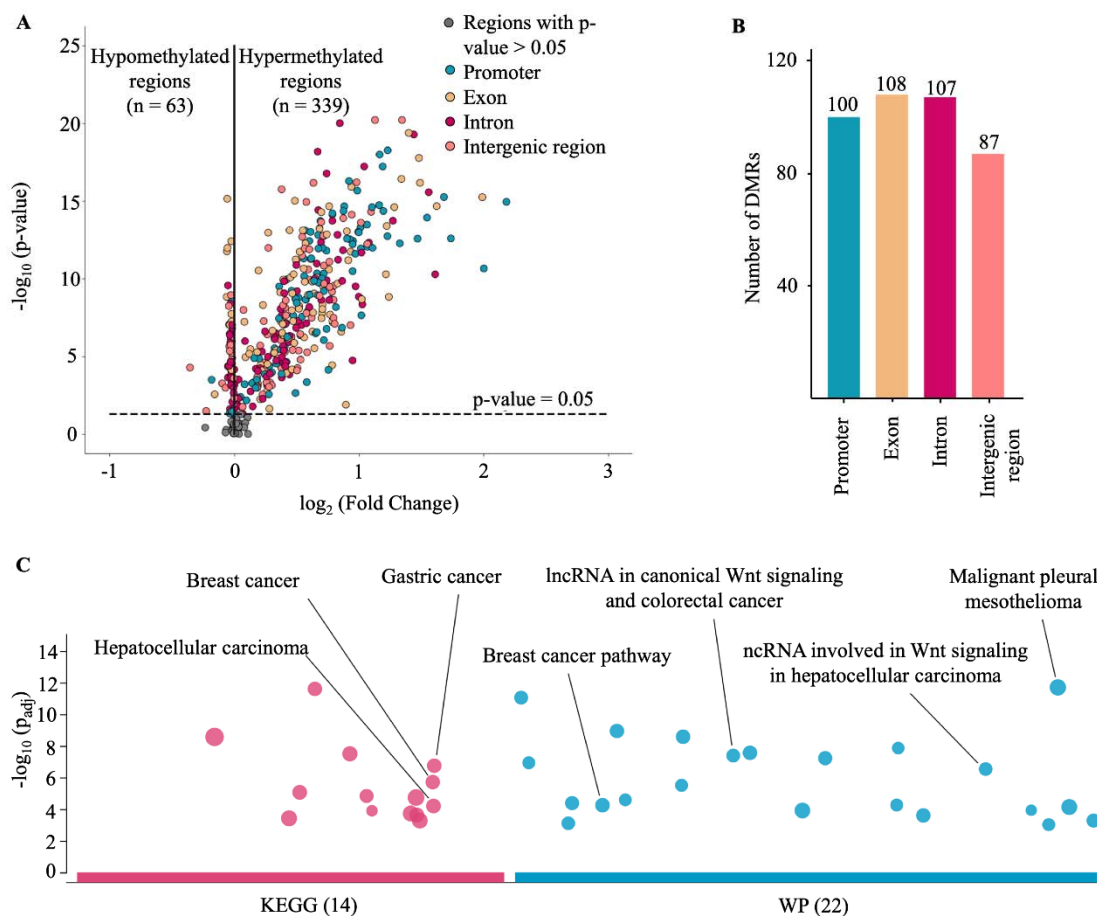
Step 2: Target and whole-genome fraction separation and sequencing



Step 3: Analysis of cfDNA signatures and model construction

941
942 **Figure 1. Workflow of SPOT-MAS assay for multi-cancer detection and localization.** There
943 are three main steps in the SPOT-MAS assay. Firstly, cfDNA is isolated from peripheral blood,
944 then treated with bisulfite conversion and adapter ligation to make whole-genome bisulfite
945 cfDNA library. Secondly, whole-genome bisulfite cfDNA library is subjected to hybridization by
946 probes specific for 450 target regions to collect the target capture fraction. The whole-genome

