

1 **Germline cancer gene expression quantitative trait loci influence local and global tumor**
2 **mutations**

3 Yuxi Liu^{1,2,3}, Alexander Gusev³, Peter Kraft^{1,2,4,*}

4 1. Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA,
5 02115, USA

6 2. Program in Genetic Epidemiology and Statistical Genetics, Harvard T.H. Chan School of
7 Public Health, Boston, MA, 02115, USA

8 3. Department of Medical Oncology, Dana-Farber Cancer Institute and Harvard Medical
9 School, Boston, MA, 02215, USA

10 4. Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA,
11 02115, USA

12 * **Corresponding author:** Peter Kraft, PhD. Program in Genetic Epidemiology and Statistical
13 Genetics, Harvard T.H. Chan School of Public Health, 655 Huntington Avenue, Boston, MA,
14 02115. Email: pkraft@hsph.harvard.edu

15 **Running title:** Germline cancer gene eQTL influence tumor mutations

16 **Keywords:** germline variant, somatic mutation, cancer driver gene, TMB, gene expression

17 **Conflict of interest:** The authors declare no potential conflicts of interest.

18

19 Abstract

20 Somatic mutations drive cancer development and are relevant to patients' response to treatment.
21 Emerging evidence shows that variations in the somatic genome can be influenced by the
22 germline genetic background. However, the mechanisms underlying these germline-somatic
23 associations remain largely obscure. We hypothesized that germline variants can influence
24 somatic mutations in a nearby cancer gene ("local impact") or a set of recurrently mutated cancer
25 genes across the genome ("global impact") through their regulatory effect on gene expression.
26 We integrated tumor targeted sequencing from 12,413 patients across 11 cancer types in the
27 Dana-Farber Profile cohort with germline cancer gene expression quantitative trait loci (eQTL)
28 from the Genotype-Tissue Expression Project. We identified variants that upregulate *ATM*
29 expression which are also associated with a decreased risk of having somatic *ATM* mutations
30 across 8 cancer types ($P = 3.43 \times 10^{-5}$). We also identified *GLI2*, *WRN*, and *CBFB* eQTL that are
31 associated with global tumor mutational burden of cancer genes in ovarian cancer, glioma, and
32 esophagogastric carcinoma, respectively ($P < 3.45 \times 10^{-6}$). An *EPHA5* eQTL was associated
33 with the number of mutations in cancer genes specific to colorectal cancer, and eQTL associated
34 with expression of *APC*, *WRN*, *GLII*, *FANCA*, and *TP53* were associated with mutations in
35 genes specific to endometrial cancer ($P < 1.73 \times 10^{-5}$). Our findings provide evidence for the
36 germline-somatic associations mediated through expression of specific cancer genes and open
37 new avenues for research on the underlying biological processes, especially those related to
38 immunotherapy responses.

39 Introduction

40 Cancer is a genetic disease driven by somatic events occurring in the genome over time.
41 Identifying genes carrying driver mutations (cancer driver genes) and elucidating their roles in
42 the related signaling pathways have become primary goals in cancer genomic research because
43 of the contribution of these genetic changes to abnormal and uncontrolled cell growth and
44 transformation which drive the development of a malignant tumor (1-4). Many of these driver
45 genomic alterations have been found to be clinically actionable drug or therapeutic targets for
46 precision medicine. With the advancement of low-cost, high-throughput next-generation
47 sequencing (NGS) technologies, genomic profiling of tumors using targeted NGS panels is
48 becoming part of routine cancer care (5-8).

49 Different cancers are characterized by different patterns of somatic mutations (9,10). Even
50 patients with the same cancer may have substantial heterogeneity in the overall tumor mutational
51 burden (TMB), mutation patterns characterized by mutational signatures, or the cancer genes and
52 oncogenic signaling pathways altered (4,11-14). These heterogeneities in the somatic mutational
53 profile can lead to differential cancer progression, prognosis, and treatment response (15,16). A
54 well-known example is the predictive association of TMB and response to immunotherapy (17).
55 Mounting evidence suggests that somatic variations in tumors can have a germline genetic basis
56 (12,18-23). This germline-somatic relationship has been established at different levels, from the
57 impact of a single germline variant on somatic mutation rate of a cancer gene (e.g., rs25673 at
58 19p13.3 with *PTEN* alterations that involved in the mTOR signaling pathway) (20), to the
59 associations between germline polygenic risk scores (PRS) and somatic mutational signatures
60 (e.g., germline PRS of inflammatory bowel disease with APOBEC signatures in breast cancer)
61 (23). Emerging evidence also shows an interactive effect of germline and somatic variations on

62 clinical outcomes (24). However, the study of germline-somatic interactions is still at an early
63 stage and the mechanisms responsible for these observed associations are still largely uncovered.

64 Germline variants may affect somatic mutations through gene expression (19,22,25). In the well-
65 established example of the APOBEC mutational process, rs17000526-A allele in the *APOBEC3B*
66 region is associated with higher expression of this gene, which contributes to somatic
67 mutagenesis of APOBEC signatures in bladder tumor (19). Chen et al. systematically assessed
68 the impacts of expression level of putative cancer-susceptibility genes on mutational signatures
69 and TMB and identified a wide range of associations across six cancer types (25). Many
70 underlying mechanisms may co-exist, but an intuitive and interpretable hypothesis would be that
71 the germline cancer gene expression quantitative trait loci (eQTL) alter the propensity of
72 acquiring somatic mutations in those specific genes or globally through their regulatory effect on
73 gene expression. Although prior studies included gene expression information in the analysis of
74 germline-somatic interactions, this is no systematic study focusing on both the local and global
75 impact of eQTL on somatic mutations in cancer genes across multiple cancers. Many latent
76 associations and mechanisms may thus have been missed.

77 Here, we performed a pan-cancer analysis of the germline genetic impacts on both the local and
78 global tumor mutations, incorporating regulatory information of germline variants on gene
79 expression. Specifically, we evaluated the associations between germline cancer gene eQTL and
80 i) somatic mutation status of those cancer genes or any hotspot mutation in those genes, ii) tumor
81 mutation counts (TMC) of all recurrently mutated cancer genes for a cancer type, and iii) TMB
82 of all targeted cancer genes from the OncoPanel sequencing platform across 11 cancer types in
83 the Dana-Farber Profile cohort. Clinical targeted sequencing cohorts are well suited for such
84 germline-somatic analysis because the tumor sequencing specifically targets those actionable

85 cancer drivers and the cancer patient population is usually large, unselected, and has extensive
86 clinical data. Our results demonstrate evidence for germline-somatic associations that are
87 potentially mediated through cancer gene expression and provide insights into the mechanisms of
88 mutagenesis in somatic cells.

89 **Materials and Methods**

90 **Study population**

91 The Dana-Farber Profile, initiated in 2011, is a cohort study of unselected cancer patients who
92 presented at the Dana-Farber Cancer Institute, Brigham and Women's Hospital, or Boston
93 Children's Hospital, received genomic profiling and consented to participate. Tumor specimens,
94 mainly formalin-fixed paraffin-embedded tissues, were retrieved from all consented patients for
95 targeted sequencing. Comprehensive clinical and pathologic data were collected along with the
96 genomic data (6,26). The study protocol was approved by the institutional review board (IRB) of
97 Dana-Farber/Partners Cancer Care Office for the Protection of Research Subjects (11-104/17-
98 000). Secondary analyses of previously collected data were approved by the Dana-Farber IRB
99 (19-033/19-025).

100 **Tumor targeted sequencing**

101 A workflow of the full data generating and processing pipeline is present in **Fig. 1**. All collected
102 tumor samples were sequenced on OncoPanel, a targeted NGS platform designed for detecting
103 somatic variations in a panel of actionable cancer genes. There are three versions of the panel
104 targeting the exon and/or intron regions of 304, 326, and 447 genes, respectively; each patient in
105 the cohort was sequenced on one of the panels (Supplementary Table S1). All targeted genes
106 were previously identified oncogenes or tumor suppressor genes involved in cancer-related
107 signaling pathways (27). Sequencing was performed using an Illumina HiSeq 2500 with 2×100
108 paired-end reads followed by somatic mutation calling using MuTect (for single-nucleotide
109 variants) (28) and Indelocator (for indels; <http://www.broadinstitute.org/cancer/cga/indelocator>)

110 from reads aligned to the targeted genome regions with $> 50\times$ reads (“On-target reads”). More
111 details about the tumor sequencing pipeline can be found in prior studies (6,27).
112 We collected somatic mutation data from the tumor sequences of 18,472 primary cancer samples
113 spanning over 60 cancer types and subtypes. Some tumors exhibit microsatellite instability
114 (MSI) with high mutational burden; the germline-somatic relationship for those hypermutable
115 subtypes might be substantially different from the microsatellite stable (MSS) tumors. We thus
116 further classified each sample as MSI or MSS using MSIDetect (29). Cancer types with > 500
117 samples were selected; for each selected cancer, we removed those rare subtypes with < 3
118 samples. The remaining 12,413 samples across 11 cancer types were included in the downstream
119 analysis (Supplementary Table S2).

120 **Germline imputation from tumor sequences**

121 Details of inferring common germline variants from the OncoPanel tumor sequencing data are
122 described elsewhere (26) and briefly summarized here. Tumor targeted sequencing generated
123 both high-coverage “on-target reads” aligned to the targeted regions and low-coverage “off-
124 target reads” aligned to the rest of the genome (**Fig. 1**). Common germline variants with $> 1\%$
125 frequency in the European population were imputed from these tumor sequences (mainly relied
126 on off-target reads) using linkage disequilibrium (LD) information with the 1000 Genomes Phase
127 3 release as the haplotype reference panel. Imputation accuracies from several algorithms
128 designed for imputing germline variants from low coverage data were evaluated by comparing
129 the imputed allele dosage to the gold standard germline data generated from genotyping array.
130 The STITCH algorithm (30) yielded the highest overall accuracy and the resulting imputed
131 germline data were used for the downstream analysis. The imputed variants were subsequently

132 restricted to an imputation INFO score of greater than 0.4, which produced a mean imputation
133 correlation of 0.86 between tumor imputed and germline SNP array variants (26).

