



28 better performance for classification of the two groups and combination of genetic  
29 and epigenetic features of cfDNA along with serum protein marker further  
30 improved the classification accuracy. We also identified novel methylation-based  
31 prognostic markers and showed that an integrated model that combined cfDNA  
32 mutational status and methylation-based prognostic markers improved prediction  
33 for lung cancer survival. Our results highlight the potential of the multi-analyte  
34 assay for non-invasive lung cancer diagnosis and prognosis.

35

## 36 **Introduction**

37 Lung cancer (LC), with the highest incidence and mortality rates among cancers  
38 worldwide, is the leading cause of death in many countries including China [1].  
39 The stage at which lung cancer is diagnosed has a significant impact on the  
40 prognosis of this disease. A study showed that the 5-year overall survival rate was  
41 57.4% for localized lung and bronchus cancers and merely 5.2% for distant ones  
42 [2]. However, timely detection of lung cancer remains difficult since patients are  
43 often asymptomatic at an early stage of the disease.

44

45 Low-dose computed tomography (LDCT), as a replacement of chest radiography,  
46 is the most extensively recommended lung cancer screening method currently [3].  
47 Its effectiveness has been proved by the National Lung Screening Trial (NLST)  
48 which demonstrated a relative reduction of 20.0% in lung cancer mortality with  
49 this approach [4]. However, LDCT as a screening method poses radiation risk. The  
50 cumulative radiation exposure of a participant following the current lung cancer  
51 screening protocols over 30 years could reach 420 mSv, which exceed those  
52 among nuclear power workers as well as atomic bomb survivors [5]. Additionally,  
53 the false-positive rate of LDCT can be up to 50%, while the positive predictive  
54 value could be as low as 2.4% [6]. This is due to the difficulty in distinguishing

55 between malignant and benign lung nodules by CT scans [7]. The resultant  
56 overdiagnosis and overtreatment could potentially lead to adverse medical events  
57 [8]. Therefore, new screening technologies for overcoming these drawbacks are  
58 required.

59  
60 Derived from tumor cells, circulating tumor DNA (ctDNA) in plasma of cancer  
61 patients provides valuable information about cancer and also holds great promise  
62 for non-invasive early cancer detection [9-13]. However, since ctDNA is diluted by  
63 circulating cell-free DNA (cfDNA) of noncancerous origins, its detection poses  
64 significant challenges especially during early stages of cancer when the tumor  
65 mass is small [14,15]. Notably, ctDNA contains both genetic and epigenetic  
66 information that may derive from the tumor, including but are not limited to  
67 mutation spectrum, copy number variation (CNV), changes in genomic  
68 methylation level, and fragmentation patterns [13,16-18]. Therefore, it is an  
69 attractive hypothesis that simultaneous analysis of multiple features may improve  
70 ctDNA detection. Nevertheless, previous studies on early cancer detection have  
71 mostly focused on a single feature of the ctDNA, such as cancer driver gene  
72 mutations or alterations in the methylome [14,19-21].

73  
74 In this study, we have developed a set of experimental and computational tools to  
75 measure both genetic and epigenetic signals from plasma cfDNA of lung cancer  
76 (LC) patients as well as patients bearing benign lung lesions (BLN) using high-  
77 throughput sequencing, aiming to explore the potential utility of blood-based  
78 biomarkers for lung cancer diagnosis and for prediction tumor recurrence risk.

79

## 80 **Methods**

### 81 **Patients enrolled and samples collected in this study**

82 Between December 2013 and December 2018, 128 LC and 94 BLN patients were  
83 enrolled in this study at the Peking University People's Hospital, Beijing, China,  
84 with the informed consent form signed by every participant. This study was  
85 approved by the Ethics Committee of Peking University People's Hospital  
86 (No.2017PHB106-01). The histopathological classification was based on the 2015  
87 World Health Organization classification [22]. 4-8 mL blood was collected from  
88 the participants before surgery into 10 mL K2EDTA tubes (BD, 366643) and  
89 stored at room temperature. Plasma separation was performed within 4 hours after  
90 collection by centrifugation at 1,600×g for 10 minutes and then at 16,000×g for  
91 another 10 minutes at room temperature. Separated plasma was stored at -80 °C  
92 until DNA extraction. 25 pairs of lung cancer tissues and adjacent normal tissues  
93 were collected during surgery at stored at -80 °C.

94

#### 95 **DNA Extraction and Quality Control**

96 Plasma cfDNA extraction was conducted by MagPure Circulating DNA Maxi Kit  
97 (Magen, 12917PC-100) following the manufacturer's instructions with some  
98 modifications. The concentration of cfDNA was measured using the Qubit™  
99 dsDNA HS Assay Kit (Thermo Fisher Scientific, Q32854). The quality of cfDNA  
100 was analyzed by Agilent High Sensitivity DNA Kit (Agilent Technologies, 5067-  
101 4626) and Agilent 2100 Bioanalyzer (Agilent Technologies). cfDNA samples with  
102 excessive high molecular weight nucleic acids were considered as contaminated by  
103 white blood cell genomic DNA (WBC gDNA) and were excluded from further  
104 analysis. gDNA was extracted from WBC, lung cancer tissues, and normal tissue  
105 adjacent to the tumor (NAT) using MagPure Buffy Coat DNA Midi KF Kit  
106 (Magen, D3537-02) per manufacturer's instruction, and DNA concentration was  
107 measured by Qubit™ dsDNA HS Assay Kit.

108

109 **Capture panel design for targeted ultra-deep Next Generation Sequencing**  
110 **(NGS)**

111 We used a 139-gene pan-cancer panel for targeted ultra-deep sequencing. Targeted  
112 genes and exons were selected based on mutation frequency in the The Cancer  
113 Genome Atlas (TCGA) database [23] and the COSMIC database of somatic  
114 mutations in cancer[24], prioritizing cancer driver genes [25]), and exons with  
115 TCGA or COSMIC hotspot mutations.

116

117 **Library preparation for targeted ultra-deep NGS**

118 To reduce noises that may derive from PCR and/or sequencing errors, we used a  
119 duplex unique molecular identifier (UMI) strategy in library preparation, adapted  
120 from a previous study [26]. Briefly, cfDNA was end-repaired and ligated to  
121 sequencing adapters, and index PCR was performed followed by purification by  
122 Agencourt AMPure XP beads (Beckman Coulter, A63882). WBC gDNA was  
123 processed in the same way except for it was fragmented by sonication before  
124 library preparation.

125

126 Target capture reactions were performed using xGen® Lockdown® Reagents (IDT  
127 technologies) per manufacturer's instruction. Captured Libraries were amplified in  
128 a 50 µL PCR mix composed of 25 µL 2× KAPA HiFi Hot Start Ready Mix, 5 µL  
129 PCR primer pair (10 µM) and 20 µL beads suspensions with the following cycling  
130 conditions: 45s at 98°C, followed by 13 cycles of 98°C for 15 s, 60°C for 30 s, and  
131 72°C for 30 s; final extension was performed at 72°C for 1min. Libraries were  
132 purified by Agencourt AMPure XP beads, quantified by Qubit™ dsDNA HS  
133 Assay Kit, and sequenced on MGISEQ-2000 (MGI Tech) using 2×100 paired-end  
134 sequencing.

135

### 136 **Library preparation for targeted bisulfite sequencing**

137 To improve the quality of cfDNA whole-genome bisulfite sequencing (WGBS)  
138 libraries, we adopted a single-stranded DNA (ssDNA) library preparation strategy.  
139 Briefly, bisulfite conversion was performed on input DNA using EZ DNA  
140 Methylation-Gold™ Kit (Zymo Research, D5006) per manufacturer's instructions.  
141 Next, bisulfite-converted ssDNA was ligated to sequencing adaptors as described  
142 previously [27]. gDNA extracted from lung cancer or normal tissues was  
143 fragmented by sonication before library preparation.

144 Targeted capture reactions of the WGBS libraries were performed using SeqCap  
145 Epi CpGiant Probes (Roche) following the manufacturer's instruction. Captured  
146 libraries were amplified and sequenced on MGISEQ-2000 using 2×100 paired-end  
147 sequencing.

148

### 149 **Variant analysis**

150 Targeted sequencing data from cfDNA libraries were processed as follows: UMI  
151 sequences were trimmed from fastq data using in-house scripts and were adapter  
152 trimmed and quality trimmed using SOAPnuke-2.0.3 [28]. Reads were aligned  
153 against the human reference genome (hg19) using BWA-MEM (version 0.7.17)  
154 [29]. Candidate mutations were identified from the aligned reads using a two-step  
155 procedure: Firstly, hotspot mutations (defined as point mutations, small insertions  
156 and deletions represented in COSMIC database  
157 (<https://cancer.sanger.ac.uk/cosmic>, version 85) with  $\geq 20$  cancer cases) were  
158 identified using the in-house script and filtered using an allele fraction cutoff of  
159 0.05% (except for indels, which were not filtered). Secondly, non-hotspot  
160 mutations were identified using freebayes (version 1.1.0) [30] and filtered using an

161 allele fraction cutoff of 0.05%. These two sets of variants were combined and  
162 filtered for potential germline variants (with allele fraction  $\geq 25\%$ ) [14]. Variants  
163 were further filtered for germline mutations using a custom germline database  
164 derived from the ExAC germline variants data [31] and 1000 Genome data [32], as  
165 well as a custom false-positive database. Remaining variants were then annotated  
166 using VEP (version 95.2-0) [33]. For cfDNA samples, variants were further  
167 filtered using the following set of criteria: variants were first filtered to exclude  
168 intronic and silent mutations. For the remaining hotspot variants, only those with at  
169 least 3 supporting UMI families and at least one supporting duplex UMI family  
170 were retained (except for indels). For the remaining non-hotspot variants, only  
171 ones with at least 8 supporting UMI families and at least one supporting duplex  
172 UMI family, or ones with at least 6 supporting UMI families and at least two  
173 duplex UMI families, were retained. Non-hotspot mutations with a SIFT prediction  
174 of "tolerated" and a PolyPhen prediction of "benign" were excluded. Finally,  
175 within the remaining non-hotspot variants, only those with a SIFT score  $\leq 0.02$   
176 and a PolyPhen score  $\geq 0.95$ , or a PolyPhen score of 1, or a SIFT score of 0, were  
177 retained. For WBC samples, no further filtering was applied. To derive the final set  
178 of variants for plasma sample, cfDNA variants were filtered with variants  
179 identified from the matched WBC sample.

180

### 181 **Mutation scoring system**

182 Variants were classified and weighted according to the following arbitrarily  
183 defined tiered scoring system: COSMIC hotspots with more than 500 cancer cases  
184 were given a score of 8; TCGA hotspot variants [34] or COSMIC hotspots with  
185 more than 100 cancer cases and not in the former class were given a score of 4;  
186 COSMIC hotspots with more than 20 cancer cases and not in the former class were  
187 given a score of 2; the rest of variants were given a score of 1.



188

## 189 **Methylation data analysis**

190 Targeted bisulfite sequencing data were processed as follows. First, low-quality  
191 reads and 3' sequencing adapters were trimmed by fastp (version 0.19.7) [35].  
192 Then, pair-end reads were aligned to the hg19 reference genome using  
193 BitMapperBS (version 1.0.0.8) [36]. Only reads mapped in proper pair to a unique  
194 genomic position and spanning an insert size between 30 bp and 500 bp were  
195 retained. Next, duplicates were marked with sambamba (v0.6.8) [37]. Finally,  
196 methylation rates were calculated as  $\#C/(\#C+\#T)$  for individual CpG sites with at  
197 least 4x coverage using MethylDackel (<https://github.com/dpryan79/MethylDackel>,  
198 version 0.3.0).

199

## 200 **Identification of differentially methylated regions (DMRs)**

201 A Bayesian hierarchical model was used to detect the differential methylated loci  
202 between 25 lung cancer tissues and 25 matched normal tissues ( $p < 0.001$  and  
203  $\text{delta} > 0.2$ ) [38]. To account for the spatial correlation of methylation ratio,  
204 smoothing was applied to combine the information from proximal CpG sites to  
205 identify differentially methylated regions (DMRs). DMRs were defined as the  
206 regions satisfying the following criteria:  $\geq 50$ bp, containing  $\geq 3$  CpG sites within the  
207 region, and  $\geq 80\%$  CpG sites with significant p-values. Only hypermethylated  
208 DMRs were used in the subsequent analysis.

209

## 210 **Predictive model construction**

211 Regional methylation ratio was calculated per DMR for each cfDNA sample  
212 sequenced by targeted bisulfite sequencing and processed as features by dividing  
213 the sum of methylated cytosine by the sum of depth in the DMR. Ten-fold cross-  
214 validation was performed to validate random forest models for classifying plasma



215 cfDNA of lung cancer patients from that of patients bearing benign lung nodules  
216 using the python package scikit-learn [39].

