

1 An Integrative Analysis of the Age-Associated Genomic, Transcriptomic and Epigenetic 2 Landscape across Cancers

3 Kasit Chatsirisupachai¹, Tom Lesluyes², Luminita Paraoan³, Peter Van Loo², João Pedro de
4 Magalhães^{1*}

5 ¹Integrative Genomics of Ageing Group, Institute of Life Course and Medical Sciences,
6 University of Liverpool, Liverpool L7 8TX, UK.

7 ²The Francis Crick Institute, London NW1 1AT, UK.

8 ³Department of Eye and Vision Science, Institute of Life Course and Medical Sciences,
9 University of Liverpool, Liverpool L7 8TX, UK.

10 *email: jp@senescence.info

11 Abstract

12 Age is the most important risk factor for cancer, as cancer incidence and mortality
13 increase with age. However, how molecular alterations in tumours differ among patients of
14 different age remains largely unexplored. Here, using data from The Cancer Genome Atlas, we
15 comprehensively characterised genomic, transcriptomic and epigenetic alterations in relation
16 to patients' age across cancer types. We showed that tumours from older patients present an
17 overall increase in genomic instability, somatic copy-number alterations (SCNAs) and somatic
18 mutations. Age-associated SCNAs and mutations were identified in several cancer-driver
19 genes across different cancer types. The largest age-related genomic differences were found in
20 gliomas and endometrial cancer. We identified age-related global transcriptomic changes and
21 demonstrated that these genes are controlled by age-associated DNA methylation changes. This
22 study provides a comprehensive view of age-associated alterations in cancer and underscores
23 age as an important factor to consider in cancer research and clinical practice.

24

25 **Keywords:** ageing, carcinoma, brain cancer, geriatric oncology, single nucleotide variants

26 **Introduction**

27 Age is the biggest risk factor for cancer, as cancer incident and mortality rates increase
28 exponentially with age in most cancer types¹. However, the relationship between ageing and
29 molecular determinants of cancer remains to be characterised. Cancer arises through the
30 interplay between somatic mutations and selection, in a Darwinian-like process^{2,3}. Thus, apart
31 from the mutation accumulation with age⁴⁻⁶, microenvironment changes during ageing could
32 also play a role in carcinogenesis^{2,7,8}. We therefore hypothesise that, due to the differences in
33 selective pressures from tissue environmental changes with age, tumours arise from patients
34 across different ages might harbour different molecular landscapes; consequently, some
35 molecular changes might be more or less common in older or younger patients.

36 Recently, several studies have investigated the molecular differences in the cancer
37 genome in relation to clinical factors, including gender^{9,10} and race^{11,12}. These studies
38 demonstrated gender- and race-specific biomarkers, actionable target genes and provided clues
39 to understanding the biology behind the disparities in cancer incidence, aggressiveness and
40 treatment outcome across patients from different backgrounds. Although the genomic
41 alterations in childhood cancers and the differences with adult cancers have been systematically
42 characterised^{13,14}, the age-related genomic landscape across adult cancers remains elusive.
43 Specific age-associated molecular landscapes have been reported in the cancer genome of
44 several cancer types, for example, glioblastoma¹⁵, prostate cancer¹⁶ and breast cancer¹⁷.
45 However, these studies focused mainly on a single cancer type and only on some molecular
46 data types.

47 Here, using data from The Cancer Genome Atlas (TCGA), we systematically
48 investigated age-related differences in genomic instability, somatic copy number alterations
49 (SCNAs), somatic mutations, pathway alterations, gene expression, and DNA methylation
50 landscape across various cancer types. We show that, in general, genomic instability and

51 mutations frequency increase with age. We identify several age-associated genomic alterations
52 in cancers, particularly in low-grade glioma and endometrial carcinoma. Moreover, we also
53 demonstrate that age-related gene expression changes are controlled by age-related DNA
54 methylation changes and that these changes are linked to numerous biological processes.

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76 **Results**

77 **Association between age and genomic instability, loss of heterozygosity, and whole-** 78 **genome duplication**

79 To gain insight into the role of patient age into the somatic genetic profile of tumours,
80 we evaluated associations between patient age and genomic features of tumours in TCGA data
81 (Table 1, Supplementary Table 1). Using multiple linear regression adjusting for gender, race,
82 and cancer type, we found that genomic instability (GI) scores increase with age in pan-cancer
83 data (adj. R-squared = 0.35, p-value = 5.98×10^{-7}) (Fig. 1a). We next applied simple linear
84 regression to investigate the relationship between GI scores and age for each cancer type.
85 Cancer types with a significant association (adj. p-value < 0.05) were further adjusted for
86 clinical variables. We found a significant positive association between age and GI score in
87 seven cancer types (adj. p-value < 0.05) (Fig. 1b, Supplementary Fig. 1a and Supplementary
88 Table 2). Cancer types with the strongest significant positive association were low-grade
89 glioma, ovarian cancer, endometrial cancer, and sarcoma. This result indicates that the level of
90 genomic instability increases with the age of cancer patients in several cancer types.

91 The genomic loss of heterozygosity (LOH) refers to the irreversible loss of one parental
92 allele, causing an allelic imbalance, and priming the cell for another defect at the other
93 remaining allele of the respective genes¹⁸. To investigate whether there is an association
94 between patients' age and LOH, we quantified percent genomic LOH. By using simple linear
95 regression, we found a significant positive association between age and pan-cancer percent
96 genomic LOH (p-value = 1.20×10^{-21}). However, this association was no longer significant in
97 a multiple linear regression analysis (adj. R-squared = 0.32, p-value = 0.289) (Fig. 1c). Thus,
98 it is likely that this association might be cancer type-specific. We then performed a linear
99 regression between age and percent genomic LOH for each cancer type. Six cancer types
100 showed a positive association between age and percent genomic LOH (adj. p-value < 0.05)

101 (Fig. 1d, Supplementary Fig. 1b, and Supplementary Table 3). The strongest positive
102 associations were found in low-grade glioma and endometrial cancer (adj. p-value < 0.05),
103 corroborate with the increase in GI score with age. On the other hand, lung adenocarcinoma,
104 oesophageal and liver cancer demonstrated a negative correlation between percent genomic
105 LOH and age (adj. p-value < 0.05).

106 Whole-genome duplication (WGD) is important in increasing the adaptive potential of
107 the tumour and has been linked with a poor prognosis¹⁹⁻²¹. We investigated the relationship
108 between age and WGD using logistic regression. For the pan-cancer analysis, we found an
109 increase in the probability that WGD occurs with age, using multiple logistic regression
110 accounting for gender, race, and cancer type (odds ratio per year (OR) = 1.0066, 95%
111 confidence interval (CI) = 1.0030-1.0103, p-value = 3.84×10^{-4}) (Fig. 1e). For the cancer-
112 specific analysis, a significant positive association was found in ovarian and endometrial
113 cancer (adj. p-value < 0.05, OR = 1.0320 and 1.0248, 95%CI = 1.0151-1.0496 and 1.0024-
114 1.0483, respectively) (Fig. 1e and Supplementary Table 4), indicating that tumours from older
115 patients are more likely to have doubled their genome. Taken together, the findings indicate
116 that tumours from patients with an increased age tend to harbour a more unstable genome and
117 a higher level of LOH in several cancer types. Notably, the strongest association between age
118 and an increase in genome instability, LOH, and WGD was evident in endometrial cancer,
119 suggesting the potential disparities in cancer genome landscape with age in this cancer type.

120

121 **Age-associated somatic copy-number alterations**

122 We used GISTIC2.0 to identify recurrently altered focal- and arm-level SCNAs²². We
123 calculated the SCNA score, as a representation of the level of SCNA occurring in a tumour^{12,23}.
124 For each tumour, the SCNA score was calculated at three different levels: focal-, arm- and
125 chromosome-level, and the overall score calculated from the sum of all three levels. We used

126 simple linear regression to identify the association between age and overall SCNA scores.
127 Cancer types that displayed a significant association were further adjusted for clinical
128 variables. Consistent with the GI score results described above, the strongest positive
129 association between age and overall SCNA scores was found in low-grade glioma, ovarian and
130 endometrial cancers. Other cancer types for which a positive association between age and
131 overall SCNA score was observed were thyroid cancer and clear cell renal cell carcinoma (adj.
132 p-value < 0.05). On the other hand, lung adenocarcinoma is the only cancer type exhibiting a
133 negative association between overall SCNA score and age (Fig. 2a, Supplementary Fig. 2a,
134 and Supplementary Table 5). The different SCNA classes (focal- and chromosome/arm-level)
135 may arise through different biological mechanisms^{12,21}, therefore we separately analysed the
136 association between age and focal- and chromosome/arm-level SCNA scores. Most cancers
137 that showed a significant relationship between age and overall SCNA score also had an
138 association between age and both chromosome/arm-level and focal-level SCNA scores (Fig.
139 2b-c, Supplementary Fig. 2b-c, and Supplementary Table 5). The only exception was in
140 cervical cancer, with a significant association between age and chromosome/arm-level but not
141 with focal-level and overall SCNA scores.

142 We next identified the chromosomal arms that tend to be gained and lost more often
143 with age, for 25 cancer types with sufficient samples (at least 100 tumours, Table 1). We
144 conducted the logistic regression on the significant recurrently gained and lost arms that were
145 identified by GISTIC2.0 for each cancer type. The significant association between age and
146 chromosomal arm gains and losses are shown in Fig. Fig. 2d, e, respectively (adj. p-value <
147 0.05) (Supplementary Fig. 3, Supplementary Table 6). The gain of chromosome 7p, 7q, 20p,
148 and 20q significantly increased with age in several cancer types including two types of gliomas,
149 low-grade glioma and glioblastoma. On the other hand, the gain of chromosome 10p decreased
150 with increased age in gliomas (Fig. 2d and 2f). For the arm losses, there was an increased

151 occurrence of loss in 11 arms with advanced age in endometrial cancer (Fig. 2e and 2g),
152 consistent with a higher genomic instability and LOH with age in this cancer type. Low-grade
153 glioma and ovarian cancer, two other cancer types for which we found the highest significant
154 association between age and SCNA scores, also exhibited a significant increase or decrease in
155 losses with age in multiple arms (Fig. 2e-f, Supplementary Fig. 3). We also observed that the
156 losses of chromosome 10p and 10q increased with age in gliomas. Recurrent losses of
157 chromosome 10 together with the gain of chromosome 7 are important features in IDH-wild-
158 type (IDH-WT) gliomas²⁴. This type of gliomas was more common in older patients, whereas
159 IDH-mutated gliomas were predominantly found in younger patients.

160 We further examined age-associated recurrent focal-level SCNAs. Applying a similar
161 logistic regression, we identified recurrent focal SCNAs associated with the age of the patients
162 for each cancer type. In total, we found 113 significant age-associated regions, including 67
163 gain regions across 10 cancer types and 46 loss regions across 9 cancer types (adj. p-value <
164 0.05) (Fig. 3a, Supplementary Table 7). In accordance with the arm-level result, the highest
165 number of significant regions was found in endometrial cancer (23 gain and 25 loss regions),
166 followed by ovarian cancer (13 gain 2 loss regions) and low-grade glioma (9 gain and 5 loss
167 regions) (Fig. 3b-c, Supplementary Fig. 4).