134 Genetic ancestry was inferred by projecting the imputed germline genetic data into the genetic
135 ancestry principal components using weights derived for European, African, and Asian
136 populations from the 1000 Genomes Project reference data (31). We further restricted our
137 analysis to samples with < 10% inferred non-European ancestry.

138 **Identifying recurrently mutated cancer genes and hotspot mutations**

139 We identified recurrently mutated cancer genes, defined as genes with $\geq 5\%$ carriers of missense
140 mutations, for each selected cancer type from the somatic data. Not all panel genes were
141 sequenced on every sample (multiple panel versions exist); we thus further excluded those
142 identified gene-cancer pairs with < 500 sequenced samples. We included additional genes that
143 were identified as highly significantly mutated or significantly mutated genes among known
144 cancer genes for each selected cancer type from the TumorPortal (<http://www.tumorportal.org/>)
145 (32). A total of 135 cancer genes and 342 gene-cancer pairs were identified, with the mutation
146 frequency ranging from 0.0036 to 0.73 (**Fig. 2A**; Supplementary Table S3). Mutation status for
147 each sample and gene is defined as whether this sample carries at least one functional mutation
148 (frame_shift_del, frame_shift_ins, frameshift, initiator_codon, missense and splice_region,
149 missense_mutation, nonsense_mutation, protein_altering, splice_site, start_lost, stop_lost, and
150 translation_start_site) in this gene and is considered to capture the “local” tumor mutation
151 (mutation in one cancer gene).

152 For each selected cancer type, we further identified specific mutations with $\geq 5\%$ carriers in the
153 somatic data as hotspot mutations. Seven of the 11 cancer types harbor at least one hotspot

154 mutation. A total of 17 hotspot mutations and 25 mutation-cancer pairs were identified, with the
155 mutation frequency ranging from 0.051 to 0.33 (**Fig. 2B**; Supplementary Table S4). A binary
156 variable of the local mutation status was created to indicate whether a sample carries a specific
157 hotspot mutation.

158 **Quantifying TMB of all panel genes and TMC of recurrently mutated cancer genes**

159 TMB is defined as the total number of missense mutations per megabase based on the targeted
160 sequencing data of all panel genes (**Fig. 2C**). It captures the total mutations in all targeted cancer
161 genes and is considered as a refined “global” mutational burden restricting to a set of cancer-
162 related genes rather than the genome-wide mutational burden. In addition to TMB, we also
163 calculated TMC for each sample, which is defined as the count of recurrently mutated cancer
164 genes (specific to each cancer type) that harbor at least one missense mutation. The number of
165 identified recurrently mutated cancer genes varies across cancer types (**Fig. 2D**). Compared to
166 TMB, TMC is a more refined measure of the mutational burden in likely driver genes for a
167 cancer. Moreover, by counting the genes instead of the mutations, the TMC analysis would be
168 less sensitive to hypermutable outliers.

169 **Identifying eQTL from the Genotype-Tissue Expression (GTEx) Project for all selected** 170 **genes**

171 We obtained the eQTL and gene expression association results in normal tissue for all selected
172 genes from the meta-analyzed multi-tissue eQTL results using METASOFT (33) from the GTEx
173 Analysis V8 release. We selected those genome-wide significant eQTL with $P < 5 \times 10^{-8}$ from
174 any of the fixed effect (FE), random effect (RE), or Han and Eskin's random effect (RE2)
175 models. Variants with minor allele frequency $< 1\%$ were further removed. A total of 28,486

176 eQTL for 114 genes with imputed germline data available were included in the analysis. We
177 performed LD clumping with $r^2 = 0.3$ on the final list of eQTL for each gene to identify
178 independent loci, which was used to determine the number of effective tests in the association
179 analyses (34).

180 **Assessing the associations of cancer gene eQTL with TMB and TMC**

181 We assessed the association between each selected cancer gene eQTL and TMB of all panel
182 genes for each cancer by fitting a linear model adjusting for age, gender (if applicable), panel
183 version, and tumor purity. MSI status was also adjusted as a covariate for the models of
184 colorectal and endometrial cancer where a substantial proportion of the cases display
185 hypermutability (35,36). TMB was Winsorized to 98% within each cancer type to reduce the
186 impact of potential outliers on the association results. The associations between cancer gene
187 eQTL and TMC were evaluated for recurrently mutated cancer genes for each cancer type by
188 fitting a negative binomial model with the same covariates as the TMB models. Sensitivity
189 analysis was performed to assess the impacts of potential TMB or TMC outliers on the
190 association results by varying the Winsorization thresholds and using standardized TMB. For
191 TMC, we further evaluated the impacts of using count of missense mutations instead of count of
192 mutated genes on the germline-somatic associations.

193 **Assessing the associations of cancer gene eQTL with recurrently mutated cancer genes and** 194 **hotspot mutations**

195 The local impact of each cancer gene eQTL on the risk of having somatic mutations in that gene
196 or a nearby hotspot mutation was assessed using logistic regression. These analyses further
197 adjusted for TMB along with all the covariates included in the TMB or TMC models. Meta-

198 analysis was performed to evaluate the broad impact of a cancer gene eQTL on the mutation
199 status of one gene or mutation across cancers.

200 **Data availability statement**

201 The individual-level data used in this study are not publicly available due to patient privacy
202 requirements. Other unidentifiable data generated in this study are available within the article
203 and its supplementary data files.

204 Results

205 Germline cancer gene eQTL influence global tumor mutations

206 We analyzed 28,486 eQTL for 114 cancer genes and assessed their associations with TMB of all
207 cancer genes sequenced on the panel across cancers. There were 1,317 independent eQTL ($r^2 <$
208 0.3) after LD clumping. We identified 22 significant eQTL-TMB associations representing 3
209 independent gene-cancer pairs that passed the Bonferroni correction threshold accounting for the
210 number of effective tests ($P < 3.45 \times 10^{-6}$; Supplementary Table S5). **Table 1** summarizes the
211 results for the most significant association at each locus. There exists heterogeneity in the effects
212 of these eQTL on TMB across cancers (Supplementary Table S6). Sensitivity analysis on the
213 impacts of potential outliers showed that the association of the *GLI2* eQTL and TMB in ovarian
214 cancer was sensitive to the changing Winsorization threshold (Supplementary Table S7). This
215 association also became non-significant if we use standardized TMB as the outcome (beta =
216 0.26 , $P = 0.43$) while the other two top associations remained nominally significant (beta =
217 -2.33 , $P = 1.57 \times 10^{-3}$ for rs139944315 (*WRN*) and TMB in glioma; beta = -0.23 , $P = 0.04$ for
218 rs11075646 (*CBFB*) and TMB in esophagogastric carcinoma).

219 To further investigate the relationship between the observed germline-somatic associations and
220 gene expression, we compared our results with the association results between the identified top
221 eQTL and expression level of the specific cancer genes in normal tissue in GTEx (**Table 1; Fig.**
222 **3**). The T allele of rs1530578 was associated with elevated TMB in ovarian cancer and reduced
223 expression of *GLI2* across tissues (**Fig. 3A,D**). The largest effect of rs1530578 on *GLI2*
224 expression was observed in ovary with beta = -0.55 and $P = 3.93 \times 10^{-5}$ (**Fig. 3A**). rs139944315
225 was associated with TMB in glioma and expression of *WRN* across tissues in a consistent

226 direction (**Fig. 3B,E**). While the largest effect of this variant on *WRN* expression was observed in
227 subcutaneous adipose tissue, there was also an association in putamen of basal ganglia with beta
228 = -0.51 for the T allele and $P = 0.05$ (**Fig. 3B**). Finally, we found that the C allele of rs11075646
229 was associated with decreased TMB in esophagogastric carcinoma and slightly increased
230 expression of *CBFB* across tissues (**Fig. 3C,F**). This variant had a nominally significant impact
231 on *CBFB* expression in both gastroesophageal junction (beta = 0.10 and $P = 0.02$, **Fig. 3C**) and
232 mucosa of esophagus (beta = 0.08 and $P = 0.01$) while the most significant effect was observed
233 for thyroid (**Fig. 3F**). None of these three top variants or variants in high LD with them have
234 been linked to cancer incidence in genome-wide association studies (GWAS) from GWAS
235 Catalog (<https://www.ebi.ac.uk/gwas/home>) (37).

236 We next assessed the impacts of cancer gene eQTL on TMC, which quantifies the mutational
237 burden of a set of genes that are recurrently mutated in the specific cancer. There were 145
238 significant eQTL-TMC associations after Bonferroni correction ($P < 1.73 \times 10^{-5}$; Supplementary
239 Table S8), representing six independent gene-cancer pairs (**Table 1**). Sensitivity analysis showed
240 that all top TMC associations were robust to a wide range of Winsorization thresholds
241 (Supplementary Table S9). Replacing count of mutated genes with count of mutations also
242 yielded similar results compared to the main analysis (Supplementary Table S10). Given that all
243 top eQTL-TMC associations were identified for colorectal and endometrial cancer, we further
244 performed a stratified analysis by MSI status. There was no substantial deviation in the effect
245 estimates for MSS or MSI subgroup from the main analysis though the subgroup results were
246 less significant, especially for MSI samples, which was likely due to the reduced sample sizes
247 (Supplementary Table S11). Finally, we compared these top eQTL-TMC associations to the
248 previous eQTL-TMB results and found that all these top germline variants were associated with

249 TMB in the corresponding cancers in consistent directions with TMC with nominal significance
250 (Supplementary Table S12).