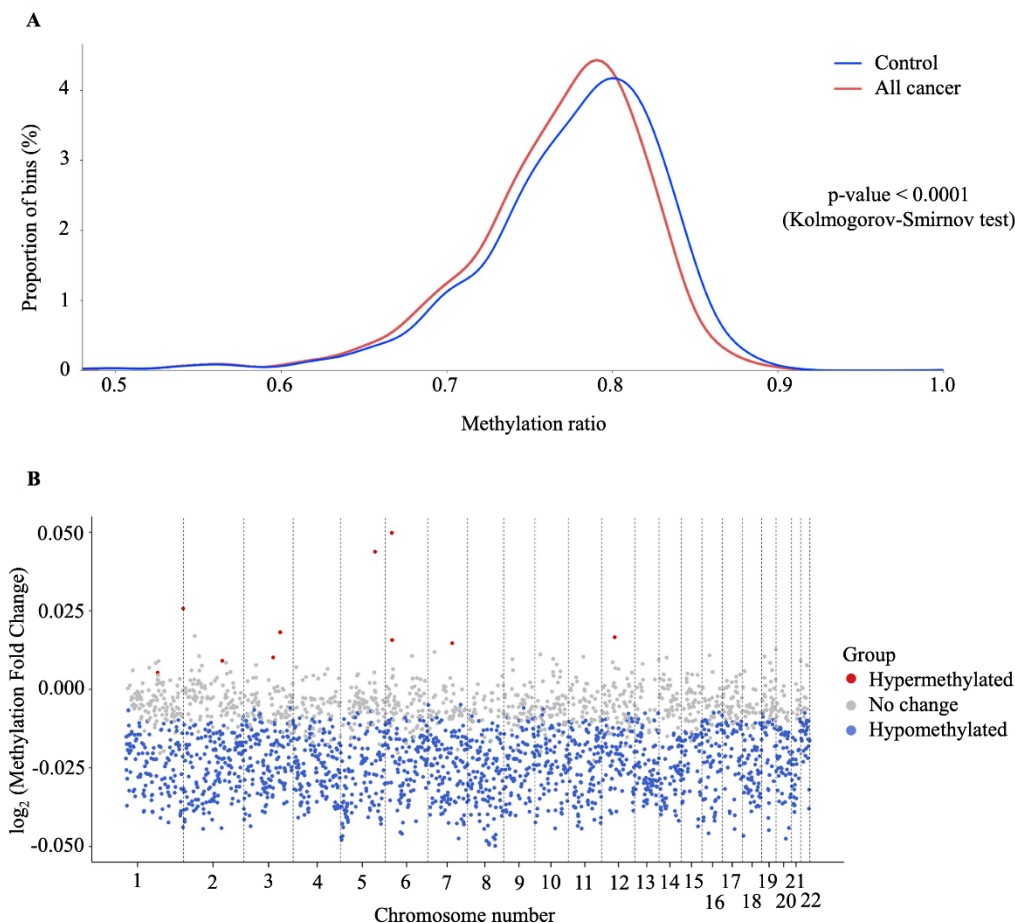
947 fraction was retrieved by collecting the ‘flow-through’ and hybridized with probes specific for
948 adapter sequences of DNA library. Both the target capture and whole-genome fractions were
949 subjected to massive parallel sequencing and the resulting data were pre-processed into five
950 different features of cfDNA: Target methylation (TM), genome-wide methylation (GWM),
951 fragment length profile (Flen), DNA copy number (CNA) and end motif (EM). Finally, machine
952 learning models and graph convolutional neural networks are adopted for classification of cancer
953 status and identification tissue of origin.



