

1 **Proteogenomic analysis of air-pollution-associated lung cancer**
2 **reveals prevention and therapeutic opportunities**

3 *Honglei Zhang^{1,*,#}, Chao Liu^{2,3,*}, Shuting Wang^{2,*}, Qing Wang^{4,*}, Xu*
4 *Feng¹, Huawei Jiang⁵, Yong Zhang^{6,#}, Xiaosan Su^{1,#} and Gaofeng Li^{2,#}*

5 ¹Center for scientific research, Yunnan University of Chinese Medicine, Kunming, Yunnan,
6 650500, China;

7 ²Department of Thoracic Surgery II, Third Affiliated Hospital of Kunming Medical University,
8 Yunnan Cancer Hospital, Kunming 650106, China;

9 ³Department of Nuclear Medicine, Third Affiliated Hospital of Kunming Medical University,
10 Yunnan Cancer Hospital, Kunming 650106, China;

11 ⁴Department of Oncology, Qijing First People's Hospital, Qijing, 202150, China;

12 ⁵Department of Plastic Surgery, First Affiliated Hospital of Kunming Medical University,
13 Kunming 650032, China;

14 ⁶Department of Nephrology, Institutes for Systems Genetics, Frontiers Science Center for Disease-
15 Related Molecular Network, West China Hospital, Sichuan University, Chengdu 610041, China.

16 *These authors contributed equally to this work.

17 #Corresponding authors. Yong Zhang (nankai1989@foxmail.com), Xiaosan Su
18 (suxs163@163.com); Gaofeng Li (ligaofeng@kmmu.edu.cn); Honglei Zhang
19 (hlzhang2014@163.com).

20

21

22

23

24

25 NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

26 **Abstract**

27 Air pollution significantly impact lung cancer progression, but there is a
28 lack of a comprehensive molecular characterization of clinical samples
29 associated with air pollution. Here, we performed a proteogenomic
30 analysis of lung adenocarcinoma (LUAD) in 169 female never-smokers
31 from the Xuanwei area (XWLC cohort), where coal smoke is the primary
32 contributor to the high lung cancer incidence. Genomic mutation analysis
33 revealed XWLC as a distinct subtype of LUAD separate from cases
34 associated with smoking or endogenous factors. Mutational signature
35 analysis suggested that Benzo[a]pyrene (BaP) is the major risk factor in
36 XWLC. The BaP-induced mutation hotspot, EGFR-G719X, was present
37 in 20% of XWLC which endowed XWLC with elevated MAPK pathway
38 activations and worse outcomes compared to common *EGFR* mutations.
39 Multi-omics clustering of XWLC identified four clinically relevant
40 subtypes. These subgroups exhibited distinct features in biological
41 processes, genetic alterations, metabolism demands, immune landscape,
42 tumor microbiota composition and radiomic features. Finally, *MAD1* and
43 *TPRN* were identified as novel potential therapeutic targets in XWLC.
44 Our study provides a valuable resource for researchers and clinicians to
45 explore prevention and treatment strategies for air-pollution-associated
46 lung cancers.

48 **Introduction**

49 Lung cancer is the leading cause of cancer deaths globally[1]. Though the
50 most common cause of lung cancer is tobacco smoking, studies estimate
51 that approximately 25% of lung cancers worldwide occur in individuals
52 who have never smoked[2]. Recently, lung cancer in never smokers
53 (LCINS) were molecular profiled and new genomic features were
54 revealed[3-9]. For now, further stratification of LCINS based on different
55 risk factors would be helpful to reveal the oncogenic mechanisms and
56 develop more targeted therapies. Air pollutants, which can directly affect
57 the pulmonary airway, play crucial roles in promoting lung
58 adenocarcinoma[10-12]. More than 20 environmental and occupational
59 agents are lung carcinogens[13, 14] and amount of studies have been
60 made to investigate molecular mechanisms in tumor progression of air
61 pollution chemicals or components using cell lines or mouse/rat
62 models[15-17]. However, a comprehensive molecular characterization of
63 clinical lung cancer samples associated with air pollution is still lacking.

64 The Xuanwei area has the highest rate of lung cancer in China, and
65 extensive research has established a strong link between lung cancer and
66 exposure to domestic coal smoke[18-23]. Specifically, etiologic link
67 between smoky coal burning and cancer was epidemiologically
68 established[18, 19] and association between household stove
69 improvement and lower risk of lung cancer was observed[23]. Moreover,

70 genomic evidence of lung carcinogenesis associated with coal smoke in
71 Xuanwei area, China was provided in our previous study[22]. Thus, lung
72 cancer in Xuanwei areas exemplifies the ideal disease to study
73 characteristics of lung cancers associated with air pollution. In recent
74 years, the molecular features of Xuanwei lung cancer have been gradually
75 revealed[24-26] [22]. A large sample size with multi-omic molecular
76 profiling is urgent needed to explicit the air pollution chemicals and
77 furthermore propose more targeted therapies.

78 To better understand the molecular mechanisms and heterogeneity of
79 XWLC and to advance precision medicine, we expanded the sample size
80 of our next-generation sequencing dataset to 169 sample size and
81 performed proteomic and phosphoproteomic profiling to 112 samples.
82 Furthermore, we integrated 107 radiomic features derived from X-ray
83 computed tomography (CT) scans to 115 samples to non-invasively
84 distinguish molecular subtypes. This allowed us to identify potential
85 major risk factor, distinguish the genomic features and establish clinically
86 relevant molecular subtypes. Our study provides an exceptional resource
87 for future biological, diagnostic, and drug discovery efforts in the study
88 of lung cancer related to air pollution.

89

90

91

92 **Results**

93 To investigate unique biological features of LUAD associated with air
94 pollution, three previous LUAD datasets related to different carcinogens
95 were used for comparison (Fig. 1a). CNLC is the subset of lung
96 adenocarcinoma from non-smoking patients in Chinese Human Proteome
97 Project (CNHPP project) [8] (n=77). TSLC is the subset of lung
98 adenocarcinoma from smoking female in TCGA-LUAD project [27]
99 (n=168). TNLC is the subset of lung adenocarcinoma from non-smoking
100 female in TCGA-LUAD project[27] (n=102). The clinicopathological
101 characteristics of patients from CNLC, TSLC and TNLC cohorts were
102 supplied in Supplementary Table 1a and 1b.

103 **Proteogenomic landscape in Xuanwei lung cancer (XWLC)**

104 The present study prospectively collected primary samples of lung
105 adenocarcinoma (LUAD) from 169 never-smoking women from the
106 Xuanwei area in China (Supplementary Table 1c). The XWLC cohort had
107 a median age of 56 years (Fig. 1b), and the majority of tissue samples
108 were in the early stages of the disease (145 were stage I/II, and 24 were
109 stage III/IV, Fig. 1c). A total of 135, 136, 102 and 102 tumor samples
110 were profiled with whole-exome sequencing (WES), RNA-seq, label-free
111 protein quantification, and label-free phosphorylation quantification,
112 respectively (Fig.1d-e and Extended Fig.1a and 1b). Analysis of the WES
113 data from the paired tumor and normal tissue samples revealed 37,149

114 somatic mutations, including 1,797 InDels, 32,972 missense mutations,
115 2,345 nonsense mutations, and 35 nonstop mutations (Supplementary
116 Table 2). Copy number analysis showed 140,396 gene-level
117 amplifications and 67,605 deletions across 40 cytobands (Supplementary
118 Table 3). The mRNA-seq data characterized the transcription profiles of
119 19,182 genes (Supplementary Table 4). The label-free global proteomics
120 identified 9,152 proteins (encoded by 6,864 genes) with an average of
121 6,457 proteins per sample (Supplementary Table 5). The label-free
122 phosphoproteomics identified 24,990 highly reliable phosphosites from
123 5,832 genes with an average of 10,478 phosphosites per sample
124 (Supplementary Table 6). The quality and reproducibility of the mass
125 spectrometry data were maintained throughout the study (Extended Data
126 Fig. 1c-e).

127 **The air pollutant Benzo[a]pyrene (BaP) primarily contributes to the**
128 **mutation landscape of XWLC**

129 To infer the primary risk factor responsible for the progression of XWLC,
130 we used SomaticSignatures[28] to identify mutational signatures from
131 single nucleotide variants. Mutational signatures were identified in each
132 cohort and a cosine similarity analysis was performed against mutational
133 signatures in COSMIC mutational signatures[29, 30] and environmental
134 agents mutational signatures[31] allowing for inference of the underlying
135 causes (Fig. 1f-i and Extended Fig. 2). Generally, exposure tobacco

136 smoking carcinogens (COSMIC signature) and chemicals such as BaP
137 (environmental agent signature) were identified as the most significant
138 contributing factors in both the XWLC and TSLC cohorts (Fig. 1f and
139 1g). In contrast, defective DNA mismatch repair and endogenous
140 mutational processes initiated by spontaneous deamination of 5-
141 methycytosine were inferred as the major causes in the TNLC and CNLC
142 cohorts (Fig. 1h and 1i). Therefore, the XWLC and TSLC cohorts are
143 more explicitly influenced by environmental carcinogens, while the
144 TNLC and CNLC cohorts may be more affected by age or endogenous
145 risk factors. BaP, a representative compound of polycyclic aromatic
146 hydrocarbons (PAHs), is found in both cigarette smoke and coal smoke
147 and is recognized as a major environmental risk factor for lung
148 cancer[32-34]. Upon metabolism, BaP forms the carcinogenic metabolite
149 7 β ,8 α -dihydroxy-9 α ,10 α -epoxy-7,8,9,10-tetrahydrobenzo[a]pyrene
150 (BPDE), which creates DNA adducts leading to mutations and malignant
151 transformations. This process involves two key regulators: *CYP1A1* and
152 *AhR*. CYP1A1 plays a crucial role in BaP epoxidation at the 7,8 positions,
153 which is the most critical step in BPDE formation[35]. AhR is a ligand-
154 activated transcription factor that responds to various chemicals,
155 including chemical carcinogens, and is activated by BaP[36].
156 Accordingly, our results demonstrated significantly higher expression of
157 CYP1A1 and AhR in tumor samples compared to normal samples (Fig. 1j

158 and 1k). To validate the association between BaP and the elevated
159 incidence of lung cancer in the Xuanwei area, we assessed the serum
160 BPDE levels in individuals residing in this region (Methods). Our
161 findings revealed that the serum BPDE content in individuals from the
162 Xuanwei area averaged around 0.85 nm/ml, and there was no difference
163 between younger and older cases (Fig. 1l and Supplementary Table 7).
164 Furthermore, we treated hiPSC (human induced pluripotent stem cell
165 (hiPSC) with BaP in vitro and found that sites such as Y1173 and Y1068
166 of EGFR were more phosphorylated in BaP treated cells whereas the total
167 abundance of EGFR showed no significant differences (Fig. 1m). All
168 these results provided further evidence for the involvement of BaP and its
169 metabolite in the development of lung cancer.

170 Though coal-smoke related lung cancer (XWLC cohort) and
171 cigarette-smoke related lung cancer (TSLC cohort) showed similar
172 environmental carcinogens, we found that downstream pathway
173 activation and therapeutic targeted potential showed distinctive features.
174 Firstly, the correlation of genomic mutations between XWLC and TSLC
175 was found to be low (Fig. 2a). Secondly, there was a remarkable
176 difference in the fraction of samples affected by pathway mutations
177 between two cohorts (Fig. 2b). Notably, the TSLC cohort exhibited a
178 higher fraction of samples affected by oncogenic pathways comparing to
179 XWLC cohort. Thirdly, mutation frequencies, such as *EGFR*, *TP53*,

180 *RBM10*, and *KRAS* (Fig. 2c), as well as the distribution of amino acid
181 changes in *EGFR* and *TP53*, showed significant differences between the
182 XWLC and other cohorts (Fig. 2d). Specifically, the XWLC cohort
183 exhibited a higher mutation rate in G719C/A/D/S and E746_A750del
184 within the *EGFR* gene compared to the other three cohorts. For the *TP53*
185 gene, frame shift mutations including D9Gfs*3 (n=1), S15Ifs*28 (n=1),
186 D22Efs*61 (n=1), and V34Wfs*49 (n=2) were exclusively detected in the
187 XWLC cohort, whereas tetramer domain mutations were only found in
188 the other three cohorts (Fig. 2d). Finally, there was a noticeable disparity
189 in the percentage of samples with actionable targets among the cohorts.
190 (Fig. 2e). Actionable targets in XWLC cohort were mainly focus on
191 *EGFR* mutations including pG719S, pG719C, pG719A, pL858R,
192 p.E746_A750del whereas TSLC cohort had more actionable targets in
193 *CHEK2* p.K373E, *KRAS* p.G12V/D/C/A (Fig. 2e).

194 Taken together, we found that the XWLC and TSLC cohorts, which
195 are smoke-related lung adenoma groups, demonstrated distinct etiology
196 compared to the TNLC and CNLC cohorts which may be influenced by
197 endogenous risk factors to a greater extent. The blood serum assay
198 provided support for BaP as a promising risk factor specifically in the
199 XWLC cohort. Additionally, significant disparities were observed
200 between XWLC and TSLC in terms of downstream pathway activations
201 and specific oncogene loci. Consequently, we conclude that air pollution-

202 associated lung cancer represents a distinct subtype within LUAD.