251 There exists substantial heterogeneity in the associations with gene expression level across
252 tissues for many of the top variants in the TMC associations (**Fig. 4**). Two of the tissue-specific
253 associations have both $P < 0.05$ and $m\text{-value} > 0.8$: rs10031417 and *EPHA5* expression in
254 sigmoid colon and rs7201264 and *FANCA* expression in uterus (**Fig. 4A,E**). The A allele of
255 rs10031417 was associated with lower somatic mutational burden in recurrently mutated cancer
256 genes in colorectal cancer and slightly higher expression of *EPHA5* across tissues (**Fig. 4A,G**);
257 this positive effect on *EPHA5* expression was larger in sigmoid colon with $\beta = 0.17$ and $P =$
258 1.11×10^{-3} (**Fig. 4A**). It is worth noting that a variant that is in LD with rs10031417
259 (rs13104357, $r^2 = 0.18$) has also been reported to be associated with *EPHA5* expression in
260 colorectal tumor samples in The Cancer Genome Atlas (TCGA) (38); the direction of this
261 association in tumor was consistent with in normal tissue. rs7201264-C allele was associated
262 with both increased TMC in endometrial cancer and decreased *FANCA* expression across tissues
263 (**Fig. 4E,K**); it had a specific significant impact on *FANCA* expression in uterus (**Fig. 4E**; $\beta =$
264 -0.28 for the C allele, $P = 0.02$). rs78378222, that is in LD with the top variant identified for
265 TMC in endometrial cancer and *TP53* expression (rs17884306, $r^2 = 0.21$ with rs78378222), has
266 been previously associated with the risk of uterine fibroids and several cancers in but not the risk
267 of endometrial cancer specifically (39,40).

268 **Local impacts of germline eQTL on somatic mutations in cancer genes**

269 Investigation of the local impacts of eQTL for a cancer gene on somatic mutations in that gene is
270 also of interest as it may point to a direct and testable mechanism of how germline variations

271 modify the susceptibility to somatic events. None of the individual associations between somatic
272 mutation status for recurrently mutated genes and their eQTL passed the Bonferroni correction
273 threshold ($P < 1.73 \times 10^{-5}$). The most significant association observed was between a *TSC2*
274 eQTL and somatic *TSC2* mutation status in endometrial cancer (beta = -1.81 for rs12918530-C
275 allele, $P = 1.56 \times 10^{-4}$; Supplementary Table S13). Looking across all cancers, there was a
276 significant ($P < 6.91 \times 10^{-5}$) association between an *ATM* eQTL (lead SNP: rs4753834 at
277 11q22.3) and somatic *ATM* mutations from a meta-analysis of 8 cancers (**Fig. 5**; Supplementary
278 Table S14). The G allele of rs4753834 was associated with a lower risk of having somatic
279 mutations in *ATM* (beta = -0.35, $P = 3.43 \times 10^{-5}$ across cancers from FE model) and increased
280 expression of *ATM* in normal tissues (beta = 0.05, $P = 1.03 \times 10^{-20}$ across tissues from RE
281 model). This variant also had specific effects on *ATM* expression in many tissues related the 8
282 cancers, including mammary tissue (beta = 0.06), sigmoid colon (beta = 0.09), hypothalamus
283 (beta = 0.12), lung (beta = 0.07), and prostate (beta = 0.11), all with $P < 0.05$ and m-value > 0.9 .
284 Moreover, variants that are in LD with rs4753834 have also been associated with *ATM*
285 expression in tumor samples of breast cancer (rs673281, $r^2 = 0.21$, beta = -0.08 for the T allele,
286 $P = 1.98 \times 10^{-4}$) and glioma (rs1003623, $r^2 = 0.21$, beta = -0.11 for the T allele, $P = 4.56 \times 10^{-4}$)
287 (38); the directions were also consistent with those in normal tissues. We additionally tested the
288 associations of *ATM* eQTL and TMB or TMC of cancer genes and found that variants in LD with
289 rs4753834 (lead SNP: rs672964, $r^2 = 0.21$ with rs4753834) were associated with TMB (beta =
290 -0.69 for rs672964-C, $P = 2.97 \times 10^{-5}$) and TMC (beta = -0.07 for rs672964-C, $P = 0.02$) in
291 non-small cell lung cancer in the consistent direction with *ATM* mutation status. No association
292 with cancer risk was found for rs4753834 or its tagging SNPs in GWAS Catalog.

293 We also identified nominal associations between eQTL for cancer genes identified in the global
294 tumor mutation analysis with the somatic mutation status of that gene in the corresponding
295 cancer. We found that rs1897693 ($r^2 = 0.42$ with rs10031417) was associated with both the
296 expression of *EPHA5* in normal tissues (beta = 0.03 for the C allele, $P = 0.03$ across tissues from
297 RE model) and the somatic mutation status of *EPHA5* in colorectal cancer (beta = -0.66 for the
298 C allele, $P = 0.01$). Another variant rs55671402 was associated with *FANCA* expression in
299 normal tissues (beta = -0.13 for the C allele, $P = 9.67 \times 10^{-13}$ across tissues from RE model; beta
300 = -0.54 for the C allele, m-value = 0.98, $P = 1.35 \times 10^{-3}$ in uterus) and somatic mutations in
301 *FANCA* in endometrial tumors (beta = -1.23 for the C allele, $P = 8.61 \times 10^{-3}$).

302 We further assessed the impacts of eQTL for a cancer gene on each identified hotspot mutation
303 in that gene. None of the associations passed the Bonferroni correction threshold ($P < 3.40 \times$
304 10^{-4}) with the most significant association observed for rs1867930 with p.S249C in *FGFR3* in
305 bladder cancer (beta = 0.60 for the G allele, $P = 3.54 \times 10^{-3}$; Supplementary Table S15). Only
306 one nominally significant ($P < 0.05$) association from the meta-analysis across cancers was
307 found for rs11047823 with p.G12D in *KRAS* across colorectal cancer, endometrial cancer, non-
308 small cell lung cancer, and pancreatic cancer (beta = 0.24 for the G allele, $P = 0.01$ across
309 cancers from the FE model), though it still did not pass the Bonferroni correction threshold for
310 significance ($P < 5 \times 10^{-3}$).

311 Discussion

312 In this study, we systematically evaluated the influence of germline variants that are associated
313 with cancer gene expression on somatic mutations in specific cancer genes across 11 cancer
314 types, leveraging large-scale clinical targeted panel sequencing data, germline data imputed from
315 tumor sequences, and cancer gene eQTL data from GTEx. Our analysis revealed novel
316 associations of germline eQTL for well-established cancer genes with local mutation status of a
317 single cancer gene or the global mutational burden. These findings provide the initial evidence
318 for the hypothesis that germline variants can influence local and global tumor mutations by
319 altering the expression level of specific cancer genes. The underlying molecular mechanisms of
320 the identified associations can be further investigated through functional analysis and in cancer
321 cell lines.

322 Although our findings are consistent with the putative mechanism that germline variants affect
323 somatic mutations through gene expression, there are also other possible scenarios that can yield
324 the same results (**Fig. 6**). First, given that there exists a causal impact of eQTL on somatic
325 mutations, we still cannot conclude that this is only mediated by the transcript abundance of the
326 specific eQTL gene. The germline eQTL may regulate the expression of other genes which
327 contribute to somatic mutagenesis, or they might be associated with somatic mutations through
328 other pathways that are not related to gene expression (**Fig. 6A**). Finding an eQTL signal in the
329 cancer-related tissue can provide further support that gene expression plays a role in the
330 germline-somatic relationship. Second, we are studying somatic mutations in developed tumor
331 (S') rather than in normal or precancerous tissue (S) (**Fig. 6A**). S' can serve as a proxy for S,
332 though it was measured after tumorigenesis and might be further influenced by other factors such
333 as the tumor microenvironment (41). Here, we are studying mutations in cancer genes that have

334 been identified as potential drivers for carcinogenesis. Even if some mutations in those genes
335 occurred after cancer initiation, our results could still inform us of the role of germline variants
336 in inducing somatic mutations during cancer progression. Finally, even when there is no direct
337 causal effect of germline variants on somatic mutations, we may still observe this association
338 among cancer patients. Consider the three possible scenarios in **Fig. 6B** given that a germline-
339 somatic association was observed: germline variants may influence somatic mutation and they
340 may or may not have an effect on cancer diagnosis through other pathways; however, under the
341 situation that the germline variants only influence cancer diagnosis through other pathways and
342 there is no causal effect on somatic mutations, we may still observe this germline-somatic
343 association among cancer patients due to collider bias (**Fig. 6B**). We are unable to distinguish
344 between these three scenarios based on our data, but we can leverage information from other
345 sources (e.g., association results of the germline variants with cancer incidence from GWAS) to
346 weigh these possible scenarios for each identified association.

347 Most of the germline-somatic associations identified here were consistent with prior evidence,
348 and many of them may be involved in the biological mechanisms that underlie patients' response
349 to immunotherapy. Among all the identified eQTL genes, *APC*, *ATM*, *CBFB*, and *TP53* have
350 been predicted as pan-cancer tumor suppressor genes across 33 cancer types in TCGA (1). We
351 observed that the germline variants associated with reduced expression of these tumor suppressor
352 genes were associated with increased tumor mutations, except for *APC* where the eQTL
353 association with gene expression was close to null across tissues (but still significant) with no
354 effect in uterus (**Table 1; Fig. 3-5**). The *APC* gene encodes the adenomatous polyposis coli
355 protein which plays an important role in the Wnt signaling pathway (42) and interacts with E-
356 cadherin, which regulates cell adhesion (43). Mutations that inactivate APC lead to disruption of

357 β -catenin degradation, resulting in its translocation into the nucleus and activation of the
358 transcription of multiple genes, which triggers cancer development, including endometrial
359 carcinogenesis (44). Active β -catenin signaling has been linked to resistance to anti-PD-L1/anti-
360 CTLA-4 monoclonal antibody therapy in melanoma (45). A recent study found that germline
361 pathogenic variants in *APC* are associated with elevated TMB (46). In our work, the minor allele
362 of the lead SNP is also associated with higher TMC of recurrently mutated cancer genes, but the
363 direction of its association with *APC* expression is not clear (**Table 1; Fig. 4**). Intuitively, we
364 would assume a variant that downregulates the expression of a tumor suppressor gene to be
365 associated with elevated risk of cancer and somatic mutational burden, but this assumption might
366 be oversimplified as the oncogenic or tumor suppressive effect of a gene on carcinogenesis and
367 on somatic mutational burden would depend on the signaling pathway that the gene involved in
368 and may vary substantially across cancer types (47). Here, the major allele of rs397768 slightly
369 downregulates *APC* expression across tissues, if this indicates activation of β -catenin signaling in
370 endometrial carcinogenesis, then it should be associated with resistance to immunotherapy and
371 reduced tumor mutations as we observed. However, this interpretation depends upon many
372 variable components involved in this complex biological process; further study is needed to
373 elucidate the molecular mechanisms underlying these associations.

