**The mutational features of aristolochic acid-induced mouse and human liver cancers**

Zhao-Ning Lu[1#], Qing Luo[1#], Li-Nan Zhao[1], Yi Shi[1], Xian-Bin Su[1], Ze-Guang Han[1*]

[1]Key Laboratory of Systems Biomedicine (Ministry of Education), Shanghai Centre for Systems Biomedicine, Shanghai Jiao Tong University, Shanghai, 200240, China

[#]These authors contribute equally to this work

[*]To whom correspondence should be addressed.

Ze-Guang Han, Key Laboratory of Systems Biomedicine (Ministry of Education), Shanghai Center of Systems Biomedicine, Shanghai Jiao Tong University, 800 Dongchuan Road, Shanghai 200240, China. Email: hanzg@sjtu.edu.cn

## Abstract

Aristolochic acid (AA) derived from traditional Chinese herbal remedies has recently been statistically associated with human liver cancer; however, the causal relationships between AA and liver cancer and the underlying evolutionary process of AA-mediated mutagenesis during tumorigenesis are obscure. Here, we subjected mice, including *Pten*-deficient ones, to aristolochic acid I (AAI) alone or a combination of AAI and carbon tetrachloride ($CCl_4$), which may induce liver injury. Significantly, AAI promoted the development of liver cancer, including hepatocellular carcinoma and intrahepatic cholangiocarcinoma, in a dose-dependent manner, and it increased the incidence of liver cancer, together with $CCl_4$ or *Pten* deficiency. AAI could lead to DNA damage and AAI-DNA adducts that initiate liver cancer via characteristic A>T transversions, as indicated by the comprehensive genomic analysis, which revealed recurrent mutations in *Hras* and some genes encoding components of the Ras/Raf, PI3K, Notch, Hippo, Wnt, DNA polymerase family and the SWI/SNF complex, some of which are also often found in human liver cancer. Mutational signature analysis across human cancer types revealed that the AA-related dominant signature was especially implicated in liver cancer in China, based on very stringent criteria derived from the animal cancer form, in which mutations of *TP53* and *JAK1* are prone to be significantly enriched. Interestingly, AAI-mediated characteristic A>T mutations were

36   the earliest genetic event driving malignant subclonal evolution in mouse

37   and human liver cancer. In general, this study provides documented

38   evidence for AA-induced liver cancer with featured mutational processes

39   during malignant clonal evolution, laying a solid foundation for the

40   prevention and diagnosis of AA-associated human cancers, especially liver

41   cancer.

## Introduction

Aristolochic acid (AA) is present in plants in the genera *Aristolochia, Bragantia*, *Asarum* and others[1], which have been widely used in traditional Chinese herbal remedies. AA is one of the most potent carcinogens known to man, belonging to the Group I human carcinogens classified the by International Agency for Research on Cancer (IARC). Aristolochic acid I (AAI) and II (AAII) are the major components of the AA mixture contained in the plant extract of *Aristolochia* species[2]. AA is a genotoxic carcinogen because its metabolite can bind purines to form AA-DNA adducts, aristolactam (AL)-DNA adducts (dA-AL and dG-AL), which are specific markers of exposure to aristolochic acids and induce DNA mutations with characteristic adenine-to-thymine (A>T) transversions *in vitro* and *in vivo*[2,3].

The dA-AL-I (7-(deoxyadenosin-N6-yl) aristolactam I) adducts induced by AAI can show long-term persistence in renal tissue[4], which may have led to the occurrence of aristolochic acid nephropathy (AAN) in Belgian women who had taken weight-reducing pills containing *Aristolochia fangchi*[5] and Balkan endemic nephropathy (BEN) through dietary contamination with *Aristolochia clematitis* seeds[6]. Both nephropathies are associated with urothelial carcinoma because AA-DNA adducts have been found in kidney tissue and urothelial tumor tissues of patients with AAN or BEN[3,5,6]. In Taiwan, approximately one-third of the

64   people consume Chinese herbal remedies containing AA, which could be

65   associated with the highest incidence of upper urinary tract cancers (UTUC)

66   in the world[7]. The genome-wide mutational signature of characteristic A>T

67   transversions (COSMIC signature 22), specifically reflecting AA-

68   implicated mutagenesis, is frequently found in Taiwanese UTUC[7,8].

69       Recently, AA has been statistically associated with human liver cancer.

70   We found, for the first time, that the characteristic A:T to T:A transversions

71   were significantly enriched in 4 of 10 (40%) hepatitis B virus (HBV)-

72   associated hepatocellular carcinoma (HCC) specimens from China,

73   indicative of AA exposure in HCC tumorigenesis[9]. A survey in larger

74   cohorts of HCC patients indicated that the AA-implicated mutational

75   signature was discovered in some Asian HCC patients, especially in more

76   than 75% of Taiwanese HCC cases[10].

77       There has been no direct evidence that AA can induce liver cancer until

78   now, although AA-DNA adducts have been detected in many organs,

79   including the liver, in experimental animals exposed to AA during a

80   relatively short peroid[11], and even in the livers of some nephropathy

81   patients with known AA exposure[12,13]. To confirm whether AA can directly

82   induce liver cancer, including HCC, here we subjected mice, including

83   *Pten*-deficient ones, to AAI alone or a combination of AAI and carbon

84   tetrachloride ($CCl_4$), a well-documented liver injury agent. Significantly,

85   AAI administration alone increased the incidence of liver cancer in a dose-

86    dependent manner, and the combination of AAI and CCl$_4$ also led to a

87    higher incidence of mouse liver cancer. Interestingly, the types of liver

88    cancer included HCC, intrahepatic cholangiocarcinoma (ICC), and

89    combined hepatocellular and intrahepatic cholangiocarcinoma (cHCC-

90    ICC). Genome-wide analysis of the AAI-induced liver cancer showed the

91    characteristic mutational signature and process during clonal evolution,

92    providing new insights into the pathogenesis of AA-induced liver cancer.

93    **Results**

94    **AAI can induce mouse liver cancer**

95    To validate whether AA could directly induce liver cancer, especially HCC,

96    we first subjected C57BL/6 male mice to AAI administration alone by

97    intraperitoneal injection. Based on previous research[11,14], AAI was

98    administered at a lower dose (2.5 and 5 mg/kg body weight) for injection

99    of 3, 7, and 14 times, respectively (see **Methods**). Moreover, considering

100   the possibility that AAI could enhance tumorigenesis due to liver injury or

101   a genetic defect, we designed the combination of AAI and CCl$_4$ to treat

102   mice, in which CCl$_4$ can induce liver injury, compensatory proliferation,

103   inflammation, and fibrosis[15]; we also subjected liver-specific *Pten*-

104   deficient mice, who frequently develop liver cancer by 74–78 weeks of

105   age[16], to AAI administration. In general, a total of eight experimental

106   groups of mice subjected to AAI administration (**Fig. 1a** and

107   **Supplementary Fig. 1a**) included the following: (I) "AAI (3x)",

108    administration of AAI at a dose of 2.5 mg/kg every other day for 3 doses

109    at 2 weeks after birth; (II) "AAI (14x)", administration of AAI at a dose of

110    2.5 mg/kg/day for 14 days at 2 weeks of age; (III) "AAI (high 3x)",

111    administration of AAI at a dose of 5 mg/kg every other day for 3 doses at

112    2 weeks of age; (IV) "AAI (high 14x)", administration of AAI at a dose of

113    5 mg/kg/day for 14 days at 2 weeks of age; (V) "AAI (3x) + $CCl_4$",

114    administration of $CCl_4$ three times per week for 4 weeks at 2 months after

115    AAI injection; (VI) "AAI (14x) + $CCl_4$", administration of $CCl_4$ once per

116    week for 10 weeks at 4 weeks after AAI injection; (VII) "AAI (7x)",

117    administration of AAI at a dose of 2.5 mg/kg every other day for 7 doses

118    at 1 week of age, to observe the effect of AAI on younger fetal livers; (VIII)

119    "AAI (high 14x, $Pten^{LKO}$)", administration of AAI at a dose of 5 mg/kg/day

120    for 14 days at 2 weeks of age in liver-specific $Pten$-deficient mice. In

121    addition, male mice were injected with $CCl_4$ alone, which was administered

122    once or three times per week at a dose of 0.5 ml/kg body weight, and

123    vehicle as the control group (**Supplementary Fig. 1a**).

124        Interestingly, liver cancer occurred in all eight experimental groups of

125    AAI administration (**Fig. 1b-g** and **Supplementary Fig. 1b-p**). AAI

126    administration alone significantly promoted the development of liver

127    cancer in a dose-dependent manner. A lower dosage of AAI ("AAI (3x)")

128    led to liver cancer development in 2 (20%) out of 10 mice at 11.5 months

129    after the first AAI administration, demonstrating a statistically increased

130  incidence compared with the control group without AAI treatment ($P =$

131  0.038) (**Supplementary Fig. 1b**). Significantly, a greater number of

132  injections ("AAI (14x)" *vs.* "AAI (3x)") was associated with earlier

133  occurrence (8.5 M *vs.* 11.5 M) and larger tumor sizes (11.5 M, mean: 5.62

134  mm *vs.* 0.22 mm, $P = 0.048$) of liver cancer (**Fig. 1d**). Under the same AAI

135  administration durations, the larger the dosage ("AAI (high 3x)" *vs.* "AAI

136  (3x)"), the larger was the tumor size (11.5 M, mean: 4.63 mm *vs.* 0.22 mm,

137  $P = 0.048$) (**Fig. 1d**). However, many mice in the "AAI (high 14x)" group

138  died during the experimental observation period (**Supplementary Fig. 1q**),

139  and all 4 surviving mice at 8.5 or 11.5 months after the first AAI

140  administration developed liver cancer (**Supplementary Fig. 1j**). In

141  addition, the mice in the "AAI (7x)" injection group at the age of 1 week

142  also displayed an increment in tumor incidence, number and size,

143  compared to the "AAI (14x)" group, although this difference was not

144  statistically significant (**Supplementary Fig. 1b-d, o**).

145      Compared with AAI administration alone, the combination of both AAI

146  and $CCl_4$ led to a significantly higher incidence of mouse liver cancer. The

147  combined models ("AAI (3x) + $CCl_4$" *vs.* "AAI (3x)"; "AAI (14x) + $CCl_4$"

148  *vs.* "AAI (14x)") resulted in an earlier tumor occurrence (8.5 M *vs.* 11.5 M;

149  5.5 M *vs.* 8.5 M), higher tumor incidence (11.5 M, 100% *vs.* 20%, $P =$

150  0.007; 8.5 M, 100% *vs.* 71.4%, $P = 0.2$), greater number of tumor nodules

151  (11.5 M, mean: 4.3 *vs.* 0.2, $P = 0.00048$; 8.5 M, mean: 3.75 *vs.* 0.86, $P =$

152  0.0013) and larger tumor sizes (11.5 M, mean: 6.53 mm *vs*. 0.22 mm, *P* =

153  0.00051; 8.5 M, mean: 8.25 *vs*. 2.86, *P* = 0.014) (**Fig. 1b-d**).

154  Moreover, compared with the same genetic background mice as a

155  control, the liver-specific *Pten*-deficient mice (*Pten^{LKO}*) treated with AAI

156  alone developed liver cancer (6 M, 100% *vs.* 0%, *P* = 0.002)

157  (**Supplementary Fig. 1b-d**), along with obvious bile duct hyperplasia in

158  adjacent liver tissues (**Supplementary Fig. 1p**).

159  Among the examined 84 livers from the above mice that received AAI

160  administration, 60 mice (71.4%) developed liver cancer. We checked these

161  tumors based on the microscopic morphology and immunohistochemistry

162  staining for Ki67, a proliferative index; α-fetoprotein (AFP), a well-known

163  HCC marker; and cytokeratin 19 (CK19), a cholangiocyte marker.

164  Interestingly, 55 (91.7%) of 60 tumors were observed to be HCCs, which

165  exhibited expansive growth, hyperchromatic and enlarged nuclei, an

166  increased nuclear-to-cytoplasmic ratio, a high Ki67 proliferative index, an

167  absence of normal liver architecture, and focal expression of AFP (**Fig. 1e,**

168  **g** and **Supplementary Fig. 1b-p**). Additionally, 4 (6.7%) of 60 tumors

169  were classified as combined HCC and intrahepatic cholangiocarcinoma

170  (cHCC-ICC) because both AFP and CK19-positive cells were present in

171  the same tumors (**Fig. 1f, g**), which were obtained from different groups

172  treated with AAI alone, a combination of AAI and CCl$_4$, and liver-specific

173  *Pten*-deficient mice, respectively (**Fig. 1g and Supplementary Fig. 1h, i,**

174   **n, p**). Interestingly, 1 (1.7%) *Pten^{LKO}* mouse developed ICC with CK19-

175   positive cells in an examined tumor nodule (**Fig. 1f, g**); however, the HCC

176   nodule was also observed in the same liver (**Supplementary Fig. 1p**). It

177   should be pointed out that no visible liver tumors were detected in any of

178   the control groups.