203 **The EGFR-G719X mutation, which is a hotspot associated with BaP**
204 **exposure, possesses distinctive biological features**

205 Notably, the XWLC cohort displayed a distinguishable mutation pattern
206 in specific *EGFR* mutation sites compared to the other cohorts (Fig.3a
207 and Fig.2d). In particular, the G719C/A/D/S (G719X) mutation was the
208 most prevalent *EGFR* mutation in the XWLC cohort (20%), while it was
209 rarely found in the other three cohorts (CNLC: 1.9%; TSLC: 1.9%;
210 TNLC: 0) (Fig. 3b). Notably, we found it was a hot spot associated with
211 BaP exposure (Fig. 3c and 3d). Specifically, **GGC** is the 719 codon, the
212 first G can be converted to T (pG719C, n=13) or A (pG719S, n=5), the
213 second G can be converted to A (pG719D, n=1) or C (pG719C, n=8).
214 Thus, pG719C was the most detected mutation type (Fig. 3c). G>T/C>A
215 transversion can be induced by several compounds such as BaP or
216 dibenz(a,h)anthracene (DBA), unlike other compounds, the tallest peak
217 induced by BaP occurs at GpGpG, reflecting how their DNA adducts are
218 formed principally at N²-guanine[31]. Our result showed that the most
219 frequently detected pG719C Achange was correspond to
220 **GGGC>GTGG** transversion (Fig. 3d). Thus, pG719C is a hot spot
221 associated with BaP exposure.

222 We conducted further investigations into the biological characteristics
223 of samples carrying the G719X mutations. Notably, we observed a

224 moderate to high expression of MAPK signaling components, MAP2K2
225 (MEK), and MAPK3 (ERK1), in tumors harboring the EGFR-G719X
226 mutation compared to other EGFR statuses (Fig. 3e-3h). Utilizing
227 hallmark capability analysis and RNA-seq-based estimation of immune
228 cell infiltration, we found that tumors with G719X mutations exhibited
229 similarities to those with L858R mutations (Extended Data Fig. 3a-b).
230 However, patients with G719X mutations were notably younger than
231 those with L858R mutations, indicating a higher occurrence rate of
232 G719X in younger female patients (Fig. 3i). Analysis of overall survival
233 and progression-free interval (PFI) revealed that patients with the G719X
234 mutation had worse outcomes compared to other EGFR mutation
235 subtypes (Fig. 3j and 3k). Furthermore, there were no significant
236 differences in mutation burden or the number of neoantigens between
237 tumors with G719X mutations and tumors with other *EGFR* mutation
238 statuses (Extended Data Fig. 3c).

239 To explore the heterogeneity of signaling pathways activated by
240 different *EGFR* mutation statuses, we conducted Kinase-Substrate
241 Enrichment Analysis (KSEA) [37, 38] based on the XWLC
242 phosphoproteomics dataset. Our analysis of the phosphoproteome across
243 various *EGFR* mutation types revealed distinct activation patterns of
244 kinases. Specifically, the G719X mutation was associated with the
245 activation of PRKCZ, CDK2, AURKB, CSNK1A1, CDK4, and HIPK2.

246 The L858R mutation showed activation of PRKCZ, MAPK7, MAPK12,
247 HIPK2, and CSNK2A1. The Exon19del mutation exhibited activation of
248 CHUK, TTK, PRKCZ, PLK1, NEK2, MAP2K2, CDK2, PRKDC, and
249 MAP2K6. Other EGFR mutations were associated with the activation of
250 AURKB, NEK2, TTK, PLK1, PRKACB, and PRKACG. EGFR-WT
251 mutations showed activation of CSNK1E, PRKCZ, AURKB, CDK2,
252 AURKC, CDK1, CSNK1A1, PRKDC, and CSNK2A1 (Fig. 3l). In
253 Extended Data Fig. 3d, we provide a list of FDA-approved drugs that
254 target the activated kinases in tumors harboring the G719X mutation.
255 Currently, afatinib is widely regarded as a first-line therapy for patients
256 with the G719X mutation[39-42]. However, reports indicate that 80% of
257 patients with this mutation may develop resistance to afatinib, even in the
258 absence of T790M[43], underscoring the need for a deeper understanding
259 of the downstream pathways associated with the G719X mutation.
260 Therefore, a promising approach to overcome resistance in tumors with
261 this mutation could involve combining afatinib, which targets activated
262 EGFR, with FDA-approved drugs that specifically target the activated
263 kinases associated with G719X.

264 **Clinically relevant Subtyping in XWLC**

265 To uncover the inherent subgroups within air-pollution-associated tumors,
266 we employed unsupervised Consensus Clustering[44] on integrated RNA,
267 protein, and phosphoprotein profiles of XWLC tumor samples. This

268 analysis led to the identification of four distinct intrinsic clusters, denoted
269 as MC-I, II, III, and IV (Fig. 4a, Extended Data Fig. 4a and Methods).
270 Further survival analysis demonstrated that patients belonging to the MC-
271 IV group exhibited the poorest overall survival compared to the other
272 three subgroups, thus indicating the prognostic potential of multi-omic
273 clustering (Fig. 4b). Notably, there were no significant differences in
274 clinical features such as age and stage observed among the four
275 subgroups (Fig. 4a). Since BaP has the potential to induce the expression
276 of CYP1A1 and AhR, we utilized the expression levels of CYP1A1 and
277 AhR as indicators of BaP activity. Consequently, we compared the
278 expression levels of CYP1A1 and AhR among the four subtypes to
279 identify the tumor type most strongly associated with air pollution. Our
280 findings revealed that the MC-II subtype exhibited higher expression of
281 CYP1A1 (Fig. 4c) and moderately elevated expression level of AhR
282 compared to the other three subtypes (Fig. 4d) Moreover, the MC-II
283 possessed more G719X mutations (MC-I:0.39, MC-II:0.42, MC-III: 0.20,
284 MC-IV: 0.08). These results collectively suggest a stronger association
285 between the MC-II subtype and BaP exposure. Notably, there was a
286 significant correlation between CYP1A1 and EGFR expression (Fig. 4e),
287 with EGFR being more highly expressed in the MC-II subtype (Fig. 4f).
288 These results indicated that MC-II was more associated with air-
289 pollution.

290 Through Kinase Substrate Enrichment Analysis (KSEA) of the
291 phosphoproteome in tumor samples compared to normal adjacent tissues
292 (NATs), we identified specific kinase activations within the four
293 subgroups. In MC-I samples, kinase activations included PRKDC,
294 PRKCZ, CSNK1A1, NEK2, GSK3A, and ROCK1. MC-II samples
295 showed activations of CDK2, CDK1, AURKA, TTK, CDK6, and CHUK.
296 MC-III samples exhibited activations of AKT1, AKT3, RPS6KB1,
297 CSNK2A1, and PAK2. Finally, CDK2 and ROCK1 were activated in
298 MC-IV samples (Fig. 4g). Particularly noteworthy is the enrichment of
299 CDK1/2/6 kinases, which regulate cell cycle checkpoints, in the MC-II
300 subtype, indicating its high proliferation capabilities. These findings
301 imply that distinct kinase pathways are activated within each subgroup,
302 suggesting the presence of specific therapeutic targets for each subgroup.
303 Consequently, we proceeded to explore therapeutic strategies for each
304 subgroup as outlined below:

305 The **MC-IV** subtype exhibited the poorest overall survival compared
306 to the other three subtypes (Fig. 4b). Given the crucial role of epithelial-
307 mesenchymal transition (EMT) in malignant progression, our first
308 evaluation focused on the EMT process across the four subtypes. We
309 observed higher expression levels of mesenchymal markers such as *VIM*,
310 *FNI*, *TWIST2*, *SNAI2*, *ZEB1*, *ZEB2*, and others in the MC-IV subtype
311 (Fig. 5a). To comprehensively assess the EMT capability of the MC-IV

312 subtype, we calculated EMT scores using the ssGSEA enrichment
313 method based on protein levels and GSEA hallmark gene set
314 (M5930)[45]. The results confirmed the elevated EMT capability of the
315 MC-IV subtype at the protein level (Fig. 5b). Furthermore, Fibronectin
316 (FN1), an EMT marker that promotes the dissociation, migration, and
317 invasion of epithelial cells, was found to be highly expressed in the MC-
318 IV subtype at the protein level (Fig. 5c). Additionally, β -Catenin, a key
319 regulator in initiating EMT, was highly expressed in the MC-IV subtype
320 at the protein level (Fig. 5d). Collectively, our findings demonstrate that
321 the MC-IV subtype possesses an enhanced EMT capability, which may
322 contribute to the high malignancy observed in this subtype.

323 The **MC-II** subtype demonstrated the second-worst outcome and was
324 found to be more strongly associated with air pollution (Fig. 4). This
325 subtype exhibited dysregulation of cell cycle processes, including cell
326 division, glycolysis, and cell cycle biological processes (Fig. 4a). The
327 KSEA analysis revealed that the CDK1 and CDK2 pathways, which are
328 closely linked to cell cycle regulation, were predominantly activated in
329 the MC-II subtype (Fig. 5e). Consistently, we observed higher expression
330 levels of CDK1 and CDK2 at both the protein and phosphoprotein levels
331 in the MC-II subtype, indicating specific elevation of the G2M phase in
332 the cell cycle (Fig. 5e). The cell cycle and glycolysis processes are tightly
333 coordinated, allowing cells to synchronize their metabolic state and

334 energy requirements with cell cycle progression to ensure proper cell
335 growth and division[46-48]. In line with this, we found that key enzymes
336 involved in glycolysis regulation, such as Hexokinase 1 (HK1),
337 Glyceraldehyde-3-phosphate dehydrogenase (GAPDH), and
338 Glyceraldehyde-3-phosphate dehydrogenase-like protein (GPL), were
339 highly expressed in the MC-II subtype (Fig. 5f). Additionally, the MC-II
340 subtype was enriched with EGFR mutations (MC-II vs. others: 18/24 vs.
341 51/110; Fisher's exact $p = 0.013$) and TP53 mutations (MC-II vs. others:
342 14/24 vs. 35/110; Fisher's exact $p = 0.019$), consistent with the
343 characteristic loss of control over cell proliferation. In summary, the MC-
344 II subtype exhibited dysregulated cell cycle processes accompanied by an
345 elevated glycolysis capability, indicating a distinct metabolic and
346 proliferative phenotype.

347 The **MC-I** subtype exhibited enrichment in various biological
348 processes including angiogenesis, the cAMP signaling pathway,
349 complement and coagulation cascades, PDL1 expression, the PD-1
350 checkpoint pathway, leukocyte transendothelial migration, and actin
351 cytoskeleton processes (Fig. 4a). In-depth exploration of key components
352 involved in angiogenesis revealed that vascular endothelial growth factor
353 A (VEGFA), a growth factor crucial for both physiological and
354 pathological angiogenesis, was highly expressed in the MC-I subtype (Fig.
355 5g). Additionally, phosphorylation of vascular endothelial growth factor

356 receptor 1 (VEGFR1), a receptor tyrosine kinase essential for
357 angiogenesis and vasculogenesis, was also highly expressed in the MC-I
358 subtype (Fig. 5h). The angiogenesis scores, calculated using the ssGSEA
359 method based on protein levels and the hallmark gene set (M5944), were
360 relatively high in the MC-I and MC-IV subtypes (Fig. 5i). Furthermore,
361 the relationship between the Notch signaling pathway and angiogenesis is
362 well-established[49]. Notch signaling plays a role in multiple aspects of
363 angiogenesis, including endothelial cell sprouting, vessel branching, and
364 vessel maturation [50, 51]. In the MC-I subtype, the expression of Notch
365 receptors (Notch1-4) and ligands (DLL1, DLL4, JAG1, and JAG2) was
366 highly elevated, indicating increased activation of Notch signaling (Fig.
367 5j). KEGG pathview analysis demonstrated that key regulators of the
368 VEGF signaling pathway were highly expressed in the MC-I subtype
369 (Extended Fig.5a). Therefore, manipulating Notch signaling could
370 potentially serve as a strategy to regulate angiogenesis and control
371 pathological angiogenesis in the MC-I subtype.