374 *ATM* germline and somatic mutations have been linked to multiple cancers. Activated ATM
375 protein kinase phosphorylates a few key proteins which activates DNA damage checkpoint,
376 leading to its main tumor suppressive effect of cell-cycle arrest and apoptosis (48). A study of
377 pathogenic germline variants in cancer identified two-hits events for *ATM* where a germline
378 variant in *ATM* is coupled with a somatic mutation in the other copy of the gene in multiple
379 cancers (49). They also found that *ATM* pathogenic variant carriers had lower *ATM* expression,

380 which is in line with our finding that the minor allele of rs4753834 is associated with lower
381 expression of *ATM* but higher risk of having somatic mutation in the gene (**Fig. 5**). Recent
382 studies also reported that *ATM* mutations were associated with improved response to immune
383 checkpoint blockade therapy (50,51). We observed this inverse relationship of *ATM* expression
384 with both somatic *ATM* mutations across cancers and TMB in non-small cell lung cancer, which
385 may support the potential role of *ATM* as a therapeutic target for promoting the response to
386 cancer immunotherapy.

387 *TP53*, which encodes protein p53, is one of the most frequently mutated genes in cancer. Genetic
388 alterations in the p53 stress response pathway can affect the tumor suppressive role of *TP53* and
389 promote tumorigenesis (52). Results from a recent study demonstrated evidence for the
390 interactive effects of a germline cancer risk variant, rs78378222, and somatic mutation status of
391 *TP53* on cancer risk, prognosis, and drug responses (24). The C allele of rs78378222 has been
392 linked to lower expression level of wild-type *TP53* in both normal tissue and tumor, which in
393 turns reduce p53 cellular activity and lead to poorer overall survival of patients. In our analysis,
394 we found that the minor allele of rs17884306, which is correlated with the C allele of
395 rs78378222, was associated with higher TMB and lower *TP53* expression (**Table 1; Fig. 3**). One
396 study highlighted the predictive value of somatic *TP53* mutations for benefit from anti-PD-
397 1/PD-L1 immunotherapy in lung cancer (53). Our results may provide further insights into how
398 inherited genetic predisposition can influence patients' response to immunotherapy through its
399 effect on *TP53* expression and somatic mutational burden.

400 Increased expression of *GLI2* in the hedgehog signaling pathway has been found to induce PD-
401 L1 expression in gastric cancer and promote tumor resistance to immunotherapy (54). We
402 identified a germline eQTL at 2q14.2 that upregulates *GLI2* and is associated with lower TMB in

403 ovarian cancer; nominally significant associations were also found in esophagogastric carcinoma
404 and glioma in the same directions (**Table 1; Fig. 3**; Supplementary Table S6). These findings
405 may shed light on the underlying mechanism of the link between TMB and response to
406 immunotherapy in these specific cancers.

407 Reduced *EPHA5* expression has been linked to lymph node metastasis, advanced TNM stage,
408 and poor survival outcome in colorectal cancer, supporting its tumor suppressive role in this
409 cancer (55). Recent work showed that having somatic *EPHA5* mutations is positively associated
410 with TMB and response to immune checkpoint inhibitor therapy in lung cancer (56). We also
411 identified consistent associations of an *EPHA5* eQTL at 4q13.2 with both somatic *EPHA5*
412 mutations and the global tumor mutations in colorectal cancer. This eQTL influences *EPHA5*
413 expression in colorectal cancer and normal colon tissue (**Table 1; Fig. 3**); the allele that was
414 associated with reduced expression was also associated with increased somatic mutations.
415 Further studies are needed to characterize the potential interactive effect of these identified
416 germline variants, *EPHA5* expression, and somatic *EPHA5* mutations in colorectal cancer.

417 Our study has several limitations. First, as mentioned above, we cannot easily distinguish
418 between several possible scenarios of the causal relationships that may be consistent with the
419 observed associations between germline eQTL and tumor mutations. We suggest future studies
420 to further investigate these associations in normal tissue or precancerous lesions and
421 incorporating haplotype-level information. Experimental validation is also necessary to confirm
422 the putative mechanisms through gene expression for the identified associations. Second, the use
423 of germline data imputed from off-target reads in tumor sequencing provides only a probabilistic
424 estimate of the imputed variant. Although the validation analysis of imputed common germline
425 variants against SNP array yielded high accuracy (26), it would still be important to validate

426 these findings using direct germline genotyping. Third, our analysis focused on somatic
427 mutations in the tumor, but we only included eQTL identified from normal tissue, which may
428 miss tumor-specific eQTL effects (19). However, using normal tissue eQTL as the genetic
429 instrument is more consistent with our hypothesis that eQTL alter gene expression in normal
430 tissue contributes to somatic mutagenesis and tumor initiation. Where available, we also cross-
431 referenced the eQTL results to those in corresponding tumor tissue and found the results had
432 consistent direction with those in normal tissue. Finally, we only focused on missense and a few
433 other functional mutations; future studies can further investigate the germline impact on somatic
434 copy number alteration or structural variation through gene regulation.

435 In conclusion, we systematically investigated the impacts of germline cancer gene eQTL on
436 somatic mutations in cancer genes across 11 cancer types. Our results indicate that germline
437 variants that regulate the expression of cancer genes also influence local and global tumor
438 mutations. These findings provide further evidence for the important role of gene expression
439 regulation in germline-somatic associations and open avenues for future research on the
440 molecular mechanisms underlying these associations that confer cancer risk and sensitize cancer
441 to immunotherapy.

442 **Authors' Disclosures**

443 No disclosures were reported.

444 **Acknowledgements**

445 This work was supported by National Cancer Institute grants R01CA227237 and R01CA244569

446 (to A. Gusev), and R01CA260352 (to P. Kraft), Phi Beta Psi Sorority, and Emerson Collective.

447 The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the

448 Office of the Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA,

449 NIMH, and NINDS. The data used for the analyses described in this manuscript were obtained

450 from the GTEx Portal on 08/04/2021.

451

452 References

- 453
- 454 1. Bailey MH, Tokheim C, Porta-Pardo E, Sengupta S, Bertrand D, Weerasinghe A, *et al.*
455 Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell*
456 **2018**;173:371-85 e18
 - 457 2. Martinez-Jimenez F, Muinos F, Sentis I, Deu-Pons J, Reyes-Salazar I, Arnedo-Pac C, *et al.*
458 A compendium of mutational cancer driver genes. *Nat Rev Cancer* **2020**;20:555-72
 - 459 3. Pon JR, Marra MA. Driver and passenger mutations in cancer. *Annu Rev Pathol*
460 **2015**;10:25-50
 - 461 4. Sanchez-Vega F, Mina M, Armenia J, Chatila WK, Luna A, La KC, *et al.* Oncogenic
462 Signaling Pathways in The Cancer Genome Atlas. *Cell* **2018**;173:321-37 e10
 - 463 5. Malone ER, Oliva M, Sabatini PJB, Stockley TL, Siu LL. Molecular profiling for precision
464 cancer therapies. *Genome Med* **2020**;12:8
 - 465 6. Sholl LM, Do K, Shivdasani P, Cerami E, Dubuc AM, Kuo FC, *et al.* Institutional
466 implementation of clinical tumor profiling on an unselected cancer population. *JCI*
467 *Insight* **2016**;1:e87062
 - 468 7. Zehir A, Benayed R, Shah RH, Syed A, Middha S, Kim HR, *et al.* Mutational landscape of
469 metastatic cancer revealed from prospective clinical sequencing of 10,000 patients. *Nat*
470 *Med* **2017**;23:703-13
 - 471 8. Wagle N, Berger MF, Davis MJ, Blumenstiel B, Defelice M, Pochanard P, *et al.* High-
472 throughput detection of actionable genomic alterations in clinical tumor samples by
473 targeted, massively parallel sequencing. *Cancer Discov* **2012**;2:82-93
 - 474 9. Greenman C, Stephens P, Smith R, Dalgliesh GL, Hunter C, Bignell G, *et al.* Patterns of
475 somatic mutation in human cancer genomes. *Nature* **2007**;446:153-8
 - 476 10. Watson IR, Takahashi K, Futreal PA, Chin L. Emerging patterns of somatic mutations in
477 cancer. *Nat Rev Genet* **2013**;14:703-18
 - 478 11. Sha D, Jin Z, Budczies J, Kluck K, Stenzinger A, Sinicrope FA. Tumor Mutational Burden as
479 a Predictive Biomarker in Solid Tumors. *Cancer Discov* **2020**;10:1808-25
 - 480 12. Consortium ITP-CAoWG. Pan-cancer analysis of whole genomes. *Nature* **2020**;578:82-93
 - 481 13. Lawrence MS, Stojanov P, Polak P, Kryukov GV, Cibulskis K, Sivachenko A, *et al.*
482 Mutational heterogeneity in cancer and the search for new cancer-associated genes.
483 *Nature* **2013**;499:214-8
 - 484 14. Alexandrov LB, Kim J, Haradhvala NJ, Huang MN, Tian Ng AW, Wu Y, *et al.* The repertoire
485 of mutational signatures in human cancer. *Nature* **2020**;578:94-101
 - 486 15. Chan TA, Yarchoan M, Jaffee E, Swanton C, Quezada SA, Stenzinger A, *et al.*
487 Development of tumor mutation burden as an immunotherapy biomarker: utility for the
488 oncology clinic. *Ann Oncol* **2019**;30:44-56
 - 489 16. Dancey JE, Bedard PL, Onetto N, Hudson TJ. The genetic basis for cancer treatment
490 decisions. *Cell* **2012**;148:409-20
 - 491 17. Goodman AM, Kato S, Bazhenova L, Patel SP, Frampton GM, Miller V, *et al.* Tumor
492 Mutational Burden as an Independent Predictor of Response to Immunotherapy in
493 Diverse Cancers. *Mol Cancer Ther* **2017**;16:2598-608