179       In addition to liver cancer, AAI also promoted liver fibrosis in a dose-

180   dependent manner. The minimum amount of AAI administration ("AAI

181   (3x)" group) also led to fibrillar collagen deposition, as detected by Sirius

182   red staining in noncancer livers, compared to the controls without AAI

183   treatment (Sirius red area: 0.79% *vs.* 0.42%, $P = 9.9 \times 10^{-10}$) (**Fig. 1g** and

184   **Supplementary Fig. 1r**). The mice in the "AAI (14x)" and "AAI (high 3x)"

185   groups displayed a profound increment of fibrosis compared with those in

186   the "AAI (3x)" group (Sirius red area: 1.81% *vs.* 0.79%, $P = 3.2 \times 10^{-15}$;

187   1.12% *vs.* 0.79%, $P = 1.44 \times 10^{-4}$, respectively) (**Fig. 1h** and

188   **Supplementary Fig. 1r**). Interestingly, the degree of liver fibrosis

189   paralleled the incidence of liver cancer with AAI administration alone.

190       Moreover, we also checked the other organs of these mice treated with

191   AAI. Hydronephrosis or renal cysts were found in the mice that were

192   administered AAI (**Supplementary Fig. 1s**). However, no visible tumors

193   were found in other organs, such as lung, spleen, stomach, ureter, bladder

194   and testis.

195       The collective data indicated that AAI could result in liver cancer,

196    including HCC, ICC and cHCC-ICC, in a dose-dependent fashion; when

197    the liver was injured or displayed a genetic defect such as *Pten* deficiency,

198    AAI could synergistically promote liver cancer tumorigenesis. The above

199    data also implied that AAI could trigger genetic lesions in liver progenitor

200    cells with bipotent potential towards hepatocytes or cholangiocytes, which

201    further develop three subtypes of liver cancer under the genetic

202    differentiation program.

203    **AAI causes DNA damage and dA-AL-I adducts in liver**

204    It is known that AA is a genotoxic agent that can form DNA adducts such

205    as dA-AL and dG-AL; however, whether AA causes genomic DNA

206    damage in liver cells is unclear. We first examined the phosphorylated

207    histone γ-H2AX, a biomarker for DNA double-strand breaks, in the livers

208    of mice with AAI administration alone ("AAI (14x)"), through

209    immunofluorescence staining. Interestingly, the phosphorylated γ-H2AX

210    level was obviously increased in the liver at 1 month of age (four days after

211    completion of AAI administration) (**Fig. 2a**); however, during the

212    subsequent 2-12 months, the γ-H2AX level was reduced in liver

213    (**Supplementary Fig. 2a**). Excluding the phosphorylated γ-H2AX level,

214    we further evaluated the p53 level and its downstream target molecule Bax

215    as a cellular response to DNA damage, in these livers via Western blotting,

216    which revealed that both γ-H2AX and p53 levels were markedly increased

217    in these livers (**Fig. 2b**), along with a slight upregulation of Bax. These

218    data suggested that AAI could give rise to DNA damage.

219        To confirm whether AAI could indeed lead to DNA damage in liver,

220    we employed the alkaline comet assay to directly detect DNA strand breaks

221    in livers 3 hours after AAI administration with 2.5 mg and 5 mg/kg dosages,

222    respectively. The data demonstrated that AAI could cause DNA strand

223    breaks in mouse livers in a dose-dependent fashion (**Fig. 2c, d**), in parallel

224    with the increased phosphorylated γ-H2AX level via immunofluorescence

225    staining (**Fig. 2e, f**), which was positively correlated with the incidence of

226    liver cancer. Except for the increased γ-H2AX level in an exposure time-

227    depend manner (**Supplementary Fig. 2b, c**), phosphorylated ATR, a

228    molecule that responds to DNA damage, was upregulated in liver at 12 h

229    after AAI administration (**Supplementary Fig. 2d**). These data revealed

230    that AAI directly triggered DNA damage in mouse livers.

231        AA is known to form AA-DNA adducts that are further processed to

232    form somatic mutations through infidelity DNA repair system, which could

233    be critical step in tumorigenesis. We thus examined the AAI-mediated dA-

234    AL-I adduct in mouse livers after AAI administration by mass

235    spectrometry, using the identified synthetic dA-AL-I as a reference

236    (**Supplementary Fig. 2e, f**). Significantly, we could detect dA-AL-I

237    adducts in all examined noncancerous livers from "AAI (14x)" mice (**Fig.**

238    **2g**), and the quantity of the adduct in these livers gradually decreased along

239    the different time points after AAI administration (**Fig. 2g, h**), while the

240   quantity of dA-AL-I in kidneys was generally higher than in livers from

241   the same mice (**Supplementary Fig. 2g**). However, we could not detect

242   the dA-AL-I adducts in three matched liver cancers (**Fig. 2h**), implying

243   that, within these tumor cells, the activated DNA repair system had

244   removed the adduct or the adduct could be diluted by repeated DNA

245   replication via cell cycle progression.

246   　The above data indicated that AAI caused DNA damage, including

247   DNA double-strand breaks, and dA-AL-I adducts in liver cells, which

248   triggered the cellular response and DNA repair system. This process could

249   further lead to genomic instability and somatic mutations that contribute to

250   tumorigenesis.

251   **AAI leads to the characteristic mutational signature of A to T**

252   **transversions**

253   To survey the genomic instability and somatic mutations triggered by AAI,

254   we performed whole-genome sequencing (WGS) for DNA copy number

255   variations (CNVs), whole-exome sequencing (WES) and transcriptome

256   analysis of 11 AAI-induced liver tumor nodules, three matched adjacent

257   noncancerous livers, three livers prior to the occurrence of tumors (from

258   the "AAI (3x)" group) and two livers from mice treated with $CCl_4$ alone,

259   in which their corresponding mouse tails for sequencing were used as the

260   reference controls (**Supplementary Table 1**). Among the 11 tumor nodules,

261   3 were respectively resected from 3 mice of the "AAI (14x)" group (labeled

262    the AAI group), and the other 8 tumors from another 3 mice in the "AAI

263    (3x) + CCl$_4$" group (labeled as combination group) (**Supplementary Table**

264    **1** and **Supplementary Fig. 1g, l**), where 3 and 4 discrete tumor nodules

265    were resected respectively from two mice (**Supplementary Fig. 3a**).

266        CNV analysis of the 11 tumor nodules by WGS at the depth of about

267    3-fold, compared to their corresponding tail tissues as references, showed

268    that the AA-induced tumors barely had obvious CNV alterations

269    (**Supplementary Table 2**). WES at the average depth of 267-fold for all

270    examined 11 tumor nodules and 8 nontumor livers from the AAI, CCl$_4$ and

271    combination groups, in which WES data of 62-fold for their corresponding

272    tail tissues were references, identified a total of 8107 single-nucleotide

273    variants (SNVs) and 704 small insertions and deletions (indels) in the 11

274    tumor nodules (**Supplementary Tables 1 and 3**).

275        The somatic SNVs and indels of tumors from the AAI group were

276    significantly more abundant than those in the combination group (mean:

277    1555 *vs.* 518, $P = 0.012$) (**Supplementary Table 1**), possibly because of

278    the larger AAI dosage in "AAI (14x)" group than in the combination group.

279    Interestingly, somatic mutations were also found in the paratumor livers

280    and noncancerous livers with or without tumors (mean: 130 *vs.* 76, $P =$

281    0.036) (**Supplementary Table 1**) of mice treated with AAI or CCl$_4$ alone,

282    respectively, suggesting that the increased somatic mutations in livers

283    could be prerequisite to liver cancer development triggered by AAI.

284   Significantly, the tumors exhibited remarkably high proportions (69%)

285 of A>T transversions, whereas the nontumor liver tissues, except for one

286 (M4P: 33%), did not show such feature (10%) (**Supplementary Fig. 3b, c**

287 and **Supplementary Table 4**). Notably, the load of A>T mutations in the

288 AAI group was larger than those in the combination group (mean: 1043 *vs*.

289 305, $P = 0.012$) (**Supplementary Table 4**), which was consistent with the

290 observation that the total applied AAI amount was higher in the AAI group.

291   The mutational profile of each tumor nodule was depicted (**Fig. 3a** and

292 **Supplementary Fig. 3d-g**) and showed that the trinucleotide context of the

293 highest proportion of A>T mutations was CTG (or CAG on the

294 complementary strand). The pentanucleotide context of the highest

295 proportion of A>T mutations was CCTGT (or ACAGG on the

296 complementary strand) (**Supplementary Fig. 3h**). In the series of

297 mutational deciphering analysis, 7 COSMIC mutational signatures were

298 detected (**Fig. 3b**), of which signature 22 related to AA was obviously

299 dominant in all tumor nodules. Except for signature 22, signatures 1, 5, 6,

300 17 and 23 were detected in these mouse tumor nodules, of which signature

301 1 related to deamination of 5-methylcytosine, and 5 related to aging, have

302 been found in all human cancer types and most cancer samples, while

303 signature 6 is associated with defective DNA mismatch repair, and the

304 etiologies of signatures 5, 17 and 23 remains unknown. Tobacco-associated

305 signature 4 was surprisingly detected in two mouse tumors.

306    However, it was noteworthy that two tumors, M4T1 and M6T,

307    presented higher levels of T > G mutations (marked with a prominent peak

308    at ATG > AGG) (**Supplementary Fig. 3e**), the etiology of which remains

309    unknown. Subsequently, the cosine similarities between the mutational

310    spectra of these tumors and signature 22 (typical AA signature) were

311    calculated. Except for M4T1 and M6T, which showed a somewhat lower

312    similarity (0.59 and 0.64) to signature 22 due to distortion of the higher T >

313    G mutations, the mutational profiles of the tumors were nearly identical to

314    signature 22, having a cosine similarity larger than 0.9 (**Fig. 3c** and

315    **Supplementary Table 5**).

316    Interestingly, the mutational spectrum of one paratumor tissue (M4P,

317    A>T, 33%) also exhibited a similar feature to signature 22 (cosine

318    similarity = 0.69) (**Fig. 3a**), which was obviously higher than the other

319    nontumor liver tissues (average cosine similarity = 0.14) (**Fig. 3c** and

320    **Supplementary Table 5**). This result indicated that M4P could be

321    associated with a precancerous process, which was consistent with the

322    pathology of hyperplasia (**Supplementary Fig. 3a**), along with higher

323    somatic mutations (189) than the other two paratumor livers from the same

324    group (mean: 100.5), albeit being significantly lower than those of the

325    tumor nodules from the same group (mean: 518) (**Supplementary Table

326    1**).

327    Except for M4P, other nontumor liver tissues exhibited C>T (average

328  36%) rather than A>T (average 10%) as the predominant mutation

329  category (**Supplementary Fig. 3c** and **Supplementary Table 4**). Their

330  mutational spectra exhibited higher similarities to signatures 5 and 6

331  (**Supplementary Fig. 3i** and **Supplementary Table 5**).

332      Previous studies have revealed that AA-induced mutations are likely to

333  be transcriptionally strand biased[8]. Here, the calculated average ratio of

334  A>T mutations on the nontranscribed strand versus the transcribed-strand

335  was 2.02 ($P < 0.001$) in tumors (**Supplementary Fig. 3j** and

336  **Supplementary Table 6**), indicating the existence of a transcription-

337  coupled repairing (TCR) mechanism to fix the AA-mediated mutations. To

338  further validate the influence of the transcription history on the

339  asymmetries of the A>T strand distribution, we investigated the A>T

340  mutation counts on both strands in the five defined gene expression

341  categories, from low to high expressions, according to the gene expression

342  profiles of the 11 tumor nodules (**Supplementary Table 7**). Next, the

343  mutation counts on the transcribed versus nontranscribed strands were

344  analyzed within each defined gene category, showing that the strand bias

345  of A>T mutations was indeed positively correlated with the gene

346  expression levels (**Fig. 3d**).

347      The collective data revealed that the characteristic mutational signature

348  of A to T transversions and the COSMIC signature 22 induced by AAI

349  were involved in tumorigenesis and could be necessary and critical for the

350    development of liver cancer.

351    **Affected driver genes and signaling pathways**

352    The AAI-mediated characteristic A>T mutations could damage the driver

353    genes that could initiate liver cancer. To identify the driver mutations of

354    these genes in AAI-induced mouse tumors, we searched for genes that were

355    mutated more frequently than expected given the average observed

356    mutation frequency. Interestingly, we found 1919 genes with

357    nonsynonymous mutation in the 11 tumor nodules (**Supplementary**

358    **Tables 1**), of which 98 genes with a total of 123 nonsynonymous somatic

359    mutations belong to the Cancer Gene Census as known driver genes (Tier

360    1 for 77 genes), or those with strong indications for a role in cancer but

361    with less extensive available evidence (Tier 2 for 21 genes) in human

362    cancers (**Supplementary Table 8**). Interestingly, 92 (75%) of the 123

363    nonsynonymous mutations were A>T mutations.