217  
218 Feature selections were performed on the training data only, using a feature  
219 importance cutoff of 0.008. Random forest models were fitted using the selected  
220 DMRs with the parameters: number of trees=60, depth=5. The fitted models were  
221 then applied in the validation set from which the sensitivity, specificity, and area  
222 under the curve (AUC) were calculated. Multi-omics prediction models were  
223 trained and validated similarly, except that feature selections were applied to the  
224 DMR features only.

225

## 226 **Identification and validation of prognostic markers**

227 To identify methylation-based prognostic markers, samples were randomly divided  
228 into a training set and testing set using a 60/40 split. We applied the following  
229 procedure to select the potential methylation-related prognostic factors and to fit  
230 prognosis model in the training set: we first removed DMRs with a standard  
231 deviation < 0.03 from the identified lung cancer DMRs as mentioned above since  
232 less variant features provided limited information; we then used the selected DMRs  
233 to fit a LASSO Cox proportional hazard model on OS. Through 10-fold cross-  
234 validation, we chose the tuning parameter  $\lambda$  when the partial likelihood deviance  
235 reached the lowest, from which DMRs were further filtered and the coefficients of  
236 the each DMR were obtained. We calculated the methylation-based prognostic  
237 score (MPS) for each individual as the sum of the products of the DMR  
238 methylation level and its coefficient and combined the mutation score (wSUMAF)  
239 with the MPS as the multi-omics score. We then assessed the association of lung  
240 cancer prognosis with mutation score and multi-omics score separately in the  
241 training set and testing set. Kaplan–Meier curves were plotted for each analysis.

242 Finally, two separate multivariate Cox proportional hazard models were built on  
243 wSUMAF only and both wSUMAF and MPS with adjustment of age, stage,  
244 histological type and smoking status in the testing set. To avoid information loss  
245 through categorization, both wSUMAF and MPS were analyzed as continuous  
246 variables in the multivariate Cox regression. To compare the performance of two  
247 models, the incident cases / dynamic controls ROC curve was plotted [40]. The R  
248 package of glmnet, survival, survminer ([https://CRAN.R-](https://CRAN.R-project.org/package=survminer)  
249 [project.org/package=survminer](https://CRAN.R-project.org/package=survminer)), risksetROC were used. The analyses procedure is  
250 summarized in the Supplementary Figure 19.

251

## 252 **Results**

### 253 **Study design and patients**

254 In this study, we performed a comprehensive analysis of sequence alterations and  
255 methylation pattern of plasma cfDNA as well as levels of serum protein markers  
256 from lung cancer patients and patients bearing benign lung nodules, in order to  
257 explore the possibility of using these features to non-invasively distinguish  
258 between malignant and benign lung nodules (Figure 1A). Blood samples were  
259 collected from 128 lung cancer (LC) and 94 benign lung nodule (BLN) patients  
260 (Table 1). As expected, LC patients had significantly higher mean plasma cfDNA  
261 level ( $20.53 \pm 1.04$ ng/ml) than BLN patients ( $13.78 \pm 1.14$  ng/ml,  $p=2.14E-05$ ,  
262 student's t-test) (Supplementary Figure1), in accordance with the previous report  
263 [41].

264

### 265 **Targeted ultra-deep NGS detected distinct mutational spectrum of plasma** 266 **cfDNA and WBC gDNA**

267 To profile sequence alterations carried by cfDNA, we performed targeted ultra-  
268 deep NGS on plasma cfDNA extracted from 111 LC patients and 78 BLN patients

269 (Supplementary Table 1) using a panel covering exons of 139 cancer driver genes  
270 selected based on TCGA and COSMIC databases (Supplementary Table 2, see  
271 Methods for panel design). We achieved an average raw target sequencing depth  
272 over 50,000× and an average deduped sequencing depth over 5,000×  
273 (Supplementary Figure 2). We designed a set of stringent thresholds to identify the  
274 most reliable variants, based on the number of supporting UMI families and duplex  
275 UMI families, the allele fractions, and function predictions (see Methods for  
276 details). Potential germline variants were also removed before downstream  
277 analysis. To test the limit of detection (LOD) and evaluate the accuracy of our  
278 method, we first performed spike-in experiments using a reference standard  
279 containing 8 single-nucleotide variants (SNVs) and cfDNA from two healthy  
280 individuals by the method previously reported [42]. Results indicated that our  
281 targeted ultra-deep NGS method could efficiently detect mutations with variant  
282 allele frequencies (VAFs) of 0.1% and 0.25%, with a sensitivity of 91.7% (22/24)  
283 and 95.5% (21/22) respectively (Figure 1B and Supplementary Figure 3). The  
284 VAFs of identified sequence alterations using this method ranged from 0.03% to  
285 6.82% with a median of 0.16% for LC patients and from 0.05% to 2.00% with a  
286 median of 0.22% for BLN patients (Figure 1C). In total, 193 and 46 mutations  
287 were detected in 75 (out of 111, 68%) LC patients and 33 (out of 78, 42%) BLN  
288 plasma cfDNA, respectively (Supplementary Figure 4 and 5). As expected, cfDNA  
289 of LC patients appears to harbor more sequence alterations than that of BLN  
290 patients (Supplementary Figure 6).

291  
292 Recent studies suggest that some of the variants found in cfDNA may derive from  
293 the process of clonal hematopoiesis and confound the analysis [43]. To address this,  
294 genomic DNA (gDNA) of white blood cell (WBC) from cfDNA mutation-positive  
295 participants were also sequenced with ultra-deep targeted sequencing (see

296 Methods). Non-synonymous variants were detected in WBC of 73 (out of 75, 97%)  
297 LC patients and 33 (out of 33, 100%) BLN patients respectively (Supplementary  
298 Figure 7 and 8). The AFs of variants observed in WBC samples were mostly less  
299 than 1%, ranging from 0.04% to 7.10% (Figure 1D). Within these WBC-shared  
300 cfDNA variants, the most frequently mutated genes included *TP53*, *CBL*, *APOB*  
301 and *CSMD3* for LC plasma, and *CBL*, *CSMD3*, and *STAT3* for BLN plasma  
302 (Supplementary Figure 9). Of these, *TP53* and *CBL* are regarded as canonical CH  
303 genes (other canonical CH genes such as *DNMT3A*, *TET2* and *ASXL1* were not  
304 included our targeted panel) [43]. Moreover, AFs of variants shared by plasma  
305 cfDNA and matched WBC samples are highly correlated (Figure 1E), suggesting  
306 that these mutations indeed originated from WBC and should be removed for  
307 downstream analysis. Notably, a number of these mutations were hotspot  
308 mutations of cancer driver genes (defined as variants with  $\geq 20$  reported cases in  
309 the COSMIC database (Supplementary Figure 10). The percentages of cfDNA  
310 variants matching corresponding WBC sample were 20.7% (40 out of 193) for LC  
311 cfDNA and 39.1% (18 out of 46) for BLN cfDNA, suggesting that a significant  
312 portion of cfDNA variants derive from clonal hematopoiesis, especially in BLN  
313 plasma ( $p=8.89E-03$ , chi-squared test).

314  
315 After filtering for variants potentially derived from clonal hematopoiesis, 153  
316 variants remained in 67 (out of 111, 60.36%) cfDNA samples from LC patients  
317 (Figure 2B and Supplementary Table 3), with AFs ranging from 0.03% to 6.00%  
318 (median was 0.13%, Figure 2A-C and Supplementary Figure 11). The most  
319 frequent variants were missense mutations ( $n=102$ , 67%), followed by nonsense  
320 mutations ( $n=33$ , 22%). SNVs ( $n=137$ , 90%) were predominantly C>T transitions  
321 ( $n=95$ , 69%) (Figure 2C), which was a feature discovered recently in East Asian  
322 LC patients [44]. Mutation frequency analysis revealed that *TP53* was the most

323 commonly mutated gene in LC plasma (mutated in 23% of LC cfDNA samples),  
324 consistent with TCGA findings [45-47]. Other frequently mutated genes included  
325 *EGFR* (8%), *PTPN11* (8%), *APC* (7%), *APOB* (7%), *KMT2C* (5%), and *KMT2D*  
326 (5%) (Figure 2A and Supplementary Figure 12), hence identifying a spectrum  
327 mostly consistent with previous reports of lung cancer mutation spectrum  
328 (Supplementary Figure 13) [45,48-53]. As for LC subtypes, cfDNA samples from  
329 LUSC patients were more frequently mutated than that from LUAD patients  
330 (Supplementary Figure 13 and Supplementary Table 4).

331  
332 After stringent QC filtering as well as filtering for WBC-matched variants , as  
333 many as 28 mutations remained in 23 (out of 78) BLN plasma cfDNA samples  
334 (Figure 2D and Supplementary Table 5), although the percentage of positive  
335 samples was much less compared to LC plasma (29.49% vs. 60.36%,  $p=2.87E-05$ ,  
336 chi-squared test). These mutations had AFs ranging from 0.05% to 1.91% (Figure  
337 2A). The median AF (0.13%) was the same as that of LC plasma cfDNA, yet the  
338 highest AF was much less (1.91% vs. 6.00%). The most frequently mutated genes  
339 in BLN plasma were *KRAS* (5%), *CSMD3* (4%), *APC* (3%), *ATM* (3%), *KMT2D*  
340 (3%), and *TP53* (3%) (Figure 2D). Notably, 39.3% (11 out of 28) of these were  
341 COSMIC hotspot mutations (such as variants affecting the *KRAS* G12 residue)  
342 (Supplementary Table 5). These results suggest that BLN cfDNA harbored  
343 common cancer driver mutations, a phenomenon consistent with previous and  
344 recent reports that cancer driver mutations are prevalent among normal tissues  
345 [43,54,55]. The mutation spectrum of BLN plasma cfDNA is notably different  
346 from that of LC plasma: the most frequent mutations found in BLN plasma cfDNA  
347 were *KRAS* (5%) and *CSMD3* (4%). *EGFR* variants were never detected in BLN  
348 plasma. *KRAS* variants, however, had almost identical frequency in both groups.

349 *TP53* mutations were also observed in BLN plasma cfDNA albeit at a much lower  
350 frequency (3%).

351  
352 **A predictive model based on somatic mutations to distinguish LC from BLN**  
353 Next, we asked whether it would be possible to differentiate LC and BLN plasma  
354 based on their different cfDNA mutation signatures and AF distribution. To  
355 quantify the cfDNA mutational burden, we constructed a mutation score for each  
356 cfDNA sample as either a simple summation of the allele fractions of all variants  
357 identified therein (SUMAF) or a weighted sum of the allele fractions based on a set  
358 of pre-defined weights (weighted SUMAF, or wSUMAF), weighing more on  
359 TCGA hotspot cancer driver mutations and COSMIC hotspot mutations, and less  
360 on other variants (see Methods for details). We found both scoring methods  
361 produced modest classification accuracy for distinguishing LC from BLN plasma:  
362 the wSUMAF model generated an area under curve (AUC) value of 0.68, with a  
363 sensitivity of 59.5% and a specificity of 71.8% (Figure 2E and Supplementary  
364 Figure 14) and the SUMAF model had a similar performance.

365  
366 **Classification of LC and BLN plasma based on cfDNA methylation data**

367 To identify lung cancer-specific epigenetic changes, such as abnormalities in 5-mC  
368 methylome, we performed whole-genome bisulfite sequencing (WGBS) on 25  
369 pairs of LC tissue and normal tissue adjacent to the tumor (NAT) among which 21  
370 pairs were from LUAD, 3 from LUSC and 1 from LCSC (Figure 3A) . 315  
371 differentially methylated regions (DMRs) were identified using a cutoff of delta  $\beta$   
372 value greater than 0.2 and p-value less than 0.001 (see Methods), including 293  
373 hyper DMRs and 22 hypo DMRs (Figure 3B). There were a lot more hyper DMRs  
374 than hypo DMRs, consistent with the belief that genomic regulatory regions such  
375 as promoters of potential tumor suppressor genes undergo remarkable



376 hypermethylation in tumorigenesis. Gene ontology (GO) annotations revealed that  
377 the 293 hyper DMRs were significantly enriched for genes encoding DNA-binding  
378 domains and homeobox domains, as well as genes involved in the developmental  
379 and transcriptional regulation process (Figure 3C), consistent with the possibility  
380 that these genes may be involved in lung cancer development by regulating cell  
381 differentiation, and when silenced by promoter methylation, may cause cell  
382 transformation.