954

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

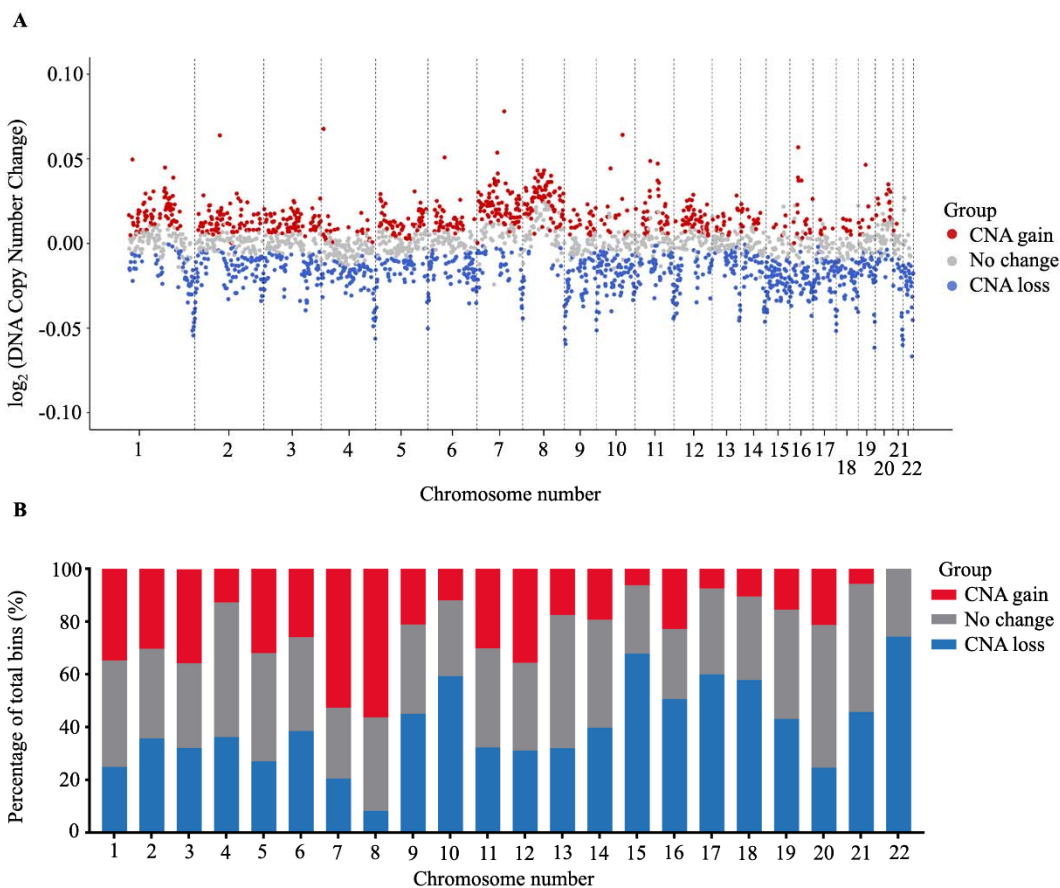
962



963

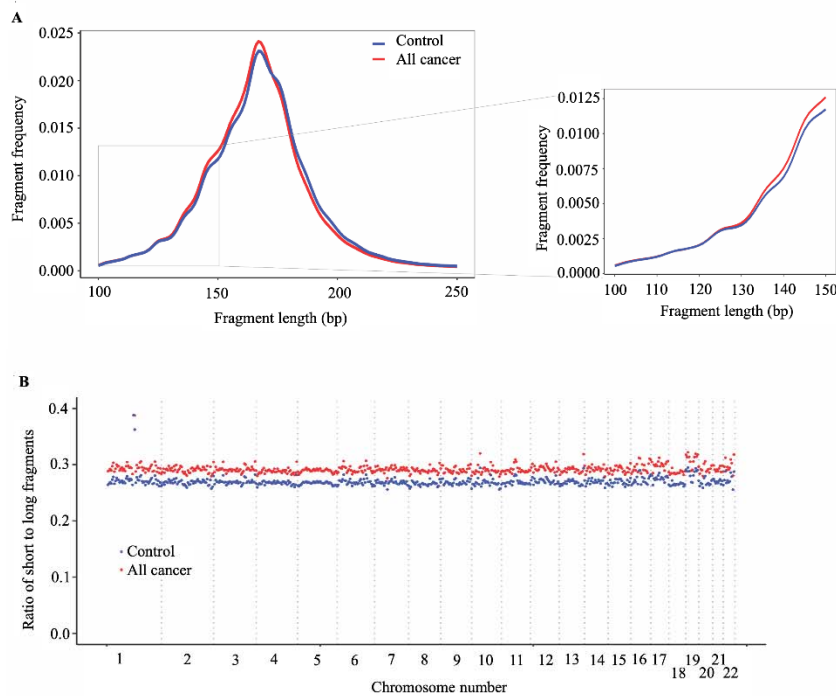
964 **Figure 3. Genome-wide methylation changes in cfDNA of cancer patients.** (A) Density plot
965 showing the distribution of genome-wide methylation ratio for all cancer patients (red curve, n=
966 499) and healthy participants (blue curve, n= 1,076). The left-ward shift in cancer samples
967 indicates global hypomethylation in the cancer genome ($p < 0.0001$, two-sample Kolmogorov-
968 Smirnov test). (B) \log_2 fold change of methylation ratio between cancer patients and healthy
969 participants in each bin across 22 chromosomes. Each dot indicates a bin, identified as
970 hypermethylated (red), hypomethylated (blue), or no significant change in methylation (grey).