- 494 18. Nik-Zainal S, Wedge DC, Alexandrov LB, Petljak M, Butler AP, Bolli N, *et al.* Association of
495 a germline copy number polymorphism of APOBEC3A and APOBEC3B with burden of
496 putative APOBEC-dependent mutations in breast cancer. *Nat Genet* **2014**;46:487-91
- 497 19. Middlebrooks CD, Banday AR, Matsuda K, Udquim KI, Onabajo OO, Paquin A, *et al.*
498 Association of germline variants in the APOBEC3 region with cancer risk and enrichment
499 with APOBEC-signature mutations in tumors. *Nat Genet* **2016**;48:1330-8
- 500 20. Carter H, Marty R, Hofree M, Gross AM, Jensen J, Fisch KM, *et al.* Interaction Landscape
501 of Inherited Polymorphisms with Somatic Events in Cancer. *Cancer Discov* **2017**;7:410-
502 23
- 503 21. Srinivasan P, Bandlamudi C, Jonsson P, Kemel Y, Chavan SS, Richards AL, *et al.* The
504 context-specific role of germline pathogenicity in tumorigenesis. *Nat Genet*
505 **2021**;53:1577-85
- 506 22. Sun X, Xue A, Qi T, Chen D, Shi D, Wu Y, *et al.* Tumor Mutational Burden Is Polygenic and
507 Genetically Associated with Complex Traits and Diseases. *Cancer Res* **2021**;81:1230-9
- 508 23. Liu Y, Gusev A, Heng YJ, Alexandrov LB, Kraft P. Somatic mutational profiles and
509 germline polygenic risk scores in human cancer. *Genome Med* **2022**;14:14
- 510 24. Zhang P, Kitchen-Smith I, Xiong L, Stracquadanio G, Brown K, Richter PH, *et al.* Germline
511 and Somatic Genetic Variants in the p53 Pathway Interact to Affect Cancer Risk,
512 Progression, and Drug Response. *Cancer Res* **2021**;81:1667-80
- 513 25. Chen Z, Wen W, Beeghly-Fadiel A, Shu XO, Diez-Obrero V, Long J, *et al.* Identifying
514 Putative Susceptibility Genes and Evaluating Their Associations with Somatic Mutations
515 in Human Cancers. *Am J Hum Genet* **2019**;105:477-92
- 516 26. Gusev A, Groha S, Taraszka K, Semenov YR, Zaitlen N. Constructing germline research
517 cohorts from the discarded reads of clinical tumor sequences. *Genome Med*
518 **2021**;13:179
- 519 27. Garcia EP, Minkovsky A, Jia Y, Ducar MD, Shivdasani P, Gong X, *et al.* Validation of
520 OncoPanel: A Targeted Next-Generation Sequencing Assay for the Detection of Somatic
521 Variants in Cancer. *Arch Pathol Lab Med* **2017**;141:751-8
- 522 28. Cibulskis K, Lawrence MS, Carter SL, Sivachenko A, Jaffe D, Sougnez C, *et al.* Sensitive
523 detection of somatic point mutations in impure and heterogeneous cancer samples. *Nat*
524 *Biotechnol* **2013**;31:213-9
- 525 29. Maruvka YE, Frazer R, Grimsby J, Van Seventer E, Gelincik O, Ibrahim M, *et al.* Detection
526 of tumors with microsatellite instability (MSI) using minimal sequencing of cfDNA. in
527 submission
- 528 30. Davies RW, Flint J, Myers S, Mott R. Rapid genotype imputation from sequence without
529 reference panels. *Nat Genet* **2016**;48:965-9
- 530 31. Chen CY, Pollack S, Hunter DJ, Hirschhorn JN, Kraft P, Price AL. Improved ancestry
531 inference using weights from external reference panels. *Bioinformatics* **2013**;29:1399-
532 406
- 533 32. Lawrence MS, Stojanov P, Mermel CH, Robinson JT, Garraway LA, Golub TR, *et al.*
534 Discovery and saturation analysis of cancer genes across 21 tumour types. *Nature*
535 **2014**;505:495-501
- 536 33. Sul JH, Han B, Ye C, Choi T, Eskin E. Effectively identifying eQTLs from multiple tissues by
537 combining mixed model and meta-analytic approaches. *PLoS Genet* **2013**;9:e1003491

- 538 34. Sobota RS, Shriner D, Kodaman N, Goodloe R, Zheng W, Gao YT, *et al.* Addressing
539 population-specific multiple testing burdens in genetic association studies. *Ann Hum*
540 *Genet* **2015**;79:136-47
- 541 35. Boland CR, Goel A. Microsatellite instability in colorectal cancer. *Gastroenterology*
542 **2010**;138:2073-87 e3
- 543 36. Stelloo E, Jansen AML, Osse EM, Nout RA, Creutzberg CL, Ruano D, *et al.* Practical
544 guidance for mismatch repair-deficiency testing in endometrial cancer. *Ann Oncol*
545 **2017**;28:96-102
- 546 37. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, *et al.* The
547 NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays
548 and summary statistics 2019. *Nucleic Acids Res* **2019**;47:D1005-D12
- 549 38. Gong J, Mei S, Liu C, Xiang Y, Ye Y, Zhang Z, *et al.* PancanQTL: systematic identification of
550 cis-eQTLs and trans-eQTLs in 33 cancer types. *Nucleic Acids Res* **2018**;46:D971-D6
- 551 39. Sakaue S, Kanai M, Tanigawa Y, Karjalainen J, Kurki M, Koshiya S, *et al.* A cross-
552 population atlas of genetic associations for 220 human phenotypes. *Nat Genet*
553 **2021**;53:1415-24
- 554 40. Stacey SN, Sulem P, Jonasdottir A, Masson G, Gudmundsson J, Gudbjartsson DF, *et al.* A
555 germline variant in the TP53 polyadenylation signal confers cancer susceptibility. *Nat*
556 *Genet* **2011**;43:1098-103
- 557 41. Sonugur FG, Akbulut H. The Role of Tumor Microenvironment in Genomic Instability of
558 Malignant Tumors. *Front Genet* **2019**;10:1063
- 559 42. Fearnhead NS, Britton MP, Bodmer WF. The ABC of APC. *Hum Mol Genet* **2001**;10:721-
560 33
- 561 43. Markowska A, Pawalowska M, Lubin J, Markowska J. Signalling pathways in endometrial
562 cancer. *Contemp Oncol (Pozn)* **2014**;18:143-8
- 563 44. Moreno-Bueno G, Hardisson D, Sanchez C, Sarrio D, Cassia R, Garcia-Rostan G, *et al.*
564 Abnormalities of the APC/beta-catenin pathway in endometrial cancer. *Oncogene*
565 **2002**;21:7981-90
- 566 45. Spranger S, Bao R, Gajewski TF. Melanoma-intrinsic beta-catenin signalling prevents
567 anti-tumour immunity. *Nature* **2015**;523:231-5
- 568 46. Chatrath A, Ratan A, Dutta A. Germline variants predictive of tumor mutational burden
569 and immune checkpoint inhibitor efficacy. *iScience* **2021**;24:102248
- 570 47. Jonsson P, Bandlamudi C, Cheng ML, Srinivasan P, Chavan SS, Friedman ND, *et al.*
571 Tumour lineage shapes BRCA-mediated phenotypes. *Nature* **2019**;571:576-9
- 572 48. Cremona CA, Behrens A. ATM signalling and cancer. *Oncogene* **2014**;33:3351-60
- 573 49. Huang KL, Mashl RJ, Wu Y, Ritter DI, Wang J, Oh C, *et al.* Pathogenic Germline Variants in
574 10,389 Adult Cancers. *Cell* **2018**;173:355-70 e14
- 575 50. Hu M, Zhou M, Bao X, Pan D, Jiao M, Liu X, *et al.* ATM inhibition enhances cancer
576 immunotherapy by promoting mtDNA leakage and cGAS/STING activation. *J Clin Invest*
577 **2021**;131
- 578 51. Zhang Q, Green MD, Lang X, Lazarus J, Parsels JD, Wei S, *et al.* Inhibition of ATM
579 Increases Interferon Signaling and Sensitizes Pancreatic Cancer to Immune Checkpoint
580 Blockade Therapy. *Cancer Res* **2019**;79:3940-51