364        The statistically significantly mutated genes included *Hras*, *Sfi1*, *Muc4*,

365    *Sp140*, *Vmn2r121* and *Inpp5d* (**Supplementary Table 9**). The well-studied

366    oncogenic A>T mutations led to the change of *Hras* Q61L (CAA>CTA) in

367    8 of 11 (72.7%) tumor nodules, and *Kras* Q61L (CAA>CTA) and *Braf*

368    V637E (GTG>GAG) were also identified in two other tumor nodules (**Fig.**

369    **4a**), indicating that the cancer-promoting mutations of the Ras/Raf

370    pathway were crucial in AAI-induced liver cancer. Interestingly, 4 of 11

371    (36.4%) tumors presented *Muc4* (4 A>T mutations), in which the same

372   *Muc4* (c.3869T>A) mutation was detected in two tumors (**Fig. 4a**). *Muc4*

373   as an oncogene is listed in the Cancer Gene Census Tier 2, mutations of

374   have appeared in many human cancers, including HCC[10,17,18], and are

375   associated with tumor metastasis[19]. *Sfi1* encoding a spindle assembly

376   associated protein, which was reported to be mutated in human HCC[18],

377   showed 8 mutations (A>T mutations) in 4 mouse tumors (**Fig. 4a**). *Sp140*

378   encodes a member of the SP100 family of proteins, *Inpp5d* encodes a

379   member of the inositol polyphosphate-5-phosphatase (INPP5) family, and

380   unknown functional *Vmn2r121* also showed a higher mutation frequency

381   in these mouse tumors (3/11), of which *Inpp5d* (3/11), also named *SHIP1*,

382   involved in the PI3K-AKT pathway, has been known to be mutated in

383   human cancers, including HCC[17] (**Supplementary Table 3** and **Fig. 4a**).

384       Both the Ras/Raf and PI3K-AKT pathways could participate in the

385   pathogenesis of all these liver tumors (**Fig. 4a, b**). Except for *Ras* and *Braf*,

386   five genes that regulate Ras activity, *Nf1* (2/11), *Rasal2* (1/11), *Sos2* (1/11),

387   *Rasgrp1* (1/11) and *Rasgrp4* (1/11), were also identified with A>T

388   mutations (**Supplementary Table 3**). In addition, some genes encoding

389   growth factors and receptors with tyrosine kinase, such as *Hgf* (3/11), *Egf*

390   (1/11), *Fgf13* (1/11), *Kdr* (2/11), *Pdgfra* (1/11), *Pdgfrb* (1/11), and *Fgfr3*

391   (1/11), except for *Met* (1/11), were also influenced by A>T mutations. *Hgf*

392   mutations also appear in human HCC[10,17,18]. Moreover, the mutant genes

393   were significantly enriched in the PI3K-AKT signaling pathway ($P = 6.6$

394  $\times 10^{-6}$) (**Supplementary Table 10** and **Fig. 4a, b**), which included those

395  encoding growth factors and the receptors mentioned earlier. Except for

396  *Inpp5d*, some genes encoding phosphoinositide-3-kinase (PI3K), such as

397  *Pik3cg* (2/11), and modulators of AKT activity such as *Ppp2r2d* (1/11),

398  *Ppp2r5e* (1/11), *Phlpp1* (1/11) and *Tcl1* (1/11), were also influenced by

399  A>T mutations. The tumor suppressor gene *Tsc1* (2/11), as a negative

400  regulator of mTORC1 and *Rheb* (1/11) activating the protein kinase

401  activity of mTORC1, demonstrated A>T mutations. (**Fig. 4a**).

402  Some mutations could damage development-related genes, including

403  components of the Hippo, Notch and Wnt pathways (**Fig. 4a, b**). *Fat4*

404  (2/11), *Cdh1* (1/11), *Nf2* (1/11), *Lats1* (1/11), *Mst1* (1/11), *Tead1* (1/11) and

405  *Wwc1* (1/11), belonging to the Hippo signaling pathway, had A>T

406  mutations. *Notch1* (1/11), *Notch2* (1/11), *Notch3* (2/11), *Notch4* (2/11) and

407  *Ncor2* (1/11), encoding components of the NOTCH signaling pathway, had

408  A>T mutations. Three genes involved in the WNT signaling pathway, *Apc*

409  (1/11), *Axin2* (1/11) and *Wnt1* (1/11), were mutated in 2 tumors

410  (**Supplementary Table 3**).

411  It was noticeable that these genes encoding DNA polymerases,

412  including *Polq* (2/11), *Pold1* (1/11), *Pold3* (1/11), *Pole* (1/11), *Poln* (1/11)

413  and *Rev1* (1/11), had A>T mutations in 6 of 11 (54.5%) tumors (**Fig. 4a,**

414  **b**). Excluding DNA replication, these DNA polymerases perform

415  exonucleolytic proofreading for DNA repair. It is known that defective

DNA polymerase proofreading contributes to human malignancy, and DNA polymerase mutations in the exonuclease domain have been reported in human tumors with an extremely high mutation load[20,21]. Other DNA repair-related genes, such as *Atm* (2/11), *Prkdc* (2/11), *Mcm8* (2/11) and *Trp53bp1* (1/11), were also mutated in these tumors (**Supplementary Table 3**). Here, we statistically analyzed the correlation between somatic mutations and these gene mutations, showing that the mutations of these DNA genes in tumors were positively associated with the total somatic mutations (**Fig. 4c**).

Some genes related to epigenetic regulation exhibited somatic mutations (**Fig. 4a, b**), including *Tet1* (1/11)*, Dnmt3b* (1/11) *and Dnmt3l* (1/11) for DNA methylation, *Crebbp* (2/11), *Trrap* (2/11), *Kdm6a* (1/11) and *Kmt2c* (1/11) for histone modifications and *Srcap* (3/11), *Smarca2* (2/11), *Smarca4* (1/11), *Smarcb1* (1/11) and *Arid1b* (1/11) for the chromatin remodeling SWI/SNF complex. Mutations of these genes have been described in human cancers, including liver cancer.

In addition, some genes related to ubiquitination and sumoylation were also mutated, including *Usp34* (3/11), *Trim33* (2/11), *Birc6* (2/11) and *Ranbp2* (2/11) (**Fig. 4a**). *Usp34* encoding ubiquitin carboxyl-terminal hydrolase 34 can remove conjugated ubiquitin from Axin1 and Axin2, as a regulator of the Wnt signaling pathway, which is also mutated in human HCC[10,17,18].

438    To further assess the effect of these mutations on the pathogenesis of

439    liver cancer, we analyzed the transcriptome data from these tumor nodules.

440    Some target genes of important pathways were upregulated, such as the

441    Ras, PI3K-AKT, Hippo and Wnt pathways disrupted by the mutations,

442    especially Ras and Hippo (**Supplementary Fig. 3k, Fig. 4d and**

443    **Supplementary Table 11**). Interestingly, some genes, such as *Afp*, *Dlk1*,

444    *Gpc-3*, *Prom1*, *Itga6*, *Cd34* and *Igdcc4* related to stem cells/progenitor

445    cells, along with downstream target genes, such as *Fstl1*, *Dab2, Hes1* and

446    *Mycn* of the Hippo, Notch and Wnt pathways, were upregulated in liver

447    tumors, suggesting that cell differentiation arrest or dedifferentiation

448    occurred in these liver cancers (**Supplementary Fig. 3k** and

449    **Supplementary Table 11**). The transcription of cell cycle-related genes,

450    such as *Ccnd1*, *Ccne1*, and *Cdks,* were increased in tumors, possibly due

451    to activation of the Ras and PI3K-AKT pathways (**Supplementary Table**

452    **7**). Further, the activation of these signaling pathways including Ras, PI3K-

453    AKT and Hippo was verified in these AAI-induced tumors, as compared

454    to adjacent non-tumorous livers (**Fig. 4e** and **Supplementary Fig. 3l**),

455    along with the up-regulated hepatic stem cells/progenitor cell biomarkers.

456    The collective data suggested that AAI contributed to tumorigenesis of

457    liver cancer through activating the RAS pathway, in combination with

458    other deregulated important pathways such as PI3K-AKT, DNA replication

459    and repair, the chromatin remodeling SWI/SNF complex, epigenetic

460  regulation, the development-related Hippo, Notch and Wnt pathways,

461  spindle integrity, and cell adhesion.

462  **AAI-mediated mutations are the early event during malignant clonal**

463  **evolution**

464  Though it was testified that the AA signature was dominant in mouse liver

465  tumors, we had particular interest in whether the AA-mediated mutations

466  were the originating source driving tumor initiation and progression,

467  especially for tumors in the combination group composed of AAI and $CCl_4$,

468  as the application of $CCl_4$ inevitably cast doubts on the role of AA in the

469  carcinogenic processes. Therefore, we further investigated the clonal

470  architecture and AA-related mutational signature distribution in these

471  malignant subclones within tumors from the "AAI alone" and "AAI and

472  $CCl_4$ combination" groups.

473      We first performed a clonality analysis of these 11 liver tumor nodules

474  (see **Methods**). Ten of the 11 nodules exhibited multiple subclones, of

475  which 3 contained 3 subclones and 7 had 2 subclones (**Fig. 5a, b** and

476  **Supplementary Fig. 4a-f, left**).

477      Pure tumor was expected to present a high-density region with nearly

478  a 50% variant allele frequency (VAF) in the Sciclone deconvolution

479  results[22]. The estimated weights of VAF in the tumor dominant clones,

480  however, ranged from 16% to 39% (**Supplementary Table 12**), reflecting

481  a substantial immunological cell infiltration into the tumors as presented in

482    the pathological sections (**Supplementary Fig. 3a**). Alternatively, the

483    subclones could have been initiated in parallel style in the tumorigenesis

484    procedures, especially in tumors M1T, M4T1 and M6T, as their dominant

485    clone had VAF centered at 16%, 19% and 17%. If the subclones formed

486    simultaneously, they should present similar mutational signatures induced

487    by the same etiologies. However, we found that the mutational signatures

488    varied within the multiple subclones within each tumor (**Fig. 5a, b, right**

489    and **Supplementary Fig. 4a-f, right**). In addition, there was a trend

490    towards a diminished AA signature from the higher-weighted to the lower-

491    weighted clones in their signature profiles. Therefore, we considered the

492    subclones with the largest weight to be the initiating founding subclones

493    and the lower-weighted subclones to be formed sequentially in later

494    processes [22]. To determine how the AA-related mutations evolved between

495    the founding clone and subsequent subclones, we calculated the A>T

496    proportions of each subclone within the tumors. For example, after

497    characterizing the multiclonal architectures in M3T (two subclones) and

498    M5T2 (three subclones) (**Fig. 5a, b, left**), we then retrieved their

499    mutational profiles (**Fig. 5a, b, right**). Significantly, the results indicated

500    that the AAI-mediated A>T transversions were predominant in the earliest

501    founding clones and then gradually were reduced in the later subclones.

502    The other 8 tumor nodules composed of multiple subclones also exhibited

503    a similar AA-related signature distribution pattern (**Supplementary Fig.**

504  **4a-f).** This result indicated that the AAI-mediated A>T mutations were the

505  early event during the malignant clonal evolution process, contributing to

506  an average of 82% in the AAI group (M1T, 86%; M2T, 71%; M3T, 88%)

507  and an average of 75% in the combination group (M4T1, 55%; M4T2, 77%;

508  M5T1, 78%; M5T2, 91%; M5T3, 88%; M5T4, 80% and M6T, 52%). By

509  contrast, the non-A>T mutational patterns, such as C>T mutations, slightly

510  increased gradually during clonal evolution and merged into the late

511  subclones **(Fig. 5a, b and Supplementary Fig. 4 a-f, right)**. As illustrated

512  by the downward-directed lines in these 11 tumor nodules (**Fig. 5c**), there

513  was a general trend of A>T mutations that diminished along with the

514  malignant clonal progression within tumors, typically M5T1 and M5T2, in

515  which the A>T transversions almost disappeared in their late-formed

516  subclones.

517  Based on the above analysis, we speculate that, regardless of the AAI

518  alone group or the combination group, AAI-mediated A>T somatic

519  mutations are responsible for the initiation of liver cancer, and the second

520  non-A>T mutations drive malignant clonal evolution and tumor

521  progression, possibly through a synergistic effect between AAI-mediated

522  A>T and non-A>T mutations.

523  **AAI-mediated tumors exhibit diversiform evolution process**

524  The above clonality analysis depicted the intratumor clonal heterogeneity

525  within single tumor nodules, but AAI can induce multiple tumor nodules

526    in the same livers in some mice, and the phylogenetic relationship of these

527    nodules is unclear. To explore their phylogenetic relationship and

528    evolutionary process, here we investigated these tumor nodules and their

529    paratumor livers from mice in the combination group, with a total of the 8

530    discrete tumor nodules; 3 were from one mouse (M4), 4 were from another

531    (M5), and 1 was from the last mouse (M6). We reconstructed their

532    phylogenetic tree to examine the relationship of the somatic mutational

533    patterns among the discrete tumor nodules within the same mice. To

534    establish a control, we applied the reconstruction to all 11 tumor nodules

535    to generate their phylogenetic tree (**see Method**). It was seen that, despite

536    a tiny overlapping distance between M1T and M2T, the other tumor

537    nodules were categorized properly, in accordance with their mouse source

538    (**Fig. 5d**), implying that multiple nodules in M4 and M5 could have

539    evolved from identical ancestors, respectively.