383  
384 Unsupervised hierarchical clustering using the regional methylation ratio of the  
385 identified DMRs perfectly separated LC tissues and NAT with the exception of a  
386 single LC sample, highlighting the pronounced epigenetic dysregulation of lung  
387 cancer cells. We did not observe notable differences between cancer stages (Figure  
388 3B), consistent with the notion that epigenomic change is an early driver of  
389 oncogenesis that persists through later stages of cancer progression.

390  
391 We next performed comprehensive analysis of 5-mC methylation profile of plasma  
392 cfDNA for 111 LC patients and 87 BLN patients using targeted bisulfite  
393 sequencing, covering 5.6 million CpG sites located within gene regions, as well as  
394 CpG islands, shelves, and shores (Supplementary Table 1). Based on the hyper  
395 DMRs we defined using tissue WGBS, we were able to build random forest  
396 models that classify LC from BLN plasma (see Methods). To estimate  
397 classification accuracy, we performed 10-fold cross-validation (CV), for which  
398 AUC was 0.75 (Supplementary Figure 15), a performance slightly better than the  
399 mutation-based model. In order to determine whether we could effectively  
400 distinguish lung cancer plasma from healthy plasma using fewer DMR markers,  
401 we also performed further feature selection. We found that by selecting DMRs  
402 with feature importance  $> 0.008$  in each random forest model for each CV, we



403 could achieve a CV AUC of 0.72 (Figure 3D). A total of 76 DMRs were retained  
404 in the final CV model. Among cancer patients, the detection sensitivity of the final  
405 CV models was higher in LUSC patients than that in LUAD samples  
406 (Supplementary Table 6).

407

### 408 **Multi-omics analysis to differentiate LC from BLN plasma**

409 Next, we attempted to integrate multi-omics features to further improve the  
410 diagnostic power of our classification model. Indeed, on 91 LC and 71 BLN  
411 cfDNA samples that had been sequenced with both targeted deep sequencing and  
412 targeted bisulfite sequencing (Supplementary Table 1), we found that combination  
413 of methylation features (based on a total of 81 DMRs selected using the same  
414 feature selection criteria described earlier) and the SUMAF mutation score  
415 achieved an AUC of 0.75, generating a sensitivity of 78.0% and specificity of 60.5%  
416 (Figure 4C), compared to an AUC of 0.68 achieved by mutation score alone  
417 (Figure 4A), and an AUC of 0.72 achieved by methylation features alone in the  
418 same set of samples using the same CV procedures (Figure 4B).

419

420 In addition, levels of 5 serum marker, CEA, CYFRA21-1, NSE, CA19-9, and  
421 CA125, were also measured in a subset of the blood samples (Supplementary  
422 Table 1). We found that among the five protein markers, only CEA level appeared  
423 to be significantly higher in LC patients than BLN patients ( $p=0.0438$ ), producing  
424 a modest AUC of 0.66 for classifying the two groups (Supplementary Figure 16  
425 and 17). Therefore, we tried to incorporate CEA into the predictive model in  
426 addition to the DMR features and mutation score in samples with complete  
427 measurements (74 LC and 60 BLN samples). The multi-omics predictive models  
428 based on SUMAF mutation score, top DMRs (a total of 81 regions) and serum  
429 CEA level achieved an AUC of 0.79, with 75.7% sensitivity and 68.3% specificity

430 (Figure 4E), which showed further improvement compared to the model without  
431 CEA on this set of samples (AUC = 0.75) (Figure 4D). Similar to classification  
432 models based solely on mutation or methylation, higher prediction accuracy was  
433 found for LUSC patients than LUAD patients in integrated bi-omics and multi-  
434 omics models (Supplementary Table 6).

435

### 436 **cfDNA mutational burden and methylation level as prognostic factors for lung** 437 **cancer**

438 We first tested whether mutational status (wSUMAF,  $<0$  vs.  $>0$ ) was associated  
439 with lung cancer overall survival (OS) among lung cancer patients. We found that  
440 among lung cancer patients, the high mutational burden was associated with a  
441 significantly worse OS (Figure 5A). Notably, among stage I lung cancer patients, a  
442 significant association was observed between mutation score and OS  
443 (Supplementary Figure 18). To further improve the performance of model  
444 prediction on lung cancer prognosis, we attempted to identify potential  
445 methylation-based prognostic biomarkers and incorporate these features into the  
446 model (Supplementary Figure 19). We divided the LC cases into training and  
447 testing set and applied a DMR selection procedure on the training set using the  
448 penalized COX regression, which identified 12 DMRs that were potentially  
449 associated with lung cancer prognosis and obtained corresponding coefficients  
450 (Supplementary Table 7, see Methods for details on the analysis procedure). Of  
451 these, DMRs of gene *FOXG1/LINC01551*, *TMEM240*, *AKR7L*, *CBLN4*, and  
452 *GCK/MYL7* appeared to be associated with a worse lung cancer prognosis while  
453 DMRs of *PRDM11*, *LOC440028/SBF2-AS1*, *GFII1*, and *ST3GAL1* were associated  
454 with a better prognosis. We then calculated a methylation-based prognostic score  
455 (MPS) for each individual as the sum of the products of the DMR methylation  
456 level and corresponding coefficient. We thereafter experimented with combining

457 the mutation score with the MPS as the bi-omics prognosis score and tested its  
458 association with survival. Patients with a high mutational burden and a high MPS  
459 were categorized as the high prognosis score group, while other patients were  
460 categorized as the low prognosis score group. Compared to patients with a low  
461 prognosis score, patients with a high prognosis score had a significantly worse OS  
462 in the testing set (Figure 5B). Finally, to avoid information loss due to  
463 categorization, we modeled both mutation score and MPS continuously (see  
464 Methods for details) and built two separate multivariate Cox proportional hazard  
465 models on wSUMAF only, as well as on combined wSUMAF and MPS with  
466 adjustment of age, stage, histological type and smoking status. Higher AUCs were  
467 obtained using the bi-omics prognosis model than the mutation only model  
468 (Supplementary Figure 20). Taken together, these results suggest that integrated  
469 genomic features have the potential to be used as better prognostic markers for  
470 lung cancer.

471

## 472 **Discussion**

473 In this study, we applied targeted ultra-deep sequencing to plasma cfDNA of LC  
474 and BLN patients. Matched WBC DNA were sequenced in parallel, which  
475 revealed that non-synonymous variants were prevalent in WBC DNA for both the  
476 LC and BLN group. Further analyses showed that a notable portion of cfDNA  
477 variants was detected in matched WBC (20.7% for LC plasma, and 39.1% for BLN  
478 plasma) and their VAFs were well correlated (Figure 1E), suggesting that these  
479 variants most likely derived from WBCs. VAFs of the majority of WBC-matched  
480 somatic mutations detected in the cfDNA were less than 1%, hence would have  
481 been missed if WBC DNA has not been sequenced with ultra-deep sequencing.  
482 These results corroborate recent findings that WBCs, carrying variants  
483 accumulated through clonal hematopoiesis (CH), constitute an important source of

484 somatic mutations found in cfDNA [43]. Notably, some of the shared variants  
485 between cfDNA and matched WBC samples were cancer hotspot mutations  
486 (Supplementary Figure 10), suggesting that CH variants may indeed significantly  
487 confound cfDNA analysis if not analyzed in parallel. Interestingly, *TP53* variants  
488 were never detected as shared variants between BLN plasma cfDNA and matched  
489 BLN WBC sample (Supplementary Figure 9). This may indicate that the  
490 accumulation of TP53 mutations through CH may be somehow related to the  
491 cancer risk of the individual.

492

493 After removing WBC derived variants, we found that cfDNA mutations were  
494 prevalent in BLN plasma cfDNA (29.5% of samples analyzed contained at least  
495 one variant). This finding is consistent with recent studies showed that benign  
496 tumors may also harbor somatic mutations, including those in cancer driver genes  
497 [56,57]. However, because we did not obtain matched BLN tissue for these plasma  
498 samples, it remained to be determined whether the mutations found in BLN cfDNA  
499 could be attributed to mutations that may have arisen in the benign lesions of the  
500 lung. If this indeed is the case, then we would expect it to be challenging to  
501 classify malignant and benign disease solely based on cfDNA mutational status.  
502 Further study would be needed to clarify whether patients bearing benign lung  
503 lesions indeed have a higher mutational burden in their cfDNA than healthy  
504 individuals.

505

506 Not surprisingly, predictive models built on mutation score alone had limited  
507 classification ability for distinguishing between LC and BLN plasma (AUC = 0.68).  
508 Some earlier studies suggested that mutational status can be used to diagnose LC  
509 from benign lung nodules with high specificity and modest sensitivity [58,59], but

510 these conclusions may have suffered from potential bias caused by limited sample  
511 sizes used in the study. Our results were obtained from a larger sample size (128  
512 LC and 94 BLN plasma) and showed that diagnostic model based on mutational  
513 status alone had sub-optimal classification accuracy and hints that a multi-analyte  
514 approach is more likely to improve the detection of cancer signal.

515 By performing WGBS on lung cancer tissues and NATs, we identified more than  
516 300 lung cancer-specific DMRs, with the majority of them being hypermethylated  
517 DMRs, suggesting that hypermethylation of genome regulatory regions is an  
518 important event in lung cancer development. Indeed, these DMRs are enriched for  
519 genes involved in transcriptional regulation and are likely to cause profound  
520 downstream changes in gene expression and contribute to cell transformation;  
521 these genes are likely to be potential tumor suppressor genes, and many of which  
522 haven't been implicated as such previously (such as *SEC31B*, *ZNF274*, and  
523 *NXP1*). A small number of DMRs are hypomethylated, and therefore may encode  
524 potential cancer driver genes. To our knowledge, this is the first time many of  
525 these genes are implicated in epigenetic dysregulation of lung cancer. DAVID  
526 functional GO analysis of biological processes revealed that these DMR genes  
527 were enriched in skeletal system/embryonic organ development/morphogenesis  
528 (Supplementary Figure 21), indicating that the lung cancer cells may have obtained  
529 some characteristics of embryonic stem cells. Functions of these genes remain to  
530 be elucidated in further studies and may help us better understand the underlying  
531 molecular mechanisms of lung cancer development and progression.

532 Our cfDNA methylation-based classification model for LC and BLN plasma  
533 achieved a slightly better cross-validation AUC (0.72) than the mutation-based  
534 model, suggesting that LC-specific methylation changes are potentially useful  
535 markers for diagnosing lung cancers versus benign lesions. Further multi-center

536 studies with larger sample sizes will be needed to validate the utility of selected  
537 markers and the robustness of our diagnostic model. We also noted that  
538 performance of our models are slightly inferior to an earlier study which used a  
539 panel of 9 methylation markers for differentiating early-stage lung cancers from  
540 benign pulmonary nodules (AUC of 0.82 (0.70-0.93) in the test set), even though  
541 we used more methylation markers in our model. The difference in model  
542 performance could be attributed to the different study populations and cfDNA  
543 analysis methods. The smaller cohort size in the previous study may have also  
544 caused over-fitting and/or over-estimation of the model performance. The  
545 discrepancy also suggests that we need to be cautious with the development and  
546 validation of such cfDNA-based diagnostic models, considering the intrinsic  
547 technical difficulties of detecting a minute amount of cancer-derived signals in  
548 circulation and relying on machine-learning approaches to build diagnostic model,  
549 a process that can be heavily affected by batch effect as well as variations in  
550 sample characteristics, especially with single-center clinical study.

551 In our results, we found that a higher percentage of LUSC cfDNA samples were  
552 mutation positive than LUAD samples (Supplementary Table 4). We also observed  
553 higher sensitivity for detecting LUSCs than LUADs using methylation as well as  
554 multi-omics based classification models (Supplementary Table 6). These results  
555 are consistent with the notion that LUSCs are significantly more necrotic than  
556 LUADs and are more likely to shed ctDNA into circulation [60].

557 Previous studies have shown that detection of cancer driver mutations in cfDNA,  
558 when combined with serum protein markers, can be used to increase sensitivity  
559 without significantly sacrificing specificity for cancer detection [9,61]. In our study,  
560 the multi-omics model integrating mutation, methylation and serum protein marker  
561 further improved the performance of the classification model (AUC of 0.79),



562 compared to the mutation-based model or methylation-based model. To our  
563 knowledge, this is the first proof-of-concept study to demonstrate that genetic,  
564 epigenetic, and proteomic analytes could be combined to increase the performance  
565 of liquid biopsy-based diagnostic model for lung cancer. Further study with a  
566 larger size of clinical samples will be needed to validate the robustness of this  
567 approach.