- 581 52. Stracquadanio G, Wang X, Wallace MD, Grawenda AM, Zhang P, Hewitt J, *et al.* The
582 importance of p53 pathway genetics in inherited and somatic cancer genomes. *Nat Rev*
583 *Cancer* **2016**;16:251-65
- 584 53. Dong ZY, Zhong WZ, Zhang XC, Su J, Xie Z, Liu SY, *et al.* Potential Predictive Value of TP53
585 and KRAS Mutation Status for Response to PD-1 Blockade Immunotherapy in Lung
586 Adenocarcinoma. *Clin Cancer Res* **2017**;23:3012-24
- 587 54. Chakrabarti J, Holokai L, Syu L, Steele NG, Chang J, Wang J, *et al.* Hedgehog signaling
588 induces PD-L1 expression and tumor cell proliferation in gastric cancer. *Oncotarget*
589 **2018**;9:37439-57
- 590 55. Gu S, Feng J, Jin Q, Wang W, Zhang S. Reduced expression of EphA5 is associated with
591 lymph node metastasis, advanced TNM stage, and poor prognosis in colorectal
592 carcinoma. *Histol Histopathol* **2017**;32:491-7
- 593 56. Huang W, Lin A, Luo P, Liu Y, Xu W, Zhu W, *et al.* EPHA5 mutation predicts the durable
594 clinical benefit of immune checkpoint inhibitors in patients with lung adenocarcinoma.
595 *Cancer Gene Ther* **2021**;28:864-74

Table 1. Significant associations between cancer gene eQTL and global tumor mutations

Cancer type	eQTL						Association results with TMB or TMC			Association results with gene expression from GTEx ^a					
	Region	Lead SNP	Pos ^b	Effect allele	Other allele	EAf	Beta	SE	P value	Gene	Beta (FE)	P value (FE)	Beta (RE)	P value (RE)	P value (RE2)
<i>TMB of all cancer genes on OncoPanel</i>															
Ovarian Cancer	2q14.2	rs1530578	121702128	T	C	0.01	17.61	3.02	1.26E-08	<i>GLI2</i>	-0.11	2.41E-16	-0.11	2.00E-10	1.33E-16
Glioma	8p12	rs139944315	30332577	T	A	0.99	-16.36	2.65	1.21E-09	<i>WRN</i>	-0.30	6.35E-44	-0.28	4.16E-20	6.79E-45
Esophagogastric Carcinoma	16q22.1	rs11075646	66969176	C	G	0.90	-2.57	0.53	1.48E-06	<i>CBFB</i>	0.05	2.06E-12	0.05	1.34E-08	3.03E-12
<i>TMC of recurrently mutated cancer genes</i>															
Colorectal Cancer	4q13.2	rs10031417	66700879	A	G	0.55	-0.18	0.04	2.04E-06	<i>EPHA5</i>	0.03	5.63E-05	0.03	4.12E-02	6.36E-13
Endometrial Cancer	5q22.2	rs397768	112181576	G	A	0.38	0.24	0.05	1.43E-05	<i>APC</i>	0.01	7.91E-02	0.00	7.36E-01	3.77E-10
Endometrial Cancer	8p12	rs11782945	31082006	G	A	0.41	0.58	0.12	4.95E-07	<i>WRN</i>	0.06	1.06E-26	0.06	4.95E-21	3.14E-26
Endometrial Cancer	12q13.3	rs73115907	57754701	G	A	0.76	-0.31	0.06	4.61E-07	<i>GLII</i>	-0.02	2.75E-03	-0.02	1.31E-01	2.99E-08
Endometrial Cancer	16q24.3	rs7201264	90036122	C	G	0.20	0.41	0.07	1.10E-08	<i>FANCA</i>	-0.07	4.70E-18	-0.09	1.32E-07	8.17E-41
Endometrial Cancer	17p13.1	rs17884306	7572101	C	T	0.94	-0.60	0.13	4.36E-06	<i>TP53</i>	0.09	2.41E-16	0.08	4.45E-12	8.15E-16

^a Meta-analysis results of the associations between the eQTL and normalized gene expression levels across 49 tissues

^b Position based on GRCh37/hg19

Figures

Figure 1. A workflow of the germline and somatic data generation pipeline in the Profile

cohort. Tumor samples were collected from all consented patients in the Profile cohort, followed by targeted sequencing using OncoPanel. Somatic data were generated from on-target reads from the tumor sequences. Germline data were imputed using both the off-target and on-target reads generated from tumor sequencing. Four measures of the local and global tumor mutations: i) Mutation status of recurrently mutated cancer genes, ii) Mutation status of hotspot mutations, iii) TMC of recurrently mutated cancer genes, and iv) TMB of all panel genes were generated for all selected primary cancer samples across 11 cancer types from somatic data. Germline eQTL were identified from GTEx for all identified genes, followed by germline allele dosage extraction from the germline imputed data.

Figure 2. Local and global tumor mutations of 11 cancer types. A,

Mutation frequency and sample size of the identified recurrently mutated cancer genes for each cancer type. A total of 135 genes were selected for 11 cancer types (Supplementary Table S3); only genes that are recurrently mutated in ≥ 5 cancer types are shown on this figure. **B,** Mutation frequency and sample size of the identified hotspot mutations for each cancer type. There are 17 hotspot mutations in 10 genes for 7 cancer types. **C,** Distribution of TMB of all panel genes across cancers with sample sizes. **D,** Distribution of TMC of recurrently mutated cancer genes across cancers. Each dot represents a sample. The red horizontal line represents the median of TMC for each cancer type. The total number of recurrently mutated cancer genes selected for each cancer is listed on the top of the figure.

Figure 3. Associations of eQTL with TMB of all panel genes and gene expression across tissues. A-C, All selected eQTL for genes identified from the top eQTL-TMB associations are shown. Association results ($-\log_{10}(P)$) for eQTL and TMB are from linear models adjusting for age, gender (if applicable), tumor purity, and panel version. Association results (m-value, the posterior probability that an effect exists in a tissue) for eQTL and gene expression in the “matching tissue” are from GTEx; matching tissue was selected as the tissue with the largest m-value among all relevant tissues for the corresponding cancer type. Each dot represents a variant; variants that are significantly associated with both TMB and gene expression (in any meta-analysis model) are in red with the top variant marked as yellow diamond. RSID, effect allele, effect size, P value, and m-value for the top variant are annotated on the plots. The horizontal red dashed lines denote the significant threshold for TMB associations ($P = 3.45 \times 10^{-6}$) and “has an effect” threshold for gene expression associations in the matching tissue (m-value = 0.9). **D-F,** Association results of the top variants with expression level of the eQTL genes identified from the top eQTL-TMB associations by tissue from GTEx. The $-\log_{10}(P)$ are from single-tissue eQTL analysis. Each dot represents a tissue with the matching tissue for the specific cancer marked as yellow triangle. Meta-analysis results across tissues from FE and RE models are provided on the plots. See Fig. 5 for tissue annotations.

Figure 4. Associations of eQTL with TMC of recurrently mutated cancer genes and gene expression across tissues. A-F, All selected eQTL for genes identified from the top eQTL-TMC associations are shown. Association results ($-\log_{10}(P)$) for eQTL and TMC are from negative binomial models adjusting for age, gender (if applicable), tumor purity, panel version, and MSI status. Association results (m-value, the posterior probability that an effect exists in a tissue) for eQTL and gene expression in the “matching tissue” are from GTEx; matching tissue was

selected as the tissue with the largest m-value among all relevant tissues for the corresponding cancer type. Each dot represents a variant; variants that are significantly associated with both TMC and gene expression (in any meta-analysis model) are in red with the top variant marked as yellow diamond. RSID, effect allele, effect size, P value, and m-value for the top variants are annotated on the plots. The horizontal red dashed lines denote the significant threshold for TMC associations ($P = 1.73 \times 10^{-5}$) and “has an effect” threshold for gene expression associations in the matching tissue (m-value = 0.9). **G-L**, Association results of the top variants with the expression level of the eQTL genes identified from the top eQTL-TMC associations by tissue from GTEx. The $-\log_{10}(P)$ are from single-tissue eQTL analysis. Each dot represents a tissue with the matching tissue for the specific cancer marked as yellow triangle. Meta-analysis results across tissues from FE and RE models are provided on the plots. See Fig. 5 for tissue annotations.

Figure 5. rs4753834 is associated with both *ATM* expression and somatic *ATM* mutations.

A, Associations between rs4753834 and risk of having somatic mutations in *ATM* across 8 cancers. The odds ratio is associated with the G allele of rs4753834. Meta-analysis results from the fixed-effect model are shown. **B**, Sample sizes and mutation frequencies of the 8 cancer types. Note that these numbers are based on samples included in the final models. **C**, Association results of rs4753834 and *ATM* expression by tissue from GTEx. The $-\log_{10}(P)$ are from single-tissue eQTL analysis. m-value is the posterior probability that an effect exists in a tissue. Results from the FE and RE meta-analysis across tissues are also shown on the plot.

Figure 6. Hypothetical relationships between germline variants, cancer gene expression, somatic mutations, and cancer diagnosis. **A**, Complete relationships between germline eQTL, environmental factors (E), expression level of cancer genes in normal tissue (transcript

abundance, T), somatic mutations in cancer genes in normal tissue before tumor development (S), cancer diagnosis (D), and somatic mutations in cancer genes in tumor after cancer diagnosis (S'). Our hypothesis is that germline eQTL regulate the expression of cancer genes; the transcript abundance of those cancer genes modifies the propensity of acquiring somatic mutations in those genes; having somatic mutations in those cancer genes is associated with an increased risk of cancer (the path shown by red arrows). Here, we are testing the associations between the eQTL and S', which can serve as a proxy for S, among cancer patients (conditioning on D). **B**, Three possible relationships between germline variants (G), somatic mutations in normal tissue before tumor development (S), cancer diagnosis (D), and somatic mutations in tumor after cancer diagnosis (S') given that an association between G and S' is observed. Blue arrows on the graphs show the paths from G to S' through S given that only cancer patients are included in the study (conditioning on D).