540        Interestingly, the three nodules in M4 shared the specific oncogenic

541    *Hras* Q61L mutation and another 15 identical mutations, supporting the

542    assumption that the separate tumor nodules might have originated from a

543    common ancestor. Moreover, 12 passenger mutations were shared by both

544    the paratumor liver and all three tumor nodules, although they had lower

545    allelic frequencies in the paratumor tissue (e.g., Chr14: 5140816 C>G),

546    suggesting that the tumor could be initiated by the emerging driver

547    mutations, such as *Hras* Q61L, in the background originating cell with

548  passenger mutations (**Fig. 5e**). Additionally, M4T1 and M4T2 might

549  branch later than M4T3, as indicated in the phylogenetic tree (**Fig. 5d**).

550  Unlike the tumor nodules in mouse M4, the four tumor nodules in M5

551  did not share a commonly known driver mutation, although three (M5T1,

552  M5T2 and M5T4) of them shared an identical oncogenic *Hras* Q61L

553  mutation, and the other (M5T3) harbored an oncogenic *Braf* V637E

554  mutation. However, as indicated in the phylogenetic tree, all four nodules

555  could have arisen from the same ancestor and then evolved separately in

556  the late phase. The observation that M5T1, M5T2 and M5T4 were located

557  within a branch (**Fig. 5d**), rather than M5T3, revealed a closer phylogenetic

558  relationship, in accordance with their differences in the initiating driving

559  force. Interestingly, all four tumor nodules in M5 also shared 10 somatic

560  mutations with their paratumor liver, including Chr11: 3176625 G>A that

561  increased the mutation allele frequency from the paratumor liver to the

562  tumor, suggesting that these tumor nodules could be initiated from the same

563  precancerous cells through the emerging driver mutations such as *Hras*

564  Q61L or *Braf* V637E (**Fig. 5f**). Although all four tumor nodules could

565  originate from the same precancerous cells with a similar genetic

566  background, the M5T3 nodule with the *Braf* V637E mutation was

567  distinguished from the other three nodules sharing a common ancestor,

568  wherein the two malignant transformed clones had undergone parallel

569  evolution within M5 liver (**Fig. 5f**).

570    To reveal the kinship between different subclones within the separated

571    tumor nodules from M4 and M5, here we adopted the assumption that the

572    second subclone was generated dependently from the founding clones.

573    However, whether the third weighted subclone was generated dependently

574    or independently of the second-weighted clone was unclear

575    (**Supplementary Fig. 4g**). Here, the M4T2, M5T1 and M5T2 nodules were

576    composed of three subclones (**Supplementary Fig. 4d, e**), where the third

577    subclones within single nodules were depicted to emerge via parallel

578    evolution along with the second subclones by driver genes such as *Ranbp2*,

579    *Actbl2*, *Smarca2* and others (**Fig. 5e, f**). The relationship between the

580    second and third clones could also be replaced by the other model as

581    provided in supplementary Fig. 4g with the same suggested driver genes

582    and cell proportions.

583    In M6, both the paratumor liver and the corresponding tumor had 30

584    overlapping somatic mutations, among which 6 expanded their allelic

585    frequencies more than 5 times from the paratumor liver to the tumor. Next,

586    we examined the presence of mutated reads of essential genes in the

587    paratumor liver tissue in Integrated Genome Viewer (IGV). We found that

588    the paratumor liver of M6 has the same positioned *Muc4* mutation reads in

589    its corresponding tumor, albeit with a very low number (2 reads)

590    (**Supplementary Fig. 4h**). Therefore, we speculate that the *Muc4* mutation,

591    along with the other passenger mutations (like the 6 expanded mutations),

was not sufficient to trigger tumor initiation and that one of the *Muc4* mutated cells, if acquiring the oncogenic *Kras* mutation, would be malignantly transformed and then become proliferative (**Supplementary Fig. 4i**).

Together with the above clonality analysis within tumor nodules (**Fig. 5e, f** and **Supplementary Fig. 4i**), we may see that, except for Ras/Raf A>T mutations as the earliest events in the tumorigenesis of M4, M5 and M6 mice, the patterns of other driver mutations exhibit obvious heterogeneity among different nodules and subclones within single nodules. The different driver genes, such as *Polq, Fgfr3, Met*, *Asxl1*, *Pdgfrb*, *Notch3, Mllt10* and *Tet1* with A>T mutations and late emerging driver genes such as *Fat4, Smarca2* and *Inpp5d* with non-A>T mutations, synergistically facilitate malignant subclonal evaluation.

**Human liver cancer exhibits an AA-mediated mutational signature**

To explore the AA signature intensities in human cancers, we first grasped a quick estimation of the AA signature contribution through the webserver mSignatureDB (**see Methods**). Interestingly, some human cancer resources presented a possible characteristic A>T mutational signature to certain degrees (**Supplementary Fig. 5a**). However, current approaches for signature deconvolution were mostly based on nonnegative matrix factorization (NMF), which do not consider mutation counts. Based on a simulated dataset with 844 samples (see **Methods**), we noticed that, when

614    the mutation number fell within 100, the mean squared error (MSE) of

615    deconvolution increased exponentially as the mutation number decreased

616    (**Supplementary Fig. 5b**). Therefore, in an effort to reduce false positives

617    with a low number of mutation counts, we improved the traditional

618    decomposing strategy by introducing the bootstrap sampling technique to

619    make up for the shortcoming that the deconvolution originally did not

620    provide empirical $P$ values. Afterwards, we used the simulated dataset to

621    evaluate the performance of the bootstrap performance by considering its

622    accuracy, specificity, sensitivity and F1 measure. Surprisingly, the

623    evaluation revealed no false positive mistakes (specificity = 1) throughout

624    each threshold (**Supplementary Fig. 5c**). In addition, the method yielded

625    fine-tuned accuracies, sensitivities and F1 measures for detecting the AA

626    signature intensity in the interval between 10 to 90%. Therefore, we

627    decided to identify the AA signature intensities in human cancers with a

628    cutoff of both 0 and 10% ($P = 0.05$). Furthermore, we required that the AA

629    signature proportion should be larger than the MSE to balance the

630    instabilities due to low mutation counts and unpredictable noise. Using this

631    method, we found that the primary positive detection of AA signature in

632    the majority of tumors and leukemia was prone to be false.

633        However, with this rigorous method, we identified an AA signature in

634    liver cancer, including HCC and ICC (**Table 1**). We detected a

635    characteristic AA signature in 52 (20%) of 313 China (mainland) HCC

636    samples catalogued in the International Cancer Genome Consortium

637    (ICGC) project, 68 (69%) of 98 Taiwan-based samples[10], and 7 (8%) of 88

638    samples that were accepted in Hong Kong[23], as well as 6 (55%) of another

639    11 samples in mainland China[24]. In total, 133 (26%) of 510 Chinese HCCs

640    were identified with an AA signature. Moreover, we detected the AA

641    signature in 3 ($<$ 1%) of 594 HCCs from Japan[17,18], 22 (10%) of 231 HCCs

642    in Korea[25], 4 (44%) of 9 HCCs in Singapore[26], 29 (10%) of 364 HCCs in

643    the US from The Cancer Genome Atlas (TCGA) dataset, and 1 ($<$ 1%) of

644    249 HCCs from France in the ICGC. Among the TCGA HCCs, Asian

645    ethnicity patients had higher detection rate of AA signature, which is 24

646    (15%) of 160. In addition, we noticed that 11 (11%) of 103 ICC from

647    China[27] showed the AA signature, indicating that AA might play a role in

648    the etiology of human ICCs (**Table 1**).

649    We also investigated other human cancer types. Bladder cancer in

650    China[28] and kidney cancer in Europe (from ICGC) presented different

651    extents of AA signature contributions (**Supplementary Table 13)**. In

652    addition, it was noteworthy that we identified 1 case of esophagus cancer

653    in China (from ICGC) that exhibited the AA signature (**Supplementary**

654    **Tables 13, 14)**. Generally, in all the cancer types, liver cancer, including

655    HCC and ICC, presents the most disturbingly high proportions of the AA

656    signature. Additionally, we retrieved the T$>$A mutations of liver cancer in

657    COSMIC and found that the T$>$A mutation profile was highly consistent

658    with the AA signature (cosine similarity = 0.94) (**Supplementary Fig. 5d**),

659    which implicated that the AA signature operated predominantly in causing

660    the T>A mutation in human liver cancers. Next, we compared the AA

661    signature intensities in the affected human cancers, which exhibited wide

662    intertype and regional variabilities (**Fig. 6a**). Kidney cancer and liver

663    cancer were more susceptible to AA genotoxicity, as their AA intensities

664    were prominently higher. According to the proportions of affected

665    populations and their AA signature intensities, it seemed that the most

666    influenced cancer type was HCC in China.

667        To uncover the affected driver mutation by AA, next we investigated

668    the liver cancers that harbored a nonsilent or splicing site with AA

669    signature mutations in known driver genes. Here, the mutations were

670    ascribed to the signatures using a Bayesian classifier, which showed that

671    83 (16%) of the 510 HCCs from China had driver genes affected by the

672    AA signature, while 16 (4%), 1 (11%), and 11 (5%) were respectively

673    identified in the US, Singapore and Korea HCCs, as well as 4 (4%) ICCs

674    from China (**Table 1**).

675        In addition, we also applied a more stringent threshold to estimate the

676    AA contribution to human liver cancers by applying the lowest observed

677    AA signature intensities in our experimental mouse tumors, requiring the

678    95% lower confident interval of the AA signature contribution to be larger

679    than 52%, which was the lowest contribution of the bootstrapped results

680  obtained for these mouse tumors (**Supplementary Fig. 5e**). We used this

681  value as an indication of carcinogenic dosage in liver cancer. Our method

682  appeared to perform well for detecting an exposure contribution above

683  50%, as the accuracy, specificity, sensitivity and F-measure all equaled 1

684  (**Supplementary Fig. 5c**). As a result, 64 (3%) of the 1957 HCCs

685  worldwide were identified as having an AA exposure greater than 52%

686  (**Table 1** and **Fig. 6b**). Significantly, among them, 57 (89%) of the 62 were

687  from China. Additionally, 3 (3%) of the ICCs in China were identified with

688  the same standard.

689      The finding that the AA mutational signature exists in human cancer,

690  particular in Chinese liver cancer, based on the different criteria (**Table 1**),

691  strongly indicates that AA exposure in Chinese population might have been

692  one of the major risk factors for the onset of liver cancer, including HCC

693  and ICC.

694  **AA exposure could be operative in an earlier stage of human liver**

695  **cancer**

696  We hope to determine whether AA is operative in the initial stage of human

697  liver cancer, similar to its role in mouse liver cancer. Therefore, we

698  performed a clonality analysis in the TCGA-derived HCC samples with the

699  AA signature and DNA copy number profiles, as we did in the mouse

700  tumors. Here, we chose 6 HCC samples with an AA signature contribution

701  above 30% to compare their AA exposures along with clonal evolution

702    (**Fig. 6c**). It was shown that, in all 6 tumors, AA-associated mutagenesis

703    was operative in the initial subclones. In sample DD.AADF, the AA

704    signature was found to be the predominant etiology in the founding

705    subclone and decreased in later formed subclones, demonstrating a similar

706    trend to that observed in mouse liver cancers. In contrast, the AA signature

707    seemed to continue or even increase throughout the clonal evolution of the

708    other five human HCCs, possibly due to the prolonged AA exposure in

709    patients rather than acute exposure in our mouse model. In addition, among

710    3 of 9 Singapore patients with HCCs carrying the AA signature, we

711    ascribed the mutations as early (trunk) and late (branch), as the study

712    provided multisector sequencing results. It was shown that trunk mutations

713    had higher proportions of A>T transversions, indicating that AA exposure

714    was likely to play roles in the initial stage of tumorigenesis

715    (**Supplementary Fig. 5f**).

716        Moreover, we analyzed the contribution of the AA signature to known

717    driver genes in HCCs. Excluding *TP53*, *ARID2*, *ARID1A* and *AXIN1*

718    frequently harbored nonsilent A>T mutations (**Fig. 6d**), which encode

719    components of SWI/SNF complex and Wnt-β-catenin pathway that were

720    also disrupted in AAI-induced mouse liver cancer.

721        In addition, to identify the driver mutations induced by AA exposure,

722    we further investigated the cancer genomic data regarded as AA-induced

723    liver cancer according to a more stringent criterion referring to an AA

724 signature > 52%. We retrieved the genomic data from 62 HCCs meeting

725 the criterion to obtain 74045 putative mutations, and we performed

726 MutSigCV analysis, which revealed that the scattered A>T mutations in

727 the tumor suppressor gene *TP53* were significantly affected in AA

728 dominant human liver cancer. *TP53* (q < 0.1), along with eight A>T

729 transversions, was significantly mutated in these AA-related liver cancers

730 (**Fig. 6e**), where the A>T transversions led to nonsense mutations that

731 dispute the structure of TP53, or missense mutations that alter its functions

732 by the mutated DNA-binding domain.

733 We also ascribed each mutation to a specific signature and then selected

734 the AA signature A>T mutations (**Supplementary Table 15**) for analysis

735 via the oncodriveCLUST algorithm, a positional clustering method,

736 because most of the oncogenic mutations of one gene were enriched at a

737 few specific loci (aka hot-spots). Significantly, 17 genes with FDR smaller

738 than 0.1 were identified (**Supplementary Table 16**). It is noteworthy that

739 JAK1 S729C induced by a c. 2185 A>T mutation was identified (FDR = 5

740 $\times\ 10^{-4}$) as a candidate driver (**Fig. 6e** and **Supplementary Table 16**), as

741 there were four hits at the exact same locus in *JAK1*. Interestingly, the

742 mutation was found in Chinese liver cancer and validated as an oncogenic

743 driver in a previous study[23].