568  
569 We also investigated the association of prognosis of lung cancer patients with  
570 cfDNA mutation and methylation status. Lung cancer patients with any mutation  
571 were observed to have an unfavorable outcome compared with those without, in  
572 line with previous studies [62-64]. We also identified a group of potential  
573 methylation-based lung cancer prognostic markers from the pool of lung cancer  
574 tissue-specific DMRs and constructed a methylation-based prognostic score (MPS).  
575 Previously, multiple methylation-based prognostic classifiers had been reported for  
576 lung cancer, however, the reported markers were mostly inconsistent [25,65-67].  
577 The inconsistency could be explained by limited sample sizes, variations in study  
578 design, as well as different detection methods used. Further studies will be needed  
579 to validate the clinical utility of these markers including ones discovered in the  
580 current study. Finally, we found that combining continuous MPS with the mutation  
581 score could improve the prognostication model compared with the multivariate  
582 model based solely on mutation. One potential caveat to be noted here is that since  
583 the estimated coefficients of DMR markers for generating the MPS might not be  
584 accurate enough due to relatively limited sample size; therefore, additional study  
585 with larger sample size would be necessary to validate current findings. Overall,  
586 we provided proof-of-principle evidence that combination of multiple blood-  
587 derived biomarkers has the potential to improve lung cancer prognostication.



588

## 589 **Availability of Data**

590 The data reported in this study are also available in the CNGB Nucleotide Sequence  
591 Archive (CNSA: <https://db.cngb.org/cnsa>; accession number CNP 0001236).

592

## 593 **Acknowledgements**

594 This study was supported by the National Natural Science Foundation of China  
595 (No.81602001), Peking University People's Hospital Research and Development  
596 Funds (RS2019-01), and Shenzhen Engineering Laboratory for Innovative  
597 Molecular Diagnostics (DRC-SZ[2016]884).

598

## 599 **References:**

- 600 1. Bray, F., Ferlay, J., Soerjomataram, I., Siegel, R. L., Torre, L. A., & Jemal, A. (2018). Global  
601 cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36  
602 cancers in 185 countries. *CA: A Cancer Journal for Clinicians*, 68(6), 394-424.
- 603 2. Howlander N, N. A., Krapcho M, Miller D, Brest A, Yu M, Ruhl J, Tatalovich Z, Mariotto A,  
604 Lewis DR, Chen HS, Feuer EJ, Cronin KA (eds). (April 2019). SEER Cancer Statistics Review,  
605 1975-2016. Retrieved from [https://seer.cancer.gov/csr/1975\\_2016/](https://seer.cancer.gov/csr/1975_2016/)
- 606 3. Henschke Ci Fau - Yankelevitz, D. F., Yankelevitz Df Fau - Libby, D. M., Libby Dm Fau -  
607 Pasmantier, M. W., Pasmantier Mw Fau - Smith, J. P., Smith Jp Fau - Miettinen, O. S., &  
608 Miettinen, O. S. (2006). Survival of patients with stage I lung cancer detected on CT screening.  
609 *The New England journal of medicine*(355), 1763-1771.
- 610 4. National Lung Screening Trial Research, T., Aberle, D. R., Adams, A. M., Berg, C. D., Black,  
611 W. C., Clapp, J. D., . . . Sicks, J. D. (2011). Reduced lung-cancer mortality with low-dose  
612 computed tomographic screening. *The New England journal of medicine*, 365(5), 395-409.
- 613 5. McCunney, R. J., & Li, J. (2014). Radiation Risks in Lung Cancer Screening Programs.  
614 *CHEST*, 145(3), 618-624.
- 615 6. Pinsky, P. F. (2014). Assessing the benefits and harms of low-dose computed tomography  
616 screening for lung cancer. *Lung cancer management*, 3(6), 491-498.
- 617 7. Li, Q., Li, F., Shiraishi, J., Katsuragawa, S., Sone, S., & Doi, K. (2003). Investigation of new  
618 psychophysical measures for evaluation of similar images on thoracic computed tomography for  
619 distinction between benign and malignant nodules. *Medical Physics*, 30(10), 2584-2593.
- 620 8. Qian, F., Yang, W., Chen, Q., Zhang, X., & Han, B. J. J. o. T. D. (2018). Screening for early  
621 stage lung cancer and its correlation with lung nodule detection. 2018, S846-S859.
- 622 9. Cohen, J. D., Li, L., Wang, Y., Thoburn, C., Afsari, B., Danilova, L., . . . Papadopoulos, N.  
623 (2018). Detection and localization of surgically resectable cancers with a multi-analyte blood test.  
624 *Science*, 359(6378), 926-930.
- 625 10. Calapre, L., Warburton, L., Millward, M., Ziman, M., & Gray, E. S. (2017). Circulating tumour  
626 DNA (ctDNA) as a liquid biopsy for melanoma. *Cancer Letters*, 404, 62-69.

- 627 11. Ye, Q., Ling, S., Zheng, S., & Xu, X. (2019). Liquid biopsy in hepatocellular carcinoma:  
628 circulating tumor cells and circulating tumor DNA. *Molecular Cancer*, *18*(1), 114.
- 629 12. Lim, S. Y., Lee, J. H., Diefenbach, R. J., Kefford, R. F., & Rizos, H. (2018). Liquid  
630 biomarkers in melanoma: detection and discovery. *Molecular Cancer*, *17*(1), 8.
- 631 13. Di Meo, A., Bartlett, J., Cheng, Y., Pasic, M. D., & Yousef, G. M. (2017). Liquid biopsy: a  
632 step forward towards precision medicine in urologic malignancies. *Molecular Cancer*, *16*(1), 80.
- 633 14. Phallen, J., Sausen, M., Adleff, V., Leal, A., Hruban, C., White, J., . . . Velculescu, V. E.  
634 (2017). Direct detection of early-stage cancers using circulating tumor DNA. *Science*  
635 *Translational Medicine*, *9*(403), eaan2415.
- 636 15. Jamal-Hanjani, M., Wilson, G. A., McGranahan, N., Birkbak, N. J., Watkins, T. B., Veeriah,  
637 S., . . . Rosenthal, R. (2017). Tracking the Evolution of Non-Small-Cell Lung Cancer. *New*  
638 *England Journal of Medicine*.
- 639 16. Otandault, A., Anker, P., Al Amir Dache, Z., Guillaumon, V., Meddeb, R., Pastor, B., . . .  
640 Thierry, A. R. (2019). Recent advances in circulating nucleic acids in oncology. *Annals of*  
641 *Oncology*.
- 642 17. Mouliere, F., & Rosenfeld, N. (2015). Circulating tumor-derived DNA is shorter than somatic  
643 DNA in plasma. *Proceedings of the National Academy of Sciences*, *112*(11), 3178.
- 644 18. Murtaza, M., & Caldas, C. (2016). Nucleosome mapping in plasma DNA predicts cancer  
645 gene expression. *Nature Genetics*, *48*, 1105.
- 646 19. Liu, X., Liu, L., Ji, Y., Li, C., Wei, T., Yang, X., . . . Wang, X. (2019). Enrichment of short  
647 mutant cell-free DNA fragments enhanced detection of pancreatic cancer. *EBioMedicine*, *41*,  
648 345-356.
- 649 20. Shen, S. Y., Singhania, R., Fehringer, G., Chakravarthy, A., Roehrl, M. H. A., Chadwick,  
650 D., . . . De Carvalho, D. D. (2018). Sensitive tumour detection and classification using plasma  
651 cell-free DNA methylomes. *Nature*, *563*(7732), 579-583.
- 652 21. Xu, R.-h., Wei, W., Krawczyk, M., Wang, W., Luo, H., Flagg, K., . . . Zhang, K. (2017).  
653 Circulating tumour DNA methylation markers for diagnosis and prognosis of hepatocellular  
654 carcinoma. *Nature Materials*, *16*, 1155.
- 655 22. Travis, W. D., Brambilla, E., Burke, A. P., Marx, A., & Nicholson, A. G. (2015). Introduction  
656 to The 2015 World Health Organization Classification of Tumors of the Lung, Pleura, Thymus,  
657 and Heart. *Journal of Thoracic Oncology*, *10*(9), 1240-1242.
- 658 23. Kandoth, C., McLellan, M. D., Vandin, F., Ye, K., Niu, B., Lu, C., . . . Ding, L. (2013).  
659 Mutational landscape and significance across 12 major cancer types. *Nature*, *502*(7471), 333-  
660 339.
- 661 24. Tate, J. G., Bamford, S., Jubb, H. C., Sondka, Z., Beare, D. M., Bindal, N., . . . Forbes, S. A.  
662 (2018). COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Research*,  
663 *47*(D1), D941-D947.
- 664 25. Bailey, M. H., Tokheim, C., Porta-Pardo, E., Sengupta, S., Bertrand, D., Weerasinghe,  
665 A., . . . Ding, L. (2018). Comprehensive Characterization of Cancer Driver Genes and Mutations.  
666 *Cell*, *173*(2), 371-385.e318.
- 667 26. Newman, A. M., Lovejoy, A. F., Klass, D. M., Kurtz, D. M., Chabon, J. J., Scherer, F., . . .  
668 Alizadeh, A. A. (2016). Integrated digital error suppression for improved detection of circulating  
669 tumor DNA. *Nat Biotechnol*, *34*(5), 547-555.
- 670 27. Gansauge, M.-T., Gerber, T., Glocke, I., Korlevic, P., Lippik, L., Nagel, S., . . . Meyer, M.  
671 (2017). Single-stranded DNA library preparation from highly degraded DNA using T4 DNA  
672 ligase. *Nucleic acids research*, *45*(10), e79-e79.
- 673 28. Chen, Y., Chen, Y., Shi, C., Huang, Z., Zhang, Y., Li, S., . . . Chen, Q. (2017). SOAPnuke: a  
674 MapReduce acceleration-supported software for integrated quality control and preprocessing of  
675 high-throughput sequencing data. *GigaScience*, *7*(1).
- 676 29. Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows–Wheeler  
677 transform. *Bioinformatics*, *25*(14), 1754-1760.

- 678 30. Garrison, E., & Marth, G. (2012). Haplotype-based variant detection from short-read  
679 sequencing. *arXiv*, 1207.
- 680 31. McVean, G. A., Altshuler, D. M., Durbin, R. M., Abecasis, G. R., Bentley, D. R., Chakravarti,  
681 A., . . . University of, G. (2012). An integrated map of genetic variation from 1,092 human  
682 genomes. *Nature*, 491(7422), 56-65.
- 683 32. Siva, N. (2008). 1000 Genomes project. *Nature Biotechnology*, 26(3), 256-256.
- 684 33. McLaren, W., Gil, L., Hunt, S. E., Riat, H. S., Ritchie, G. R. S., Thormann, A., . . .  
685 Cunningham, F. (2016). The Ensembl Variant Effect Predictor. *Genome Biology*, 17(1), 122.
- 686 34. Bailey, M. H., Tokheim, C., Porta-Pardo, E., Sengupta, S., Bertrand, D., Weerasinghe,  
687 A., . . . Ding, L. (2018). Comprehensive Characterization of Cancer Driver Genes and Mutations.  
688 *Cell*, 173(2), 371-385.e318.
- 689 35. Chen, S., Zhou, Y., Chen, Y., & Gu, J. (2018). fastp: an ultra-fast all-in-one FASTQ  
690 preprocessor. *Bioinformatics*, 34(17), i884-i890.
- 691 36. Cheng, H., & Xu, Y. (2018). BitMapperBS: a fast and accurate read aligner for whole-  
692 genome bisulfite sequencing. *bioRxiv*, 442798.
- 693 37. Tarasov, A., Vilella, A. J., Cuppen, E., Nijman, I. J., & Prins, P. (2015). Sambamba: fast  
694 processing of NGS alignment formats. *Bioinformatics*, 31(12), 2032-2034.
- 695 38. Wu, H., Xu, T., Feng, H., Chen, L., Li, B., Yao, B., . . . Conneely, K. N. (2015). Detection of  
696 differentially methylated regions from whole-genome bisulfite sequencing data without replicates.  
697 *Nucleic Acids Research*, 43(21), e141-e141.
- 698 39. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., . . . Dubourg,  
699 V. J. J. o. m. l. r. (2011). Scikit-learn: Machine learning in Python. 12(Oct), 2825-2830.
- 700 40. Heagerty, P. J., & Zheng, Y. J. B. (2005). Survival model predictive accuracy and ROC  
701 curves. 61(1), 92-105.
- 702 41. Szpechcinski, A., Rudzinski, P., Kupis, W., Langfort, R., Orłowski, T., & Chorostowska-  
703 Wynimko, J. (2016). Plasma cell-free DNA levels and integrity in patients with chest radiological  
704 findings: NSCLC versus benign lung nodules. *Cancer Letters*, 374(2), 202-207.
- 705 42. Liu, J., Chen, X., Wang, J., Zhou, S., Wang, C. L., Ye, M. Z., . . . Qian, Z. Y. (2019).  
706 Biological background of the genomic variations of cf-DNA in healthy individuals. *Annals of*  
707 *Oncology*, 30(3), 464-470.
- 708 43. Razavi, P., Li, B. T., Brown, D. N., Jung, B., Hubbell, E., Shen, R., . . . Reis-Filho, J. S.  
709 (2019). High-intensity sequencing reveals the sources of plasma circulating cell-free DNA  
710 variants. *Nature Medicine*.
- 711 44. Chen, Y.-J., Roumeliotis, T. I., Chang, Y.-H., Chen, C.-T., Han, C.-L., Lin, M.-H., . . . Chen,  
712 Y.-J. (2020). Proteogenomics of Non-smoking Lung Cancer in East Asia Delineates Molecular  
713 Signatures of Pathogenesis and Progression. *Cell*, 182(1), 226-244.e217.
- 714 45. Cancer Genome Atlas Research, N. (2012). Comprehensive genomic characterization of  
715 squamous cell lung cancers. *Nature*, 489(7417), 519-525.
- 716 46. The Cancer Genome Atlas Research, N., Collisson, E. A., Campbell, J. D., Brooks, A. N.,  
717 Berger, A. H., Lee, W., . . . Tsao, M.-S. (2014). Comprehensive molecular profiling of lung  
718 adenocarcinoma. *Nature*, 511, 543.
- 719 47. Campbell, J. D., Alexandrov, A., Kim, J., Wala, J., Berger, A. H., Peadarallu, C. S., . . .  
720 Meyerson, M. (2016). Distinct patterns of somatic genome alterations in lung adenocarcinomas  
721 and squamous cell carcinomas. *Nature Genetics*, 48, 607.
- 722 48. Ohgaki, H., Kros, J. M., Okamoto, Y., Gaspert, A., Huang, H., & Kurrer, M. O. (2004). APC  
723 mutations are infrequent but present in human lung cancer. *Cancer Letters*, 207(2), 197-203.
- 724 49. Bentires-Alj, M., Paez, J. G., David, F. S., Keilhack, H., Halmos, B., Naoki, K., . . . Neel, B. G.  
725 (2004). Activating Mutations of the Noonan Syndrome-Associated *SHP2/PTPN11* Gene in  
726 Human Solid Tumors and Adult Acute Myelogenous Leukemia. *Cancer Research*, 64(24), 8816.
- 727 50. Zhang, Y., Wang, D. C., Shi, L., Zhu, B., Min, Z., & Jin, J. (2017). Genome analyses identify  
728 the genetic modification of lung cancer subtypes. *Seminars in Cancer Biology*, 42, 20-30.