372 The **MC-III** subtype is characterized by the upregulation of various
373 metabolic processes, including oxidative phosphorylation, peroxisome
374 function, adipogenesis, fatty acid metabolism, and xenobiotic
375 metabolism-related processes (Fig. 5k). Additionally, we conducted
376 further investigations into the immune features across the subtypes.
377 Interestingly, we observed higher expression of genes associated with

378 PD-1 signaling (GSEA, SYSTEMATIC_NAME M18810) in the MC-III
379 subtype (Fig. 5l). Since PD-1 is primarily expressed on the surface of
380 certain immune cells, particularly activated T cells, we inferred the
381 immune cell infiltration using the ssGSEA method based on immune cell-
382 specific gene sets. We found that activated CD8⁺ T cells exhibited higher
383 infiltration levels in the MC-III subtype compared to the other three
384 subtypes (Fig. 5m and Supplementary Table 8), which may explain the
385 elevated PD-1 signaling in the MC-III subtype. Furthermore, we
386 examined the expression of receptor-ligand pairs involved in both anti-
387 tumor and pro-tumor lymphocyte recruitment. Remarkably, the MC-III
388 subtype exhibited specific high expression of anti-tumor lymphocyte
389 receptors and ligands, while the expression of pro-tumor lymphocyte
390 receptors and ligands was relatively lower (Fig. 5n). In general, the MC-I
391 subtype showed the reverse expression trend in terms of anti-tumor and
392 pro-tumor receptor-ligand pairs (Fig. 5n).

393 In conclusion, our classification of lung adenocarcinoma associated
394 with air pollution resulted in the identification of four subtypes, each
395 exhibiting distinct biological pathway activation and immune features.
396 The MC-I subtype demonstrated elevated angiogenesis processes, while
397 the MC-II subtype showed a high capacity for cell division and glycolysis.
398 The MC-III subtype exhibited a notable infiltration of CD8⁺ cells, and
399 the MC-IV subtype was characterized by high EMT capability, which

400 may contribute to its poor outcome. These findings have significant
401 implications for the development of precision treatments for XWLC
402 (Extended Fig.5b).

403 **Tumor microbiota composition and radiomic features across** 404 **subtypes**

405 Host-bacterial immune interactions profoundly influence tumorigenesis,
406 cancer progression, and response to therapy[52-54]. Thus, we further
407 generated species-level resolution data allowing the identification of
408 specific bacteria across subgroups. Firstly, we extracted and mapped
409 high-quality reads of bacterial origin from RNA-Seq datasets across
410 normal and four tumor subgroups of the XWLC cohort using PathSeq[55]
411 while carefully mitigating the influences of microbial contamination
412 (Supplementary Table 9 and Methods). We found that the subgroups
413 showed distinct bacterial abundances (Fig. 6a). Principal coordinate
414 analysis (PCoA)[56] of UniFrac distances of species abundances across
415 subgroups revealed that the first PCoA axis generally captured
416 differences between tumor and adjacent normal tissues, while the second
417 axis mainly captured the differences between the MC-IV subgroup and
418 the other three subgroups (Fig. 6b). Thus, the MC-IV subgroup showed
419 distinguishable bacterial composition compared to the other three
420 subgroups. We further identified 16 bacterial species that were enriched
421 in the MC-IV subgroup, mainly belonging to the *Pseudomonas*,

422 *Stenotrophomonas*, *Acidovorax*, *Rheinheimera*, *Brevundimonas*,
423 *Sphingomonas*, and *Methyloversatilis* genera (Fig. 6c). Among them,
424 *Acidovorax* and *Brevundimonas* have been identified within lung tumor
425 sections[57]. Further investigation is warranted to determine whether and
426 how the unique microbiota composition of MC-IV contributes to its
427 malignancy.

428 Furthermore, we built a noninvasive method to distinguish MC
429 subtypes with radiomics which entails the extensive quantification of
430 tumor phenotypes by utilizing numerous quantitative image features. In
431 the initial step, we defined 107 quantitative image features that describe
432 various characteristics of tumor phenotypes, including tumor image
433 intensity, size, shape, and texture. These features were derived from X-
434 ray computed tomography (CT) scans of 155 patients with XWLC
435 (Methods). The baseline characteristics of this cohort can be found in
436 Supplementary Table 10. Firstly, all features were compared among the
437 four subtypes, and notably, eight features showed significant differences
438 between the MC-II subtype and the other three subtypes (Fig. 6d),
439 suggesting a denser image in the MC-II subtype. We further established a
440 signature using a multivariate linear regression model with five image
441 features to distinguish MC-II from the other three subgroups (Extended
442 Fig.6). The performance of the five-feature radiomic signature was
443 validated using the AUC value, which is a generation of the area under

444 the ROC curve. The radiomic signature had an AUC value of 0.94 in the
445 training set and 0.83 in the validation set (Fig. 6e). The confusion matrix
446 revealed an overall accuracy of 0.875 for sample classification using the
447 signature, indicating proficient performance. However, it exhibited
448 suboptimal performance in terms of false-negative classification (Fig. 6f).
449 Taken together, we found that MC-II showed a dense image phenotype,
450 which can be noninvasively distinguished using radiomic features.

451 **Identification of novel targets based on mutation-informed protein-** 452 **protein interface (PPI) analysis**

453 The integration of genomics and interactomics has enabled the discovery
454 of functional and biological consequences of disease mutations[58, 59].
455 To explore novel targets with the concept, we created PPI networks with
456 structural resolution using missense mutations from the XWLC, CNLC,
457 TSLC, and TNLC cohorts (Fig.7a and Methods). OncoPPIs, defined as a
458 significant enrichment of interface mutations in either of the two protein-
459 binding partners across individuals, were identified in each cohort and
460 were provided in Supplementary Tables 11. The OncoPPIs from the four
461 cohorts are named XWLC_oncoPPIs, CNLC_oncoPPIs,
462 TSLC_oncoPPIs, and TNLC_oncoPPIs, respectively (Fig. 7b and
463 Extended Fig. 7a-c). Initially, the nodes from these four OncoPPIs were
464 subjected to biological process enrichment analysis (Extended Fig. 7d).
465 The analysis revealed that biological processes such as regulation of

466 mitotic cell cycle, TGF-beta signaling pathway, and immune system were
467 predominantly enriched in the genes related to OncoPPIs. Moreover, the
468 processes disrupted by interface mutations showed a relatively higher
469 similarity between the XWLC and TSLC cohorts (Extended Fig. 7d)
470 suggesting convergent targets or pathways affected by smoke-induced
471 mutations.

472 To refine the novel targets from XWLC_oncoPPs, we performed
473 molecular dynamics simulations to predict the binding affinity change by
474 the interface-located mutations (Methods). Mitotic Arrest Deficient 1
475 Like 1 (MAD1), a crucial component of the mitotic spindle-assembly
476 checkpoint[60, 61], forms a tight core complex with MAD2, facilitating
477 the binding of MAD2 to CDC20[62], which plays a critical role in sister
478 chromatid separation during the metaphase-anaphase transition[63].
479 Specifically, MAD1 Arg558His has been identified as a susceptibility
480 factor for lung cancer[64] and colorectal cancer[65]. Here, we found that
481 MAD1 allele carrying a p.Arg558His substitution may disrupt the
482 interaction between MAD1 and MAD2 (Fig.7b-c). To assess this, we
483 performed molecular dynamics simulations and found that the binding
484 affinity between Arg558His MAD1 and MAD2 was -195.091 kJ/mol and
485 that of wild type was -442.712 kJ/mol (Fig.7c). Furthermore, on a per-
486 residue basis, the predicted binding affinity ($\Delta\Delta G$) of Arg558His
487 (Extended Fig.7e) was projected to increase by 118.319 kJ/mol relative to

488 the wild type (Extended Fig.7f), indicating that the substitution of
489 Arg558His in MAD1 perturbs the binding affinity. Thus, our findings
490 suggest that the MAD1 Arg558His attributed to lung cancer progression
491 by disrupting of the interaction between MAD1 and MAD2, which
492 showed potential to be explored as a target.

493 Notably, we identified TPRN as a novel significantly mutated gene in
494 the XWLC cohort (Extended Fig.7g and 7h) whose status was also
495 associated with patients' outcomes (Fig.7e). Previous studies have
496 reported that TPRN interacts with PPP1CA[66]. Thus, we assessed the
497 binding affinity of the TPRN-PPP1CA complex affected by the mutant
498 variant His550Gln. Our result showed that the binding affinities of the
499 complex were -694.372 kJ/mol and -877.570 kJ/mol in mutant and WT
500 cases, respectively (Fig.7d). On a per-residue basis, the predicted $\Delta\Delta G$ of
501 His550Gln compared to the wild type exhibited an increase of 96.774
502 kJ/mol (Extended Fig. 7i-j). All these results indicated that TPRN
503 His550Gln increase the binding affinity of the TPRN-PPP1CA complex.
504 To investigate the effect of TPRN His550Gln mutation on tumor
505 progression, we examined proliferation and migration capabilities in both
506 A549 and H1299 lung adenocarcinoma cell lines. CCK-8 assay showed
507 significantly enhanced cell growth after transfection of TPRN mutant
508 allele in both A549 (Fig.7f) and H1299 cells (Extended Fig. 8a).
509 Moreover, wound-healing assay showed that TPRN mutant cell had

510 achieved enhanced migration capacity (Fig.7g-h and Extended Fig. 8b-c).
511 Finally, more cell clones in TPRN His550Gln mutation cells were
512 observed in both TPRN-mutant A549 and H1299 cells (Fig.7i). All these
513 results supported that TPRN His550Gln could be explored as a target in
514 XWLC.

515 Taken together, our integrated analysis of oncoPPIs and molecular
516 dynamics simulations showed potential to explore novel therapeutic
517 vulnerabilities.

518 **Discussion**

519 In this study, we conducted proteogenomic and characterized air-
520 pollution-related lung cancers. We found that Benzo[a]pyrene (BaP)
521 influenced the mutation landscape, particularly the EGFR-G719X hotspot
522 found in 20% of cases. This mutation correlated with elevated MAPK
523 pathway activation, worse clinical outcomes, and younger patients. Multi-
524 omics clustering identified four subtypes with unique biological pathways
525 and immune cell patterns. Moreover, our analysis of protein-protein
526 interfaces unveiled novel therapeutic targets. These findings have
527 significant implications for preventing and developing precise treatments
528 for air-pollution-associated lung cancers.

529 Previously considered uncommon, the EGFR-G719X mutation was
530 detected in only 1-2% of CNLC or TCGA-LUAD cohort samples.
531 Limited knowledge exists from G719X, mostly based on isolated case

532 reports or small series studies[67-71]. In vitro experiments using G719X
533 mutant cell lines and patient-derived xenografts (PDX) demonstrated that
534 osimertinib effectively inhibits signaling pathways and cellular growth,
535 leading to sustained tumor growth inhibition[72]. However, in silico
536 protein structure analysis suggests that G719 alterations may confer
537 osimertinib resistance due to reduced EGFR binding[73]. Presently,
538 afatinib is proposed as the first-line therapy for G719X mutation
539 patients[39-42]. Unfortunately, 80% of G719X patients develop acquired
540 resistance to afatinib without detecting the T790M mutation⁴⁰. Hence,
541 further mechanistic studies are warranted for G719X. Our study reveals
542 that the G719X mutation is prevalent in the XWLC cohort, significantly
543 impacting treatment selection. Additionally, the large number of G719X
544 samples allowed us to uncover variations in biology and pathway
545 activation, which may facilitate the development of more precise targeted
546 therapies for these patients.

547 There is substantial evidence linking lung cancer in the Xuanwei area
548 to coal smoke[18, 19, 21, 22, 74, 75]. In addition, we conducted a rat
549 model study that demonstrated the induction of lung cancer by local coal
550 smoke exposure[22]. However, the specific chemical compound in coal
551 smoke responsible for causing lung cancer remains largely unknown.
552 Previous research has mainly focused on studying indoor concentrations
553 of airborne particles and BaP[18, 19, 21, 23, 76]. For instance, studies

554 have shown an association between the concentration of BaP and lung
555 cancer rates across counties[18]. Moreover, improvements in household
556 stoves have led to reduced exposure to benzopyrene and particulate
557 matter, benefiting people’s health[23, 76]. However, these studies
558 primarily relied on epidemiological data, which may be influenced by
559 confounding factors. In our study, we used clinical samples and linked the
560 mutational signatures of XWLC to the chemical compound BaP, which
561 advanced the etiology and mechanism of air-pollution-induced lung
562 cancer.

563 In summary, our proteogenomic analysis of clinical tumor samples
564 provides insights into air-pollution-associated lung cancers, especially
565 those induced by coal smoke and offers an opportunity to expedite the
566 translation of basic research to more precise diagnosis and treatment in
567 the clinic.

568 **Declarations**

569 **Ethical Approval**

570 The study protocol was reviewed and approved by the ethical committees
571 of Yunnan Cancer Hospital & The Third Affiliated Hospital of Kunming
572 Medical University (KYCS2022067) and conformed to the ethical
573 standards for medical research involving human subjects, as laid out in
574 the 1964 Declaration of Helsinki and its later amendments.

575 **Consent to participate**

576 Participants provided written informed consent prior to taking part in the
577 study.

578 **Consent for publication**

579 Written informed consent was obtained from the patient and the ethical
580 committees of Yunnan Cancer Hospital & The Third Affiliated Hospital
581 of Kunming Medical University for publication of this study.