744 **Discussion**

745 Liver cancer is the seventh most common cancer and the third leading

746  cause of cancer-related death worldwide[29]. In China, its incidence and

747  mortality rate are higher[30]. Liver cancer has several known risk factors,

748  including infection with HBV and hepatitis C virus, alcohol consumption,

749  and aflatoxin B1 contamination of food. Recently, AA has been statistically

750  associated with human liver cancer, especially in Chinese patients[8,10].

751  However, no experimental evidence supports the notion that AA can

752  directly lead the liver cancer. Significantly, our results demonstrate that

753  AAI can directly induce mouse liver cancer, including HCC and ICC, in a

754  dose-dependent manner, and increases the incidence of liver cancer when

755  the liver is injured, such as $CCl_4$ administration. This finding is consistent

756  with clinical practice in China and some Asia countries, because some

757  Chinese patients, including hepatitis patients, often take traditional Chinese

758  herbal remedies that could contain AA. In fact, in human liver cancer, we

759  found the characteristic AA signature in HCC and ICC patients based on

760  very stringent criteria, indicating that AA exposure was the leading cause

761  in some liver cancers, especially in Chinese patients. Therefore, our animal

762  experiments and analyses of human liver cancer strongly indicate that AA

763  can directly lead to liver cancer and should be listed as major risk factor

764  for liver cancer.

765      AA exposure could be prone to trigger some driver mutations by

766  characteristic A>T transversions, which lead to a growth advantage of the

767  malignant transformed clones. Significantly, AAI-mediated mutations are

768  found to be the early event during malignant clonal evolution in mouse and

769  human liver cancer. AAI-DNA adducts could be detected not only in livers

770  from mice exposed to AA but also in multiple heterogeneous subclones

771  within the same tumor nodules, and the AAI-mediated characteristic A>T

772  mutational signature was found in the founding subclones in both mouse

773  and human liver cancers, further supporting the critical nature of AA

774  exposure in some forms of liver cancer.

775      AA could prefer to damage different driver genes in different species,

776  although these mutations share similar A>T transversions. A previous

777  study had shown that the *ras* family, including *Hras*, *Kras* and *Nras*, had

778  the same activating mutation — Q61L (CAA to CTA) — in oral

779  administration AA-induced rat tumors[14]. In this study, the same activating

780  mutation in *Hras* and *Kras* also occurred in most AA-induced HCC

781  samples (*Hras*, 8/11; *Kras*, 1/11), indicating that the Ras pathway is crucial

782  in AAI-induced mouse liver cancer. Interestingly, activating mutations of

783  *KRAS* are frequent in human ICC[31], although relatively lower in human

784  HCC.

785      More interestingly, AA-mediated mutations also alter other genes that

786  can lead to the deregulation of some signaling pathways, such PI3K-AKT,

787  the chromatin remodeling SWI/SNF complex, epigenetic regulation, and

788  the development-related Hippo, Notch and Wnt pathways (**Figure 4**),

789  which are often associated with human HCC and ICC. Like human HCC,

790 AAI-induced mouse HCCs also express hepatic stem or progenitor cell-

791 related biomarkers such as *Afp*, *Gpc3*, *Dlk1* and *Prom1* (**Supplementary**

792 **Fig. 3k**). In addition, it should be pointed out that, although *Tp53* and *Jak1*

793 mutations were not found in the AAI-induced mouse HCCs, some point

794 mutations of *TP53* and *JAK1*, especially *JAK1* S729C, could be considered

795 as candidate biomarkers for AA exposure, similar to TP53 R249S for

796 aflatoxin B1 contamination.

797 In conclusion, this study provides documented evidence indicating that

798 AA can directly induce mouse liver cancers, including HCC and ICC,

799 similar to the genetic pathogenesis of human liver cancers. In light of the

800 animal model, AA exposure is considered as a major risk factor for some

801 human liver cancers, especially among Chinese patients. The featured

802 mutational process during malignant clonal evolution in AA-induced liver

803 cancer reveals that AAI-mediated characteristic mutations are the earliest

804 genetic event in tumorigenesis. Our data lay a solid foundation for the

805 prevention and diagnosis of AA-associated human cancers, especially liver

806 cancer.

807 **URLs.** ICGC data portal, https://dcc.icgc.org/; TCGA data portal,

808 https://portal.gdc.cancer.gov/; COSMIC mutation signatures,

809 http://cancer.sanger.ac.uk/cosmic/signatures; COSMIC cancer census

810 genes, http://cancer.sanger.ac.uk/census; HMMcopy (v1.22.0),

811 http://www.bioconductor.org/packages/release/bioc/html/HMMcopy.html;

812 MsigDB, http://software.broadinstitute.org/gsea/msigdb; R package

813 pracma (v2.1.4), https://cran.r-

814 project.org/web/packages/pracma/index.html; mSignatureDB,

815 http://tardis.cgu.edu.tw/msignaturedb/Browse/.

816 **Methods**

817 **Mice.** The mice used in this study were on a C57BL6/J background. The

818 wild-type C57BL6/J mice were purchased from the Slaccas Company

819 (Shanghai, China). LoxP-flanked (floxed [f]) *Pten* (*Pten$^{f/f}$*) mice (The

820 Jackson Laboratory) and Alb-Cre mice (The Jackson Laboratory) were

821 crossed to generate conditional liver-specific PTEN-KO mice designated

822 as *Pten$^{LKO}$*. All animal experiments were conducted under procedural

823 guidelines and severity protocols with the approval granted by the

824 Institutional Review Board on Bioethics of Shanghai Jiao Tong University.

825 **HCC induction.** As mentioned in the results section, the male mice were

826 randomly grouped and administered with AAI (2.5 or 5 mg/kg, dissolved

827 in PBS) or a combination of AAI (2.5 mg/kg) and $CCl_4$ (0.5 ml/kg,

828 dissolved in corn oil) by intraperitoneal injection for different doses and

829　times since the age of 1 or 2 weeks. Additional groups were injected with

830　$CCl_4$ alone or vehicle as controls. The exposure timeline was presented in

831　**Fig. 1a**. Mice were sacrificed with $CO_2$ anesthesia. The visible discrete

832　tumors at mice livers were dissected and counted. Tumor sizes was

833　measured with a caliper at its largest diameter. No mice were excluded in

834　subsequent analyses.

835　**Histology, immunohistochemistry and immunofluorescence.** Formalin-

836　fixed tissues were embedded in paraffin. Sections (5 μm) were stained with

837　hematoxylin and eosin (H&E), and PicroSirius Red. For

838　immunohistochemical (IHC) staining, sections were incubated with

839　primary antibodies against AFP (polyclonal rabbit, 1:100; ab46799,

840　Abcam), Ki67 (polyclonal rabbit, 1:200; ab15580, Abcam) and CK19

841　(monoclonal rabbit, 1:400; ab52625, Abcam) overnight at 4 °C. HRP-

842　conjugated anti-rabbit secondary antibody (polyclonal goat, 1:400; A0545,

843　Sigma) and DAB (Sangon Biotech, Shanghai, China) were used to detect

844　the primary antibodies, followed by hematoxylin redyeing. For

845　immunofluorescence assay, tissues were embedded in OCT (optimal-

846　cutting-temperature compound). Cryosections (5 μm) were fixed in 4%

847　paraformaldehyde for 10 min and then permeabilized with 0.1% Triton X-

848　100. Sections were incubated with primary antibody against γ-H2AX

849　(monoclonal rabbit, 1:200; 9718, Cell Signaling) overnight at 4°C. Primary

850　antibody were detected using fluorescent-conjugated secondary antibody

851 (polyclonal Donkey, 1:1000; A21206, Invitrogen). The nuclei were stained

852 with DAPI and mounted with anti-fading mounting reagent

853 (Fluoromount™ Aqueous Mounting Medium, Sigma). Brightfield images

854 were taken using a Nikon Eclipse Ni microscope. Fluorescence images

855 were taken using fluorescent confocal microscope (Nikon A1Si). Sirius

856 Red stained area or fluorescence intensities of 10 non-overlapping fields in

857 each section were quantified using Fiji Image J at $\times$ 100 or $\times$ 400

858 magnification.

859 **Western blot analysis.** Liver tissues were lysed and the protein

860 concentrations were determined using the BCA assay (Thermo Scientific).

861 Membranes were incubated with the following primary antibodies: anti-γ-

862 H2AX (monoclonal rabbit, 1:1000; 9718, Cell Signaling), anti-p53

863 (monoclonal mouse, 1:200; sc-126, Santa Cruz), anti-Bax (polyclonal

864 rabbit, 1:200; sc-493, Santa Cruz), anti-p-ATR (Ser428, polyclonal rabbit,

865 1:1000; 2853, Cell Signaling), anti-AFP (polyclonal rabbit, 1:1000;

866 ab46799, Abcam), GPC3 (polyclonal rabbit, 1:400; ab66596, Abcam), E-

867 cadherin (polyclonal rabbit, 1:1000; 20874-1-AP, Proteintech), p-ERK (Tyr

868 204, monoclonal mouse, 1:200, sc-7383, Santa Cruz), ERK1 (polyclonal

869 rabbit, 1:400, sc-94, Santa Cruz), p-AKT (Ser473, monoclonal rabbit,

870 1:2000; 4060, Cell Signaling), AKT (polyclonal rabbit, 1:1000; 9272, Cell

871 Signaling) and YAP (polyclonal rabbit, 1:1000; 4912, Cell Signaling) and

872 anti-β-actin (monoclonal mouse, 1:5000; A2228, Sigma). HRP-conjugated

873    secondary antibodies (polyclonal goat, 1:10000; A6154 and A4416, Sigma)

874    were applied. Proteins was detected by ECL reagent (Share-bio, Shanghai,

875    China).

876    **In vivo alkaline comet assays.** The male mice (n = 4) were administered

877    with PBS (10 ml/kg) or AAI (2.5 or 5 mg/kg) by intraperitoneal injection

878    at the age of 2 weeks. Mice were anesthetized at 3 hours after

879    administration. Livers were perfused with Hanks' balanced salt solution

880    (HBSS). Then alkaline comet assay was performed with CometAssay kit

881    (4250-050-K, Trevigen) following the manufacturer's instructions.

882    Afterwards, the slides were stained with SYBR Green I (Sangon Biotech,

883    Shanghai, China). Fluorescence images were taken using fluorescent

884    confocal microscope (Nikon A1Si) at $\times$ 100 or $\times$ 400 magnification. Tail

885    DNA were analyzed using the CASP[32]. At least 100 cells were randomly

886    selected and analyzed per sample. A total of 600 cells were analyzed per

887    group.

888    **DNA and RNA extraction.** Tissues were minced and digested overnight

889    at 55 °C in 10 mM Tris-HCl (pH 8.0) containing 100 mM EDTA, 10 mM

890    NaCl, 0.1% SDS, proteinase K (0.2 mg/ml), and RNase A (0.2 mg/ml).

891    DNA was purified by phenol/$CHCl_3$. RNA was extracted from frozen tissue

892    using TRIzol according to the manufacturer's instructions.

893    **Synthesis and identification of dA-AL-I.** 7-(deoxyadenosin-N6-yl)

894    aristolactam I (dA-AL-I) was synthesized by incubating deoxyadenosine

895  (dA) with AAI according to the method described previously[33]. A mixture

896  containing dA-AL-I, AAI and dA was yielded.

897      Synthetic dA-AL-I was analyzed with an ACQUITY ultra performance

898  liquid chromatography (UPLC) system (Waters) connected to a XEVO-

899  G2XS quadrupole time-of-light (QTOF) mass spectrometer (UPLC-

900  QTOF-MS) (Waters) with electron spray ionization (ESI). Seven

901  microliters synthetic dA-AL-I was injected into an ACQUITY HSS T3

902  column (2.1 mm × 100 mm i.d., 1.8 μm particle size) (Waters) in positive

903  electrospray ionization mode at a flow rate of 0.4 ml/min. Mobile phase A

904  and B were 0.1% formic acid in water and acetonitrile, respectively. The

905  gradient program used was: 0-1 min, 1% B; 1-3 min, 1-30% B; 3-7 min,

906  30% B; 7-9 min, 30-100% B; 9-11.2 min, 100% B; 11.2-11.3 min, 100-1%

907  B and 11.3-13 min, 1% B. The ESI source was operated in positive ion

908  mode with a capillary voltage of 2 kV, cone voltage of 40 V, source

909  temperature of 115 °C, desolvation temperature of 450 °C, cone gas flow

910  of 50 l/h, and desolvation gas flow of 900 l/h. The mass spectra were

911  acquired over m/z 50-1200 in full scan mode. The secondary mass spectra

912  (MS/MS) were also acquired in the positive mode in the range of m/z 50-

913  600 with the collision energy of 10-30 eV. Data were acquired by software

914  Masslynx v 4.1 and analyzed by UNIFI 1.8.1.

915  **Mass spectrometry identification and quantitation of AAI-DNA**

916  **adducts.** Liver and renal tissue DNA (500 μg and 50 μg in 5 mM bis-tris-

HCl buffer (pH 7.1) containing 10 mM $MgCl_2$) was digested with DNase I (Worthington), nuclease $P_1$ (Sigma), alkaline phosphatase (Worthington) and phosphodiesterase I (Worthington) as described previously[34]. Protein was precipitated by adding 2 vol of chilled $C_2H_5OH$ and centrifuging. The supernatant was concentrated by vacuum centrifugation and dissolved in a solvent of 1:1 $H_2O$/DMSO (50 μl)[34].