168 To further investigate the impact of these SCNAs, we studied the correlation between
169 the SCNA level and gene expression for tumours that have both types of data using Pearson
170 correlation. In total, 81 genes in the list of previously identified cancer driver genes
171 (Supplementary Table 8) were presented in at least one significant age-associated focal region
172 in at least one cancer type and showed a significant correlation between SCNAs and gene
173 expression (adj. p-value < 0.05) (Fig. 3d). For example, regions showing an increased gain
174 with age in endometrial cancer included 1q22, where the gene *RIT1* is located in (OR = 1.0355,
175 95%CI = 1.0151-1.0571, adj. p-value = 0.0018) (Fig. 3c, e). The Ras-related GTPases *RIT1*

176 gene has been reported to be highly amplified and correlated with poor survival in endometrial
177 cancer²⁵. Therefore, an increase in the gain of the *RIT1* gene with age might relate to a poor
178 prognosis in older patients. The 16p13.3 loss increased in older endometrial cancer patients
179 (OR = 1.0335, 95%CI = 1.0048-1.0640, adj. p-value = 0.0328). This region contains the p53
180 coactivator gene *CREBBP*. The gain of 8q24.21 harbouring the oncogene *MYC* decreased with
181 patient age in low-grade glioma (OR = 0.9737, 95%CI = 0.9541-0.9927, adj. p-value = 0.0128)
182 and ovarian cancer (OR = 0.9729, 95%CI = 0.9553-0.9904, adj. p-value = 0.0063) (Fig. 3d, e).
183 In addition, in low-grade glioma, we found an increase in 9p21.3 loss with age (OR = 1.0332,
184 95%CI = 1.0174-1.0496, adj. p-value = 0.00017). This region contains the cell cycle-regulator
185 genes *CNKN2A* and *CDKN2B* (Fig. 3b, d, e). The full list of age-associated focal regions across
186 cancer types and the correlation between SCNA status and gene expression can be found in
187 Supplementary Table 7. Taken together, our analysis demonstrates the association between age
188 and SCNAs level across cancer types. We also identified age-associated arms and focal-
189 regions, and these regions harboured several cancer-driver genes. Our results suggest a possible
190 contribution of different SCNA events in cancer initiation and progression of patients with
191 different ages.

192

193 **Age-associated somatic mutations in cancer**

194 The increase in the mutational burden with age is well-established⁴⁻⁶. This age-related
195 mutation accumulation is largely explained by a clock-like mutational process, the spontaneous
196 deamination of 5-methylcytosine to thymine⁵. As expected, we confirmed the correlation
197 between age and mutation load (somatic non-silent SNVs and indels) in the pan-cancer cohort
198 using multiple linear regression adjusting for gender, race, and cancer type (adj. R-squared =
199 0.53, p-value = 1.41×10^{-37}) (Supplementary Fig. 5a). For cancer-specific analysis, 18 cancer
200 types exhibited a significant relationship between age and mutation load using linear regression

201 (adj. p-value < 0.05) (Supplementary Fig. 5, Supplementary Table 9). Only endometrial cancer
202 showed a negative correlation between mutational burden and age. We observed a high
203 proportion of hypermutated tumours (> 1,000 non-silent mutations per exome) from younger
204 endometrial cancer patients. Thirteen out of 38 tumours (34%) from the younger patients (age
205 ≤ 50) were hypermutated tumours, while there were only 42 hypermutated tumours from 383
206 tumours from older patients (11%) (Fisher's exact, p-value = 0.0003) (Fig. 4a). Microsatellite
207 instability (MSI) is a unique molecular alteration caused by defects in DNA mismatch
208 repair^{26,27}. The MSI-high (MSI-H) tumours occur as a subset of high mutational burden
209 tumours²⁸. We investigated whether high mutation loads in endometrial cancer from young
210 patients were due to the presence of MSI-H tumours. Using multiple logistic regression, we
211 found that MSI-H tumours were associated with younger endometrial cancer (OR = 0.9751,
212 95%CI = 0.9531-0.9971, p-value = 0.0264) (Fig. 4b). Another source of hypermutation in
213 cancer is the defective DNA polymerase proofreading ability by mutations in polymerase ϵ
214 (*POLE*) or polymerase δ (*POLD1*) genes^{29,30}. We showed that mutations in *POLE* (OR =
215 0.9690, 95%CI = 0.9422-0.9959, p-value = 0.0243) and *POLD1* (OR = 0.9573, 95%CI =
216 0.9223-0.9925, p-value = 0.0177) were both more prevalent in younger endometrial cancer
217 patients (Fig. 4c). Therefore, the negative correlation between age and mutation loads in
218 endometrial cancer could be explained by the presence of hypermutated tumours in younger
219 patients, which are associated with MSI-H and *POLE/POLD1* mutations. Previous studies on
220 *POLE* and MSI-H subtypes in hypermutated endometrial tumours revealed that these subtypes
221 associated with a better prognosis when compared with the copy-number high subtype³¹⁻³³.
222 Together with our SCNA results, younger UCEC patients are likely to associate with a *POLE*
223 and MSI-H subtypes, high mutation rate and better survival, whilst tumours from older patients
224 are characterized by high SCNAs and are generally associated with a worse prognosis. We
225 extended the age and MSI-H analysis to other cancer types known to have a high prevalence

226 of MSI-H tumours, including colon, rectal, and stomach cancers²⁶. Only in stomach cancer we
227 found an association between older age and the presence of MSI-H tumours (OR = 1.0392,
228 95%CI = 1.0091-1.0720, p-value = 0.01, Supplementary Fig. 6a). When we further examined
229 the association between age and mutations in *POLE* and *POLD1* in other cancers apart from
230 endometrial cancer, no significant association was observed (Supplementary Fig. 6b).

231 Although the increase in mutation load with age in cancer is well studied^{4,28}, the bias
232 of mutation in particular genes with age across cancer types is largely unclear. To better
233 understand this, we conducted logistic regression to investigate genes that are more or less
234 likely to be mutated with an increased age. To prevent the potential bias caused by
235 hypermutated tumours, we restricted the analysis to samples with < 1,000 non-silent mutations
236 per exome (Table 1). We first investigated the association between age and pan-cancer gene-
237 level mutations. Using multiple logistic regression correcting for gender, race, and cancer type,
238 mutations in *IDHI* (OR = 0.9619, 95%CI = 0.9510-0.9730, adj. p-value = 4.18×10^{-10}) and
239 *ATRX* (OR = 0.9803, 95%CI = 0.9724-0.9881, adj. p-value = 9.85×10^{-6}) showed a negative
240 association with age. On the other hand, mutations in *PIK3CA* were more common in older
241 individuals (OR = 1.0082, 95%CI = 1.0022-1.0143, adj. p-value = 4.18×10^{-10}) (Fig. 4d). We
242 next identified genes in which mutations associated with age in a cancer-specific manner in 24
243 cancers with at least 100 samples (Table 1). Using logistic regression, we identified 35
244 mutations from 13 cancers that increased or decreased with the patients' age (adj. p-value <
245 0.05) (Fig. 4e-f, Supplementary Fig. 7 and Supplementary Table 10). The most striking
246 negative associations between mutations and age in low-grade glioma and glioblastoma were
247 found in *IDHI* (OR = 0.9509 and 0.8962, 95%CI = 0.9328-0.9686 and 0.8598-0.9291, adj. p-
248 value = 4.33×10^{-7} and 1.88×10^{-9} , respectively), *ATRX* (OR = 0.9471 and 0.9120, 95%CI =
249 0.9310-0.9628 and 0.8913-0.9466, adj. p-value = 1.75×10^{-10} and 2.45×10^{-8} , respectively), and
250 *TP53* (OR = 0.9431 and 0.9736, 95%CI = 0.9274-0.9582 and 0.9564-0.9905, adj. p-value =

251 1.13×10^{-12} and , respectively). Our observation was consistent with the fact that the median age
252 of IDH-mutants is younger than IDH-WT gliomas. Patients carrying the *IDHI* mutation
253 generally had longer survival than those with IDH-WT³⁴. Previous studies also reported that
254 *IDHI* mutations often co-occurred with *ATRX* and *TP53* mutations, and mutations in these
255 three genes were more prevalent in gliomas without *EGFR* mutations^{15,35}. Indeed, we found
256 that *EGFR* mutations were more common in older low-grade glioma patients (OR = 1.0865,
257 95%CI = 1.0525-1.1258, adj. p-value = 4.35×10^{-7}) (Fig. 4f). Moreover, our SCNA analysis
258 revealed an increase in the gain of *EGFR* with age in low-grade glioma but not in glioblastoma
259 (Fig. 3d), suggesting the difference in age-associated genomic landscape between the two
260 glioma types. Together with the SCNA results, gliomas from younger patients are associated
261 with *IDHI*, *ATRX*, and *TP53* mutations, lower SCNAs, and longer survival. In contrast,
262 gliomas from older patients were more likely to be IDH-WT with *EGFR* mutations,
263 chromosome 7 gain and 10 loss, *CDKN2A* deletion and worse prognosis.

264 Mutations in cancer driver genes showed a positive or negative association with age
265 depending on cancer types. For instance, *PTEN* mutations decreased with patient's age in colon
266 (OR = 0.9347, 95%CI = 0.8935-0.9738, adj. p-value = 0.0029) and endometrial cancers (OR
267 = 0.9586, 95%CI = 0.9331-0.9840, adj. p-value = 0.0033) but increased with age in cervical
268 cancer (OR = 1.0550, 95%CI = 1.0174-1.0959, adj. p-value = 0.0067). *CDHI* mutations were
269 more frequent in younger stomach cancer patients (OR = 0.9414, 95%CI = 0.9027-0.9800, adj.
270 p-value = 0.0061) but more common in older breast cancer patients (OR = 1.0218, 95%CI =
271 1.0049-1.0392, adj. p-value = 0.0171). These results highlight cancer-specific patterns of
272 genomic alterations in relation to age. Overall, our results demonstrate that non-silent
273 mutations in cancer driver genes were not uniformly distributed across ages and we have
274 comprehensively identified, based on data available at present, genes that show age-associated

275 mutation patterns. These patterns might point out age-associated disparities in carcinogenesis,
276 molecular subtypes and survival outcome.

277

278 **Age-associated alterations in oncogenic signalling pathways**

279 As we have identified numerous age-associated alterations in cancer driver genes in
280 both SCNA and somatic mutation levels, we asked if the age-associated patterns also exist in
281 particular oncogenic signalling pathways. We used the data from a previous TCGA study,
282 which had comprehensively characterized 10 highly altered signalling pathways in cancers³⁶.
283 To make the subsequent analysis comparable to previous analyses, we restricted the analysis
284 to samples that were used in our previous analyses, yielding 8,055 samples across 33 cancer
285 types (Table 1). Using logistic regression adjusting for gender, race and cancer type, we
286 identified five out of 10 signalling pathways that showed a positive association with age (adj.
287 p-value < 0.05), indicating that the genes in these pathways are altered more frequently in older
288 patients, concordant with the increase in overall mutations and SCNAs with age (Fig. 5a,
289 Supplementary Table 11). The strongest association was found in cell cycle (OR = 1.0122,
290 95%CI = 1.0076-1.0168, adj. p-value = 1.40×10^{-6}) and Wnt signalling (OR = 1.0122, 95%CI
291 = 1.0073-1.0172, adj. p-value = 6.39×10^{-6}). We next applied logistic regression to investigate
292 the cancer-specific association between age and oncogenic signalling alterations for cancer
293 types that contained at least 100 samples. In total, we identified 28 significant associations
294 across 15 cancer types (adj. p-value < 0.05) (Fig. 5b, Supplementary Table 11). Alterations in
295 Hippo and TP53 signalling pathways significantly associated with age, both positively and
296 negatively, in five cancer types. Consistent with a pan-cancer analysis, cell cycle, Notch and
297 Wnt signalling each showed an increase in alterations with age in three cancer types. We found
298 that alterations in cell cycle pathway increased with age in low-grade glioma (OR = 1.0313,
299 95%CI = 1.0161-1.0467, adj. p-value = 0.00035). This was largely explained by the increase

300 in *CDKN2A* and *CDKN2B* deletions with age as well as epigenetic silencing of *CDKN2A* in
301 older patients (Fig. 5c). On the other hand, TP53 pathway alteration was more pronounced in
302 younger patients (OR = 0.9520, 95%CI = 0.9372-0.9670, adj. p-value = 2.63×10^{-8}), due to the
303 mutations in the *TP53* gene (Fig. 5c). In endometrial cancer, two pathways – Hippo (OR =
304 0.9681, 95%CI = 0.9459-0.9908, adj. p-value = 0.0126) and Wnt (OR = 0.9741, 95%CI =
305 0.9541-0.9946, adj. p-value = 0.0240) - showed a negative association with age, that may be
306 explained by the presence of hypermutated tumours in younger patients. Collectively, we
307 reported pathway alterations in relation to age in several cancer types, highlighting differences
308 in oncogenic pathways that might be important in cancer initiation and progression in an age-
309 related manner.