582 **Competing interests**

583 The authors declare no competing interests.

584 **Authors' contributions**

585 Honglei Zhang, Gaofeng Li, Xiaosan Su and Yong Zhang conceived the
586 study. Honglei Zhang, Chao Liu, Qing Wang and Shuting Wang
587 performed data analysis. Honglei Zhang wrote the manuscript. Xu Feng,
588 Heng Li, Zhiyong Deng, Huawei Jiang, Yunyan He, Chao Luo, Rou Qian,
589 Minjun Zhou, Zhi Li, Li Xiao, Yong Zhang interpreted the data analysis
590 and drafted the manuscript and critically revised the manuscript. All
591 authors critically revised and gave final approval of the manuscript.

592 **Funding**

593 This work was supported by National Natural Science Foundation (Nos.
594 82273501, 81960322, 82160343, 82060519, 81860171); Yunnan young
595 and middle-aged academic and technical leaders Reserve Talents Project
596 (No. 202005AC160048); Yunnan basic research program
597 (Nos.202101AZ070001-002, 202001AY070001-277, 2019FE001(-236));

598 Medical Reserve Personnel Training Program of Yunnan Provincial
599 Health Commission(No. H-2018097); “Famous Doctor” Special Project
600 of Ten Thousand People Plan of Yunnan Province(Nos.CZ0096, YNWR-
601 MY-2020-095); Medical Leading Talents Training Program of Yunnan
602 Provincial Health Commission (No.L-2019028).

603 **Acknowledgements**

604 We thank Novogene Co., Ltd. and Beijing Qinglian Biotech Co., Ltd. for
605 the analysis of mass spectrometric data.

606 **Availability of data and materials**

607 Raw sequencing data have been deposited in the Genome Sequence
608 Archive for Human (GSA-Human, <https://ngdc.cncb.ac.cn/gsa-human/>)
609 under accession codes HRA000124, HRA001481 and HRA001482.
610 Proteomics and phosphoproteomics data have been deposited in the Open
611 Archive for Miscellaneous Data (OMIX, <https://ngdc.cncb.ac.cn/omix/>)
612 under accession codes OMIX001292. The raw lung CT images used in
613 this paper are available from OMIX under accession codes OMIX002491.

614

615

616 References

- 617 1. Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A: **Global cancer statistics**
618 **2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185**
619 **countries.** *CA Cancer J Clin* 2018, **68**:394-424.
- 620 2. Parkin DM, Bray F, Ferlay J, Pisani P: **Global cancer statistics, 2002.** *CA Cancer J Clin*
621 2005, **55**:74-108.
- 622 3. Chen J, Yang H, Teo ASM, Amer LB, Sherbaf FG, Tan CQ, Alvarez JJS, Lu B, Lim JQ,
623 Takano A, et al: **Genomic landscape of lung adenocarcinoma in East Asians.** *Nat Genet*
624 2020, **52**:177-186.
- 625 4. Govindan R, Ding L, Griffith M, Subramanian J, Dees ND, Kanchi KL, Maher CA, Fulton R,
626 Fulton L, Wallis J, et al: **Genomic landscape of non-small cell lung cancer in smokers and**
627 **never-smokers.** *Cell* 2012, **150**:1121-1134.
- 628 5. Chen YJ, Roumeliotis TI, Chang YH, Chen CT, Han CL, Lin MH, Chen HW, Chang GC,
629 Chang YL, Wu CT, et al: **Proteogenomics of Non-smoking Lung Cancer in East Asia**
630 **Delineates Molecular Signatures of Pathogenesis and Progression.** *Cell* 2020, **182**:226-
631 244.e217.
- 632 6. Zhang T, Joubert P, Ansari-Pour N, Zhao W, Hoang PH, Lokanga R, Moye AL, Rosenbaum J,
633 Gonzalez-Perez A, Martinez-Jimenez F, et al: **Genomic and evolutionary classification of**
634 **lung cancer in never smokers.** *Nat Genet* 2021, **53**:1348-1359.
- 635 7. Zhang XC, Wang J, Shao GG, Wang Q, Qu X, Wang B, Moy C, Fan Y, Albertyn Z, Huang X,
636 et al: **Comprehensive genomic and immunological characterization of Chinese non-small**
637 **cell lung cancer patients.** *Nat Commun* 2019, **10**:1772.
- 638 8. Xu JY, Zhang C, Wang X, Zhai L, Ma Y, Mao Y, Qian K, Sun C, Liu Z, Jiang S, et al:
639 **Integrative Proteomic Characterization of Human Lung Adenocarcinoma.** *Cell* 2020,
640 **182**:245-261.e217.
- 641 9. Gillette MA, Satpathy S, Cao S, Dhanasekaran SM, Vasaikar SV, Krug K, Petralia F, Li Y,
642 Liang WW, Reva B, et al: **Proteogenomic Characterization Reveals Therapeutic**
643 **Vulnerabilities in Lung Adenocarcinoma.** *Cell* 2020, **182**:200-225 e235.
- 644 10. Hill W, Lim EL, Weeden CE, Lee C, Augustine M, Chen K, Kuan FC, Marongiu F, Evans EJ,
645 Jr., Moore DA, et al: **Lung adenocarcinoma promotion by air pollutants.** *Nature* 2023,
646 **616**:159-167.
- 647 11. Turner MC, Andersen ZJ, Baccarelli A, Diver WR, Gapstur SM, Pope CA, 3rd, Prada D,
648 Samet J, Thurston G, Cohen A: **Outdoor air pollution and cancer: An overview of the**
649 **current evidence and public health recommendations.** *CA Cancer J Clin* 2020.
- 650 12. Fajersztajn L, Veras M, Barrozo LV, Saldiva P: **Air pollution: a potentially modifiable risk**
651 **factor for lung cancer.** *Nat Rev Cancer* 2013, **13**:674-678.
- 652 13. Cogliano VJ, Baan R, Straif K, Grosse Y, Lauby-Secretan B, El Ghissassi F, Bouvard V,
653 Benbrahim-Tallaa L, Guha N, Freeman C, et al: **Preventable exposures associated with**
654 **human cancers.** *J Natl Cancer Inst* 2011, **103**:1827-1839.
- 655 14. Humans IWGotEoCRt: **Some non-heterocyclic polycyclic aromatic hydrocarbons and**
656 **some related exposures.** *IARC Monogr Eval Carcinog Risks Hum* 2010, **92**:1-853.
- 657 15. Shi Q, Godschalk RWL, van Schooten FJ: **Inflammation and the chemical carcinogen**
658 **benzo[a]pyrene: Partners in crime.** *Mutat Res Rev Mutat Res* 2017, **774**:12-24.

- 659 16. Saravanakumar K, Sivasantosh S, Sathiyaseelan A, Sankaranarayanan A, Naveen KV, Zhang
660 X, Jamla M, Vijayasarathy S, Vishnu Priya V, MubarakAli D, Wang MH: **Impact of**
661 **benzo[a]pyrene with other pollutants induce the molecular alternation in the biological**
662 **system: Existence, detection, and remediation methods.** *Environ Pollut* 2022, **304**:119207.
- 663 17. Abd El-Fattah EE, Abdelhamid AM: **Benzo[a]pyrene immunogenetics and immune**
664 **archetype reprogramming of lung.** *Toxicology* 2021, **463**:152994.
- 665 18. Mumford JL, He XZ, Chapman RS, Cao SR, Harris DB, Li XM, Xian YL, Jiang WZ, Xu CW,
666 Chuang JC, et al.: **Lung cancer and indoor air pollution in Xuan Wei, China.** *Science*
667 1987, **235**:217-220.
- 668 19. Barone-Adesi F, Chapman RS, Silverman DT, He X, Hu W, Vermeulen R, Ning B, Fraumeni
669 JF, Jr., Rothman N, Lan Q: **Risk of lung cancer associated with domestic use of coal in**
670 **Xuanwei, China: retrospective cohort study.** *BMJ* 2012, **345**:e5414.
- 671 20. Chapman RS, He X, Blair AE, Lan Q: **Improvement in household stoves and risk of**
672 **chronic obstructive pulmonary disease in Xuanwei, China: retrospective cohort study.**
673 *BMJ* 2005, **331**:1050.
- 674 21. Lin H, Ning B, Li J, Ho SC, Huss A, Vermeulen R, Tian L: **Lung cancer mortality among**
675 **women in Xuan Wei, China: a comparison of spatial clustering detection methods.** *Asia*
676 *Pac J Public Health* 2015, **27**:NP392-401.
- 677 22. Zhang H, Liu C, Li L, Feng X, Wang Q, Li J, Xu S, Wang S, Yang Q, Shen Z, et al: **Genomic**
678 **evidence of lung carcinogenesis associated with coal smoke in Xuanwei area, China.** *Natl*
679 *Sci Rev* 2021, **8**:nwab152.
- 680 23. Lan Q, Chapman RS, Schreinemachers DM, Tian L, He X: **Household stove improvement**
681 **and risk of lung cancer in Xuanwei, China.** *J Natl Cancer Inst* 2002, **94**:826-835.
- 682 24. Hosgood HD, Pao W, Rothman N, Hu W, Pan YH, Kuchinsky K, Jones KD, Xu J, Vermeulen
683 R, Simko J, Lan Q: **Driver mutations among never smoking female lung cancer tissues in**
684 **China identify unique EGFR and KRAS mutation pattern associated with household**
685 **coal burning.** *Respiratory Medicine* 2013, **107**:1755-1762.
- 686 25. Wang J, Duan Y, Meng QH, Gong R, Guo C, Zhao Y, Zhang Y: **Integrated analysis of DNA**
687 **methylation profiling and gene expression profiling identifies novel markers in lung**
688 **cancer in Xuanwei, China.** *PLoS One* 2018, **13**:e0203155.
- 689 26. Wang X, Li J, Duan Y, Wu H, Xu Q, Zhang Y: **Whole genome sequencing analysis of lung**
690 **adenocarcinoma in Xuanwei, China.** *Thorac Cancer* 2017, **8**:88-96.
- 691 27. Cancer Genome Atlas Research N: **Comprehensive molecular profiling of lung**
692 **adenocarcinoma.** *Nature* 2014, **511**:543-550.
- 693 28. Gehring JS, Fischer B, Lawrence M, Huber W: **SomaticSignatures: inferring mutational**
694 **signatures from single-nucleotide variants.** *Bioinformatics* 2015, **31**:3673-3675.
- 695 29. Alexandrov LB, Nik-Zainal S, Wedge DC, Campbell PJ, Stratton MR: **Deciphering**
696 **signatures of mutational processes operative in human cancer.** *Cell Rep* 2013, **3**:246-259.
- 697 30. Alexandrov LB, Nik-Zainal S, Wedge DC, Aparicio SA, Behjati S, Biankin AV, Bignell GR,
698 Bolli N, Borg A, Borresen-Dale AL, et al: **Signatures of mutational processes in human**
699 **cancer.** *Nature* 2013, **500**:415-421.
- 700 31. Kucab JE, Zou X, Morganello S, Joel M, Nanda AS, Nagy E, Gomez C, Degasperi A, Harris
701 R, Jackson SP, et al: **A Compendium of Mutational Signatures of Environmental Agents.**
702 *Cell* 2019, **177**:821-836 e816.