971



972

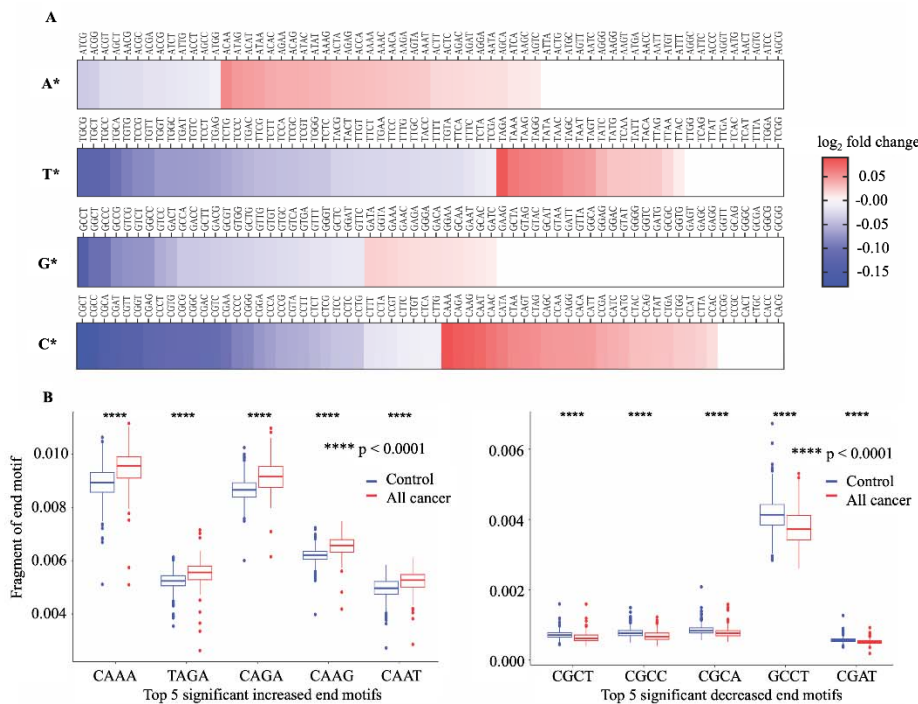
973 **Figure 4. Analysis of copy number aberration (CNA) in cfDNA.** (A) \log_2 fold change of
974 DNA copy number in each bin across 22 autosomes between 499 cancer patients and 1,076
975 healthy participants in the discovery cohort. Each dot represents a bin identified as gain (red),
976 loss (blue) or no change (grey) in copy number. (B) Proportions of different CNA bins in each
977 autosomes.



978

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

984



985

986 **Figure 6. Differences in 4-mer end motif between cancer and healthy cfDNA.** (A) Heatmap