Figure 1

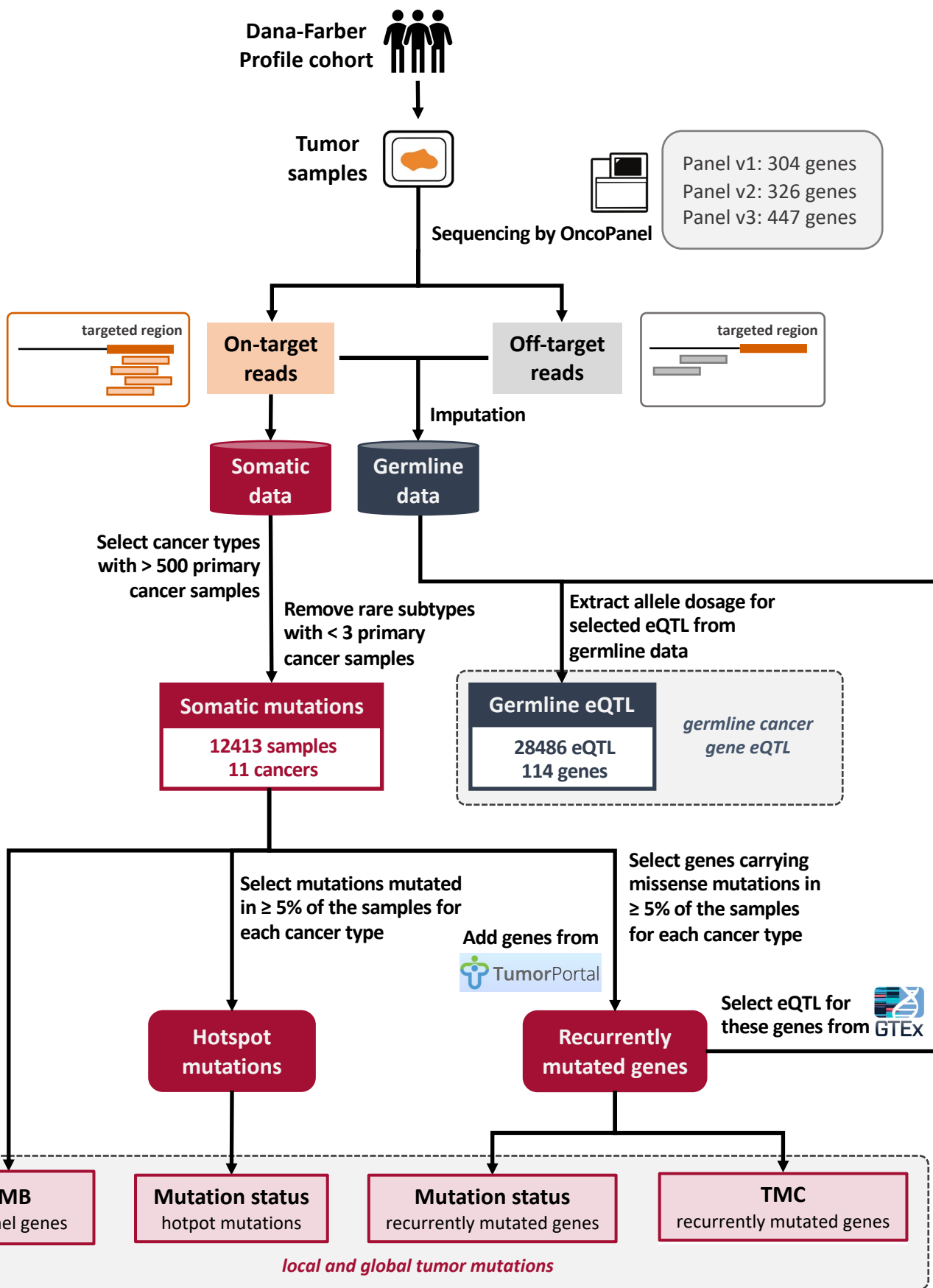
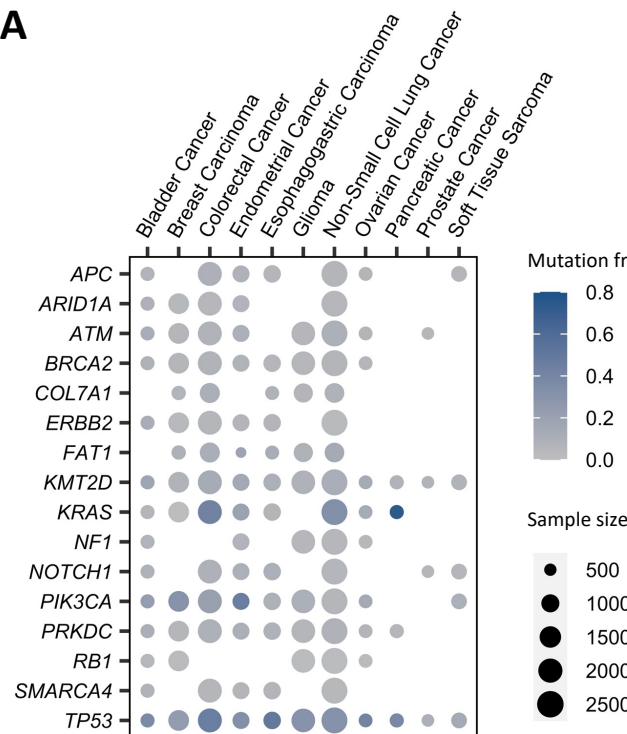
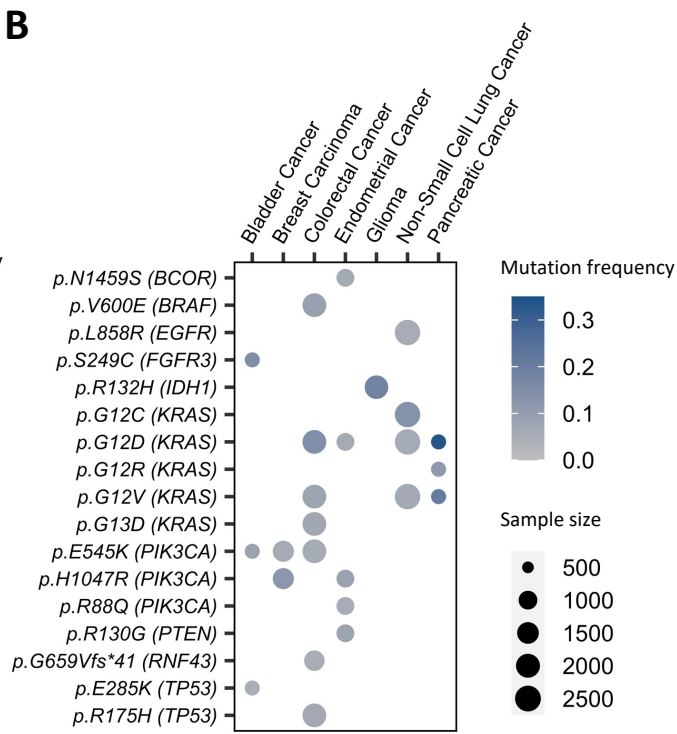