- 703 32. Petit P, Maitre A, Persoons R, Bicout DJ: **Lung cancer risk assessment for workers exposed**
704 **to polycyclic aromatic hydrocarbons in various industries.** *Environment International*
705 2019, **124**:109-120.
- 706 33. Widziewicz K, Rogula-Kozłowska W, Majewski G: **Lung Cancer Risk Associated with**
707 **Exposure to Benzo(A)Pyrene in Polish Agglomerations, Cities, and Other Areas.**
708 *International Journal of Environmental Research* 2017, **11**:685-693.
- 709 34. Mangal D, Vudathala D, Park JH, Lee SH, Penning TM, Blair IA: **Analysis of 7,8-dihydro-8-**
710 **oxo-2'-deoxyguanosine in cellular DNA during oxidative stress.** *Chem Res Toxicol* 2009,
711 **22**:788-797.
- 712 35. Chung JY, Kim JY, Kim YJ, Jung SJ, Park JE, Lee SG, Kim JT, Oh S, Lee CJ, Yoon YD, et al:
713 **Cellular defense mechanisms against benzo[a]pyrene in testicular Leydig cells:**
714 **implications of p53, aryl-hydrocarbon receptor, and cytochrome P450 1A1 status.**
715 *Endocrinology* 2007, **148**:6134-6144.
- 716 36. Hidaka T, Fujimura T, Aiba S: **Aryl Hydrocarbon Receptor Modulates Carcinogenesis and**
717 **Maintenance of Skin Cancers.** *Front Med (Lausanne)* 2019, **6**:194.
- 718 37. Wiredja DD, Koyuturk M, Chance MR: **The KSEA App: a web-based tool for kinase**
719 **activity inference from quantitative phosphoproteomics.** *Bioinformatics* 2017, **33**:3489-
720 3491.
- 721 38. Casado P, Rodriguez-Prados JC, Cosulich SC, Guichard S, Vanhaesebroeck B, Joel S, Cutillas
722 PR: **Kinase-substrate enrichment analysis provides insights into the heterogeneity of**
723 **signaling pathway activation in leukemia cells.** *Sci Signal* 2013, **6**:rs6.
- 724 39. Ettinger DS, Wood DE, Aisner DL, Akerley W, Bauman JR, Bharat A, Bruno DS, Chang JY,
725 Chirieac LR, D'Amico TA, et al: **Non-Small Cell Lung Cancer, Version 3.2022, NCCN**
726 **Clinical Practice Guidelines in Oncology.** *J Natl Compr Canc Netw* 2022, **20**:497-530.
- 727 40. Janning M, Suptitz J, Albers-Leischner C, Delpy P, Tufman A, Velthaus-Rusik JL, Reck M,
728 Jung A, Kauffmann-Guerrero D, Bonzheim I, et al: **Treatment outcome of atypical EGFR**
729 **mutations in the German National Network Genomic Medicine Lung Cancer (nNGM).**
730 *Ann Oncol* 2022, **33**:602-615.
- 731 41. Yang JC, Schuler M, Popat S, Miura S, Heeke S, Park K, Marten A, Kim ES: **Afatinib for the**
732 **Treatment of NSCLC Harboring Uncommon EGFR Mutations: A Database of 693**
733 **Cases.** *J Thorac Oncol* 2020, **15**:803-815.
- 734 42. Cho JH, Lim SH, An HJ, Kim KH, Park KU, Kang EJ, Choi YH, Ahn MS, Lee MH, Sun JM,
735 et al: **Osimertinib for Patients With Non-Small-Cell Lung Cancer Harboring Uncommon**
736 **EGFR Mutations: A Multicenter, Open-Label, Phase II Trial (KCSG-LU15-09).** *J Clin*
737 *Oncol* 2020, **38**:488-495.
- 738 43. Harada T, Futamura S, Inoue Y, Sawada R, Okuda T, Kagawa K: **P2.03-13 Acquired**
739 **Resistance to Afatinib in Non-Small Cell Lung Cancer with EGFR G719X Mutation.**
740 *Journal of Thoracic Oncology* 2019, **14**:S687.
- 741 44. Wilkerson MD, Hayes DN: **ConsensusClusterPlus: a class discovery tool with confidence**
742 **assessments and item tracking.** *Bioinformatics* 2010, **26**:1572-1573.
- 743 45. Liberzon A, Birger C, Thorvaldsdottir H, Ghandi M, Mesirov JP, Tamayo P: **The Molecular**
744 **Signatures Database (MSigDB) hallmark gene set collection.** *Cell Syst* 2015, **1**:417-425.
- 745 46. Vander Heiden MG, Cantley LC, Thompson CB: **Understanding the Warburg effect: the**
746 **metabolic requirements of cell proliferation.** *Science* 2009, **324**:1029-1033.

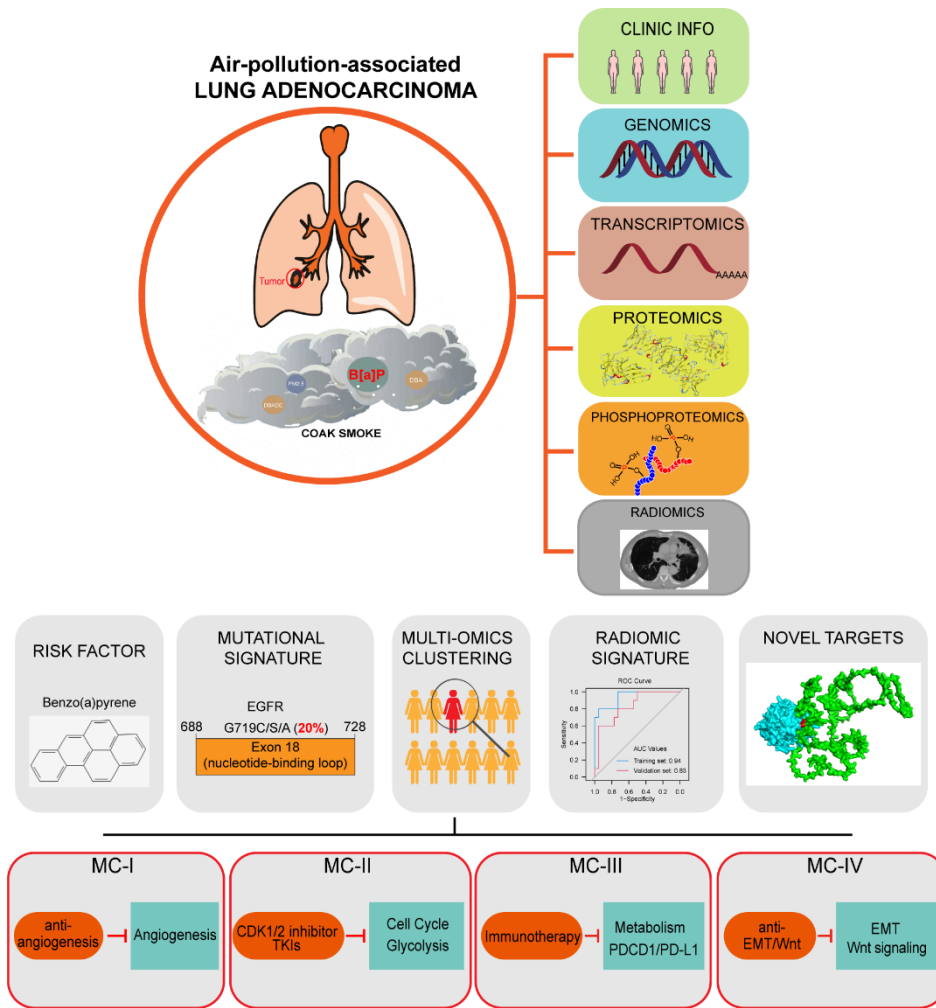
- 747 47. DeBerardinis RJ, Thompson CB: **Cellular metabolism and disease: what do metabolic**
748 **outliers teach us?** *Cell* 2012, **148**:1132-1144.
- 749 48. Cairns RA, Harris IS, Mak TW: **Regulation of cancer cell metabolism.** *Nat Rev Cancer*
750 2011, **11**:85-95.
- 751 49. Carmeliet P: **Angiogenesis in life, disease and medicine.** *Nature* 2005, **438**:932-936.
- 752 50. Pitulescu ME, Schmidt I, Giaimo BD, Antoine T, Berkenfeld F, Ferrante F, Park H, Ehling M,
753 Biljes D, Rocha SF, et al: **Dll4 and Notch signalling couples sprouting angiogenesis and**
754 **artery formation.** *Nat Cell Biol* 2017, **19**:915-927.
- 755 51. Gridley T: **Notch signaling in vascular development and physiology.** *Development* 2007,
756 **134**:2709-2718.
- 757 52. Davar D, Dzutsev AK, McCulloch JA, Rodrigues RR, Chauvin JM, Morrison RM, Deblasio
758 RN, Menna C, Ding Q, Pagliano O, et al: **Fecal microbiota transplant overcomes resistance**
759 **to anti-PD-1 therapy in melanoma patients.** *Science* 2021, **371**:595-602.
- 760 53. Dzutsev A, Badger JH, Perez-Chanona E, Roy S, Salcedo R, Smith CK, Trinchieri G:
761 **Microbes and Cancer.** *Annu Rev Immunol* 2017, **35**:199-228.
- 762 54. Finlay BB, Goldszmid R, Honda K, Trinchieri G, Wargo J, Zitvogel L: **Can we harness the**
763 **microbiota to enhance the efficacy of cancer immunotherapy?** *Nat Rev Immunol* 2020,
764 **20**:522-528.
- 765 55. Kostic AD, Ojesina AI, Peadarallu CS, Jung J, Verhaak RG, Getz G, Meyerson M: **PathSeq:**
766 **software to identify or discover microbes by deep sequencing of human tissue.** *Nat*
767 *Biotechnol* 2011, **29**:393-396.
- 768 56. Lozupone C, Knight R: **UniFrac: a new phylogenetic method for comparing microbial**
769 **communities.** *Appl Environ Microbiol* 2005, **71**:8228-8235.
- 770 57. Greathouse KL, White JR, Vargas AJ, Bliskovsky VV, Beck JA, von Muhlinen N, Polley EC,
771 Bowman ED, Khan MA, Robles AI, et al: **Interaction between the microbiome and TP53 in**
772 **human lung cancer.** *Genome Biol* 2018, **19**:123.
- 773 58. Sahni N, Yi S, Taipale M, Fuxman Bass JI, Coulombe-Huntington J, Yang F, Peng J, Weile J,
774 Karras GI, Wang Y, et al: **Widespread macromolecular interaction perturbations in**
775 **human genetic disorders.** *Cell* 2015, **161**:647-660.
- 776 59. Cheng F, Zhao J, Wang Y, Lu W, Liu Z, Zhou Y, Martin WR, Wang R, Huang J, Hao T, et al:
777 **Comprehensive characterization of protein-protein interactions perturbed by disease**
778 **mutations.** *Nat Genet* 2021, **53**:342-353.
- 779 60. Wang Y, Skibbe JR, Hu C, Dong L, Ferchen K, Su R, Li C, Huang H, Weng H, Huang H, et al:
780 **ALOX5 exhibits anti-tumor and drug-sensitizing effects in MLL-rearranged leukemia.**
781 *Sci Rep* 2017, **7**:1853.
- 782 61. Tsukasaki K, Miller CW, Greenspun E, Eshaghian S, Kawabata H, Fujimoto T, Tomonaga M,
783 Sawyers C, Said JW, Koeffler HP: **Mutations in the mitotic check point gene, MAD1L1, in**
784 **human cancers.** *Oncogene* 2001, **20**:3301-3305.
- 785 62. Yang M, Li B, Liu CJ, Tomchick DR, Machius M, Rizo J, Yu H, Luo X: **Insights into mad2**
786 **regulation in the spindle checkpoint revealed by the crystal structure of the symmetric**
787 **mad2 dimer.** *PLoS Biol* 2008, **6**:e50.
- 788 63. Sironi L, Mapelli M, Knapp S, De Antoni A, Jeang KT, Musacchio A: **Crystal structure of**
789 **the tetrameric Mad1-Mad2 core complex: implications of a 'safety belt' binding**
790 **mechanism for the spindle checkpoint.** *EMBO J* 2002, **21**:2496-2506.

- 791 64. Guo Y, Zhang X, Yang M, Miao X, Shi Y, Yao J, Tan W, Sun T, Zhao D, Yu D, et al:
792 **Functional evaluation of missense variations in the human MAD1L1 and MAD2L1 genes**
793 **and their impact on susceptibility to lung cancer.** *J Med Genet* 2010, **47**:616-622.
- 794 65. Zhong R, Chen X, Chen X, Zhu B, Lou J, Li J, Shen N, Yang Y, Gong Y, Zhu Y, et al:
795 **MAD1L1 Arg558His and MAD2L1 Leu84Met interaction with smoking increase the risk**
796 **of colorectal cancer.** *Sci Rep* 2015, **5**:12202.
- 797 66. Ferrar T, Chamousset D, De Wever V, Nimick M, Andersen J, Trinkle-Mulcahy L, Moorhead
798 GB: **Taperin (c9orf75), a mutated gene in nonsyndromic deafness, encodes a vertebrate**
799 **specific, nuclear localized protein phosphatase one alpha (PP1alpha) docking protein.**
800 *Biol Open* 2012, **1**:128-139.
- 801 67. Chiu CH, Yang CT, Shih JY, Huang MS, Su WC, Lai RS, Wang CC, Hsiao SH, Lin YC, Ho
802 CL, et al: **Epidermal Growth Factor Receptor Tyrosine Kinase Inhibitor Treatment**
803 **Response in Advanced Lung Adenocarcinomas with G719X/L861Q/S768I Mutations.** *J*
804 *Thorac Oncol* 2015, **10**:793-799.
- 805 68. Kunishige M, Ichihara S, Kadota N, Okano Y, Machida H, Hatakeyama N, Naruse K,
806 Shinohara T, Takeuchi E: **Non-small cell lung cancer with EGFR (L858R and E709X) and**
807 **CNNB1 mutations responded to afatinib.** *Thorac Cancer* 2022.
- 808 69. Yang JC, Sequist LV, Geater SL, Tsai CM, Mok TS, Schuler M, Yamamoto N, Yu CJ, Ou SH,
809 Zhou C, et al: **Clinical activity of afatinib in patients with advanced non-small-cell lung**
810 **cancer harbouring uncommon EGFR mutations: a combined post-hoc analysis of LUX-**
811 **Lung 2, LUX-Lung 3, and LUX-Lung 6.** *Lancet Oncol* 2015, **16**:830-838.
- 812 70. Massarelli E, Johnson FM, Erickson HS, Wistuba, II, Papadimitrakopoulou V: **Uncommon**
813 **epidermal growth factor receptor mutations in non-small cell lung cancer and their**
814 **mechanisms of EGFR tyrosine kinase inhibitors sensitivity and resistance.** *Lung Cancer*
815 2013, **80**:235-241.
- 816 71. Han SW, Kim TY, Hwang PG, Jeong S, Kim J, Choi IS, Oh DY, Kim JH, Kim DW, Chung
817 DH, et al: **Predictive and prognostic impact of epidermal growth factor receptor**
818 **mutation in non-small-cell lung cancer patients treated with gefitinib.** *J Clin Oncol* 2005,
819 **23**:2493-2501.
- 820 72. Floc'h N, Lim S, Bickerton S, Ahmed A, Orme J, Urosevic J, Martin MJ, Cross DAE, Cho
821 BC, Smith PD: **Osimertinib, an Irreversible Next-Generation EGFR Tyrosine Kinase**
822 **Inhibitor, Exerts Antitumor Activity in Various Preclinical NSCLC Models Harboring**
823 **the Uncommon EGFR Mutations G719X or L861Q or S768I.** *Mol Cancer Ther* 2020,
824 **19**:2298-2307.
- 825 73. Yang Z, Yang N, Ou Q, Xiang Y, Jiang T, Wu X, Bao H, Tong X, Wang X, Shao YW, et al:
826 **Investigating Novel Resistance Mechanisms to Third-Generation EGFR Tyrosine Kinase**
827 **Inhibitor Osimertinib in Non-Small Cell Lung Cancer Patients.** *Clin Cancer Res* 2018,
828 **24**:3097-3107.
- 829 74. Wong JYY, Downward GS, Hu W, Portengen L, Seow WJ, Silverman DT, Bassig BA, Zhang
830 J, Xu J, Ji BT, et al: **Lung cancer risk by geologic coal deposits: A case-control study of**
831 **female never-smokers from Xuanwei and Fuyuan, China.** *Int J Cancer* 2019, **144**:2918-
832 2927.
- 833 75. Vermeulen R, Downward GS, Zhang J, Hu W, Portengen L, Bassig BA, Hammond SK, Wong
834 JYY, Li J, Reiss B, et al: **Constituents of Household Air Pollution and Risk of Lung**