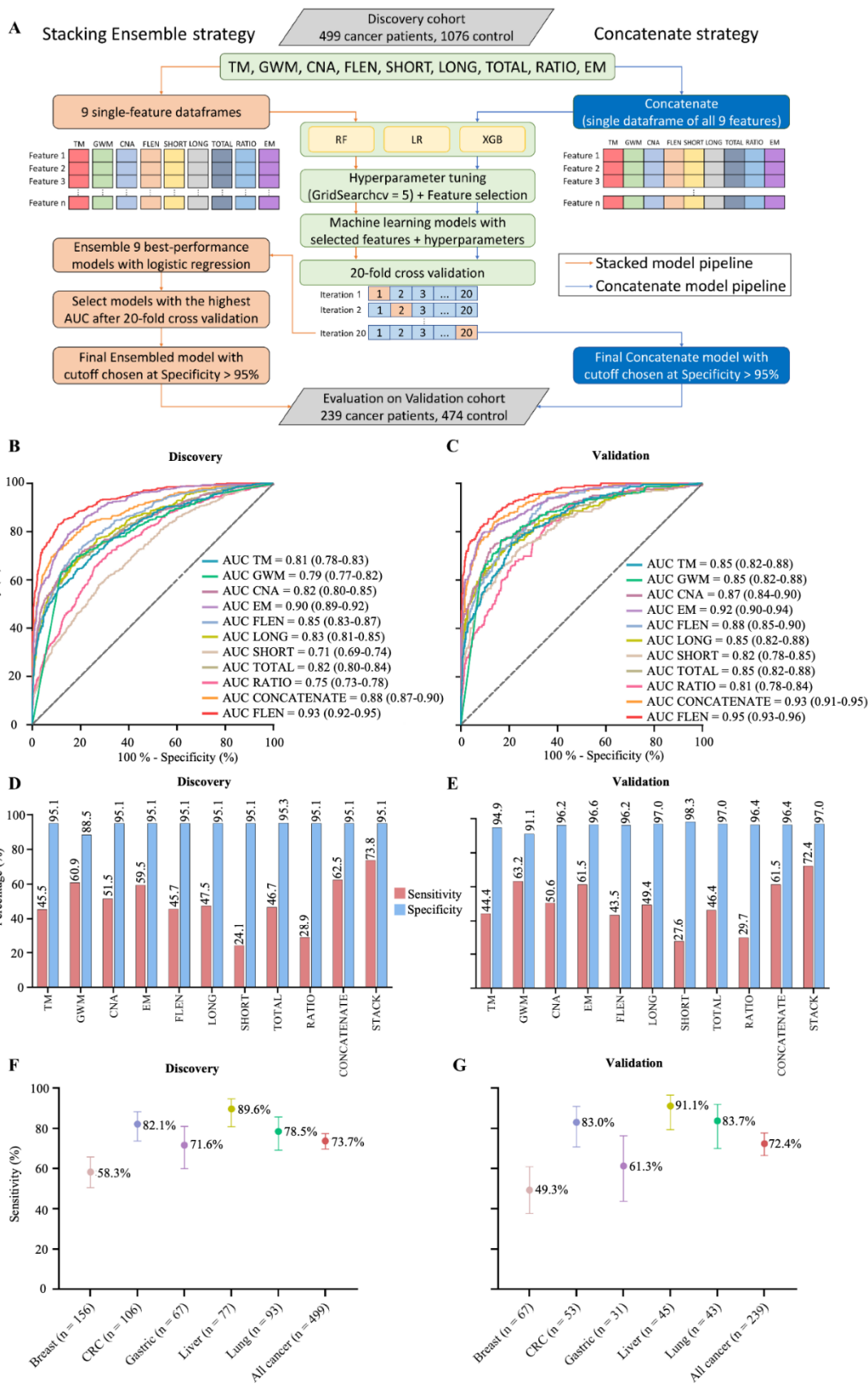
987 shows log₂ fold change of 256 4-mer end motifs in cancer patients (n=499) compared to healthy

988 controls (n=1,076). (B) Box plots showing the top ten motifs with significant differences in