310

311 **Age-associated gene expression and DNA methylation changes**

312 Apart from the genomic differences with age, we investigated age-associated
313 transcriptomic and epigenetic changes across cancers. We separately performed multiple linear
314 regression analyses on gene expression data and methylation data of 24 cancer types that
315 contained at least 100 samples in both types of data (Table 1). We noticed that, across all genes,
316 the regression coefficient of age on gene expression negatively correlated with the regression
317 coefficient of age on methylation in all cancer types (Supplementary Fig. 8), suggesting that
318 the global changes of gene expression and methylation with age are in the opposite direction.
319 This supports the established role of DNA methylation in suppressing gene expression.
320 Numbers of significant differentially expressed genes with age (age-DEGs) (adj. p-value <
321 0.05, Supplementary Table 12) varied from nearly 5,000 up- and down-regulated genes in low-
322 grade glioma to no significant gene in 5 cancers. Similarly, we also identified significant
323 differentially methylated genes with age (age-DMGs, Supplementary Table 13) (adj. p-value
324 < 0.05), the number of age-DEGs and age-DMGs were consistent for most cancer types (Fig.

325 6a). We next focused our analysis on 10 cancer types that contained at least 150 age-DEGs and
326 150 age-DMGs, including low-grade glioma, breast cancer, endometrial cancer, oesophageal
327 cancer, papillary renal cell carcinoma, ovarian cancer, liver cancer, acute myeloid leukaemia,
328 melanoma, and prostate cancer. We identified overlapping genes between age-DEGs and age-
329 DMGs and found that most of them, from 84% (37/44 genes) in ovarian cancer to 100% in
330 acute myeloid leukaemia (57 genes) and prostate cancer (7 genes), were genes that presented
331 increased methylation and decreased expression with age and genes that had decreased
332 methylation and increased expression with age (Fig. 6b-c, Supplementary Fig. 9,
333 Supplementary Table 14). We further examined the correlation coefficient between
334 methylation and expression comparing between 4 groups of genes 1) overlap genes between
335 age-DMGs and age-DEGs (age-DMGs-DEGs), 2) age-DMGs only, 3) age-DEGs only, and 4)
336 other genes. We found that age-DMGs-DEGs had the most negative correlation between DNA
337 methylation and expression when comparing with other groups of genes (Fig. 6d,
338 Supplementary Fig. 10, Supplementary Table 15), highlighting that age-associated gene
339 expression changes in cancer are repressed, at least in part, by DNA methylation.

340 We next performed Gene Set Enrichment Analysis (GSEA) to gain biological insights
341 into the expression and methylation changes with age. We identified various significantly
342 enriched Gene Ontology (GO) terms across cancers (Fig. 6e, Supplementary Fig. 11,
343 Supplementary Table 16). Notably, several GO terms were enriched in both expression and
344 methylation changes, in the opposite direction. The enriched terms in breast cancer included
345 several signalling, metabolism, and developmental pathways. The Wnt signaling pathway,
346 which was altered more frequently in older breast cancer patients (Fig. 5b), showed a decrease
347 in gene expression and increase in methylation with age. In low-grade glioma, interestingly,
348 mitochondrial terms were enriched in the gene expression of younger patients. Mitochondrial
349 dysfunction is known to be important in glioma pathophysiology³⁷, thus the different levels of

350 mitochondrial aberrations might contribute to the disparities in the aggressiveness of gliomas
351 in patients of different age. We also identified numerous immune-related terms enriched across
352 several cancer types, including oesophageal, papillary renal cell, liver, and prostate cancers
353 (Supplementary Fig. 11, Supplementary Table 16). Previous studies suggested alterations in
354 immune-related gene expression and immune cell abundance changes with age in cancers^{38,39}.
355 In the present study, we have systematically characterised the transcriptome and methylation
356 in relation to age across cancer types. Our results suggest that gene expression changes with
357 age in cancer are controlled, at least in part, by DNA methylation. These changes reflect
358 differences in biological pathways that might be important in tumour development.

359

360 **Discussion**

361 Although age is an important risk factor for cancer, how age impacts the molecular
362 landscape of cancer is not well understood. In this study, we provide a comprehensive overview
363 of the age-associated molecular landscape in cancer, including genomic instability, LOH,
364 WGD, SCNAs, somatic mutations, pathway alterations, gene expression, and DNA
365 methylation. We confirmed the known increase in mutation load^{4,5} and found an increase in
366 genomic instability, LOH and WGD with age in several cancer types. We identified several
367 age-related pan-cancer and cancer-specific alterations. The highest age-related differences
368 were evident in low-grade glioma and endometrial cancer.

369 Cancer develops through the accumulation of genetic and epigenetic alterations.
370 Mutation accumulation with age is thought to be a cause of cancer and a substantial portion of
371 mutations arise before cancer initiation⁶. The age-associated mutation accumulation has been
372 demonstrated in both cancer^{4,5} and normal tissues⁴⁰⁻⁴², providing a better understanding of an
373 early carcinogenesis event. Our results show that, in addition to mutations, SCNAs, LOH and
374 WGD increase with age in several cancers, in particular low-grade glioma, endometrial and

375 ovarian cancers. Recent evidence suggests that SCNA burden is a prognostic factor associated
376 with both recurrence and death⁴³, thus, an increased SCNA level with age might relate to poor
377 prognosis in the elderly.

378 The negative association between age and mutation in *IDH1* and *ATRX* in glioma points
379 towards the difference of patient age at diagnosis between the *IDH*-mutant and *IDH*-WT
380 subtypes. *IDH*-mutant tumours are observed in the majority of low-grade glioma and show
381 favourable prognosis. *IDH*-WT low-grade gliomas, on the other hand, more resemble
382 glioblastomas and have poorer survival. In glioblastoma, although *IDH*-mutants are a minority
383 of tumours, they are also associated with younger age⁴⁴. The present study together with
384 others^{34,45}, therefore indicates that glioma shows unique age-associated subtypes. However,
385 more research is needed to understand how age influences the evolution of glioma subtypes.

386 Our results highlighted substantial age-associated differences in the genome of
387 endometrial cancer. Younger endometrial tumours associate with a *POLE* and MSI-H
388 subtypes, leading to an enrichment of hypermutated tumours, while tumours from older
389 patients tend to harbour a higher SCNA level and lower mutation load. Previous studies have
390 classified endometrial cancer into four subtypes: *POLE*, MSI-H, copy-number low and copy-
391 number high subtypes. The *POLE* subtype and MSI-H subtype are dominated by the *POLE*
392 and defective mismatch repair mutational signatures, respectively³³. Conversely, the copy-
393 number low and copy-number high subtypes had a dominant ageing-related mutational
394 signature³¹. The *POLE* and MSI-H subtypes have a favourable prognosis, while the copy-
395 number high subtype is associated with poor survival. Therefore, endometrial cancer from
396 younger patients is associated with *POLE* mutations, mismatch repair defects, high mutation
397 load and better survival outcomes. Older endometrial cancer, however, is related to extensive
398 SCNAs and worse prognosis. Importantly, apart from low-grade glioma and endometrial
399 cancer, we demonstrate that other cancer types also present an age-associated genomic

400 landscape in cancer driver genes and oncogenic signalling pathways. These results highlight
401 the impact of age on the molecular profile of cancer.

402 Having identified these age-related differences in the molecular landscapes of various
403 cancers, the obvious question is what drives these differences. Accumulating evidence has
404 underscored the importance of tissue environment changes with ageing in cancer initiation and
405 progression^{7,8,39,46}. We reason that tissue environment changes during ageing and might
406 provide different selective advantages for tumours harbouring different molecular alterations
407 in turn directing the tumours to different evolutionary routes. Therefore, cancer with different
408 genomic alterations might thrive better in younger or older patients. Gene expression and
409 epigenetic changes related to ageing have been studied and linked to cancer^{8,38,47,48}. Here, we
410 identified numerous age-associated gene expression and corresponding DNA methylation in a
411 broad range of cancers. Indeed, age-DMGs-DEGs are those with the strongest negative
412 correlation between methylation and expression when comparing with other groups, indicating
413 that differentially expressed genes with age in cancer are partly regulated by methylation.
414 Expression and methylation changes with age link to several biological processes, showing that
415 cancer from patients with different ages present different phenotypes. We also noticed that
416 cancer in female reproductive organs including breast, ovarian and endometrial cancers are
417 among those with the highest number of age-DEGs and age-DMGs. These cancers tend to have
418 a higher mass-normalised cancer incidence, which may reflect evolutionary trade-offs
419 involving selective pressures related to reproduction⁴⁹. The age-associated hormonal changes
420 could also be responsible for this age-related expression differences in cancer⁵⁰. The limitation
421 of this analysis is that although we have already included tumour purity in our linear model, it
422 is not possible to account for the different tumour-constituent cell proportions and thus fully
423 exclude the influence from gene expression of non-tumor cells such as infiltrating immune

424 cells³⁹. Further studies are required to provide mechanistic understanding of the impact of an
425 ageing microenvironment in shaping tumour evolution.

426 During the preparation of our manuscript, a study based on a similar concept has been
427 released by Li et al⁵¹. In this work, Li et al. used TCGA and the recent pan-cancer analysis of
428 whole genomes (PCAWG) data to study the age-associated genomic differences in cancer.
429 Results from the two studies are consistent on several points. Firstly, both studies indicate the
430 increase in mutations and SCNA levels with age. Next, despite using slightly different
431 statistical cutoffs and models, several age-associated genomic features are identified by both
432 studies, for example, the higher frequency of *IDH1* and *ATRX* mutations in younger glioma
433 patients. Li et al. explored mutational timing and signatures, which suggested the possible
434 underlying mechanisms for age-associated genomic differences. Our study, however, has also
435 featured an age-related genomic profile in endometrial cancer. We have investigated cancer-
436 specific associations between age and LOH, WGD and oncogenic signalling. Furthermore, we
437 have analysed age-related global transcriptomic and DNA methylation changes. Our study are
438 complementary with the Li et al. study, both studies thus serve as a foundation for
439 understanding age-related differences and effects on the cancer molecular landscape and
440 emphasise the importance of age in cancer genomic research that is particularly valuable in the
441 clinical practice.

442

443

444

445

446

447

448

449 **Methods**

450 **Data acquisition**

451 Publicly available copy-number alteration seg files (nocnv_hg19.seg), normalized mRNA
452 expression in RSEM (.rsem.genes.normalized_results TCGA files from the legacy archive,
453 aligned to hg19), and clinical data (XML files) from TCGA were downloaded using
454 *TCGAbiolinks* (version 2.14.1)⁵². The mutation annotation format (MAF) file was downloaded
455 from the TCGA MC3 project⁵³ (<https://gdc.cancer.gov/about-data/publications/mc3-2017>).
456 The somatic alterations in 10 canonical oncogenic pathways across TCGA samples were
457 obtained from a previous study by Sanchez-Vega et al³⁶. The TCGA Illumina
458 HumanMethylation450K array data (in β -values) was downloaded from Broad GDAC
459 Firehose (<http://gdac.broadinstitute.org/>). The allele-specific copy number, tumour ploidy,
460 tumour purity were estimated using ASCAT (version 2.4.2)⁵⁴ on hg19 SNP6 arrays with
461 penalty=70 as previously described^{55,56}. We restricted our subsequent analyses to samples that
462 have these profiles available. WGD duplication was determined using fraction of genome with
463 LOH and ploidy information. Genomic instability (GI) scores have been computed as fraction
464 of genomic regions that are not in 1+1 (for non WGD tumours) or 2+2 (for WGD tumours)
465 statuses. For each data type and each cancer type, the summary of the numbers of TCGA
466 samples included in the analysis, alongside clinical variable analysed are presented in the
467 Supplementary Table 1.