- 835 **Cancer among Never-Smoking Women in Xuanwei and Fuyuan, China.** *Environ Health*
836 *Perspect* 2019, **127**:97001.
- 837 76. Kim C, Chapman RS, Hu W, He X, Hosgood HD, Liu LZ, Lai H, Chen W, Silverman DT,
838 Vermeulen R, et al: **Smoky coal, tobacco smoking, and lung cancer risk in Xuanwei,**
839 **China.** *Lung Cancer* 2014, **84**:31-35.
- 840 77. Charoentong P, Finotello F, Angelova M, Mayer C, Efremova M, Rieder D, Hackl H,
841 Trajanoski Z: **Pan-cancer Immunogenomic Analyses Reveal Genotype-Immunophenotype**
842 **Relationships and Predictors of Response to Checkpoint Blockade.** *Cell Rep* 2017,
843 **18**:248-262.
- 844
- 845

846

Graphical Abstract



847

848

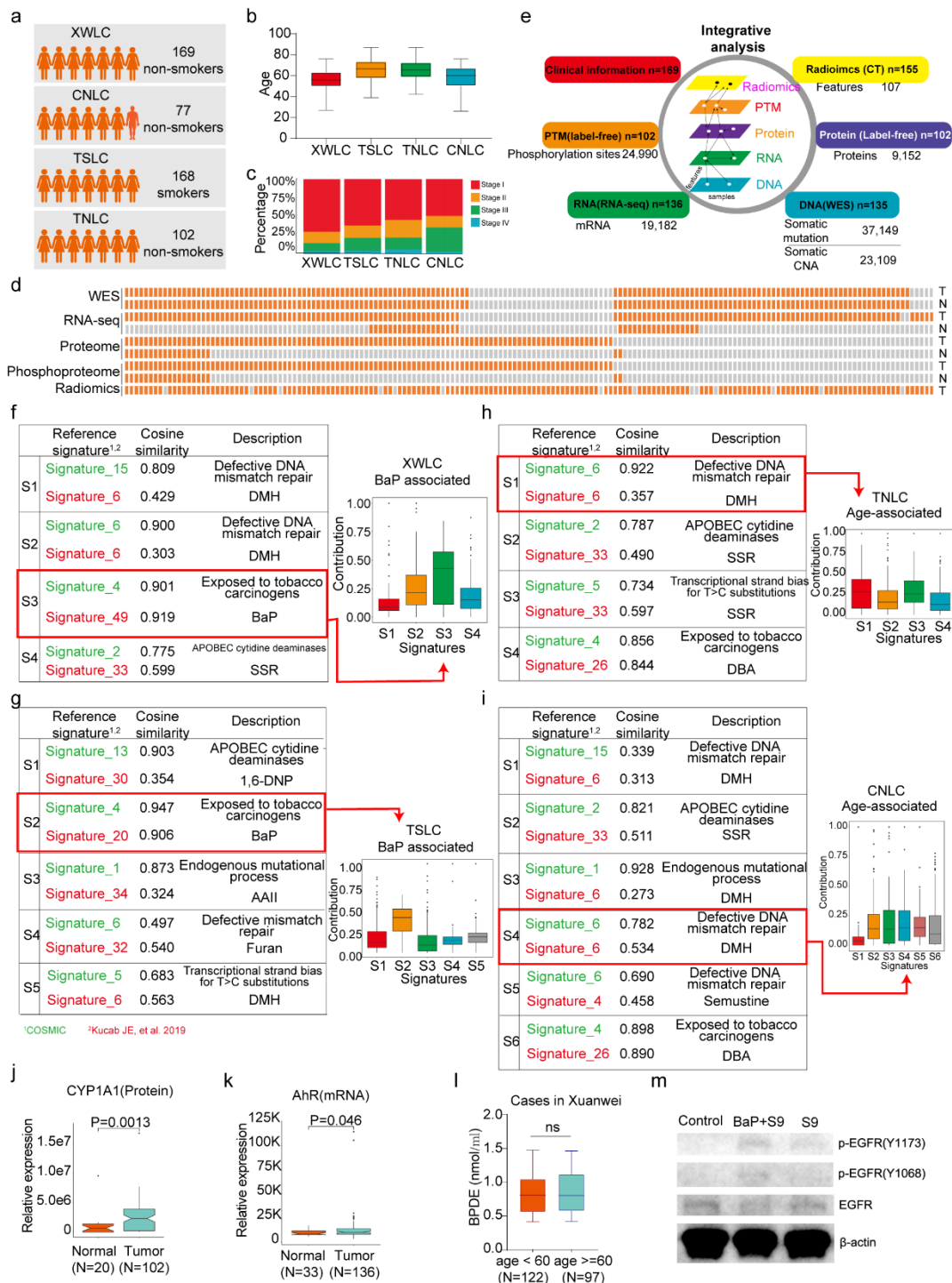
849 Highlights:

- 850 ● We conducted comprehensive multi-omic profiling of air-pollution-
- 851 associated LUAD.
- 852 ● Our study revealed the significant roles of the air pollutant BaP and
- 853 its induced hot mutation G719X in lung cancer progression.
- 854 ● Multi-omic clustering enabled the identification of personalized
- 855 therapeutic strategies.
- 856 ● Through mutation-informed interface analysis, we identified novel
- 857 targets for therapeutic intervention.

858

859

860 Figure legends



861

862 Fig 1. | Proteogenomic profiling and mutational signatures in XWLC

863 a. Four cohort datasets used in this study: XWLC (Lung adenocarcinoma

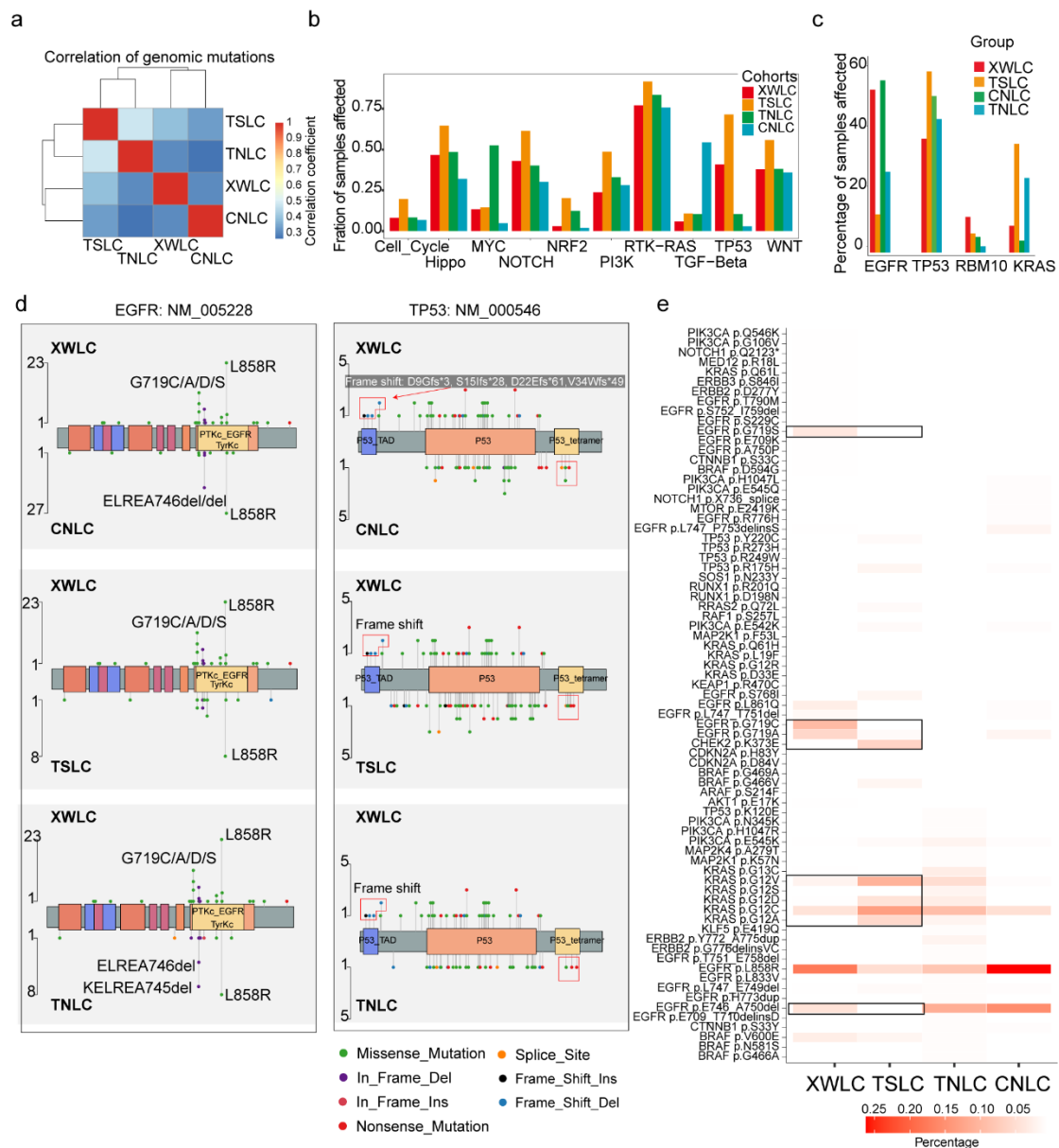
864 adenocarcinoma from non-smoking females in Xuanwei area), CNLC (subset of lung

865 adenocarcinoma from non-smoking patients in Chinese Human Proteome

866 Project), TSLC (subset of lung adenocarcinoma from smoking females in

867 TCGA-LUAD project), TNLC (subset of lung adenocarcinoma from non-
868 smoking females in TCGA-LUAD project); b. Age distribution of
869 patients at the time of operation in the four cohorts; c. Distribution of
870 tumor stages across the cohorts; d. Data availability for the XWLC
871 datasets. Each bar represents a sample, with orange bars indicating data
872 availability and gray bars indicating data unavailability; e. Summary of
873 data generated from the XWLC cohort; f-i. Mutational signatures
874 identified in XWLC (f), TSLC (g), TNLC (h), and CNLC (i) cohorts.
875 Cosine similarity analysis of the signatures compared to well-established
876 COMIC signatures (in green) and Kucab et al. signatures (in red).
877 Contribution of signatures in each cohort provided on the right; j. Protein
878 abundance of CYP1A1 in tumor and normal samples within the XWLC
879 cohort; k. Expression levels of the AhR gene in tumor and normal
880 samples within the XWLC cohort; l. Comparison of serum BPDE content
881 in individuals from the Xuanwei area, categorized by young cases and
882 older cases. m. western blotting of EGFR-Y1173, EGFR-Y1068 and
883 EGFR abundance in hiPSC (Control), BaP+S9 treated hiPSC (BaP+S9),
884 and S9 treated hiPSC (S9). Liver S9 is a variety of biological sources
885 represent the post-mitochondrial supernatant fraction from homogenized
886 liver and is known to be a rich source of drug metabolizing enzymes
887 including P-450. Two-tailed Wilcoxon rank sum test used to calculate p-
888 values in j-l.