- 729 51. Rao, R. C., & Dou, Y. (2015). Hijacked in cancer: the KMT2 (MLL) family of  
730 methyltransferases. *Nature Reviews Cancer*, *15*, 334.
- 731 52. Chen, K.-Z., Lou, F., Yang, F., Zhang, J.-B., Ye, H., Chen, W., . . . Wang, J. (2016).  
732 Circulating Tumor DNA Detection in Early-Stage Non-Small Cell Lung Cancer Patients by  
733 Targeted Sequencing. *Scientific Reports*, *6*(1), 31985.
- 734 53. Hoadley, K. A., Yau, C., Hinoue, T., Wolf, D. M., Lazar, A. J., Drill, E., . . . Laird, P. W.  
735 (2018). Cell-of-Origin Patterns Dominate the Molecular Classification of 10,000 Tumors from 33  
736 Types of Cancer. *Cell*, *173*(2), 291-304.e296.
- 737 54. Yizhak, K., Aguet, F., Kim, J., Hess, J. M., Kübler, K., Grimsby, J., . . . Getz, G. (2019). RNA  
738 sequence analysis reveals macroscopic somatic clonal expansion across normal tissues.  
739 *Science*, *364*(6444).
- 740 55. Risques, R. A., & Kennedy, S. R. (2018). Aging and the rise of somatic cancer-associated  
741 mutations in normal tissues. *PLoS genetics*, *14*(1), e1007108-e1007108.
- 742 56. Makinen, N., Mehine, M., Tolvanen, J., Kaasinen, E., Li, Y., Lehtonen, H. J., . . . Aaltonen, L.  
743 A. (2011). MED12, the mediator complex subunit 12 gene, is mutated at high frequency in  
744 uterine leiomyomas. *Science*, *334*(6053), 252-255.
- 745 57. Lim, W. K., Ong, C. K., Tan, J., Thike, A. A., Ng, C. C., Rajasegaran, V., . . . Teh, B. T.  
746 (2014). Exome sequencing identifies highly recurrent MED12 somatic mutations in breast  
747 fibroadenoma. *Nat Genet*, *46*(8), 877-880.
- 748 58. Ye, M., Li, S., Huang, W., Wang, C., Liu, L., Liu, J., . . . Liang, W. (2018). Comprehensive  
749 targeted super-deep next generation sequencing enhances differential diagnosis of solitary  
750 pulmonary nodules. *J Thorac Dis*, *10*(Suppl 7), S820-s829.
- 751 59. Peng, M., Xie, Y., Li, X., Qian, Y., Tu, X., Yao, X., . . . Tian, G. (2019). Resectable lung  
752 lesions malignancy assessment and cancer detection by ultra-deep sequencing of targeted  
753 gene mutations in plasma cell-free DNA. *Journal of Medical Genetics*, *56*(10), 647.
- 754 60. Abbosh, C., Birkbak, N. J., Wilson, G. A., Jamal-Hanjani, M., Constantin, T., Salari, R., . . .  
755 The, T. c. (2017). Phylogenetic ctDNA analysis depicts early-stage lung cancer evolution.  
756 *Nature*, *545*(7655), 446-451.
- 757 61. Cohen, J. D., Li, L., Wang, Y., Thoburn, C., Afsari, B., Danilova, L., . . . Papadopoulos, N.  
758 (2018). Detection and localization of surgically resectable cancers with a multi-analyte blood test.  
759 *Science*, *359*(6378), 926.
- 760 62. Isaksson, S., George, A. M., Jonsson, M., Cirenajwis, H., Jonsson, P., Bendahl, P. O., . . .  
761 Planck, M. (2019). Pre-operative plasma cell-free circulating tumor DNA and serum protein  
762 tumor markers as predictors of lung adenocarcinoma recurrence. *Acta Oncol*, *58*(8), 1079-1086.
- 763 63. Kang, J., Luo, Y., Wang, D., Men, Y., Wang, J., Che, Y.-Q., & Hui, Z. (2019). Tumor  
764 Mutation Load: A Novel Independent Prognostic Factor in Stage IIIA-N2 Non-Small-Cell Lung  
765 Cancer. *Disease Markers*, *2019*, 3837687.
- 766 64. Corradetti, M. N., Torok, J. A., Hatch, A. J., Xanthopoulos, E. P., Lafata, K., Jacobs, C., . . .  
767 Nixon, A. B. (2019). Dynamic Changes in Circulating Tumor DNA During Chemoradiation for  
768 Locally Advanced Lung Cancer. *Adv Radiat Oncol*, *4*(4), 748-752.
- 769 65. Sandoval, J., Mendez-Gonzalez, J., Nadal, E., Chen, G., Carmona, F. J., Sayols, S., . . .  
770 Esteller, M. (2013). A prognostic DNA methylation signature for stage I non-small-cell lung  
771 cancer. *J Clin Oncol*, *31*(32), 4140-4147.
- 772 66. Li, Y., Gu, J., Xu, F., Zhu, Q., Ge, D., & Lu, C. (2019). Novel methylation-driven genes  
773 identified as prognostic indicators for lung squamous cell carcinoma. *Am J Transl Res*, *11*(4),  
774 1997-2012.
- 775 67. Wang, Y., Deng, H., Xin, S., Zhang, K., Shi, R., & Bao, X. (2019). Prognostic and Predictive  
776 Value of Three DNA Methylation Signatures in Lung Adenocarcinoma. *Front Genet*, *10*, 349.

777

778

779

		LC (N=128)		BLN (N=94)	
		Number	Percentage	Number	Percentage
<b>Gender</b>	<b>Female</b>	53	41%	48	51%
	<b>Male</b>	75	59%	46	49%
<b>Age</b>	<b>Median±SD (Range)</b>	63.00±11.58 (30-86)		55.00±10.49 (18-79)	
<b>Histology</b>	<b>LUAD</b>	97	76%		
	<b>LUSC</b>	23	18%		
	<b>LCC</b>	3	2%		
	<b>SCLC</b>	5	4%		
<b>Stage</b>	<b>0</b>	2	2%		
	<b>IA</b>	54	42%		
	<b>IB</b>	29	23%		
	<b>II</b>	17	13%		
	<b>III</b>	19	15%		
	<b>IV</b>	7	5%		
<b>Smoking History</b>	<b>Smokers</b>	20	16%	23	24%
	<b>Non-smokers</b>	33	26%	70	74%
	<b>Unknown</b>	75	59%	1	1%

780

781 **Table 1: Clinicopathological characteristics of the patients enrolled in this**  
 782 **study.** LUAD: lung adenocarcinoma. LUSC: lung squamous cell carcinoma. LCC:  
 783 large cell carcinoma. SCLC: small cell lung carcinoma.

784



785 **Figure Legends**

786 **Figure 1: Study design and variants detected by targeted ultra-deep**

787 **sequencing in cfDNA and WBC gDNA.** (A) Schematic view of the study design.

788 See Methods for additional details. (B) Spike-in experiments using Multiplex I

789 cfDNA Reference Standard Set which contains 8 SNVs to test the LOD of our

790 targeted ultra-deep sequencing method. The sensitivity was calculated as the

791 number of detected SNVs divided by the number of total spiked-in SNVs in all the

792 replicates for each condition. (C) Allele fractions (x-axis, log scale) of mutations

793 detected in plasma cfDNA of BLN patients (blue) and LC patients (red). (D) AF

794 (x-axis, log scale) distribution of WBC gDNA variants from BLN and LC patients.

795 (E) Pearson correlation of AF in cfDNA (x-axis, log scale) and AF in WBC gDNA

796 (y-axis, log scale). Each point represents one variant detected in matched cfDNA

797 and WBC gDNA samples from the same patient.

798 **Figure 2: Predictive model based on variants detected in plasma cfDNA after**  
799 **filtering with matched WBC sample for shared variants.** (A) AFs (x-axis, log

800 scale) of cfDNA variants of LC patients (red) and BLN patients (blue). (B)

801 Oncoplot showing the 153 mutations detected in 67 out of 111 (60.36%) LC

802 samples. 45 LC samples without any mutation detected were not shown. Each

803 column represents a sample and each row a different gene. The upper barplot

804 represents the frequency of mutations for each sample, and the right barplot

805 represents the frequency of mutations for each gene. Samples are ordered by the

806 most mutated genes. (C) Summary plot of the 153 mutations detected in LC

807 samples. Upper panel from left to right: Variant classification, Variant Type, and

808 SNV Class. Lower panel from left to right: Variants per sample and Variant

809 classification summary. (D) Oncoplot of the 28 mutations detected in 23 out of 78

810 (29.49%) BLN samples. 55 BLN samples without any mutation detected were not

811 shown. (E) Predictive models to distinguish LC from BLN based on mutations

812 detected. SUMAF (green): sum of AF model. Weighted\_SUMAF (red): sum of

813 weighted AF model (see Methods). The AUC of SUMAF model is 0.67 with 55.9%

814 sensitivity and 76.9% specificity. The AUC of weighted\_SUMAF model is 0.68

815 with 59.5% sensitivity and 71.8% specificity.

816 **Figure 3: Diagnostic model for Distinguishing LC from BLN plasma by**  
817 **analyzing cfDNA methylation levels.** (A) Differentially methylated regions  
818 (DMRs) discovered by WGBS of LC tumor and  
819 normal tissue adjacent to the tumour (NAT). Red points: Hypermethylated DMRs  
820 in LC tissues. Blue points: hypomethylated DMRs in LC tissues. From outer to  
821 inner circle, the first circle is overview of DMRs, the second circle is the area  
822 statistics of hypermethylation regions (methy.diff>0.2), and the third circle is the  
823 area statistics of hypomethylation regions (methy.diff<-0.2). (B) Heatmap of the  
824 DMRs with hierarchical clustering. Block color represents the methylation  $\beta$  value  
825 and black represents N.A. (C) Functional annotation of the genes associated with  
826 the 293 hypermethylated DMRs by gene ontology (GO) terms using DAVID. (D)  
827 Predictive models to distinguish LC from BLN based on cfDNA methylation level  
828 with selected features (feature importance  $\geq 0.008$ ). A total of 76 DMRs were  
829 retained in the final model. The AUC is 0.72 with 80.5% sensitivity and 57.5%  
830 specificity.