468

469 **Statistical analysis and visualisation**

470 Simple linear regression and multiple linear regression adjusting for clinical variables were
471 performed using the *lm* function in R to access the relationship between age and continuous
472 variables of interest. Simple logistic regression to investigate the association between age and
473 binary response (e.g. mutation as 1 and wild-type as 0) and multiple logistic regression

474 adjusting for covariates were carried out using the *glm* function in R. In pan-cancer analyses,
475 gender, race and cancer type were variables included in the linear model. Clinical variables
476 used in cancer-specific analyses included gender, race, pathologic stage, neoplasm histologic
477 grade, smoking status, alcohol consumption and cancer-specific variables such as oestrogen
478 receptor (ER) status in breast cancer. To avoid the potential detrimental effect caused by
479 missing data, we retained only variables with missing data less than 10% of samples used in
480 the somatic copy number alteration analysis (Supplementary Table 1). To account for the
481 difference in the proportion of cancer cells in each tumour, tumour purity (cancer cell fraction)
482 estimated from ASCAT was included in the linear model. When necessary, to avoid the
483 separation problem that might occur due to the sparse-data bias⁵⁷, *logistf* function from the
484 *logistf* package (version 1.23)⁵⁸ was used to perform multivariable logistic regression with
485 Firth's penalization⁵⁹. Effect sizes from logistic regression analyses were reported as odds ratio
486 per year and 95% confidence intervals. P-values from the analyses were accounted for
487 multiple-hypothesis testing using Benjamini–Hochberg procedure⁶⁰. Statistical significance
488 was considered if adj. p-value < 0.05, unless specifically indicated otherwise.

489 All statistical analyses were carried out using R (version 3.6.3)⁶¹. Plots were generated
490 using *ggplot2* (version 3.3.2)⁶², *ggrepel* (version 0.8.2)⁶³, *ggpubr* (version 0.4.0)⁶⁴,
491 *ComplexHeatmap* (version 2.2.0)⁶⁵, and *VennDiagram* (version 1.6.20)⁶⁶.

492

493 **GI score analysis**

494 GI score was calculated as a genome fraction (percent-based) that does not fit the estimated
495 tumour ploidy, 2 for normal diploid, and 4 for tumours that have undergone the WGD process.
496 Simple linear regression was performed to identify the association between age and GI score.
497 For pan-cancer analysis, multiple linear regression was used to adjust for gender, race, and
498 cancer type. For cancer-specific analysis, multiple linear regression accounting for clinical

499 variables was conducted on the cancer types that had a significant association between age and
500 GI score from the simple linear regression analysis (adj. p-value < 0.05). The complete set of
501 results is presented in Supplementary Table 2.

502

503 **Percentage genomic LOH quantification and analysis**

504 To quantify the percent genomic LOH for each tumour, we used allele-specific copy number
505 profiles from ASCAT. X and Y chromosome regions were discarded from the analysis. The
506 LOH segments were segments that harbour only one allele. The percent genomic LOH was
507 defined as 100 times the total length of LOH regions / length of the genome.

508 Simple linear regression and multiple linear regression adjusting for gender, race, and
509 cancer types were conducted to investigate the relationship between age and the percent
510 genomic LOH in the pan-cancer analysis. For cancer-specific analysis, simple linear regression
511 was performed followed by multiple linear regression accounting for clinical factors for
512 cancers with a significant association in simple linear regression analysis (adj. p-value < 0.05).
513 The complete set of results is in Supplementary Table 3.

514

515 **WGD analysis**

516 WGD status for each tumour was obtained from fraction of genome with LOH and tumour
517 ploidy. To investigate the association between age and WGD across the pan-cancer dataset, we
518 performed simple logistic regression and multiple logistic regression correcting for gender,
519 race, and cancer type. For cancer-specific analysis, simple logistic regression was performed
520 to assess the association between age and WGD on tumours from each cancer type. Cancer
521 types with a significant association between age and WGD (adj. p-value < 0.05) were further
522 subjected to the multiple logistic regression accounting for the clinical variables. The complete
523 set of results is in Supplementary Table 4.

524

525 **List of known cancer driver genes**

526 We compiled a list of known cancer driver genes from (1) the list of 243 COSMIC classic
527 genes from COSMIC database version 91⁶⁷ (downloaded on 1st July 2020), (2) the list of 260
528 significantly mutated genes from Lawrence et al⁶⁸, and (3) the list of 299 cancer driver genes
529 from the TCGA Pan-Cancer study⁶⁹. In total, we obtained 505 cancer genes and focused on the
530 mutations and focal-level SCNAs on these genes in our study. The full list of cancer driver
531 genes is available in Supplementary Table 8.

532

533 **Recurrent SCNA analysis**

534 Recurrent arm-level and focal-level SCNAs of each cancer type were identified using
535 GISTIC2.0²². Segmented files (nocnv_hg19.seg) from TCGA, marker file and CNV file,
536 provided by GISTIC2.0, were used as input files. The parameters were set as follows: ‘-
537 genegistic 1 -smallmem 1 -qvt 0.25 -ta 0.25 -td 0.25 -broad 1 -brlen 0.7 -conf 0.95 -armpeel 1
538 -savegene 1’. Based on these parameters, broad events were defined as the alterations happen
539 in more than 70% of an arm. The log₂ ratio thresholds for copy number gains and deletions
540 were 0.25 and -0.25, respectively. The confidence level was set as 0.95 and the q-value was
541 0.25.

542 To investigate the association between age and arm-level SCNAs for each cancer type,
543 simple logistic regression was performed for each chromosomal arm that was identified as
544 recurrent SCNA in a cancer type. Only cancer types with more than 100 samples were included
545 in this analysis (Table 1). Arms with a significant association (adj. p-value < 0.05) were further
546 adjusted for clinical variables using multiple logistic regression. The complete set of results is
547 in Supplementary Table 6. Similarly, simple and multiple logistic regression was conducted on
548 the focal-level SCNAs for each cancer type. Regions that are not overlapped with centromeres

549 or telomeres were removed from the analysis. The complete set of results is in Supplementary
550 Table 7.

551 To confirm the impact of SCNAs on gene expression, we investigated the correlation
552 between GISTIC2.0 score and RNA-seq based gene expression ($\log_2(\text{normalised RSEM} + 1)$)
553 for tumours that have both types of data using Pearson correlation. The correlation was
554 considered significant if the p-value corrected for multiple-hypothesis testing using the
555 Benjamini-Hochberg procedure < 0.05 . The complete set of results is in Supplementary Table
556 7.

557

558 **SCNA score quantification and analysis**

559 Previous studies have developed the SCNA score representing the SCNA level of a tumour^{12,23}.
560 We applied the methods described by Yuan et al¹² to calculate SCNA scores. Using SCNA
561 profiles from GISTIC2.0 analysis, SCNA scores for each tumour were derived at three different
562 levels (chromosome-, arm-, and focal-level). For each tumour, each focal-event \log_2 copy
563 number ratio from GISTIC2.0 was classified into the following score: 2 if the \log_2 ratio ≥ 1 , 1
564 if the \log_2 ratio < 1 and ≥ 0.25 , 0 if the \log_2 ratio < 0.25 and ≥ -0.25 , -1 if the \log_2 ratio $< -$
565 0.25 and ≥ -1 , and -2 if the \log_2 ratio < -1 . The $|\text{score}|$ from each focal event in a tumour was
566 then summed into a focal score of a tumour. Thereafter, the rank-based normalisation
567 (rank/number of tumours in a cancer type) was applied to focal scores from all tumours within
568 the same cancer type, resulting in normalized focal-level SCNA scores. Therefore, tumours
569 with high focal-level SCNAs will have focal-level SCNA scores close to 1, while tumours with
570 low focal-level SCNAs will have scores close to 0. For the arm- and chromosome-level SCNA
571 scores, a similar procedure was applied to the broad event \log_2 copy number ratio from
572 GISTIC2.0. An event was considered as a chromosome-level if both arms have the same \log_2
573 ratio, otherwise it was considered as an arm-level. Similar to the focal-level SCNA score, each

574 arm- and chromosome-event log₂ copy number ratio was classified into the 2, 1, 0, -1, -2 scores
575 using the threshold described above. The |score| from all arm-events and chromosome-events
576 for a tumour were then summed into an arm score and chromosome score, respectively. For
577 each cancer type, the rank-based normalisation was applied to arm scores and chromosome
578 scores from all tumours to derive normalised arm-level SCNA scores and normalised
579 chromosome-level SCNA scores, respectively. An overall SCNA score for a tumour was
580 defined as the sum of focal-level, arm-level, and chromosome-level SCNA scores. A
581 chromosome/arm-level SCNA score for a tumour was defined as the sum of chromosome-level
582 and arm-level SCNA scores.

583 The association between age and overall, chromosome/arm-level, and focal-level
584 SCNA scores for each cancer type was investigated using simple linear regression. Cancer
585 types with a significant association (adj. p-value < 0.05) were then subjected to multiple linear
586 regression analysis adjusting for the clinical variables. The complete set of results is included
587 in Supplementary Table 5.

588

589 **Analysis of age-associated somatic mutation in cancer genes**

590 We obtained the mutation data from the MAF file from the recent TCGA Multi-Center
591 Mutation Calling in Multiple Cancers (MC3) project⁵³. In the MC3 effort, variants were called
592 using seven variant callers. We filtered the variants to keep only non-silent SNVs and indels
593 located in gene bodies, retaining only “Frame_Shift_Del”, “Frame_Shift_Ins”,
594 “In_Frame_Del”, “In_Frame_Ins”, “Missense_Mutation”, “Nonsense_Mutation”,
595 “Nonstop_Mutation”, “Splice_Site” and Translation_Start_Site in the
596 “Variant_Classification” column. We focused only on mutations in the cancer genes from our
597 compiled list of cancer driver genes. To prevent the bias that might cause by hypermutated
598 tumours, we restricted the analysis to tumours with < 1,000 mutations per exome. For pan-

599 cancer analysis, multiple logistic regression accounting for gender, race and cancer type was
600 performed to investigate the association between age and mutations in 20 cancer genes that are
601 mutated in > 5% of samples (Supplementary Table 10). For cancer-specific analysis, simple
602 logistic regression was used to identify cancer genes that the mutations in these genes are
603 associated with the patient's age. Only genes that are mutated in > 5% of samples from each
604 cancer type were included in the analysis. The significant associations (adj. p-value < 0.05)
605 were further investigated using multiple logistic regression accounting for clinical variables.
606 The complete set of results is in Supplementary Table 10.

607

608 **Analysis of mutational burden, MSI-H status, and *POLE/POLD1* mutations**

609 A mutational burden was defined as the total non-silent mutations in an exome. The mutational
610 burden for each tumour was log-transformed before using it in the subsequent analysis. To
611 investigate the relationship between age and mutational burden in pan-cancer, multiple linear
612 regression adjusting for gender, race and cancer type was conducted. For cancer-specific
613 analysis, simple linear regression was performed. Cancer types with a significant association
614 between age and mutational burden in simple linear regression analysis (adj. p-value < 0.05)
615 were further examined using multiple linear regression accounting for clinical factors. The
616 complete set of results is in Supplementary Table 9.

617 Microsatellite instability status for COAD, READ, STAD, and UCEC were
618 downloaded from TCGA using *TCGAbiolinks*. To study the association between the presence
619 of high microsatellite instability (MSI-H) and age, tumours were divided into binary groups:
620 MSI-H = TRUE and MSI-H = FALSE. Multiple logistic regression adjusting for clinical
621 variables was then performed. Similarly, *POLE* and *POLD1* mutation status were in a binary
622 outcome (mutated and not mutated). Multiple logistic regression was used to investigate the

623 association between age and *POLE/POLD1* mutations in cancer types that contained
624 *POLE/POLD1* mutations in > 5% of samples.

625

626 **Oncogenic signalling pathway analysis**

627 We used the list of pathway-level alterations in ten oncogenic pathways (cell cycle, Hippo,
628 Myc, Notch, Nrf2, PI-3-Kinase/Akt, RTK-RAS, TGF β signaling, p53 and β -catenin/Wnt) for
629 TCGA tumours comprehensively compiled by Sanchez-Vega et al³⁶. Member genes in the
630 pathways were accessed for SCNAs, mutations, epigenetic silencing through promoter DNA
631 hypermethylation and gene fusions. We retained only the pathway alteration data of samples
632 that were presented in our SCNA analysis. For the pan-cancer analysis, we employed multiple
633 logistic regression adjusting for the patient's gender, race and cancer type to demonstrate the
634 relationship between pathway-level alteration and age. To investigate the association between
635 age and cancer-specific pathway alterations, we performed simple logistic regression. Cancer
636 types with a significant association (adj. p-value < 0.05) were further examined by multiple
637 logistic regression accounting for clinical variables. The complete set of results is in
638 Supplementary Table 11.