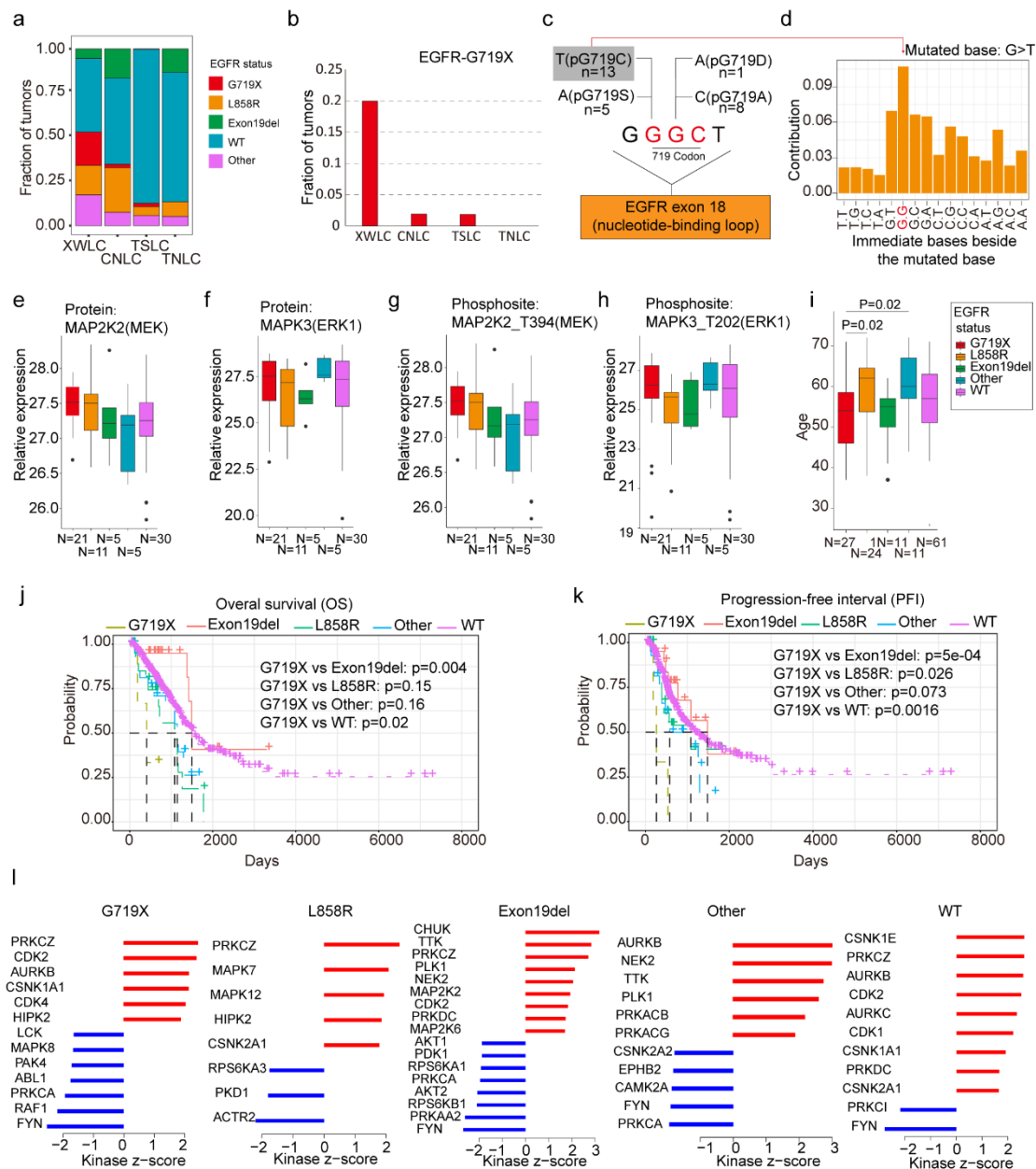
889



890

891 **Fig 2. | Genomic and genetic features in XWLC cohort**

892 a. Correlation of genomic mutations among cohorts, determined using the
 893 Pearson correlation coefficient; b. Comparison of oncogenic pathways
 894 affected by mutations in each cohort; c. Comparison of mutation
 895 frequency of four key genes across cohorts; d. Lollipop plot illustrating
 896 differences in mutational sites within EGFR (left) and TP53 (right) across
 897 XWLC/CNLC, XWLC/TSLC, and XWLC/TNLC pairs; e. Analysis of
 898 the percentage of samples with actionable alterations, with a focus on
 899 significant variations between XWLC and TSLC cohorts, highlighted by
 900 black boxes.



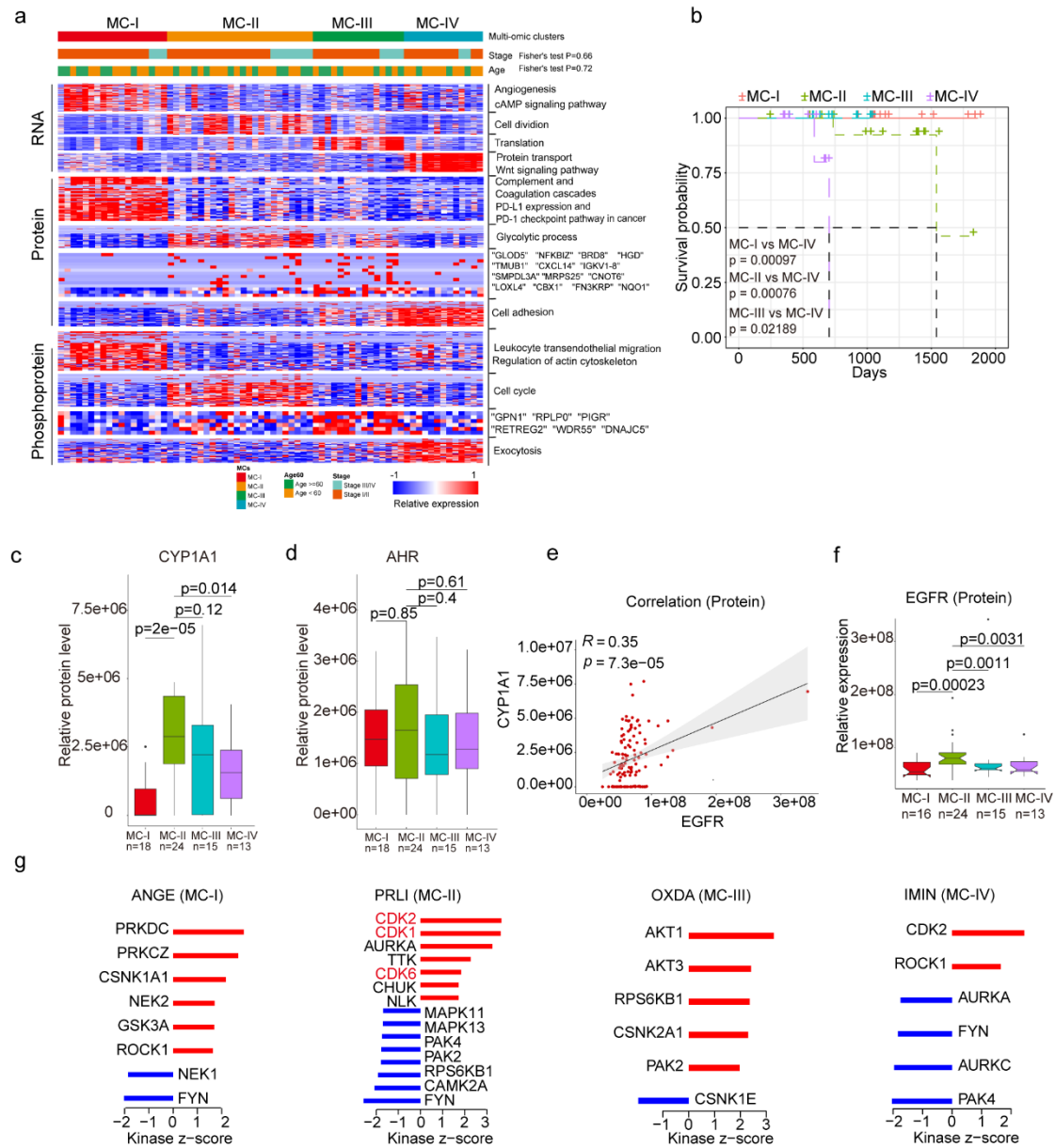
902

903 **Fig 3. | EGFR-G719X in the XWLC cohort**

904 a. Distribution of different EGFR mutation statuses across the four
 905 cohorts; b. Comparison of the fraction of G719X mutations across the
 906 four cohorts; c. Detailed information on pG719X (pG719/A/D/C/S)
 907 mutations in the XWLC cohort. The number of each mutation type is
 908 labeled; d. Distribution of nucleotide pairs surrounding the most common
 909 G>T transversion site in the XWLC cohort. The x-axis represents the
 910 immediate bases surrounding the mutated base. For BaP, the tallest G>T
 911 peak occurs at GpGpG; e-h. Comparison of activation levels of key

912 components in the MAPK pathway across different EGFR mutation
913 statuses in the XWLC cohort; i. Comparison of patient ages across
914 different EGFR mutation statuses in the XWLC cohort; j-k. Presentation
915 of overall survival (OS, j) and progression-free interval (PFI, k) analysis
916 across different EGFR mutation statuses in the TCGA-LUAD cohort; l.
917 Evaluation of kinase activities by KSEA in tumors across different EGFR
918 mutation statuses in the XWLC cohort.

919



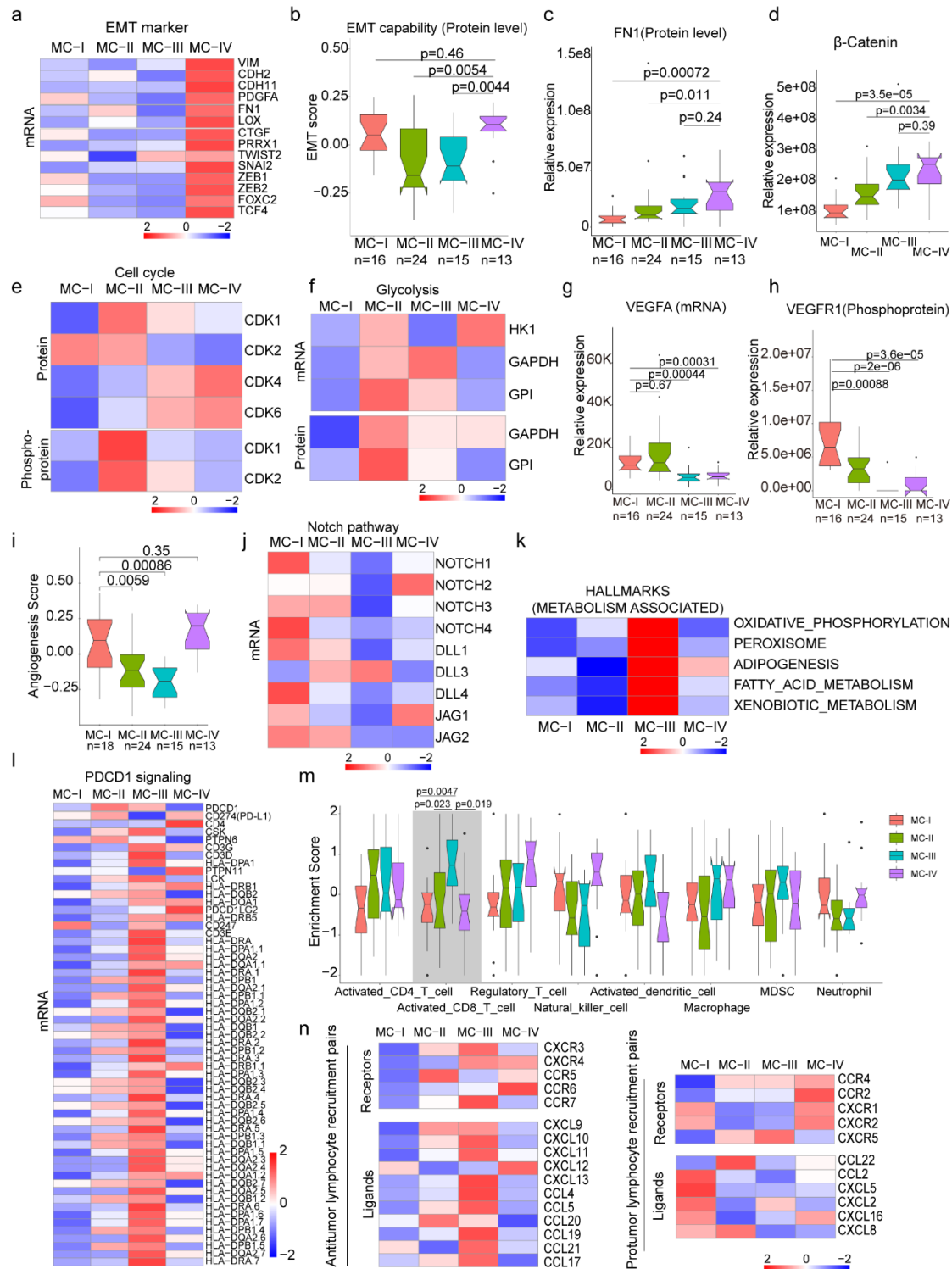
920

921 Fig 4. Subtyping of XWLC

922 a. Integrative classification of tumor samples into four
 923 ConsensusClusterPlus-derived clusters (MC-I to MC-IV). The heatmap
 924 displays the top 50 features, including mRNA transcripts, proteins, and
 925 phosphoproteins, for each multi-omic cluster. The features are annotated
 926 with representative pathways or genes. If a cluster has fewer than 50
 927 features, all features are shown. If no significant GO biological processes
 928 are associated with cluster features, all features are displayed; b.
 929 Comparison of overall survival between MC-IV and the other three
 930 subtypes; c-d. Protein abundance comparison of CYP1A1 (c) and AHR

931 (d) across subtypes; e. Protein-level correlation between CYCP1A1 and
932 EGFR; f. Protein-level comparison of EGFR across subtypes; g.
933 Evaluation of kinase activities by KSEA in tumors across subtypes in the
934 XWLC cohort.