Ultra-high-performance liquid chromatography/triple quadrupole mass spectrometry (UHPLC/QQQ MS) was used for identification and quantitation of AAI-DNA adducts. The chromatographic column and mobile phases used with this system were the same as those that were used with the UPLC-QTOF-MS system. The injection volume of synthetic dA-AL-I or DNA digestion product was two microliters. The QQQ MS system (Waters) was operated in positive ion mode with a capillary voltage of 0.5 kV, cone voltage of 30 V, source temperature of 150 °C, desolvation temperature of 500 °C, cone gas flow of 150 l/h, and desolvation gas flow of 1000 l/h. Detection of the ion pairs were performed by multiple reaction monitoring (MRM) mode. The MRM ion pairs for dA-AL-I were 543.16/427.12, 543.16/395, and 543.16/292, and the corresponding collision energies were 25, 30 and 35 eV, respectively. The quadrupoles were set at unit resolution. The analytical data was processed by Masslynx.

**WGS and copy number analysis.** Whole-genome DNA libraries were created with Illumina Truseq Nano DNA HT Sample Prep Kit following

939    the manufacturer's instructions. The libraries were sequenced on Illumina

940    Hiseq platform and 150 bp paired-end reads were generated. We

941    investigated the copy number patterns of the samples applying the suite of

942    HMMcopy (v1.22.0) on the WGS data. Briefly, the coverages were

943    initially corrected with the GC and mappability bias of the reference

944    genome. Then the corrected signals were segmented using Hidden Markov

945    Model to yield an estimate of the copy number events.

946    **WES and somatic mutation calling.** Whole-exome capture was done with

947    Agilent SureSelect Mouse All Exon V1 kit according to the manufacturer's

948    instructions. The libraries were sequenced on Illumina Hiseq platform and

949    150 bp paired-end reads were generated. Raw sequencing reads of WES

950    were aligned against mouse reference build GRCm38 using bwa

951    (v0.7.11)[35]. The duplicates were removed by Picard Tools (v1.4.5)

952    (http://broadinstitute.github.io/picard/). The base quality was recalibrated

953    using the Genome Analysis Toolkit (GATK 3.7-0-gcfedb67)[36]. Mutect

954    (v2.0)[37] was employed to predict somatic mutations of liver tumor and

955    adjacent tissue with the corresponding tail tissue being the control. The

956    mutations were removed if it was in mouse dbsnp. Filtered somatic

957    mutations were functionally annotated by ANNOVAR[38], using the

958    RefGene database. Nonsynonymous, stop-loss, stop-gain and splice-site

959    SNVs (based on RefGene annotations) were considered to be functional.

960    SNPEFF (v4.3s)[39] were used to predict functional influences of the somatic

961    mutations. Bam files were visualized in Integrated Genome Viewer

962    (IGV)[40]. Nonsynonymous mutant genes in AAI-induced mouse liver

963    cancer were performed with KEGG pathway enrichment analysis[41].

964    **RNA sequencing, analysis and annotation.** RNA-seq libraries were

965    generated using NEBNext® Ultra[TM] RNA Library Prep Kit for Illumina®

966    according to the manufacturer's instructions. Then they were sequenced on

967    an Illumina Hiseq platform to generate 150 bp paired-end reads. Sequenced

968    reads were mapped to the GRCm38 UCSC annotated transcripts via Tophat

969    (v2.1.0)[42]. Transcripts were then assembled and counted with the Cufflinks

970    suit (v2.2.1)[43]. Differentiated expressed genes were analyzed by Cuffdiff.

971    **Gene set enrichment analysis.** Gene set enrichment analysis (GSEA

972    v3.0)[44] was performed using the normalized expression values generated

973    by Cuffnorm between the liver and the tumors. Differential enrichment was

974    calculated using the signal-to-noise metric. FDR 0.1 was set as significant

975    in the analysis. To investigate the response of expression alterations to the

976    significantly mutated pathways, we investigated the gene sets of the target

977    genes of the activated transcription factors, such as Ets1 in the Ras pathway.

978    Therefore, the analysis was run using the 'motif'[45], together with 'KEGG',

979    and 'GO' signature collections from the Molecular Signature Database

980    (MsigDB). Differentially expressed genes in the Ras, Hippo and PI3K-

981    AKT downstream transcription factor associated gene sets (FDR < 0.1)

982    were selected and listed in **Supplementary Table 11**, along with

983    differentially expressed the liver cell stem markers[46], and Wnt

984    (http://web.stanford.edu/group/nusselab/cgi-bin/wnt/)    and    Notch[47]

985    signaling pathway target genes. The top 30 genes with higher fold change

986    were used to generate the heatmap in **Supplementary Fig. 3k**.

987    **Clonal and phylogenetic reconstruction.** With the information of copy

988    number and mutation allele frequency, Sciclone was used to characterize

989    coexisting subpopulations in the individual tumors, both in the mouse liver

990    tumors and the TCGA-derived liver cancers. The minimum depth of

991    coverage was set as 70-fold and 50-fold respectively for mouse and human

992    data. The phylogenetic tree of the 11 tumor nodules was reconstructed via

993    R package ape (v5.2)[48] with the application of the neighbor-joining

994    algorithm[49]. R package fishplot[50] was used to visualizing tumor evolutions

995    of the discrete tumor nodules within the same mouse.

996    **Mutation signature analysis.** Trinucleotide contextualized mutational

997    signature deconvolution was previously described as cocktail party

998    problem[51]. We used the least square root implemented in the R package

999    pracma to decipher the mutational signatures with the known mutational

1000   signatures inferred in the specific cancer type, from "Signatures of

1001   Mutational Processes in Human Cancer" in the COSMIC database, as

1002   recommended in the signature analyzing R package Mutational Patterns

1003   (v1.8.0)[52]. Briefly, the algorithm deciphers the set of mutational signatures

1004   that optimally explains the total trinucleotide frequencies. For the mice

1005  tumor signature deciphering, we adopted this method with a cutoff of 5%

1006  of signature contribution to avoid over fitting

1007  For human cancer signature investigation, we initially used the

1008  webserver mSignatureDB[53] to investigate the COSMIC signature

1009  contributions across 73 research programs over 15, 780 tumors

1010  documented on The Cancer Genome Atlas (TCGA) and the International

1011  Cancer Genome Consortium (ICGC) data portals using the deciphering

1012  method provided on the server. Next, we used the locally adopted least

1013  square root implementation to double check the positively detected cancer

1014  projects. To improve the deconvolution, bootstrap resampling

1015  implemented in the R package Signature Estimation[54] was employed to

1016  calculate the confident interval of signature exposures. In this way, we

1017  conducted 1000 times of randomized re-sampling in order to simulate the

1018  perturbation of the input data. Then we generated estimation of the

1019  exposures of the mutational signatures in each bootstrap sample. From the

1020  continuum of the estimated signature contributions, we retrieved the lower

1021  boundary of the 95% confidence intervals of the bootstrapped AA

1022  signature distribution to obtain a probability of 0.05 for rejection of the

1023  event that the AA signature contribution being above the specific retrieved

1024  threshold. To evaluate the performance of the deciphering procedure, we

1025  used the random sampling and permutation function to generate a

1026  simulated dataset of 1000 samples with known signature exposures and

1027 mutation counts. 156 samples were excluded for containing zero mutation.

1028 The remained 864 samples were used for evaluation of deciphering

1029 methods. There was a work reporting that performing the same method on

1030 the exome resided mutations or the genome resided mutations revealed

1031 different results; the former was stricter for AA signature detection[55]. It is

1032 to some degree due to the overfitting danger when dealing with a large load

1033 of mutations. Therefore, to avoid sequence bias between WES and WGS

1034 generated data, we only retained the exome resided mutations from WGS

1035 generated data for the mutation deconvolution analysis.

1036 To estimate similarities between tumor or clonal mutational profiles,

1037 and the COSMIC signatures, cosine similarities were calculated using the

1038 R package Mutational Patterns (v1.8.0).

1039 **A>T Transcriptional strand bias.** For each tumor, A>T transcriptional

1040 strand bias was analyzed by comparing the number of mutations occurring

1041 on the transcribing and non-transcribing strands over the genome with the

1042 Poisson distribution test. Later on, to correlate AAI mutational processes

1043 with gene transcription history, we categorized the UCSC (University of

1044 California Santa Cruz)[56] known genes into 5 categories from no to high

1045 transcriptional activities in the RNA-seq data of the 11 non-tumor liver

1046 samples and compared the transcriptional strand bias within each defined

1047 gene expression category.

1048 **Mutation assignment to the signatures.** Each mutation was firstly

1049 ascribed to a specific signature via a Bayesian inference method

1050 implemented in the R package Palimpsest [57], which was calculating the

1051 probability of each operative process for a certain mutation and then

1052 choosing the largest as the assigned signature for the specific mutation.

1053 **Driver gene analysis.** To analyze the significant A>T mutated genes in the

1054 mice, we calculated the nonsilent mutational counts per mega base for each

1055 gene to search for the genes that are mutated more frequently. The genes

1056 listed in the duplicated gene database were removed as they are easily to

1057 be falsely detected with mutations[58]. The MutSigCV[59] and

1058 oncodriveCLUST[60] analyses were performed on the human mutation data

1059 for drive gene identification. R package maftools[61] was used to plot the

1060 gene mutation points distribution on the motifs and Palimpsest was used to

1061 calculate the contribution of the operative signatures to the reported driver

1062 genes in HCC[57].

1063 **Statistics.** Statistical analyses were performed using SPSS software. All

1064 the statistical tests used were described in the relevant sections of the

1065 manuscript. $P$-values $< 0.05$ were considered statistically significant.

1066 **Data availability**

1067 The mouse next-generation sequencing data used in the manuscript can be

1068 downloaded from the database of NCBI under accession number: PRJNA

1069 507339.

## Acknowledgements

## Author contributions

Z.-G.H. initiated and supervised the project. Z.-N.L., L.-N.Z. and X.-B.S. performed animal test, other experiments and statistical analysis. Q.L. and Y.S. analyzed the mouse WGS and WES data. Q.L. did the other bioinformatics analysis. Z.-G.H., Z.-N.L, and Q.L. analyzed the data and wrote the manuscript.

## Competing interests

The authors declare no competing interests.

## Supplementary information

Supplementary Figures 1–5 and Supplementary Tables 1–16

# References

1. Yang, H.Y., Chen, P.C. & Wang, J.D. Chinese herbs containing aristolochic acid associated with renal failure and urothelial carcinoma: a review from epidemiologic observations to causal inference. *Biomed Res Int* **2014**, 569325 (2014).

2. Debelle, F.D., Vanherweghem, J.L. & Nortier, J.L. Aristolochic acid nephropathy: a worldwide problem. *Kidney Int* **74**, 158-69 (2008).

3. Grollman, A.P. *et al.* Aristolochic acid and the etiology of endemic (Balkan) nephropathy. *Proc Natl Acad Sci U S A* **104**, 12129-34 (2007).

4. Schmeiser, H.H. *et al.* Exceptionally long-term persistence of DNA adducts formed by carcinogenic aristolochic acid I in renal tissue from patients with aristolochic acid nephropathy. *International Journal of Cancer* **135**, 502-507 (2014).

5. Nortier, J.L. *et al.* Urothelial carcinoma associated with the use of a Chinese herb (Aristolochia fangchi). *N Engl J Med* **342**, 1686-92 (2000).

6. Jelakovic, B. *et al.* Aristolactam-DNA adducts are a biomarker of environmental exposure to aristolochic acid. *Kidney Int* **81**, 559-67 (2012).

7. Hoang, M.L. *et al.* Mutational signature of aristolochic acid exposure as revealed by whole-exome sequencing. *Sci Transl Med* **5**, 197ra102 (2013).

8. Poon, S.L. *et al.* Genome-wide mutational signatures of aristolochic acid and its application as a screening tool. *Sci Transl Med* **5**, 197ra101 (2013).

9. Huang, J. *et al.* Exome sequencing of hepatitis B virus-associated hepatocellular carcinoma. *Nat Genet* **44**, 1117-21 (2012).

10. Ng, A.W.T. *et al.* Aristolochic acids and their derivatives are widely implicated in liver cancers in Taiwan and throughout Asia. *Sci Transl Med* **9**(2017).

11. Arlt, V.M. *et al.* Gene expression changes induced by the human carcinogen aristolochic acid I in renal and hepatic tissue of mice. *Int J Cancer* **128**, 21-32 (2011).

12. Arlt, V.M. *et al.* Aristolochic acid (AA)-DNA adduct as marker of AA exposure and risk factor for AA nephropathy-associated cancer. *Int J Cancer* **111**, 977-80 (2004).