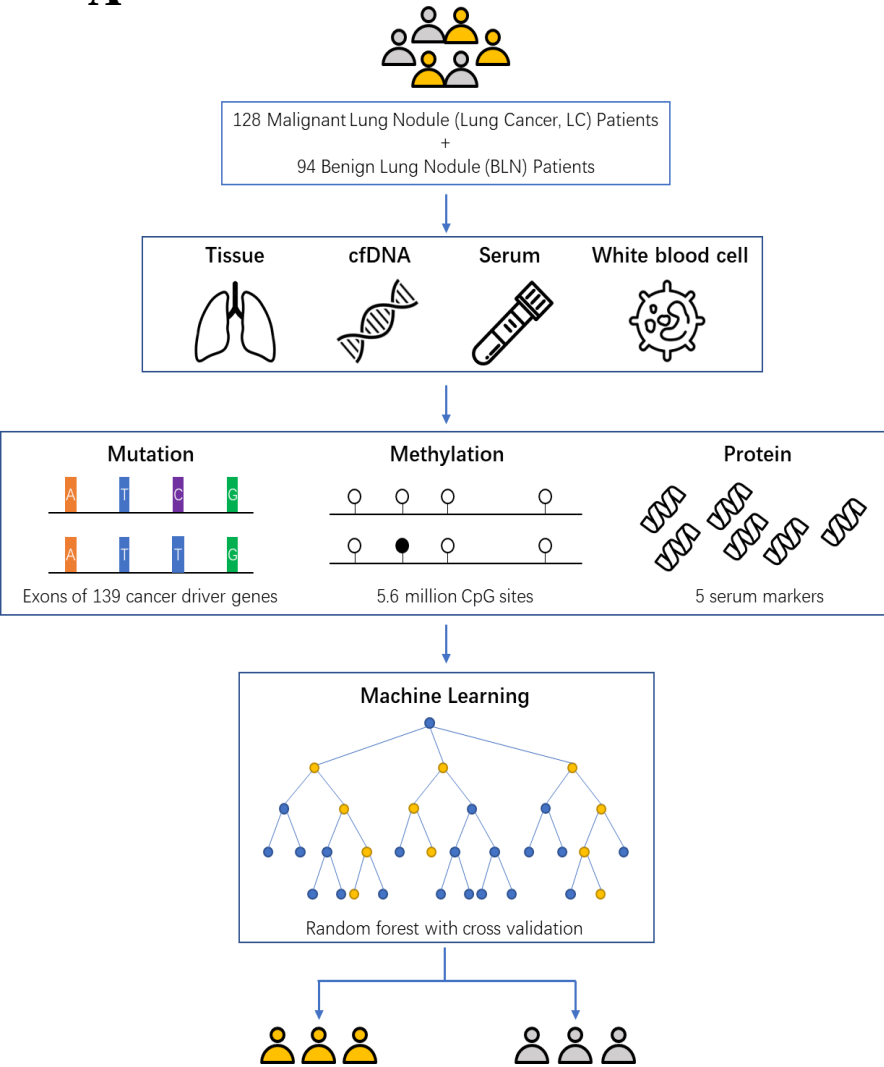
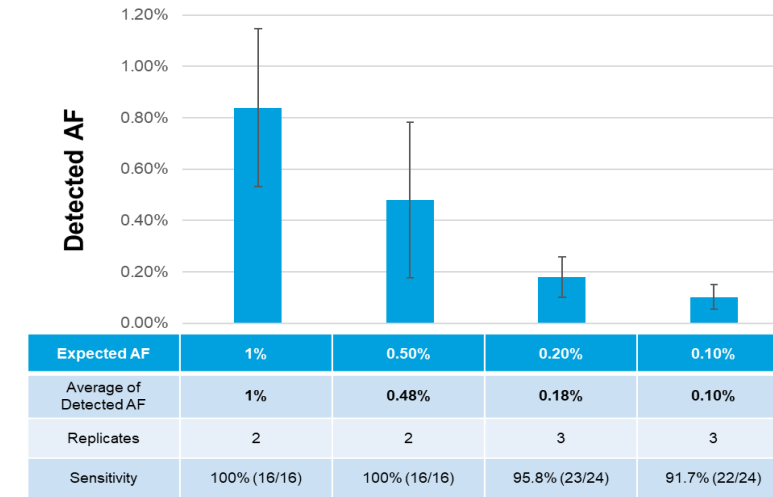
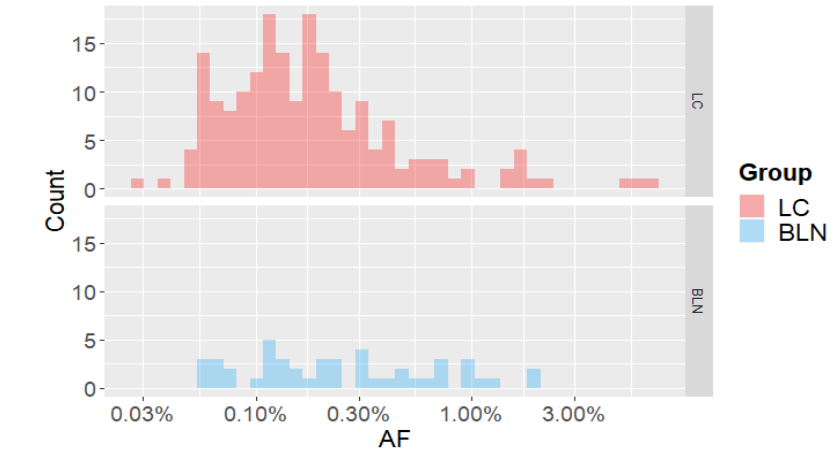
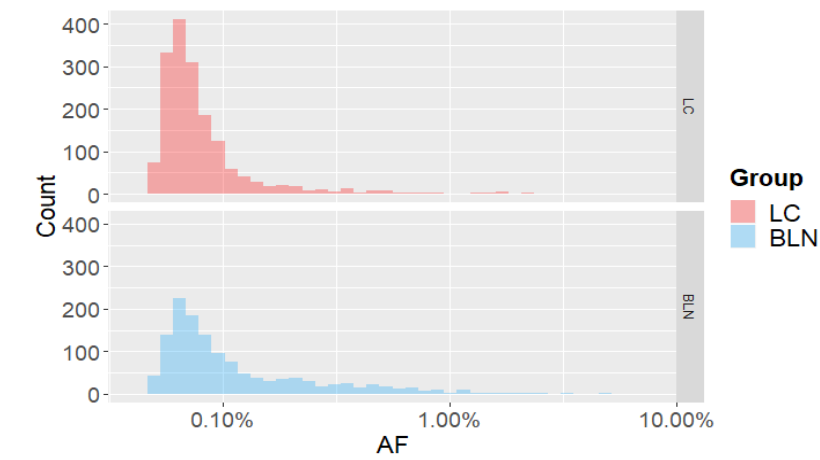
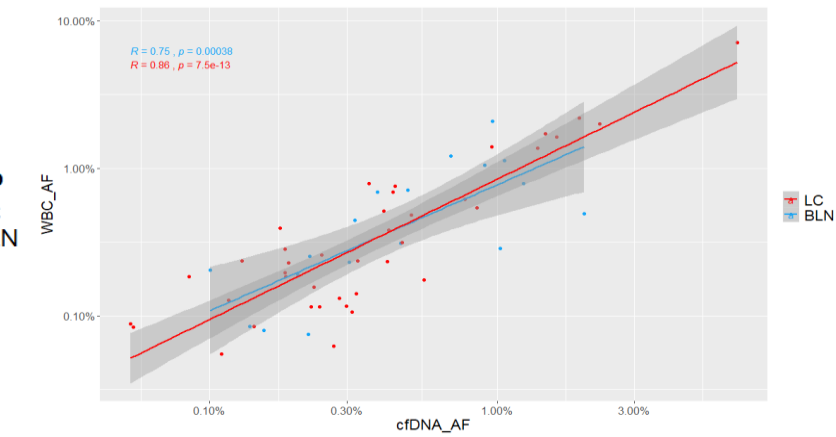
831 **Figure 4: Multiomics predictive models to distinguish LC from BLN plasma.**  
832 Bi-omics samples: cfDNA samples from 91 LC and 71 BLN patients with  
833 complete mutation and methylation data. Tri-omics samples: cfDNA samples from  
834 74 LC and 60 BLN patients with complete measurement of mutation score,  
835 methylation levels of selected DMR and serum CEA levels. (A) Classification  
836 models built based on mutation status alone in bi-omics samples; (B) Models built  
837 based on selected DMRs alone in bi-omics samples; (C) Models built based on  
838 mutation score and selected DMR in bi-omics samples; (D) Models built based on  
839 mutation score and selected DMR in tri-omics samples; (E) Models based on  
840 combined mutation score, selected DMR, and serum CEA levels in tri-omics  
841 samples.