935



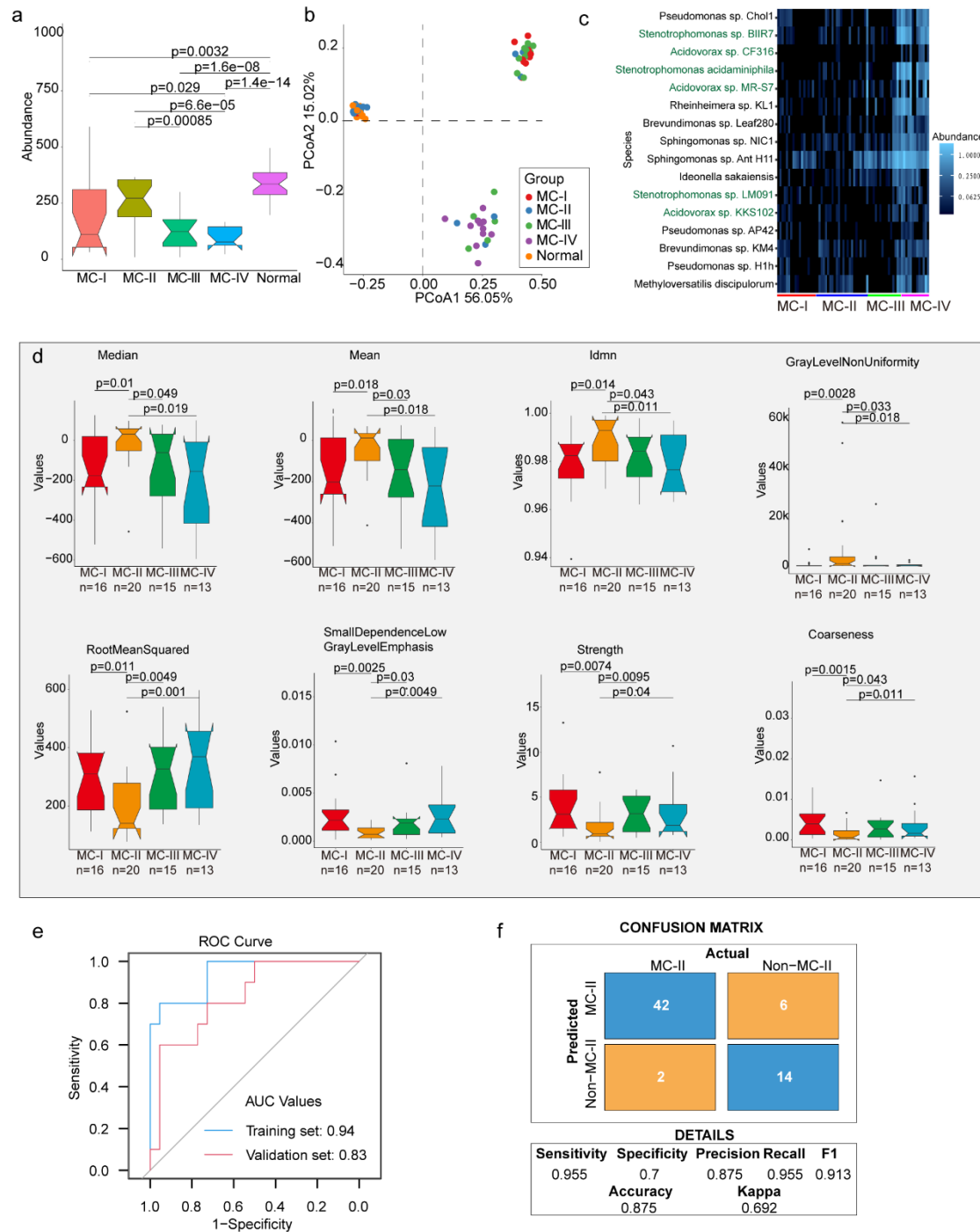
936

937 **Fig5. | Biological and immune features across MC subtypes**

938 a. Relative expression of epithelial-mesenchymal transition (EMT)
 939 markers across subtypes; b. EMT scores across subtypes using a gene set
 940 derived from MsigDB (M5930); c-d. Protein abundance comparison of
 941 FN1 (c) and β -Catenin (d) across subtypes; e. Protein and phosphoprotein

942 levels of key cell cycle kinases across subtypes; f. Expression of mRNA
943 and protein levels of glycolysis-associated enzymes; g. mRNA expression
944 of VEGFA across subtypes; h. Phosphoprotein abundance of VEGFR1
945 across subtypes; i. Angiogenesis score across subtypes using a gene set
946 derived from MsigDB (Systematic name M5944); j. Expression
947 comparison of key regulators of the Notch pathway across subtypes; k.
948 Metabolism-associated hallmarks across subtypes. Gene sets for oxidative
949 phosphorylation, peroxisome, adipogenesis, fatty acid metabolism, and
950 xenobiotic metabolism were derived from MsigDB hallmark gene sets; l.
951 Expression of PD-1 signaling-associated genes across subtypes. PD-1
952 signaling-associated genes were derived from MsigDB (Systematic name
953 M18810) ; m. Immune cell infiltration across subtypes. Gene sets for
954 each immune cell type were derived from a previous study[77]; n.
955 Expression of anti-tumor/pro-tumor lymphocyte receptors and ligands
956 across subtypes. The two-tailed Wilcoxon rank sum test was used to
957 calculate p-values in panels b, c, d, g, h, i, and m.

958



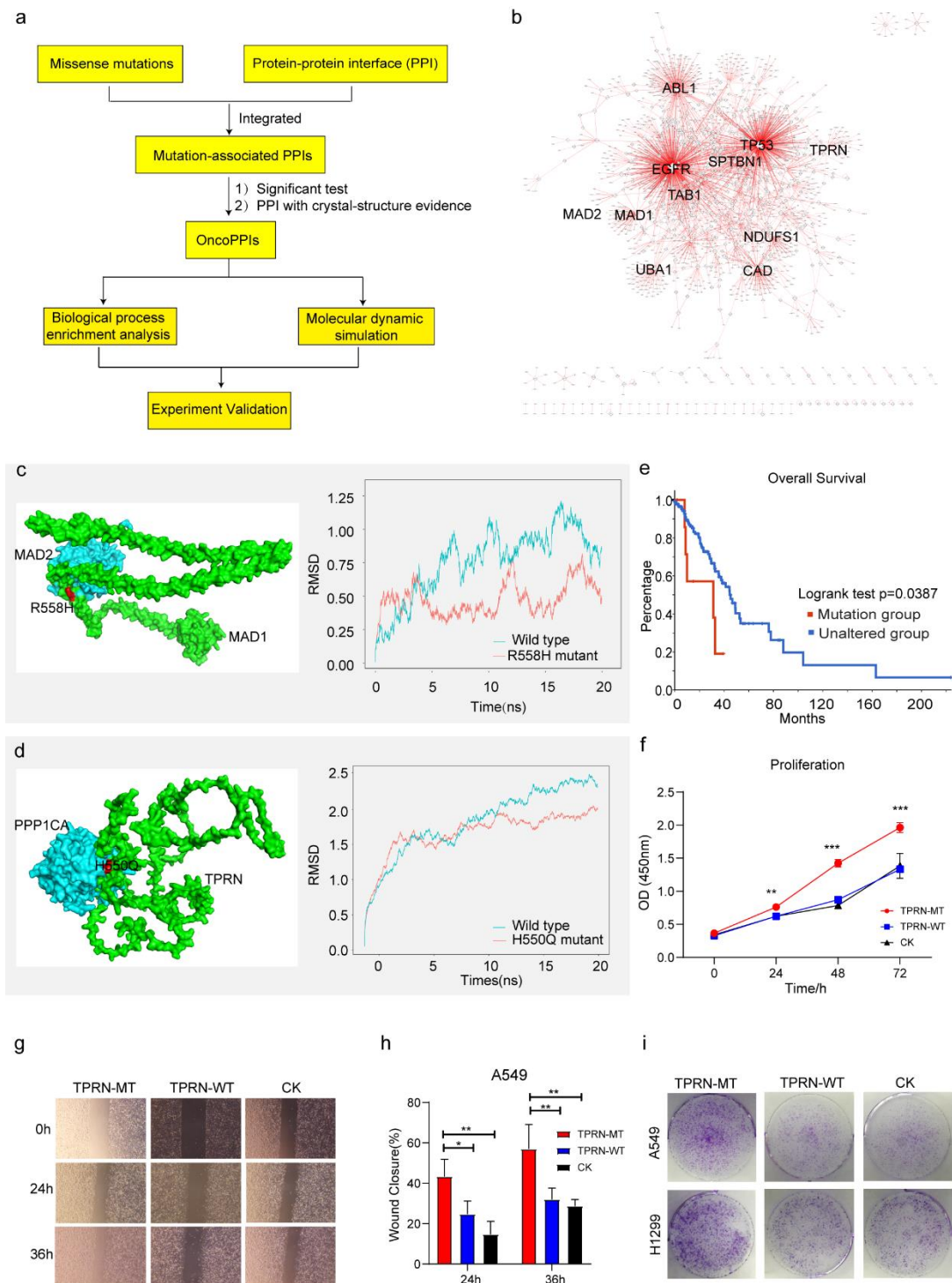
959

960 **Fig 6. | Microbiota composition and radiomic features across**
 961 **subtypes.**

962 a. Abundance of bacterial sequencing reads per million (RPM) across
 963 tumor molecular subtypes and adjacent normal tissues. Significance was
 964 determined using a paired two-sided Wilcoxon rank sum test; b. Principal
 965 coordinates analysis (PCoA) plot of mRNA data for XWLC samples
 966 shows significant variation between tumor and adjacent normal samples

967 along the first axis of variation, and variation between MC-IV subtypes
968 and the other three tumor subtypes along the second axis; c. Heatmap of
969 bacterial species' abundance highly expressed in the MC-IV subtype
970 compared to the other three subtypes in the XWLC cohort. Bacterial
971 species labeled green prefer an acidic living environment; d. Eight
972 features showing significant differences between MC-II and the other
973 three subtypes. The Wilcoxon rank sum test was used to calculate the p-
974 values; e. A receiver operating characteristic (ROC) curve was used to
975 evaluate the performance of the radiomic signature in distinguishing MC-
976 II from the other three subtypes; f. Confusion matrix allows visualization
977 of the performance of the algorithm in separating MC-II from other
978 subtypes.

979



980

981 **Fig.7|Identification of novel targets in XWLC**

982 a. Flow chart showing the integration of mutation-informed PPI analysis,
 983 molecular dynamic simulation and experiment validation to identify
 984 novel targets; b. Network visualization of XWLC_oncoPPIs. Edge
 985 thickness represents the number of missense mutations at the protein-
 986 protein interaction (PPI) interface, while node size indicates connectivity;

987 c. MAD1-MAD2 interaction model and the p.Arg558His mutation at the
988 interface (left). The complex model was generated using Zdock protein
989 docking simulation. The right distribution showing root-mean-squared
990 deviation (RMSD) during a 20 ns molecular dynamics simulation of
991 MAD1 wild type vs. MAD1 p.Arg558His in the complex; d. Model
992 showing the p.His550Gln alteration within the TPRN-PPP1CA complex
993 (left). The right distribution showing root-mean-squared deviation
994 (RMSD) during a 20 ns molecular dynamics simulation for TPRN wild
995 type vs. TPRN p.His550Gln (H550Q) in the complex; e. Survival analysis
996 of TPRN mutation group and unaltered group derived from cbiportal
997 using TCGA-LUAD cohort (<https://www.cbiportal.org/>); f. CCK8 assay
998 for TPRN-MT, TPRN-WT, and CK cell lines in A549 cells which was
999 transfected by mutant TPRN, wild-type and empty vector, respectively. g.
1000 Transwell assay for TPRN-MT, TPRN-WT, and CK after 24h and 36h in
1001 A549 cells. Magnification was set to 40x; h.Bar chart showing the
1002 statistical results of transwell assay; i. Cell colon assay for TPRN-MT,
1003 TPRN-WT and CK in A549 and H1299 cell line. The two-tailed
1004 Wilcoxon rank sum test was used to calculate p-values in f and h.
1005 *.p<0.05; **,p <0.01; ***, p<0.001;