639

640 **Gene expression and DNA methylation analysis**

641 To render the results from gene expression and DNA methylation comparable, we limited the
642 analysis to genes that are presented in both types of data. The lowly expressed genes were
643 filtered out from the analysis by keeping only genes with RSEM > 0 in more than 50 percent
644 of samples. Only protein coding genes identified using biomaRt⁷⁰ (Ensembl version 100, April
645 2020). Normalised mRNA expression in RSEM for each TCGA cancer type was log₂-
646 transformed before subjected to the multiple linear regression analysis adjusting for clinical
647 factors. RNA-seq data for colon cancer and endometrial cancer consisted of two platforms,

648 Illumina HiSeq and Illumina GA. Thus, a platform was included as another covariate in the
649 linear regression model for these two cancer types. Genes with adj. p-value < 0.05 were
650 considered significantly differentially expressed genes with age (age-DEGs) (Supplementary
651 Table 12). DNA methylation data was presented as β -values, which are the ratio of the
652 intensities of methylated and unmethylated alleles. Because multiple methylation probes can
653 be mapped to the same gene, we used the one-to-one mapping genes and probes by selecting
654 the probes that are most negatively correlated with the corresponding gene expression from the
655 files meth.by_min_expr_corr.data.txt. Similar multiple linear regression to the gene expression
656 analysis was performed on the methylation data. Genes with adj. p-value < 0.05 were
657 considered significant differentially methylated genes with age (age-DMGs). The complete set
658 of results is in Supplementary Table 13.

659 The correlation between gene expression and DNA methylation was calculated using
660 Pearson correlation. We used the Kruskal-Wallis test to investigate the differences between
661 correlation coefficients among groups (age-DMGs-DEGs, age-DMGs, age-DEGs, other
662 genes). The pairwise comparisons were carried out by Dunn's test. The complete set of results
663 is in Supplementary Table 15.

664 Gene Set Enrichment Analysis (GSEA) was performed to investigate the Gene
665 Ontology (GO) terms that are enriched in tumours from younger or older patients. The analysis
666 was done using the package *ClusterProfiler* (version 3.14.3)⁷¹. The complete list of enriched
667 GO terms is presented in Supplementary Table 16.

668

669 **Data availability**

670 TCGA data used in this study are publicly available and can be obtained from NCI's Genomic
671 Data Commons portal (<https://portal.gdc.cancer.gov/>), *TCGAbiolinks* (version 2.14.1)⁵² and
672 Broad GDAC Firehose (<http://gdac.broadinstitute.org/>).

673

674 **Code availability**

675 The custom scripts for data analysis and generate figures are available at

676 https://github.com/maglab/Age-associated_cancer_genome.

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698 References

- 699 1 de Magalhaes, J. P. How ageing processes influence cancer. *Nat Rev Cancer* **13**, 357-
700 365, doi:10.1038/nrc3497 (2013).
- 701 2 Laconi, E., Marongiu, F. & DeGregori, J. Cancer as a disease of old age: changing
702 mutational and microenvironmental landscapes. *Br J Cancer* **122**, 943-952,
703 doi:10.1038/s41416-019-0721-1 (2020).
- 704 3 Nowell, P. C. The clonal evolution of tumor cell populations. *Science* **194**, 23-28,
705 doi:10.1126/science.959840 (1976).
- 706 4 Milholland, B., Auton, A., Suh, Y. & Vijg, J. Age-related somatic mutations in the
707 cancer genome. *Oncotarget* **6**, 24627-24635, doi:10.18632/oncotarget.5685 (2015).
- 708 5 Alexandrov, L. B. *et al.* Clock-like mutational processes in human somatic cells. *Nat*
709 *Genet* **47**, 1402-1407, doi:10.1038/ng.3441 (2015).
- 710 6 Tomasetti, C., Vogelstein, B. & Parmigiani, G. Half or more of the somatic mutations
711 in cancers of self-renewing tissues originate prior to tumor initiation. *Proc Natl Acad*
712 *Sci U S A* **110**, 1999-2004, doi:10.1073/pnas.1221068110 (2013).
- 713 7 Fane, M. & Weeraratna, A. T. How the ageing microenvironment influences tumour
714 progression. *Nat Rev Cancer* **20**, 89-106, doi:10.1038/s41568-019-0222-9 (2020).
- 715 8 Chatsirisupachai, K., Palmer, D., Ferreira, S. & de Magalhaes, J. P. A human tissue-
716 specific transcriptomic analysis reveals a complex relationship between aging, cancer,
717 and cellular senescence. *Aging Cell* **18**, e13041, doi:10.1111/acer.13041 (2019).
- 718 9 Li, C. H., Haider, S., Shiah, Y. J., Thai, K. & Boutros, P. C. Sex Differences in Cancer
719 Driver Genes and Biomarkers. *Cancer Res* **78**, 5527-5537, doi:10.1158/0008-
720 5472.CAN-18-0362 (2018).
- 721 10 Yuan, Y. *et al.* Comprehensive Characterization of Molecular Differences in Cancer
722 between Male and Female Patients. *Cancer Cell* **29**, 711-722,
723 doi:10.1016/j.ccell.2016.04.001 (2016).
- 724 11 Sinha, S. *et al.* Higher prevalence of homologous recombination deficiency in tumors
725 from African Americans versus European Americans. *Nature Cancer* **1**, 112-121,
726 doi:10.1038/s43018-019-0009-7 (2020).
- 727 12 Yuan, J. *et al.* Integrated Analysis of Genetic Ancestry and Genomic Alterations across
728 Cancers. *Cancer Cell* **34**, 549-560 e549, doi:10.1016/j.ccell.2018.08.019 (2018).
- 729 13 Ma, X. *et al.* Pan-cancer genome and transcriptome analyses of 1,699 paediatric
730 leukaemias and solid tumours. *Nature* **555**, 371-376, doi:10.1038/nature25795 (2018).
- 731 14 Grobner, S. N. *et al.* The landscape of genomic alterations across childhood cancers.
732 *Nature* **555**, 321-327, doi:10.1038/nature25480 (2018).
- 733 15 Brennan, C. W. *et al.* The somatic genomic landscape of glioblastoma. *Cell* **155**, 462-
734 477, doi:10.1016/j.cell.2013.09.034 (2013).
- 735 16 Gerhauser, C. *et al.* Molecular Evolution of Early-Onset Prostate Cancer Identifies
736 Molecular Risk Markers and Clinical Trajectories. *Cancer Cell* **34**, 996-1011 e1018,
737 doi:10.1016/j.ccell.2018.10.016 (2018).
- 738 17 Liao, S. *et al.* The molecular landscape of premenopausal breast cancer. *Breast Cancer*
739 *Res* **17**, 104, doi:10.1186/s13058-015-0618-8 (2015).
- 740 18 Ryland, G. L. *et al.* Loss of heterozygosity: what is it good for? *BMC Med Genomics*
741 **8**, 45, doi:10.1186/s12920-015-0123-z (2015).
- 742 19 Lopez, S. *et al.* Interplay between whole-genome doubling and the accumulation of
743 deleterious alterations in cancer evolution. *Nat Genet* **52**, 283-293,
744 doi:10.1038/s41588-020-0584-7 (2020).
- 745 20 Bielski, C. M. *et al.* Genome doubling shapes the evolution and prognosis of advanced
746 cancers. *Nat Genet* **50**, 1189-1195, doi:10.1038/s41588-018-0165-1 (2018).

- 747 21 Van de Peer, Y., Mizrachi, E. & Marchal, K. The evolutionary significance of
748 polyploidy. *Nat Rev Genet* **18**, 411-424, doi:10.1038/nrg.2017.26 (2017).
- 749 22 Mermel, C. H. *et al.* GISTIC2.0 facilitates sensitive and confident localization of the
750 targets of focal somatic copy-number alteration in human cancers. *Genome Biol* **12**,
751 R41, doi:10.1186/gb-2011-12-4-r41 (2011).
- 752 23 Davoli, T., Uno, H., Wooten, E. C. & Elledge, S. J. Tumor aneuploidy correlates with
753 markers of immune evasion and with reduced response to immunotherapy. *Science* **355**,
754 doi:10.1126/science.aaf8399 (2017).
- 755 24 Korber, V. *et al.* Evolutionary Trajectories of IDH(WT) Glioblastomas Reveal a
756 Common Path of Early Tumorigenesis Instigated Years ahead of Initial Diagnosis.
757 *Cancer Cell* **35**, 692-704 e612, doi:10.1016/j.ccell.2019.02.007 (2019).
- 758 25 Xu, F. *et al.* Elevated expression of RIT1 correlates with poor prognosis in endometrial
759 cancer. *Int J Clin Exp Pathol* **8**, 10315-10324 (2015).
- 760 26 Bonneville, R. *et al.* Landscape of Microsatellite Instability Across 39 Cancer Types.
761 *JCO Precis Oncol* **2017**, doi:10.1200/PO.17.00073 (2017).
- 762 27 Kim, T. M., Laird, P. W. & Park, P. J. The landscape of microsatellite instability in
763 colorectal and endometrial cancer genomes. *Cell* **155**, 858-868,
764 doi:10.1016/j.cell.2013.10.015 (2013).
- 765 28 Chalmers, Z. R. *et al.* Analysis of 100,000 human cancer genomes reveals the landscape
766 of tumor mutational burden. *Genome Med* **9**, 34, doi:10.1186/s13073-017-0424-2
767 (2017).
- 768 29 Campbell, B. B. *et al.* Comprehensive Analysis of Hypermutation in Human Cancer.
769 *Cell* **171**, 1042-1056 e1010, doi:10.1016/j.cell.2017.09.048 (2017).
- 770 30 Shlien, A. *et al.* Combined hereditary and somatic mutations of replication error repair
771 genes result in rapid onset of ultra-hypermuted cancers. *Nat Genet* **47**, 257-262,
772 doi:10.1038/ng.3202 (2015).
- 773 31 Ashley, C. W. *et al.* Analysis of mutational signatures in primary and metastatic
774 endometrial cancer reveals distinct patterns of DNA repair defects and shifts during
775 tumor progression. *Gynecol Oncol* **152**, 11-19, doi:10.1016/j.ygyno.2018.10.032
776 (2019).
- 777 32 Berger, A. C. *et al.* A Comprehensive Pan-Cancer Molecular Study of Gynecologic and
778 Breast Cancers. *Cancer Cell* **33**, 690-705 e699, doi:10.1016/j.ccell.2018.03.014 (2018).
- 779 33 Cancer Genome Atlas Research, N. *et al.* Integrated genomic characterization of
780 endometrial carcinoma. *Nature* **497**, 67-73, doi:10.1038/nature12113 (2013).
- 781 34 Yan, H. *et al.* IDH1 and IDH2 mutations in gliomas. *N Engl J Med* **360**, 765-773,
782 doi:10.1056/NEJMoa0808710 (2009).
- 783 35 Cancer Genome Atlas Research, N. *et al.* Comprehensive, Integrative Genomic
784 Analysis of Diffuse Lower-Grade Gliomas. *N Engl J Med* **372**, 2481-2498,
785 doi:10.1056/NEJMoa1402121 (2015).
- 786 36 Sanchez-Vega, F. *et al.* Oncogenic Signaling Pathways in The Cancer Genome Atlas.
787 *Cell* **173**, 321-337 e310, doi:10.1016/j.cell.2018.03.035 (2018).
- 788 37 Ordys, B. B., Launay, S., Deighton, R. F., McCulloch, J. & Whittle, I. R. The role of
789 mitochondria in glioma pathophysiology. *Mol Neurobiol* **42**, 64-75,
790 doi:10.1007/s12035-010-8133-5 (2010).
- 791 38 Wu, Y. *et al.* Comprehensive transcriptome profiling in elderly cancer patients reveals
792 aging-altered immune cells and immune checkpoints. *Int J Cancer* **144**, 1657-1663,
793 doi:10.1002/ijc.31875 (2019).
- 794 39 Erbe, R. *et al.* Aging interacts with tumor biology to produce major changes in the
795 immune tumor microenvironment. *bioRxiv*,
796 doi:<https://doi.org/10.1101/2020.06.08.140764> (2020).