13. Nortier, J.L. *et al.* Invasive urothelial carcinoma after exposure to Chinese herbal medicine containing aristolochic acid may occur without severe renal failure. *Nephrol Dial Transplant* **18**, 426-8 (2003).

14. Schmeiser, H.H. *et al.* Aristolochic acid activates ras genes in rat tumors at deoxyadenosine residues. *Cancer Res* **50**, 5464-9 (1990).

15. Caviglia, J.M. & Schwabe, R.F. Mouse models of liver cancer. *Methods Mol Biol* **1267**, 165-83 (2015).

16. Horie, Y. *et al.* Hepatocyte-specific Pten deficiency results in steatohepatitis and hepatocellular carcinomas. *J Clin Invest* **113**, 1774-83 (2004).

17. Fujimoto, A. *et al.* Whole-genome mutational landscape and characterization of noncoding and structural mutations in liver cancer. *Nat Genet* **48**, 500-9 (2016).

18. Totoki, Y. *et al.* Trans-ancestry mutational landscape of hepatocellular carcinoma genomes. *Nat Genet* **46**, 1267-73 (2014).

19. Rowson-Hodel, A.R. *et al.* Membrane Mucin Muc4 promotes blood cell association with tumor cells and mediates efficient metastasis in a mouse model of breast cancer. *Oncogene* **37**, 197-207 (2018).

20. Barbari, S.R. & Shcherbakova, P.V. Replicative DNA polymerase defects in human cancers: Consequences, mechanisms, and implications for therapy. *DNA Repair (Amst)* **56**, 16-25 (2017).

21. Rayner, E. *et al.* A panoply of errors: polymerase proofreading domain mutations in cancer. *Nat Rev Cancer* **16**, 71-81 (2016).

22. Miller, C.A. *et al.* SciClone: inferring clonal architecture and tracking the spatial and temporal patterns of tumor evolution. *PLoS Comput Biol* **10**, e1003665 (2014).

23. Kan, Z. *et al.* Whole-genome sequencing identifies recurrent mutations in hepatocellular carcinoma. *Genome Res* **23**, 1422-33 (2013).

24. Lin, D.C. *et al.* Genomic and Epigenomic Heterogeneity of Hepatocellular Carcinoma. *Cancer Res* **77**, 2255-2265 (2017).

25. Ahn, S.M. *et al.* Genomic portrait of resectable hepatocellular carcinomas: implications of RB1 and FGF19 aberrations for patient stratification. *Hepatology* **60**, 1972-82 (2014).

26. Zhai, W. *et al.* The spatial organization of intra-tumour heterogeneity and evolutionary trajectories of metastases in hepatocellular carcinoma. *Nat Commun* **8**, 4565 (2017).

27. Zou, S. *et al.* Mutational landscape of intrahepatic cholangiocarcinoma. *Nat Commun* **5**, 5696 (2014).

28. Guo, G. *et al.* Whole-genome and whole-exome sequencing of bladder cancer identifies frequent alterations in genes involved in sister chromatid cohesion and segregation. *Nat Genet* **45**, 1459-63 (2013).

29. Bray, F. *et al.* Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA Cancer J Clin* **68**, 394-424 (2018).

30. Chen, W. *et al.* Cancer incidence and mortality in China, 2014. *Chin J Cancer Res* **30**, 1-12 (2018).

31. Sia, D., Villanueva, A., Friedman, S.L. & Llovet, J.M. Liver Cancer Cell of Origin, Molecular Class, and Effects on Patient Prognosis. *Gastroenterology* **152**, 745-761 (2017).

32. Konca, K. *et al.* A cross-platform public domain PC image-analysis program for the comet assay. *Mutat Res* **534**, 15-20 (2003).

33. Schmeiser, H.H., Frei, E., Wiessler, M. & Stiborova, M. Comparison of DNA adduct formation by aristolochic acids in various in vitro activation systems by 32P-post-labelling: evidence for reductive activation by peroxidases. *Carcinogenesis* **18**, 1055-62 (1997).

34. Yun, B.H. *et al.* Biomonitoring of aristolactam-DNA adducts in human tissues using ultra-performance liquid chromatography/ion-trap mass spectrometry. *Chem Res Toxicol* **25**, 1119-31 (2012).

35. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754-60 (2009).

36. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**, 1297-303 (2010).

37. Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nat Biotechnol* **31**, 213-9 (2013).

38. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res* **38**, e164 (2010).

39. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster strain w1118; iso-2; iso-3. *Fly (Austin)* **6**, 80-92 (2012).

40. Thorvaldsdottir, H., Robinson, J.T. & Mesirov, J.P. Integrative Genomics Viewer (IGV): high-

1175    performance genomics data visualization and exploration. *Brief Bioinform* **14**, 178-92 (2013).

1176    41.    Huang da, W., Sherman, B.T. & Lempicki, R.A. Systematic and integrative analysis of large gene
1177    lists using DAVID bioinformatics resources. *Nat Protoc* **4**, 44-57 (2009).

1178    42.    Trapnell, C., Pachter, L. & Salzberg, S.L. TopHat: discovering splice junctions with RNA-Seq.
1179    *Bioinformatics* **25**, 1105-11 (2009).

1180    43.    Roberts, A., Trapnell, C., Donaghey, J., Rinn, J.L. & Pachter, L. Improving RNA-Seq expression
1181    estimates by correcting for fragment bias. *Genome Biol* **12**, R22 (2011).

1182    44.    Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for
1183    interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**, 15545-50 (2005).

1184    45.    Xie, X. *et al.* Systematic discovery of regulatory motifs in human promoters and 3' UTRs by
1185    comparison of several mammals. *Nature* **434**, 338-45 (2005).

1186    46.    Su, X. *et al.* Single-cell RNA-Seq analysis reveals dynamic trajectories during mouse liver
1187    development. *BMC Genomics* **18**, 946 (2017).

1188    47.    Borggrefe, T. & Oswald, F. The Notch signaling pathway: transcriptional regulation at Notch
1189    target genes. *Cell Mol Life Sci* **66**, 1631-46 (2009).

1190    48.    Popescu, A.A., Huber, K.T. & Paradis, E. ape 3.0: New tools for distance-based phylogenetics
1191    and evolutionary analysis in R. *Bioinformatics* **28**, 1536-7 (2012).

1192    49.    Saitou, N. & Nei, M. The neighbor-joining method: a new method for reconstructing
1193    phylogenetic trees. *Mol Biol Evol* **4**, 406-25 (1987).

1194    50.    Miller, C.A. *et al.* Visualizing tumor evolution with the fishplot package for R. *BMC Genomics*
1195    **17**, 880 (2016).

1196    51.    Alexandrov, L.B., Nik-Zainal, S., Wedge, D.C., Campbell, P.J. & Stratton, M.R. Deciphering
1197    signatures of mutational processes operative in human cancer. *Cell Rep* **3**, 246-59 (2013).

1198    52.    Blokzijl, F., Janssen, R., van Boxtel, R. & Cuppen, E. MutationalPatterns: comprehensive
1199    genome-wide analysis of mutational processes. *Genome Med* **10**, 33 (2018).

1200    53.    Huang, P.J. *et al.* mSignatureDB: a database for deciphering mutational signatures in human
1201    cancers. *Nucleic Acids Res* **46**, D964-D970 (2018).

1202    54.    Huang, X., Wojtowicz, D. & Przytycka, T.M. Detecting presence of mutational signatures in
1203    cancer with confidence. *Bioinformatics* (2017).

1204    55.    Ji, X.J., Feng, G.S., Chen, G. & Shi, T.L. Lack of correlation between aristolochic acid exposure
1205    and hepatocellular carcinoma. *Science China-Life Sciences* **61**, 727-728 (2018).

1206    56.    Haeussler, M. *et al.* The UCSC Genome Browser database: 2019 update. *Nucleic Acids Res*
1207    (2018).

1208    57.    Letouze, E. *et al.* Mutational signatures reveal the dynamic interplay of risk factors and cellular
1209    processes during liver tumorigenesis. *Nat Commun* **8**, 1315 (2017).

1210    58.    Ouedraogo, M. *et al.* The duplicated genes database: identification and functional annotation
1211    of co-localised duplicated genes across genomes. *PLoS One* **7**, e50653 (2012).

1212    59.    Lawrence, M.S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-
1213    associated genes. *Nature* **499**, 214-218 (2013).

1214    60.    Tamborero, D., Gonzalez-Perez, A. & Lopez-Bigas, N. OncodriveCLUST: exploiting the positional
1215    clustering of somatic mutations to identify cancer genes. *Bioinformatics* **29**, 2238-44 (2013).

1216    61.    Mayakonda, A., Lin, D.C., Assenov, Y., Plass, C. & Koeffler, H.P. Maftools: efficient and
1217    comprehensive analysis of somatic variants in cancer. *Genome Res* **28**, 1747-1756 (2018).

1218 **Figure legends**

1219 **Figure 1. AAI can induce liver cancer.** (**a**) Simplified diagram of liver

1220 cancer induction in C57BL/6 male mice with AAI alone or a combination

1221 of AAI and $CCl_4$, where the dosages and time points of drug administration

1222 are indicated by arrows, and samples are harvested at the indicated time of

1223 sacrifice mice (Sac). (**b-d**) Tumor incidence (**b**), tumor number (**c**), and

1224 largest tumor size (**d**) of AAI-induced liver cancer. The numbers in

1225 parentheses are the numbers of mice in the corresponding group. The

1226 numbers of mice in the control groups corresponding to the first four

1227 experimental groups in the figure were 6 (5.5 M), 10 (8.5 M), and 12 (11.5

1228 M); that in the fifth group were 5 ($CCl_4$ (2 M), 8.5 M) and 6 ($CCl_4$ (2 M),

1229 11.5 M); and that in the sixth model were 6 ($CCl_4$ (1 M), 5.5 M) and 6 ($CCl_4$

1230 (1 M), 8.5 M). The asterisk directly above each group indicates a

1231 significant difference compared with the corresponding control group. The

1232 numbers marked with an asterisk indicate the numbers of surviving mice,

1233 and the initial number of mice in each group was 11. (**e, f**) Representative

1234 images of gross appearance (scale bars, 1 cm), hematoxylin and eosin

1235 (H&E) staining (scale bars, 100 μm), and immunohistochemistry (IHC)

1236 analysis with anti-AFP, Ki67 and CK19 antibodies (scale bars, 100 μm) of

1237 AAI-induced HCCs (**e**), cHCC-ICC and ICC (**f**). (**g**) Proportion of types of

1238 HCC, cHCC-ICC and ICC in all C57BL/6 male mice with liver cancers

1239 (up) and each group (down). The numbers in the column charts indicate the

1240   number of mice with the corresponding cancer types in each group. (**h**)

1241   Quantification of PicroSirius Red histochemistry staining in liver slices

1242   from the different experimental groups and control groups. The numbers

1243   of mice corresponding to each group were 12 (time after first

1244   administration, 11.5 M), 10 (11.5 M), 9 (11.5 M), 9 (11.5 M), 5 (8.5 M), 5

1245   (8.5 M), 6 (5.5 M) and 6 (5.5 M). (**b-d, h**) Values indicated by long and

1246   short horizontal lines represent the mean $\pm$ SD. Asterisks signify

1247   significant differences using the two-sided Student's *t*-test or Wilcoxon

1248   rank-sum test and Fisher's exact test. \*$P < 0.05$; \*\*$P < 0.01$; \*\*\*$P < 0.001$;

1249   NS, not significant.

1250   **Figure 2. AAI can cause liver DNA damage.** (**a**) The γ-H2AX level was

1251   measured by immunofluorescence assay in livers from the control group

1252   (PBS) and "AAI (14x)" mice at 1 month of age. (**b**) γ-H2AX, p53 and Bax

1253   levels were measured by Western blotting assay (n = 3) in mouse livers

1254   from the control group (PBS) and "AAI (14x)" group. The numbers under

1255   the Bax band are relative intensity values. (**c, d**) The DNA strand breaks

1256   were measured by the alkaline comet assay in liver cells from 2-week-old

1257   mice at 3 h once after PBS or AAI (2.5 mg/kg or 5 mg/kg) injection,

1258   including representative images (**c**) and quantitative analyses (**d**, the

1259   number of mice per group is 4; the number of nuclei per group: 600, at least

1260   100 nuclei from each mouse). The white dot, thick black bar in the center

1261   and thin black line extended from it in (**d**) stand for the median,

1262     interquartile range and 95% confidence intervals. (**e, f**) The γ-H2AX level

1263     was measured by the immunofluorescence assay in livers from the above

1264     samples, including representative images (**e**) and quantitative analyses (**f**,

1265     the number of mice per group: 3; 10 nonoverlapping fields at $\times$400

1266     magnification per mouse). Values indicated by long and short horizontal

1267     lines represent the mean ± SD. (**g, h**) The relative abundance (*m/z* 427) of

1268     dA-AL-I was measured by MS in livers of "AAI (14x)" mice at the

1269     indicated ages, including representative images (**g**) and quantitative

1270     analyses (**h**). Spots with the same colors indicate the paratumors and

1271     tumors are from the same mouse. The horizontal lines in (**h**) denote the

1272     mean. Asterisks signify significant differences using the two-sided

1273     Student's *t*-test or Wilcoxon rank-sum test. ***$P < 0.001$. Scale bars, 100

1274     μm.