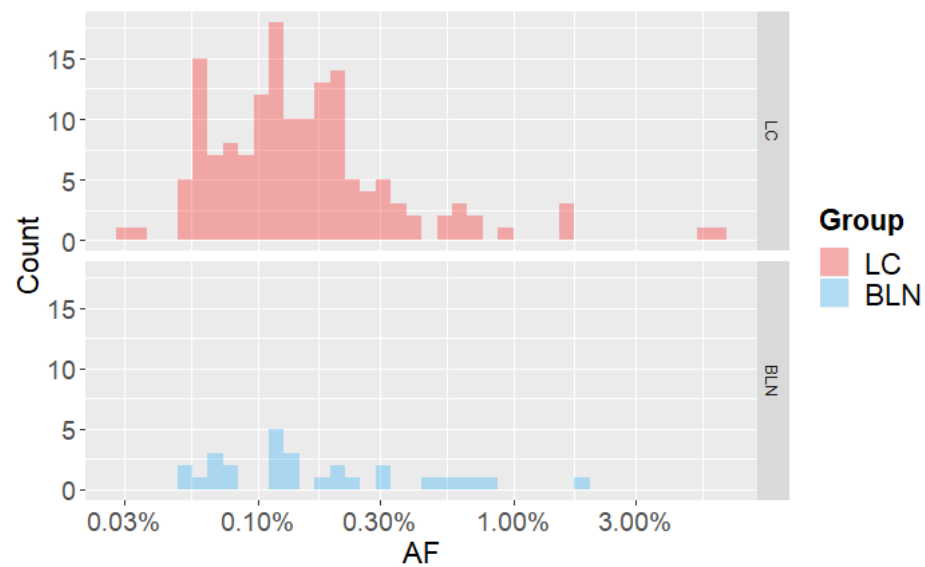
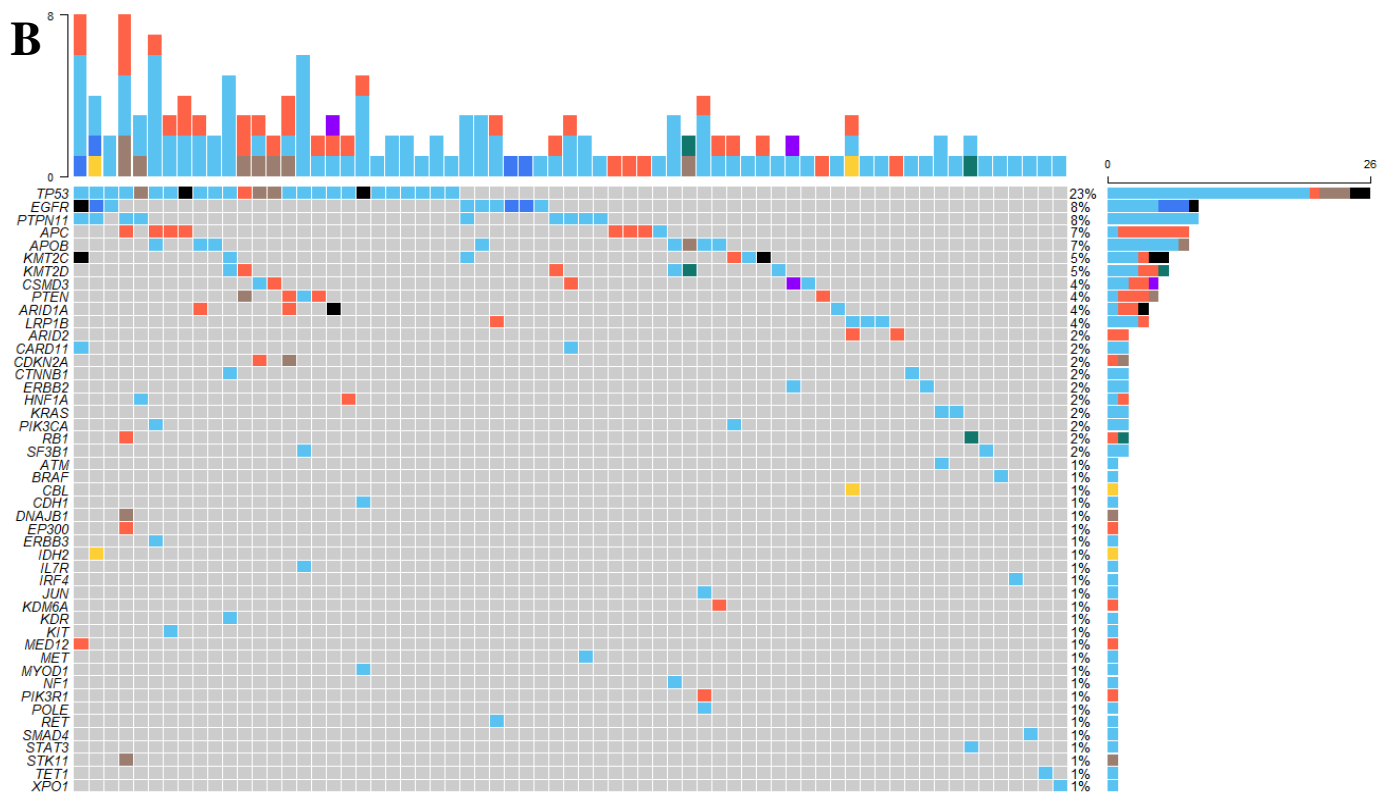
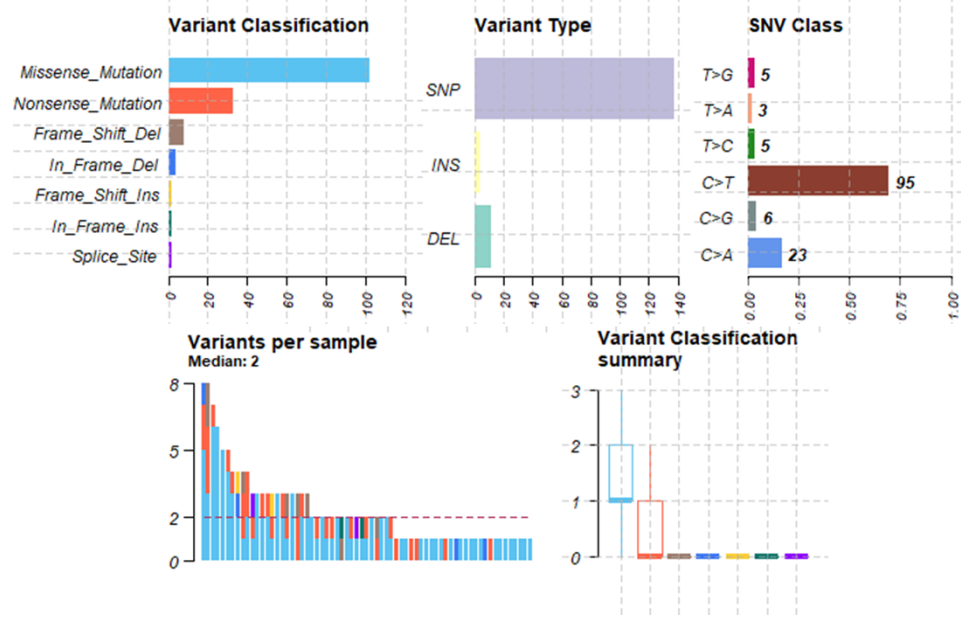
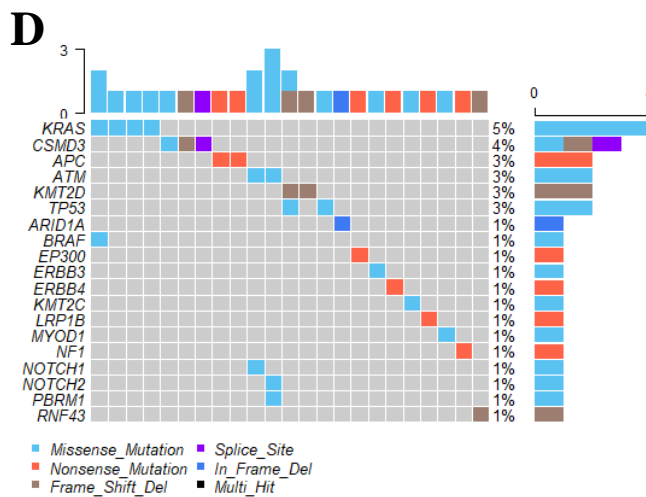
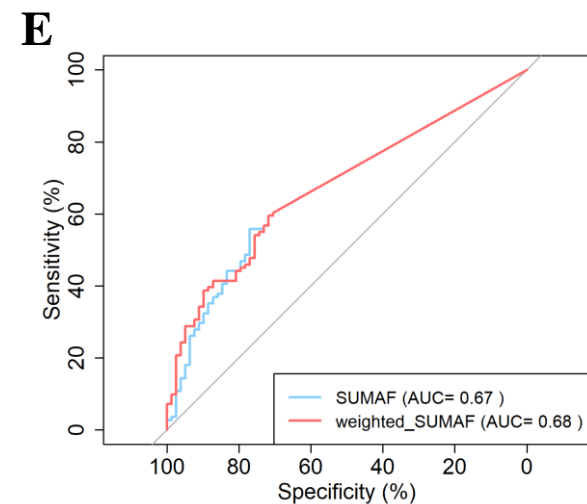
842 **Figure 5: Kaplan-Meier plot on omics-based lung cancer prognostic model in**  
843 **relation to OS.** (A) Mutation scores in the whole dataset; (B) Multi-omics scores in  
844 the testing set.

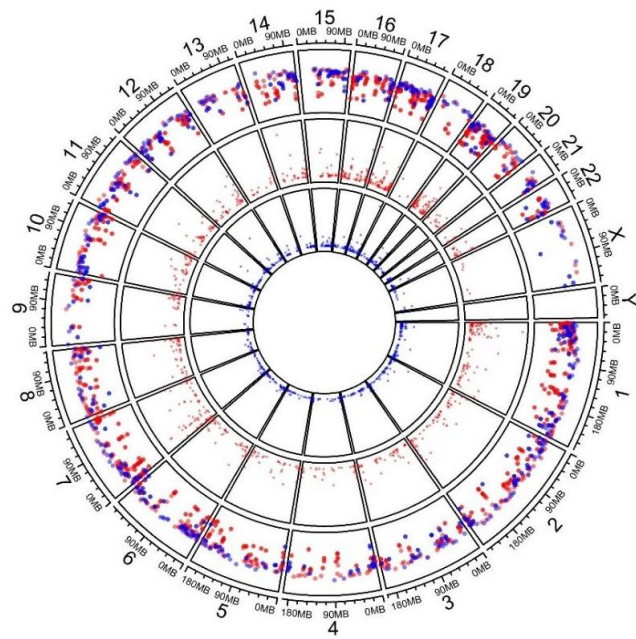
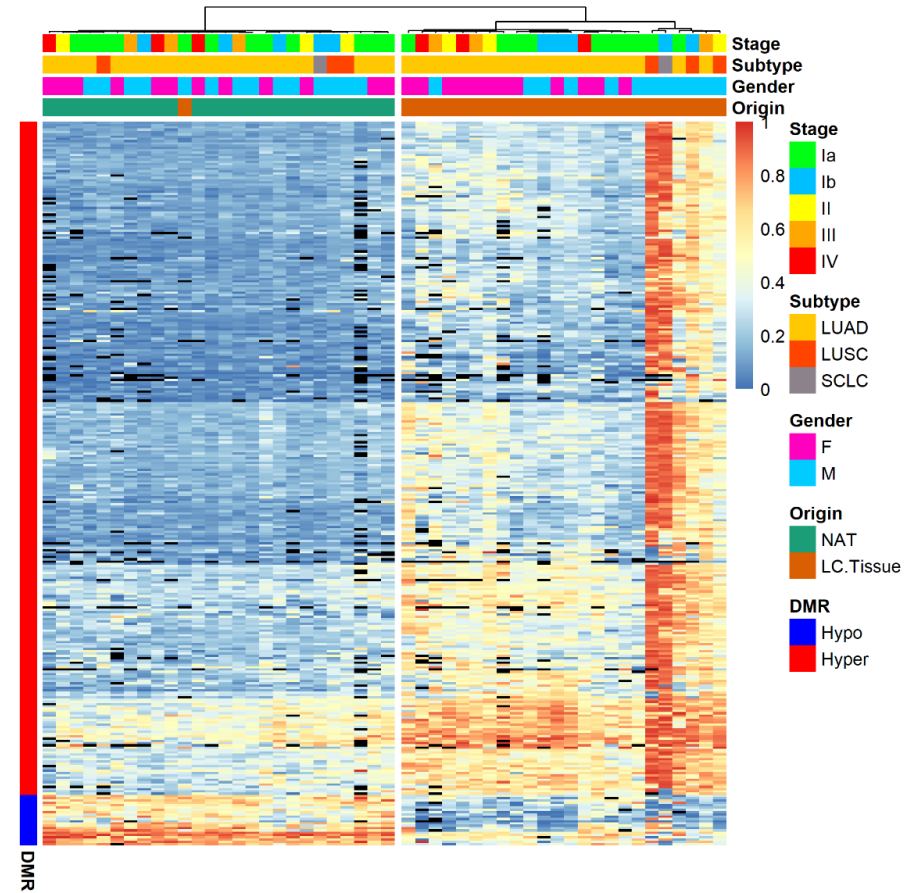
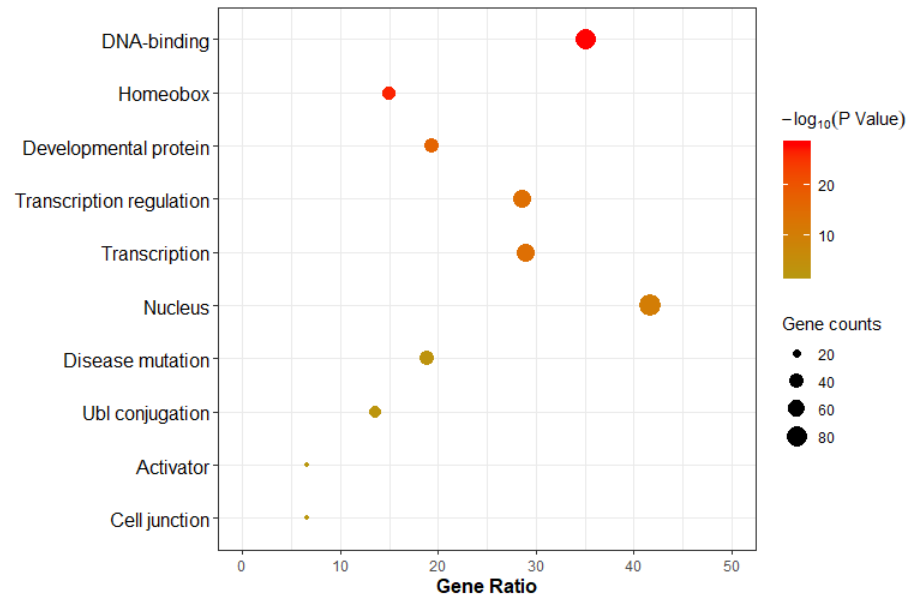
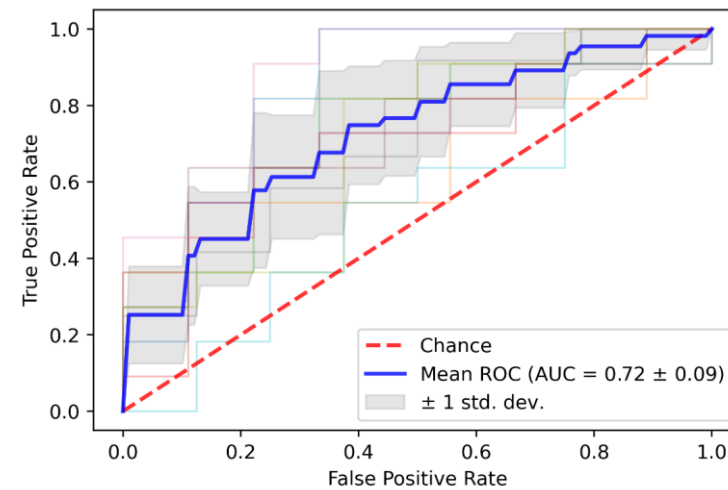
845