- 797 40 Martincorena, I. *et al.* Somatic mutant clones colonize the human esophagus with age.
798 *Science* **362**, 911-917, doi:10.1126/science.aau3879 (2018).
- 799 41 Martincorena, I. *et al.* Tumor evolution. High burden and pervasive positive selection
800 of somatic mutations in normal human skin. *Science* **348**, 880-886,
801 doi:10.1126/science.aaa6806 (2015).
- 802 42 Xie, M. *et al.* Age-related mutations associated with clonal hematopoietic expansion
803 and malignancies. *Nat Med* **20**, 1472-1478, doi:10.1038/nm.3733 (2014).
- 804 43 Hieronymus, H. *et al.* Tumor copy number alteration burden is a pan-cancer prognostic
805 factor associated with recurrence and death. *Elife* **7**, doi:10.7554/eLife.37294 (2018).
- 806 44 Mirchia, K. & Richardson, T. E. Beyond IDH-Mutation: Emerging Molecular
807 Diagnostic and Prognostic Features in Adult Diffuse Gliomas. *Cancers (Basel)* **12**,
808 doi:10.3390/cancers12071817 (2020).
- 809 45 Verhaak, R. G. *et al.* Integrated genomic analysis identifies clinically relevant subtypes
810 of glioblastoma characterized by abnormalities in PDGFRA, IDH1, EGFR, and NF1.
811 *Cancer Cell* **17**, 98-110, doi:10.1016/j.ccr.2009.12.020 (2010).
- 812 46 Rozhok, A. & DeGregori, J. A generalized theory of age-dependent carcinogenesis.
813 *Elife* **8**, doi:10.7554/eLife.39950 (2019).
- 814 47 Perez, R. F., Tejedor, J. R., Bayon, G. F., Fernandez, A. F. & Fraga, M. F. Distinct
815 chromatin signatures of DNA hypomethylation in aging and cancer. *Aging Cell* **17**,
816 e12744, doi:10.1111/accel.12744 (2018).
- 817 48 Johnson, A. A. *et al.* The role of DNA methylation in aging, rejuvenation, and age-
818 related disease. *Rejuvenation Res* **15**, 483-494, doi:10.1089/rej.2012.1324 (2012).
- 819 49 Silva, A. S. *et al.* Gathering insights on disease etiology from gene expression profiles
820 of healthy tissues. *Bioinformatics* **27**, 3300-3305, doi:10.1093/bioinformatics/btr559
821 (2011).
- 822 50 Benz, C. C. Impact of aging on the biology of breast cancer. *Crit Rev Oncol Hematol*
823 **66**, 65-74, doi:10.1016/j.critrevonc.2007.09.001 (2008).
- 824 51 Li, C. H., Haider, S. & Boutros, P. C. Age Influences on the Molecular Presentation of
825 Tumours. *bioRxiv*, doi:<https://doi.org/10.1101/2020.07.07.192237> (2020).
- 826 52 Colaprico, A. *et al.* TCGAAbiolinks: an R/Bioconductor package for integrative analysis
827 of TCGA data. *Nucleic Acids Res* **44**, e71, doi:10.1093/nar/gkv1507 (2016).
- 828 53 Ellrott, K. *et al.* Scalable Open Science Approach for Mutation Calling of Tumor
829 Exomes Using Multiple Genomic Pipelines. *Cell Syst* **6**, 271-281 e277,
830 doi:10.1016/j.cels.2018.03.002 (2018).
- 831 54 Van Loo, P. *et al.* Allele-specific copy number analysis of tumors. *Proc Natl Acad Sci*
832 *U S A* **107**, 16910-16915, doi:10.1073/pnas.1009843107 (2010).
- 833 55 Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues.
834 *Cell* **171**, 1029-1041 e1021, doi:10.1016/j.cell.2017.09.042 (2017).
- 835 56 Alexandrov, L. B. *et al.* Mutational signatures associated with tobacco smoking in
836 human cancer. *Science* **354**, 618-622, doi:10.1126/science.aag0299 (2016).
- 837 57 Greenland, S., Mansournia, M. A. & Altman, D. G. Sparse data bias: a problem hiding
838 in plain sight. *BMJ* **352**, i1981, doi:10.1136/bmj.i1981 (2016).
- 839 58 Heinze, G. & Ploner, M. logistf: Firth's Bias-Reduced Logistic Regression. (2018).
- 840 59 Heinze, G. & Schemper, M. A solution to the problem of separation in logistic
841 regression. *Stat Med* **21**, 2409-2419, doi:10.1002/sim.1047 (2002).
- 842 60 Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and
843 Powerful Approach to Multiple Testing. *J. R. Statist. Soc. B* **57**, 289-300,
844 doi:10.1111/j.2517-6161.1995.tb02031.x (1995).
- 845 61 Team, R. C. R: A Language and Environment for Statistical Computing. (2020).

- 846 62 Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag New
847 York, 2016).
- 848 63 Slowikowski, K. *ggrepel: Automatically Position Non-Overlapping Text Labels with*
849 *'ggplot2'*. (2020).
- 850 64 Kassambara, A. *ggpubr: 'ggplot2' Based Publication Ready Plots*. (2020).
- 851 65 Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in
852 multidimensional genomic data. *Bioinformatics* **32**, 2847-2849,
853 doi:10.1093/bioinformatics/btw313 (2016).
- 854 66 Chen, H. & Boutros, P. C. VennDiagram: a package for the generation of highly-
855 customizable Venn and Euler diagrams in R. *BMC Bioinformatics* **12**, 35,
856 doi:10.1186/1471-2105-12-35 (2011).
- 857 67 Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic*
858 *Acids Res* **47**, D941-D947, doi:10.1093/nar/gky1015 (2019).
- 859 68 Lawrence, M. S. *et al.* Discovery and saturation analysis of cancer genes across 21
860 tumour types. *Nature* **505**, 495-501, doi:10.1038/nature12912 (2014).
- 861 69 Bailey, M. H. *et al.* Comprehensive Characterization of Cancer Driver Genes and
862 Mutations. *Cell* **173**, 371-385 e318, doi:10.1016/j.cell.2018.02.060 (2018).
- 863 70 Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the
864 integration of genomic datasets with the R/Bioconductor package biomaRt. *Nat Protoc*
865 **4**, 1184-1191, doi:10.1038/nprot.2009.97 (2009).
- 866 71 Yu, G., Wang, L. G., Han, Y. & He, Q. Y. clusterProfiler: an R package for comparing
867 biological themes among gene clusters. *OMICS* **16**, 284-287,
868 doi:10.1089/omi.2011.0118 (2012).
- 869

870 **Acknowledgements**

871 K.C. is supported by a Mahidol-Liverpool PhD scholarship from Mahidol University,
872 Thailand, and the University of Liverpool, UK. J.P.M. is grateful to funding from the Wellcome
873 Trust (208375/Z/17/Z) and the Biotechnology and Biological Sciences Research Council
874 (BB/R014949/1). This work was supported by the Francis Crick Institute, which receives its
875 core funding from Cancer Research UK (FC001202), the UK Medical Research Council
876 (FC001202), and the Wellcome Trust (FC001202). P.V.L. is a Winton Group Leader in
877 recognition of the Winton Charitable Foundation's support towards the establishment of The
878 Francis Crick Institute. We wish to thank members of the Integrative Genomics of Ageing
879 Group for suggestions and discussion.

880

881

882

883 **Author contributions**

884 K.C., T.L., P.V.L. and J.P.M. conceived the project and designed the study. T.L. and P.V.L.
885 provided data. K.C. performed the analyses with helps from T.L. T.L., L.P., P.V.L. and J.P.M.
886 provided critical insights and were involved in data interpretation. K.C. wrote the first draft of
887 the manuscript. All authors edited and approved the manuscript.

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908 **Table 1** Summary of TCGA cancer type and number of samples used in each analysis

Cancer type	Abbreviation	GI, LOH, WGD	SCNAs	Mutations (hypermutated tumours removed)	Pathway alterations	Gene expression	DNA methylation
Adrenocortical carcinoma	ACC	89	89	89 (88)	76	77	78
Bladder Urothelial Carcinoma	BLCA	370	369	369 (364)	361	366	370
Breast invasive carcinoma	BRCA	1015	1011	954 (946)	922	1011	719
Cervical squamous cell carcinoma and endocervical adenocarcinoma	CESC	287	287	271 (263)	264	284	287
Cholangiocarcinoma	CHOL	35	35	35 (35)	35	35	35
Colon adenocarcinoma	COAD	411	411	374 (319)	323	410	278
Lymphoid Neoplasm Diffuse Large B-cell Lymphoma	DLBC	42	42	32 (32)	32	42	42
Oesophageal carcinoma	ESCA	176	176	176 (174)	165	175	176
Glioblastoma multiforme	GBM	489	489	356 (354)	116	137	259
Head and Neck squamous cell carcinoma	HNSC	489	489	472 (469)	459	481	489
Kidney Chromophobe	KICH	66	66	66 (66)	65	66	66
Kidney renal clear cell carcinoma	KIRC	496	496	343 (343)	331	493	296
Kidney renal papillary cell carcinoma	KIRP	228	228	222 (222)	215	228	213
Acute Myeloid Leukaemia	LAML	126	121	55 (54)	101	102	121
Brain Lower Grade Glioma	LGG	488	488	484 (484)	482	488	488
Liver hepatocellular carcinoma	LIHC	355	355	342 (340)	334	349	355
Lung adenocarcinoma	LUAD	460	460	456 (438)	446	456	402
Lung squamous cell carcinoma	LUSC	460	460	444 (437)	426	457	336
Mesothelioma	MESO	82	82	77 (77)	77	82	82
Ovarian serous cystadenocarcinoma	OV	556	556	397 (395)	173	288	545
Pancreatic adenocarcinoma	PAAD	133	133	130 (129)	113	127	132
Pheochromocytoma and Paraganglioma	PCPG	165	157	165 (164)	154	165	165
Prostate adenocarcinoma	PRAD	434	434	434 (432)	425	434	434
Rectum adenocarcinoma	READ	152	152	132 (127)	109	151	95
Sarcoma	SARC	229	229	213 (211)	209	227	229
Skin Cutaneous Melanoma	SKCM	434	434	432 (340)	332	433	434
Stomach adenocarcinoma	STAD	388	388	385 (345)	340	365	341
Testicular Germ Cell Tumours	TGCT	129	129	124 (124)	123	129	129
Thyroid carcinoma	THCA	260	260	249 (248)	244	259	258
Thymoma	THYM	76	76	76 (76)	73	73	76
Uterine Corpus Endometrial Carcinoma	UCEC	434	434	421 (366)	406	432	360
Uterine Carcinosarcoma	UCS	52	52	52 (51)	52	52	52
Uveal Melanoma	UVM	72	72	72 (72)	72	72	72
Total		9678	9660	8899 (8585)	8055	8946	8414

909

910 **Figure legends**

911 **Fig. 1** Association between cancer patients' age and genomic instability (GI) score, percent
912 genomic loss-of-heterozygosity (LOH) and whole-genome duplication events (WGD). **a**
913 Association between age and pan-cancer GI score. Dots are coloured by cancer type. Multiple
914 linear regression R-squared and p-value are shown in the figure. **b** Association between age
915 and cancer type-specific GI score. Linear regression coefficients and significant values are
916 shown in the figure. Cancers with a significant positive association between age and GI score
917 after using multiple linear regression (adj. p-value < 0.05) are highlighted in red. Cancers with
918 a significant association in simple linear regression but not significant after using multiple
919 linear regression are showed in black. The grey line indicates adj. p-value = 0.05. Dot size is
920 proportional to median GI score. **c** Association between age and pan-cancer percent genomic
921 LOH. Dots are coloured by cancer type. Multiple linear regression R-squared and p-value are
922 shown in the figure. **d** Association between age and cancer type-specific percent genomic
923 LOH. Simple linear regression results are shown. Cancers with a significant positive and
924 negative association between age and percent genomic LOH after using multiple linear
925 regression are highlighted in red and blue, respectively. Cancer with a significant association
926 in simple linear regression but not significant after using multiple linear regression is showed
927 in black. The grey line indicates adj. p-value = 0.05. Dot size is proportional to median percent
928 genomic LOH. **e** Association between age and WGD events in pan-cancer, OV, and UCEC.
929 Multiple logistic regression p-values were indicated in the figure.