1275     **Figure 3. Mutational signatures of AAI-induced mouse liver cancer.** (**a**)

1276     Representative trinucleotide contextualized mutational spectra in AAI-

1277     induced liver cancer. (**b**) Estimated COSMIC mutational signature

1278     contributions for each mouse liver cancer. Mutational signature

1279     decomposing with known liver cancer signatures (i.e., COSMIC signatures

1280     1, 4, 5, 6, 12, 16, 17, 22, 23 and 24) was performed using the least square

1281     root algorithm. The AA signature (or COSMIC signature 22) was dominant

1282     throughout the mouse liver cancer. (**c**) Cosine similarities of trinucleotide

1283     mutational spectra between the tumor/noncancerous liver samples and the

1284  AA signature. M1T, M2T, M3T, M4T1, M4T2, M4T3, M5T1, M5T2,

1285  M5T3, M5T4 and M6T refer to AAI-induced liver cancer; M4P, M5P and

1286  M6P are the paratumor liver tissues of the mice in the combination group;

1287  M7L, M8L and M9L are the livers from the "AAI (3x)" group. M10L and

1288  M11L are the liver tissues from the $CCl_4$ -treated group. (**d**) The mutational

1289  frequency of A>T transversions in transcribed and nontranscribed regions

1290  per megabase (Mb) in these genes as a function of the expression level.

1291  The genes with expression were divided into 4 expression quintiles

1292  according to the expression levels. NT, nontranscribed strand; Tr,

1293  transcribed strand.

1294  **Figure 4. The genes and signaling pathways affected by AAI-mediated**

1295  **mutations.** (**a**) The categories of the statistically and empirically important

1296  genes with somatic mutations in liver cancer. The genes in red or blue refer

1297  to the proto-oncogenes and tumor-suppressor genes listed in COSMIC

1298  Cancer Gene Census Tier 1, respectively. Genes in purple font refer to

1299  driver genes without a clear definition in terms of proto-oncogenes and

1300  tumor-suppressor genes. (**b**) Major signaling pathways involving genetic

1301  alterations in AAI-induced mouse liver cancer. A brown background

1302  indicates mutated genes; a white background denotes unmutated genes.

1303  Genes in red and blue refer to proto-oncogenes and tumor-suppressor genes

1304  listed in COSMIC Cancer Gene Census Tier 1, respectively. Those in

1305  purple refer to driver genes without a clear definition in terms of proto-

1306  oncogenes and tumor-suppressor genes. Percentages stand for the

1307  proportion of gene or genes in the pathway altered in liver cancer. (**c**)

1308  Correlations between the number of DNA repair-related genes with the

1309  total mutation counts in mouse liver tumors. (**d**) Gene set enrichment

1310  analysis (GSEA) plot of SRF and YAP motif target gene sets. SRF is a key

1311  regulatory transcription factor in the Ras signaling pathway, and YAP is a

1312  key regulatory transcription factor in the Hippo signaling pathway. (**e**) AFP,

1313  GPC3, E-cadherin, p-ERK, ERK1, p-AKT, AKT and YAP levels were

1314  measured by Western blotting assay (n = 3) in mouse paratumors ("P") and

1315  tumors ("T") from "AAI (3x) + CCl$_4$" (18 M) group.

1316  **Figure 5. Clonal architecture and phylogenetic reconstructions of**

1317  **AAI-induced mouse liver cancer.** (**a, b, left**) Malignant clonal

1318  architecture reconstructions within M3T1 and M5T2 tumors. Each peak

1319  indicates one subclone. The subclones lying at the right end had the largest

1320  mutational allele frequency and therefore represent the founding clones.

1321  Others are subclones. (**a, b, right**) Trinucleotide mutational spectra of the

1322  founding clone and subclones within M3T1 and M5T2. The A>T

1323  transversions were predominant in the earliest founding clones and

1324  diminished in the later formed subclones. (**c**) The downward-pointing lines

1325  of A>T mutations in different subclones within these 11 tumor nodules.

1326  The A>T mutation proportions deposited in each clone in the multiclonal

1327  tumors. (**d**) The phylogenetic tree was reconstructed using the neighbor-

1328   joining algorithm with the R package ape (v5.2). (**e**) A common ancestry

1329   evolutionary model for three discrete tumor nodules within M4 liver. The

1330   font size of the genes reflects the allele frequency in each tumor.

1331   (**f**) Common ancestry evolutionary model for three discrete tumor nodules

1332   within M5 liver. The font size reflects the genes allele frequency in each

1333   tumor.

1334   **Figure 6. AA signatures in human liver cancer and other cancers.** (**a**)

1335   The AA signature contribution in each individual of the affected human

1336   cancer types. The figure is depicted as a violin plot, in which each dot

1337   represents one human tumor. (**b**) COSMIC signature contribution

1338   according to mutational counts (upper) and proportions (bottom) in the

1339   selected human HCCs with an AA signature proportion larger than 52%.

1340   (**c**) The curves of the A>T mutation proportions in different subclones

1341   within 6 TCGA-derived human HCCs. The A>T mutations were deposited

1342   in the founding clones. Tumor DD.AAC8 presents a typical pattern,

1343   implying that AA led to tumor initiation and diminished in the later

1344   processes, as in the mouse tumors. The upward-pointing lines of the AA

1345   contributions across the clonal evolution in the other 5 tumors indicate

1346   sustained AA exposures in the patients. (**d**) Cosmic signature contributions

1347   to the essential driver genes reported in HCCs. (**e**) The AA signature caused

1348   A>T mutation sites in *TP53* and *JAK1* in the selected human HCCs with
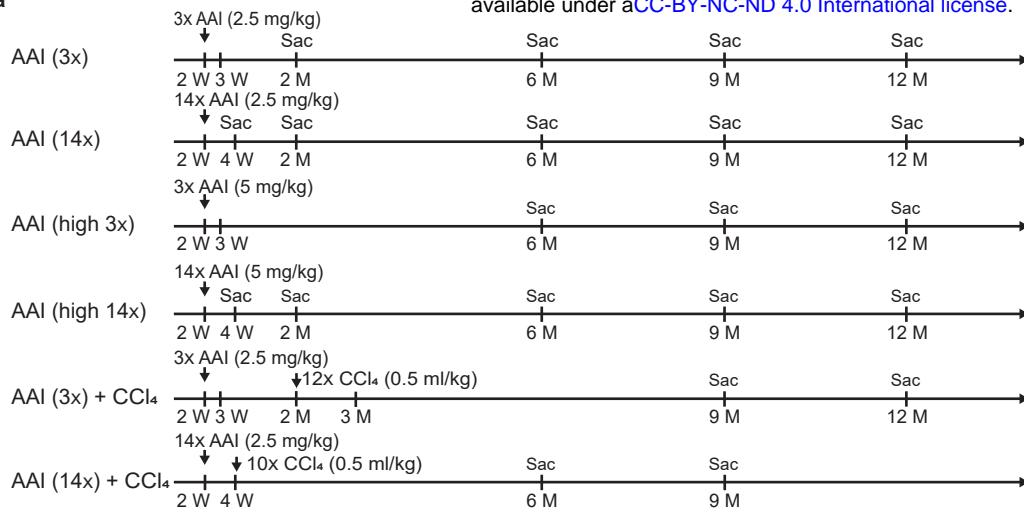
1349   an AA signature proportion larger than 52%.

**Table 1. Aristolochic acid exposure in human liver cancers**

| Cancer type | Data source | Regions | Number of patients | Numbers (%) of patients with AA exposure inferred by different criteria | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | > 0* | > 10** | Nonsilent mutations in known driver genes | > 52%*** |
| HCC | ICGC | China (mainland) | 313 | 52 (17) | 51 (17) | 27 (9) | 12 (4) |
| | Kan, et al. [23] | China (Hongkong) | 88 | 7 (8) | 4 (5) | 1 (1) | 1 (1) |
| | Lin, et al. [24] | China (mainland) | 11 | 6 (55) | 6 (55) | 2 (18) | 2 (18) |
| | Ng, et al. [10] | Taiwan | 98 | 68 (69) | 68 (69) | 53 (54) | 42 (43) |
| | Chinese HCC total | | 510 | 133 (26) | 129 (25) | 83 (16) | 57 (11) |
| | TCGA | USA (Asian) | 160 | 24 (15) | 23 (14) | 13 (8) | 4 (3) |

|  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|
|  | TCGA | USA (others) | 204 | 5 (2) | 4 (2) | 3 (1) | 0 (0) |
|  | ICGC | France | 249 | 1 (<1) | 0 (0) | 0 (0) | 0 (0) |
|  | ICGC[17,18] | Japan | 594 | 3 (<1) | 1(<1) | 2 (<1) | 0 (0) |
|  | Zhai, *et al.*[26] | Singapore | 9 | 4 (44) | 4 (44) | 1 (11) | 0 (0) |
|  | Ahn, *et al.*[25] | Korea | 231 | 22 (10) | 22 (10) | 11 (5) | 3 (1) |
|  | Worldwide HCC total |  | 1957 | 192 (10) | 183 (9) | 113 (6) | 64 (3) |
| ICC | Zou, *et al.*[27] | China (mainland) | 103 | 11 (11) | 11 (11) | 4 (4) | 3 (3) |

Note: * indicates the estimated lower boundary of 95% confidence interval of AA signature exposure lager than 0 ($P < 0.05$); ** indicates the estimated lower boundary of 95% confidence interval of AA signature exposure lager than 10% ($P < 0.05$); *** indicates the estimated lower boundary of 95% confidence interval of AA signature exposure lager than 52% ($P < 0.05$).
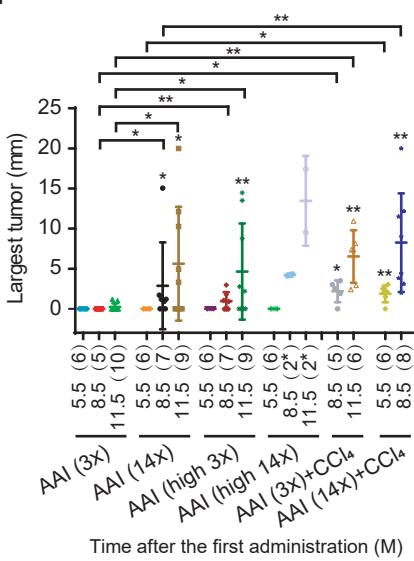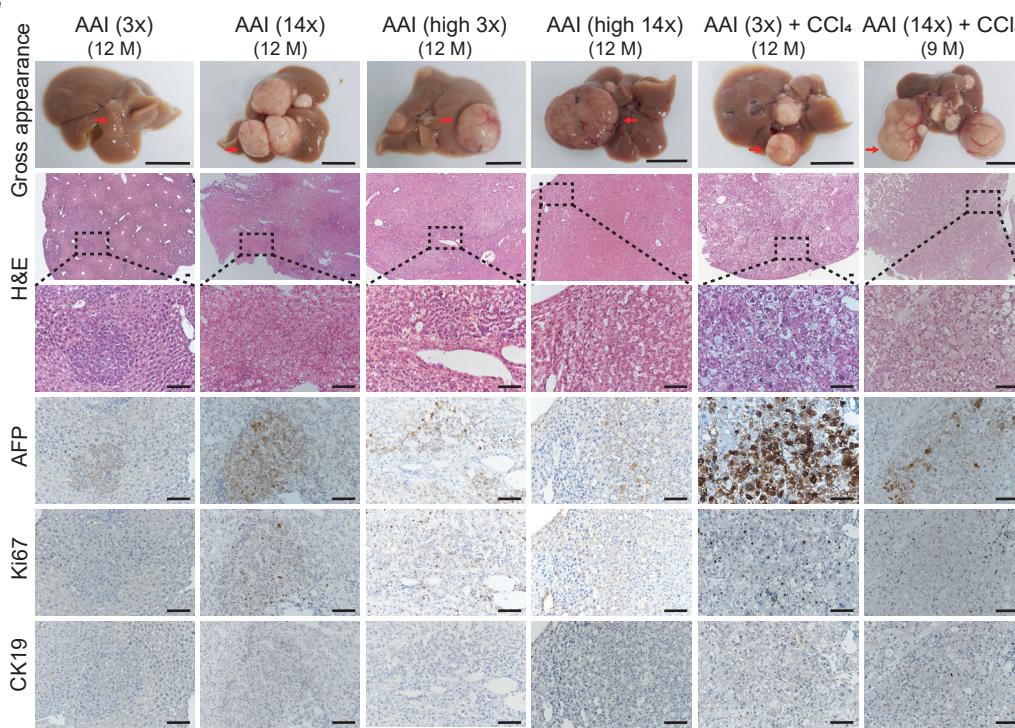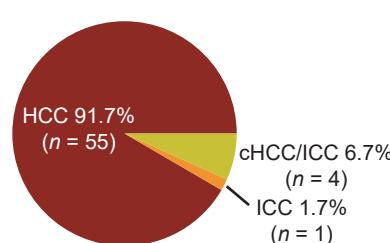
**a** "AAI (14x)" (M3T)

**b** "AAI (3x) + CCl₄" (M5T2)