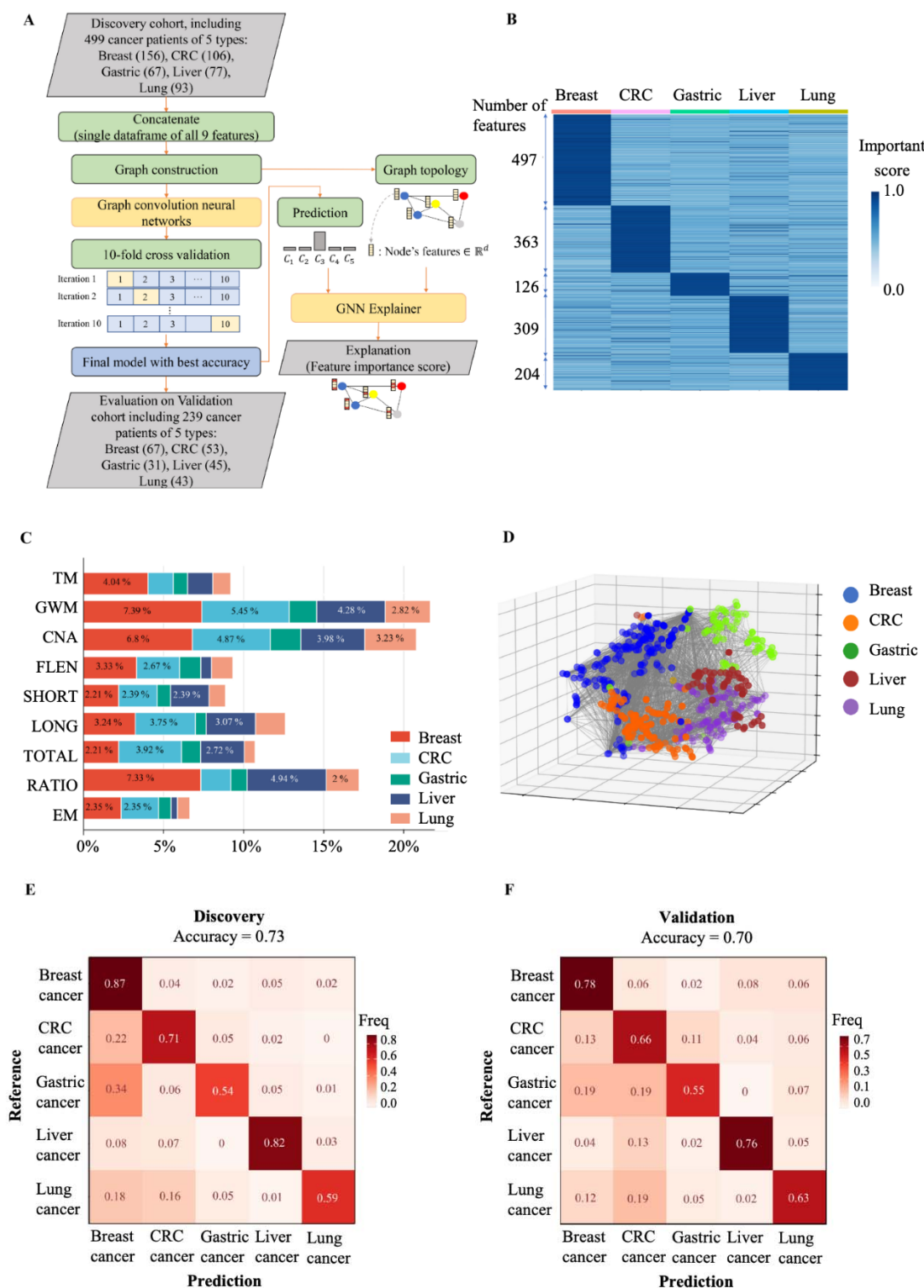
989 frequency between cancer patients (red) and healthy controls (blue) using Wilcoxon rank-sum

990 test with Bonferroni-adjusted p-value < 0.0001.

991



993 **Figure 7. Model construction and performance validation for SPOT-MAS.** (A) Two model
994 construction strategies for cancer detection. (B) and (C) ROC curves comparing the performance
995 of single-feature models, and two combination models (concatenate and ensemble stacking) in
996 the discovery (B) and validation cohorts (C). (D) and (E) Bar charts showing the specificity and
997 sensitivity of single-feature models and two combination models (concatenate and ensemble
998 stacking) in the discovery (D) and validation cohorts (E). (F) and (G) Dot plots showing the
999 sensitivity of SPOT-MAS assay in detection of 5 different cancer types in the discovery (F) and
1000 validation cohorts (G). The points and error bars represent the average sensitivity over 20 runs
1001 and 95% confidence intervals. Feature abbreviations as follows: TM – target methylation
1002 density, GWM – genome-wide methylation density, CNA – copy number aberration, EM – 4-
1003 mer end motif, FLEN – fragment length distribution, LONG – long fragment count, SHORT –
1004 short fragment count, TOTAL – all fragment count, RATIO – ratio of short/long fragment.
1005



1006

1007 **Figure 8. The performance of SPOT-MAS assay in prediction of the tissue of origin. (A)**

1008 Model construction strategy to predict tissue of origin by combining nine sets of cfDNA features

1009 using graph convolutional neural networks. (B) Heatmap shows feature important scores of five
1010 cancer types. (C) Bar chart indicates the contribution of important features for classifying five
1011 different cancers. (D) Three dimensions graph represents the classification of five cancer types.
1012 (E) and (F) Cross-tables show agreement between the prediction (x-axis) and the reference (y-
1013 axis) to predict tissue of origin in the discovery cohort (E) and validation cohort (F).

1014

1015 **Table 1.** Summary of clinical features of 738 cancer patients and 1,550 healthy controls in
 1016 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%			

1017 [#] P-values from Chi-square test; ^{##} P-values from Mann-Whitney test

1018