930

931 **Fig. 2** Association between cancer patients' age and somatic copy-number alterations
932 (SCNAs). Volcano plot representing the association between age and **(a)** overall, **(b)** focal-
933 level and **(c)** chromosome/arm-level SCNA scores. Linear regression coefficients and
934 significant values are shown. Cancers with a significant positive and negative association

935 between age and SCNA score after using multiple linear regression (adj. p-value < 0.05) are
936 highlighted in red and blue, respectively. Cancers with a significant association in simple linear
937 regression but not significant after using multiple linear regression are showed in black. The
938 grey line indicates adj. p-value = 0.05. Dot size is proportional to median SCNA score. **d, e**
939 The left and right dot plots show the association between age and arm-level copy-number gains
940 and copy-number deletions. Circle size corresponds to the significant level, red and blue
941 represent positive and negative associations, respectively. **f, g** Heatmaps represent recurrently
942 gain and deletion arms in LGG and UCEC, respectively. Samples are sorted by age. Colours
943 represent copy-number changes from GISTIC2.0, blue denotes deletion and red corresponds
944 to gain.

945

946 **Fig. 3** Association between cancer patients' age and focal-level SCNAs. **a** Number of gained
947 and deleted focal regions that showed a significant association with age per cancer type
948 (multiple logistic regression, adj. p-value < 0.05). Heatmap showing age-associated focal-level
949 SCNAs in **(b)** LGG and **(c)** UCEC. Samples are sorted by age. Colours represent copy-number
950 changes from GISTIC2.0, blue denotes deletion and red corresponds to gain. The gain_or_loss
951 legend demonstrates that the region is recurrently gained or deleted. The direction legend
952 shows whether the gain/deletion of the region increases or decreases with age. **d** Age-
953 associated SCNA changes in cancer driver genes. Cancer driver genes located in the age-
954 associated focal regions are plotted by cancer type. Colours of the dot represent the condition
955 of the focal region where the gene located in as follows: blue - decrease deletion; green -
956 increase deletion; yellow - decrease gain; and red - increase gain with age. **e** The effect of copy-
957 number changes on gene expression of *CDKN2A* in LGG, *MYC* in OV, *CREBBP* and *RIT1* in
958 UCEC. These are examples of genes with age-associated changes in SCNAs. Violin plots show

959 the $\log_2(\text{normalized expression} + 1)$ of samples grouped by their SCNA status. Pearson
960 correlation coefficient r and p -value are shown in the figures.

961

962 **Fig. 4** Association between cancer patients' age and somatic mutations. **a** The proportion of
963 hypermutated tumours ($>1,000$ mutations/exome) in young (age ≤ 50) and old (age > 50)
964 UCEC. The statistical significant (p -value) was calculated using Fisher's exact test. **b** The
965 association between age and MSI-H in UCEC. The statistical significance was calculated from
966 the multiple logistic regression adjusting for clinical variables. The p -value is shown in the
967 figure. **c** The association between age and *POLE/POLD1* mutations in UCEC. The statistical
968 significance (p -value) was calculated from the multiple logistic regression adjusting for clinical
969 variables. **d** A pan-cancer association between age and mutations. Multiple logistic regression
970 coefficient and significant values are shown. Genes with a significant positive and negative
971 association between age and somatic mutations after using multiple logistic regression (adj. p -
972 value < 0.05) are highlighted in red and blue, respectively. **e** Summary of the cancer type-
973 specific association between age and mutations. Multiple logistic regression coefficient and
974 significant values are shown. Only genes with a significant association (adj. p -value < 0.05)
975 are shown in the figure. A colour code is provided to denote the cancer type where the
976 association between age and gene mutation was found. **f** Heatmap showing age-associated
977 mutations in GBM and LGG. Samples are sorted by age. Colours represent types of mutation.
978 The right annotation legend indicates the direction of change, increase or decrease mutations
979 with age. The mutational burden of samples is presented in the dot above the heatmap.

980

981 **Fig. 5** Association between cancer patients' age and oncogenic signalling pathway alterations.
982 **a** Association between age and oncogenic pathway alterations in the pan-cancer level. Multiple
983 logistic regression coefficients and significant values are shown. Pathways with a significant

984 positive association between age and alterations (adj. p-value < 0.05) are highlighted in red. **b**
985 Cancer-specific age-associated pathway alterations. Pathways that show a significant positive
986 and negative association with age per cancer type (multiple logistic regression, adj. p-value <
987 0.05) are displayed in red and blue dots, respectively. **c** Heatmap showing age-associated
988 alterations in genes associated with TP53 and cell cycle pathways in LGG. Samples are sorted
989 by age. Colours represent types of alteration.

990

991 **Fig. 6** Age-related gene expression in cancers was controlled by age-related methylation. **a**
992 Number of age-DEGs and age-DMGs across cancer types. Red dots represent up-regulated
993 genes, while blue dots denote down-regulated genes. The dot size corresponds to the number
994 of genes. **b** Venn diagrams of the overlap between age-DEGs and age-DMGs. LGG and BRCA
995 are shown as examples. Venn diagrams of the other cancers are shown in Supplementary Fig.
996 9. **c** The distribution of overlap genes between age-DMGs and age-DEGs. The genes were
997 classified into (1) down-regulated methylation and down-regulated expression, (2) down-
998 regulated methylation and up-regulated expression, (3) up-regulated methylation and down-
999 regulated expression, and (4) up-regulated methylation and up-regulated expression. **d** Violin
1000 plots showing the distribution of the Pearson correlation coefficient between methylation and
1001 gene expression in LGG and BRCA. Genes were grouped into (1) common genes between age-
1002 DMGs and age-DEGs (age-DMGs-DEGs), (2) age-DMGs only genes, (3) age-DEGs only
1003 genes, and (4) other genes. The group comparison was performed by the Kruskal-Wallis test.
1004 The pairwise comparisons were done using Dunn's test. P-values from Dunn's test between
1005 age-DMGs-DEGs and the other groups are shown. The plots for the other cancers are shown
1006 in Supplementary Fig. 10. **e** The enriched gene ontology (GO) terms identified by GSEA in
1007 LGG and BRCA. The dot size corresponds to a significant level. A GO term was considered
1008 significantly enriched term if adj. p-value < 0.05 for gene expression and adj. p-value < 0.1 for

1009 methylation. Colours represent enrichment scores, red denotes positive score (enriched in older
1010 patients), while blue signifies negative score (enriched in younger patients). The plots for the
1011 other cancers are shown in Supplementary Fig. 11.

1012

1013 **Supplementary Fig. 1** Association between age and (a) GI score and (b) percent genomic
1014 LOH. Multiple linear regression was performed to identify the relationship between age and
1015 GI score or percent genomic LOH for each cancer type. Cancer types with a significant
1016 association (adj. p-value < 0.05) are shown together with adjusted R-squared and p-values from
1017 multiple linear regression analysis.

1018

1019 **Supplementary Fig. 2** Association between age and (a) overall SCNA score, (b)
1020 chromosome/arm-level SCNA score, and (c) focal-level SCNA score. Multiple linear
1021 regression was performed to identify the relationship between age and SCNA score for each
1022 cancer type. Cancer types with a significant association (adj. p-value < 0.05) are shown
1023 together with adjusted R-squared and p-value from multiple linear regression analysis.

1024

1025 **Supplementary Fig. 3** Heatmaps represent recurrent gain and deletion of arms across cancer
1026 types. Samples are sorted by age. Colours represent copy-number changes from GISTIC2.0,
1027 blue denotes deletion and red corresponds to gain.

1028

1029 **Supplementary Fig. 4** Heatmaps represent recurrent gain and deletion of focal-regions that
1030 showed the age-associated patterns across cancer types. Samples are sorted by age. Colours
1031 represent copy-number changes from GISTIC2.0, blue denotes deletion and red corresponds
1032 to gain. The direction legend shows whether the gain/deletion of the region increases or
1033 decreases with age.

1034

1035 **Supplementary Fig. 5** Association between age and mutational burden. Simple linear
1036 regression was performed to investigate the association between age and mutational burden.
1037 Cancer types with a significant association (adj. p-value < 0.05) were further investigated using
1038 multiple linear regression. The figure shows pan-cancer analysis and cancer type-specific
1039 analysis for every cancer that had a significant association in the simple linear regression
1040 analysis. Adjusted R-squared and p-value from multiple linear regression analysis are
1041 displayed.

1042

1043 **Supplementary Fig. 6** Association between age and (a) MSI-H in COAD, READ, and STAD,
1044 and (b) *POLE/POLD1* mutations in cancer types containing the mutations in these genes in
1045 more than 5% of the samples. The statistical significance (p-value showed) was calculated from
1046 the multiple logistic regression adjusting for clinical variables.

1047

1048 **Supplementary Fig. 7** Heatmap showing age-associated mutations in 11 cancer types.
1049 Samples are sorted by age. Colours represent types of mutation. The right annotation legend
1050 indicates the direction of change, increase or decrease mutations with age. The mutational
1051 burden of samples is presented in the dot above the heatmap.

1052

1053 **Supplementary Fig. 8** Pearson correlation between linear regression coefficient of age on
1054 DNA methylation level and linear regression coefficient of age on gene expression. The
1055 regression coefficients were obtained from the multiple linear regression analysis to investigate
1056 the association between age and DNA methylation or gene expression. Pearson correlation
1057 coefficient and p-values are shown.

1058

1059 **Supplementary Fig. 9** Venn diagrams of the overlap between age-DEGs and age-DMGs in 8
1060 cancer types. Age-DEGs were separated into genes up-regulated with age (Expr_Up) and genes
1061 down-regulated with age (Expr_Down). Age-DMGs were classified into genes with increased
1062 methylation with age (Methy_Up) and genes with decreased methylation with age
1063 (Methy_Down).

1064

1065 **Supplementary Fig. 10** Pearson correlation coefficient between methylation and gene
1066 expression. **a** Violin plots showing the distribution of the Pearson correlation coefficient
1067 between methylation and gene expression in 8 cancer types. Genes were grouped into (1)
1068 overlapping genes between age-DMGs and age-DEGs (age-DMGs-DEGs), (2) age-DMGs
1069 only genes, (3) age-DEGs only genes, and (4) other genes. The group comparison was
1070 performed by the Kruskal-Wallis test. The pairwise comparisons were done using Dunn's test.
1071 P-values from Dunn's test between age-DMGs-DEGs and the other groups are shown. **b**
1072 Density plots showing the distribution of the Pearson correlation coefficient between
1073 methylation and gene expression across 10 cancer types.

1074

1075 **Supplementary Fig. 11** The enriched age-related gene ontology (GO) terms identified by
1076 GSEA in 8 cancer types. The dot size corresponds to a significant level. A GO term was
1077 considered significantly enriched if adj. p-value < 0.05 for gene expression and adj. p-value <
1078 0.1 for methylation. Colours represent enrichment scores, red denotes positive score (enriched
1079 in older patients), while blue signifies negative score (enriched in younger patients). No
1080 enriched term was identified from the gene expression data of UCEC.