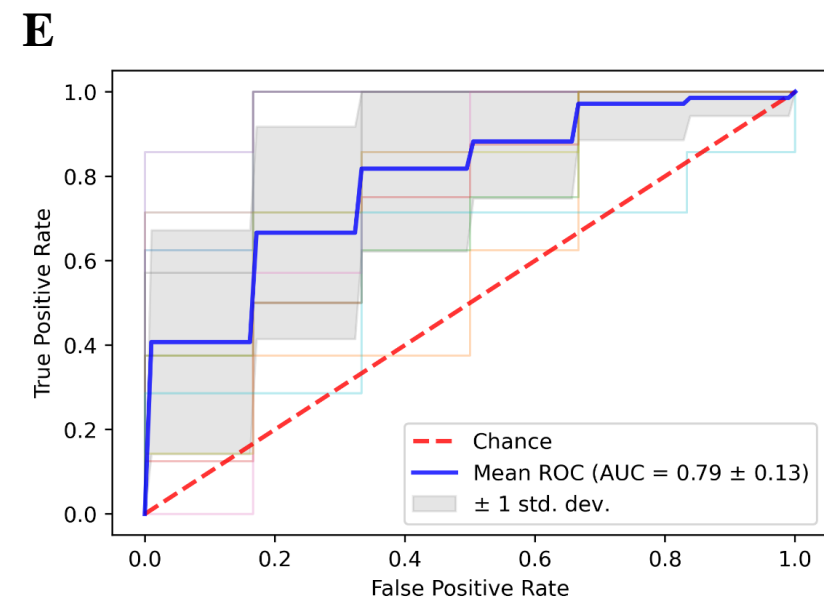
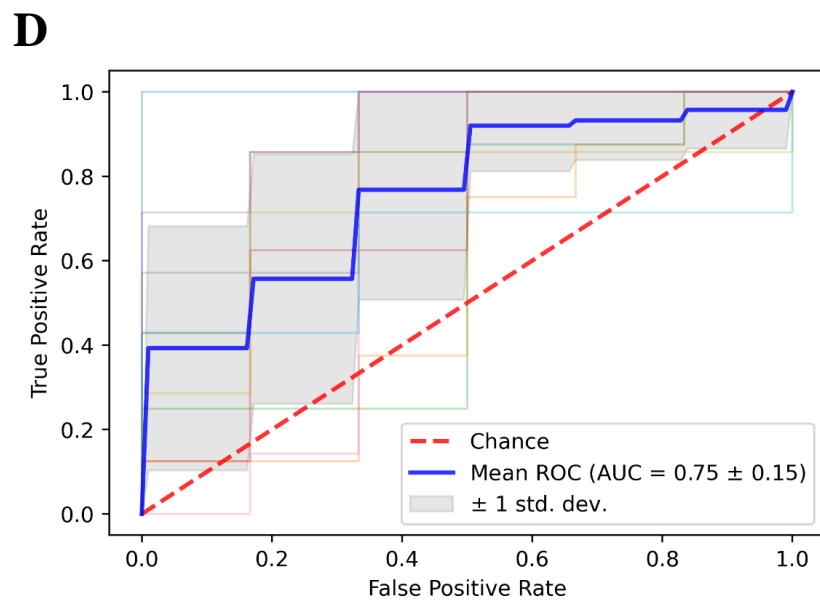
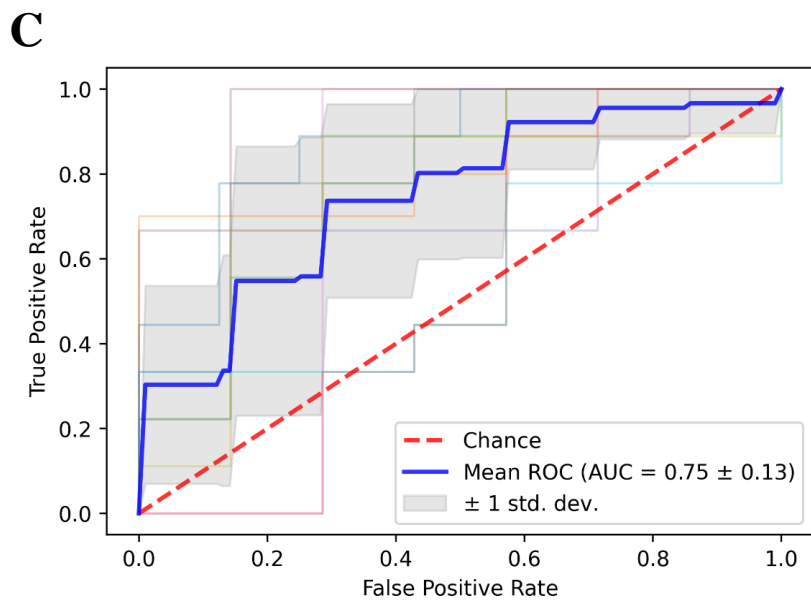
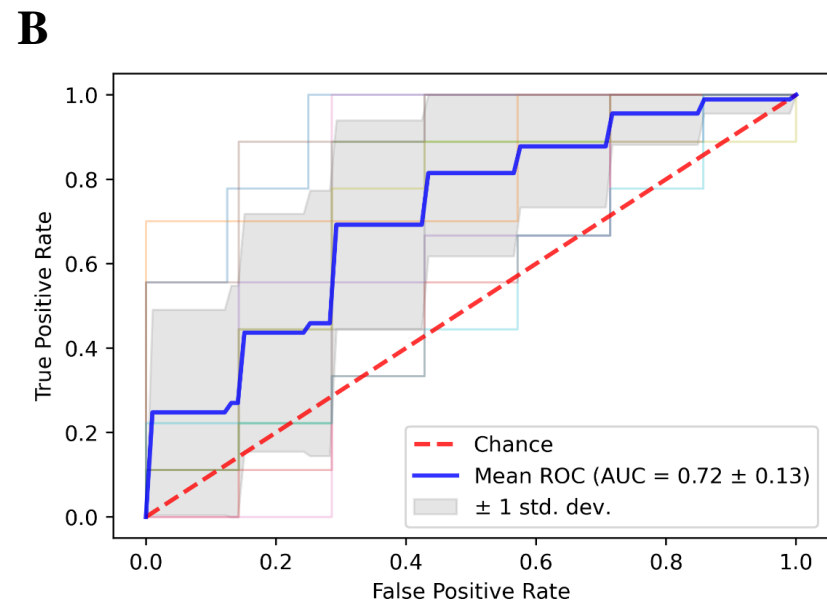
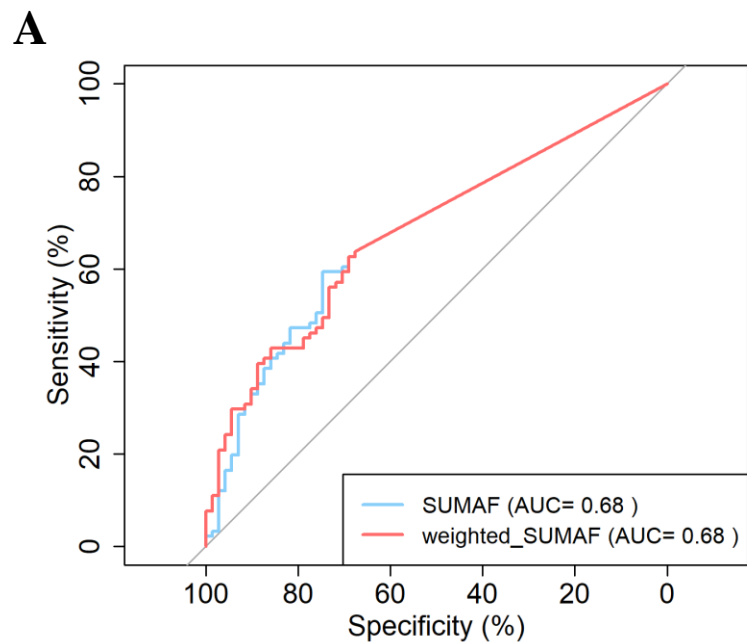


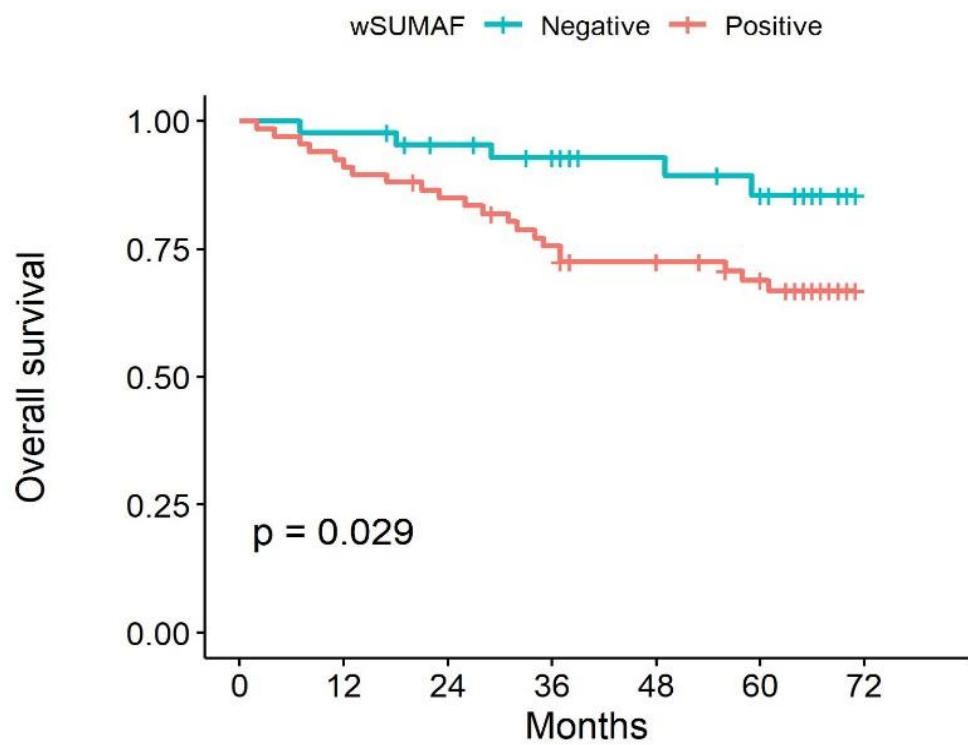
# Figure 1

**A****B****C****D****E**

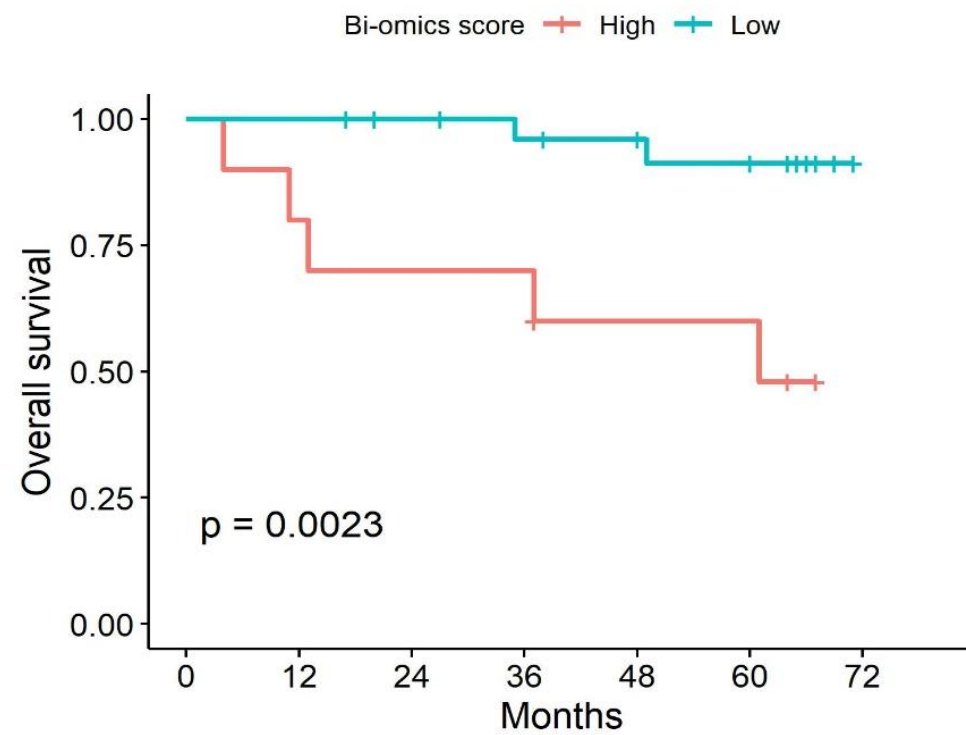
**Figure 2****A****B****C****D****E**

**Figure 3****A****B****C****D**

**Figure 4**

**Figure 5****A**

		Number at risk						
wSUMAF		0	12	24	36	48	60	72
	Negative		44	43	39	36	26	22
Positive		67	62	55	48	43	37	0

**B**

		Number at risk						
Bi-omics score		0	12	24	36	48	60	72
	High		10	8	7	7	5	5
Low		28	28	26	24	21	19	0