1081

1082

1083

1084 **List of supplementary tables**

1085 **Supplementary Table 1:** Summary of the number of samples and clinical variables used in
1086 the study

1087 **Supplementary Table 2:** Association between age and GI scores

1088 **Supplementary Table 3:** Association between age and percent genomic LOH

1089 **Supplementary Table 4:** Association between age and WGD

1090 **Supplementary Table 5:** Association between age and SCNA scores

1091 **Supplementary Table 6:** Association between age and arm-level SCNAs

1092 **Supplementary Table 7:** Association between age and focal-level SCNAs

1093 **Supplementary Table 8:** List of previously identified cancer driver genes

1094 **Supplementary Table 9:** Association between age and mutational burden

1095 **Supplementary Table 10:** Association between age and somatic mutations

1096 **Supplementary Table 11:** Association between age and oncogenic signalling pathway

1097 **Supplementary Table 12:** Gene expression changes with age

1098 **Supplementary Table 13:** DNA methylation changes with age

1099 **Supplementary Table 14:** Number of overlapping genes between age-DEGs and age-DMGs

1100 **Supplementary Table 15:** Comparison of the correlation between methylation and gene
1101 expression in age-DMGs-DEGs, age-DMGs, age-DEGs and other genes

1102 **Supplementary Table 16:** List of enriched GO terms identified using GSEA

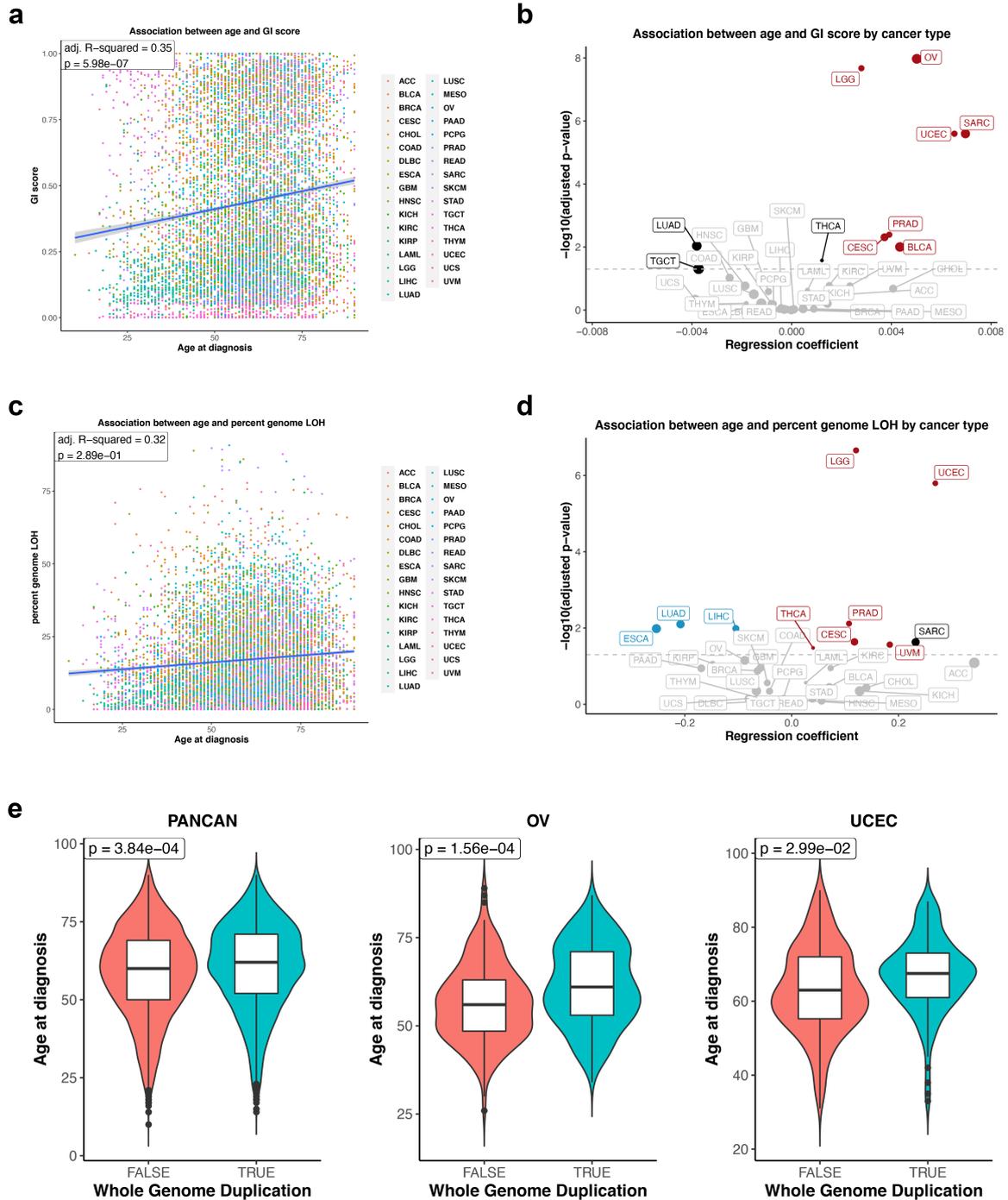


Figure 1

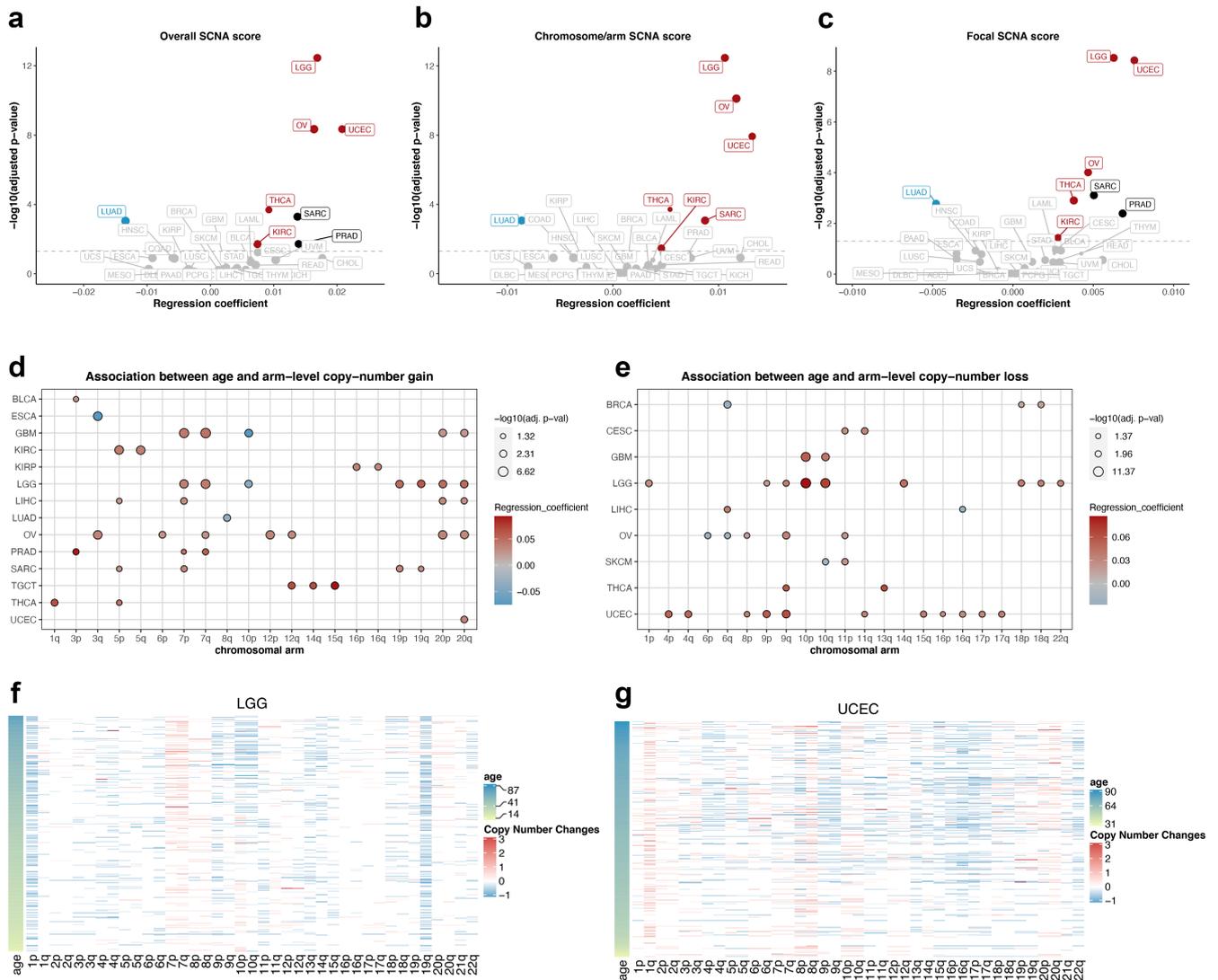


Figure 2

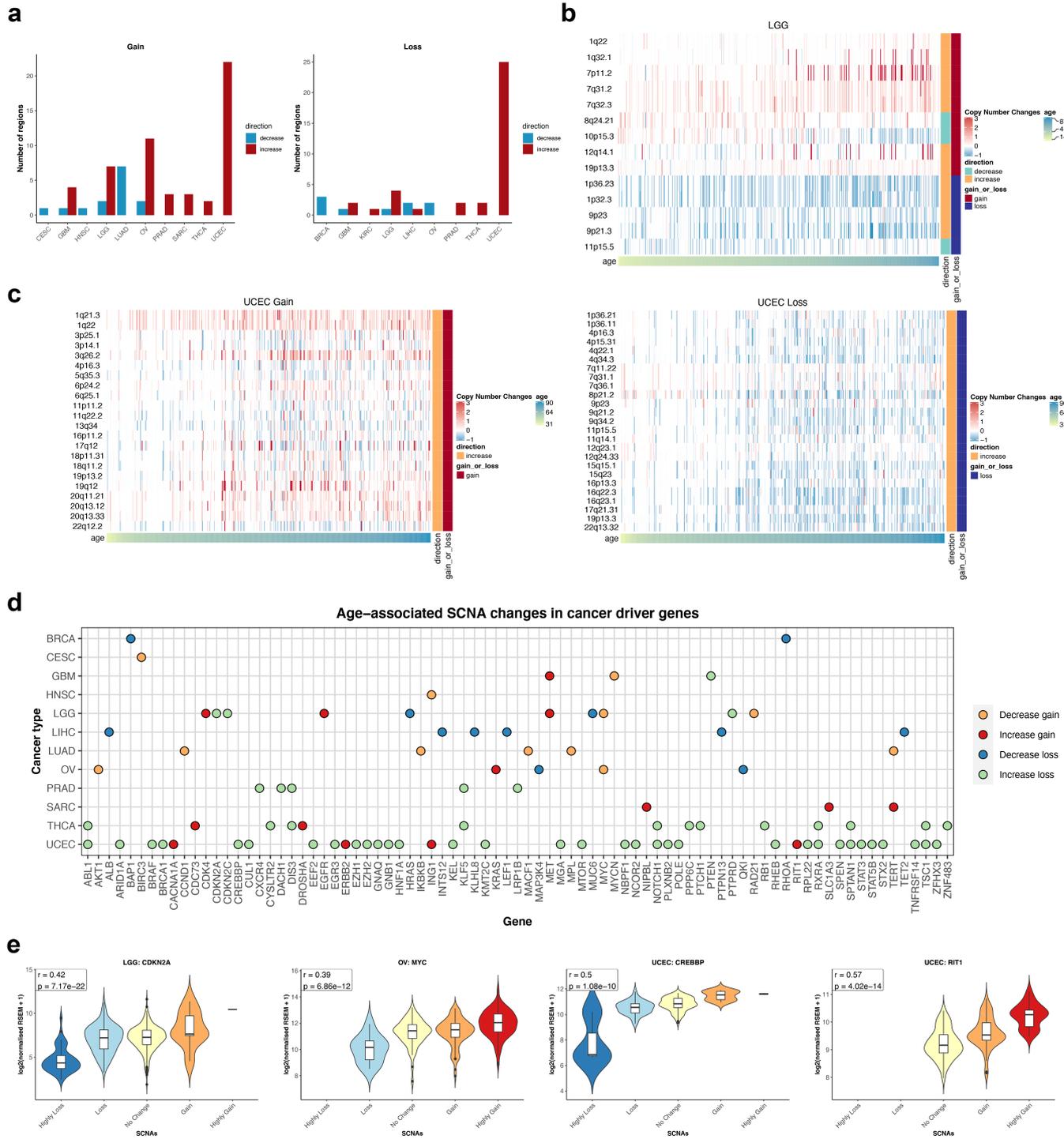


Figure 3