Figure 2

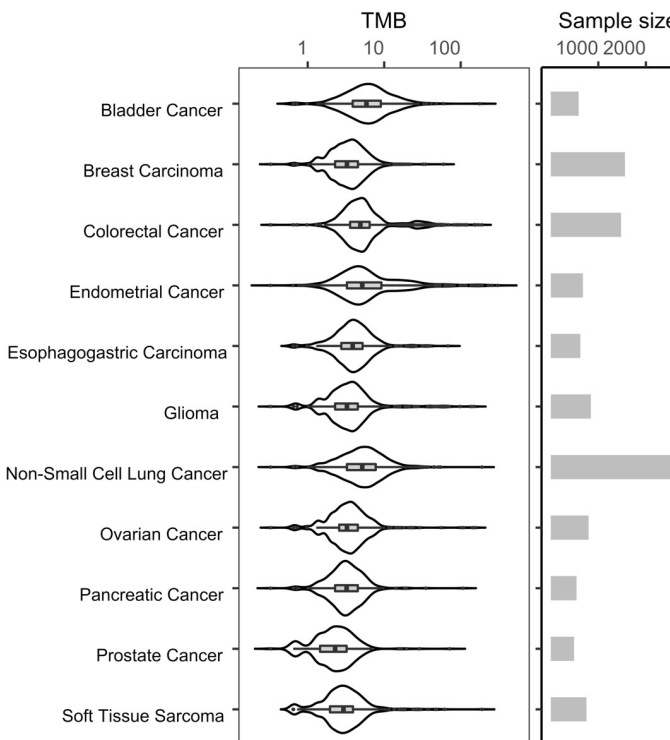
A



B



C



D

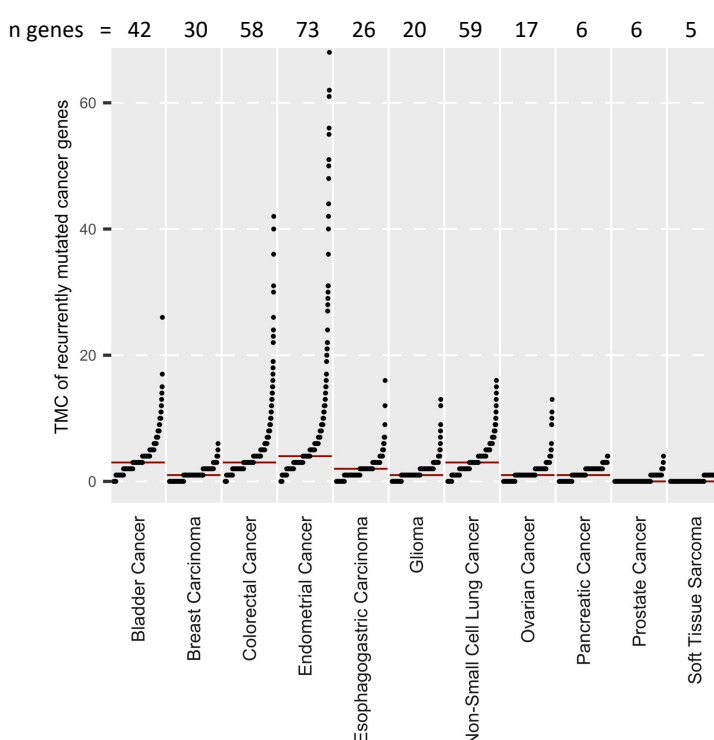


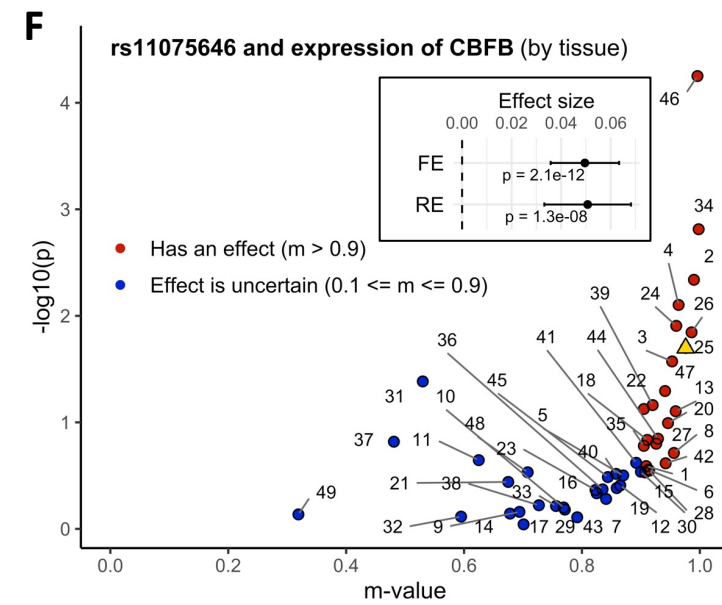
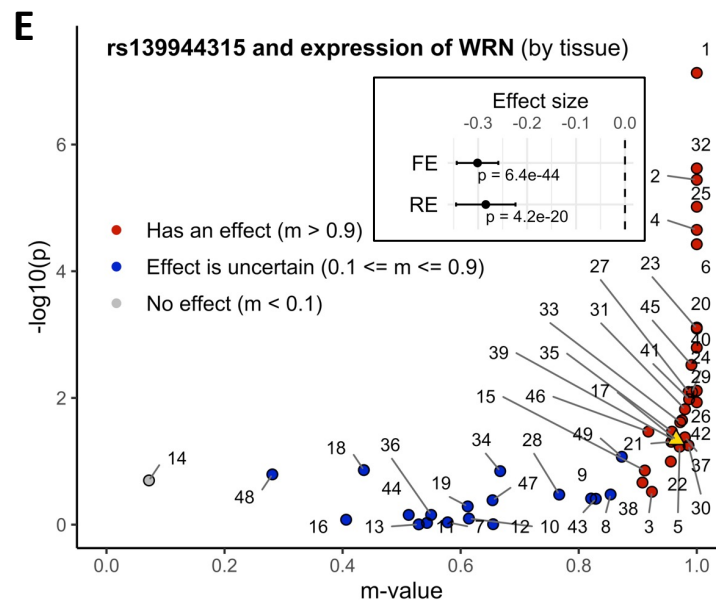
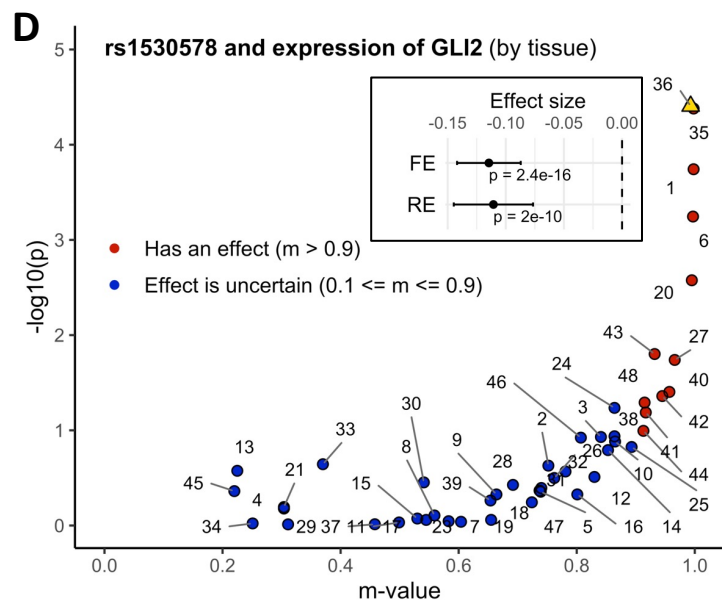
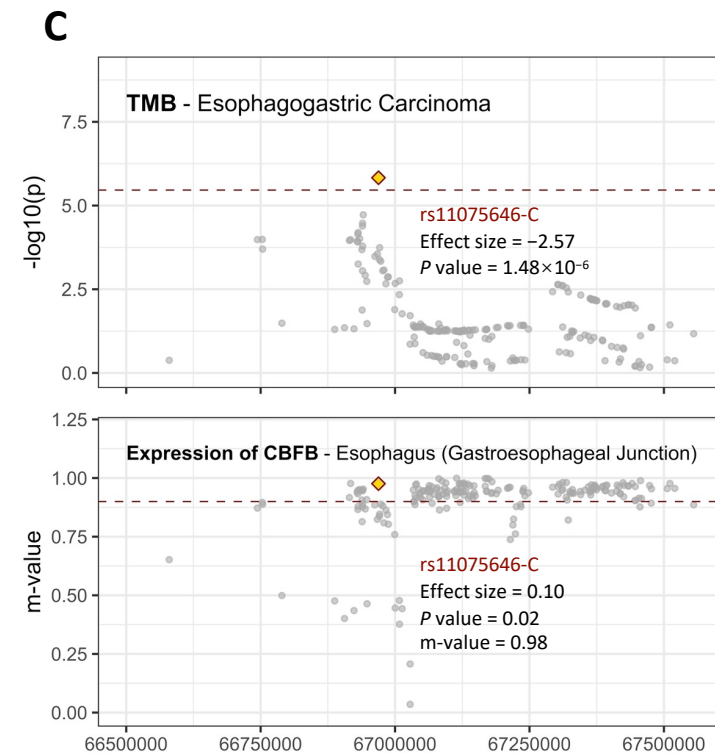
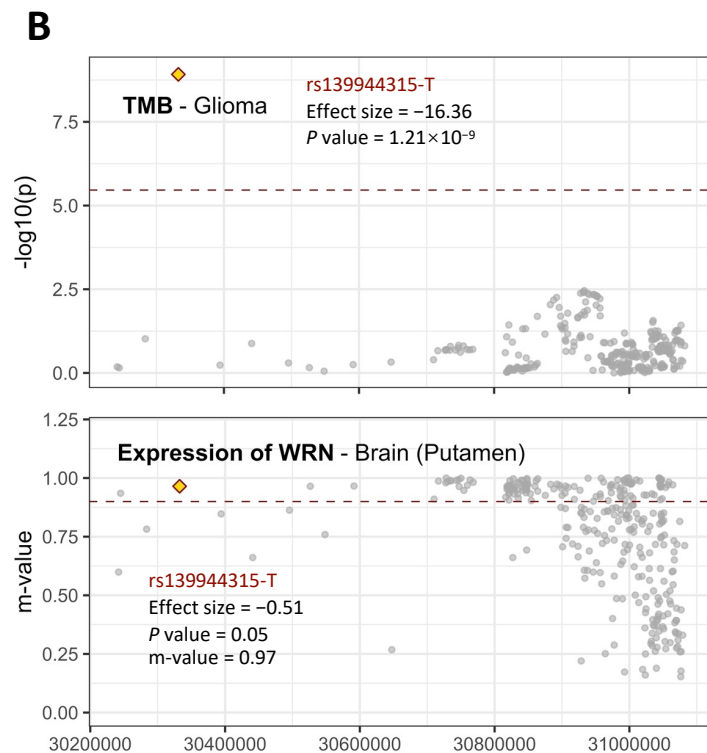
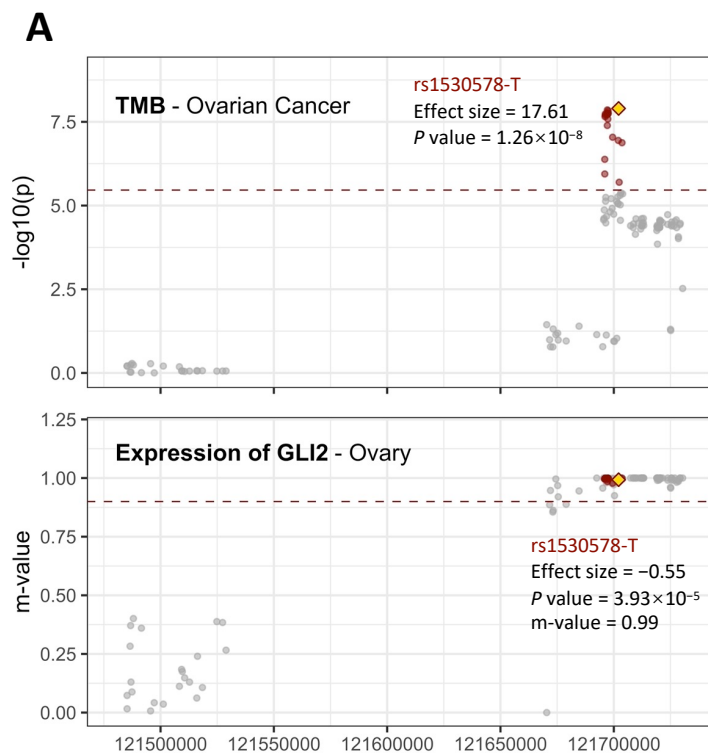
Figure 3

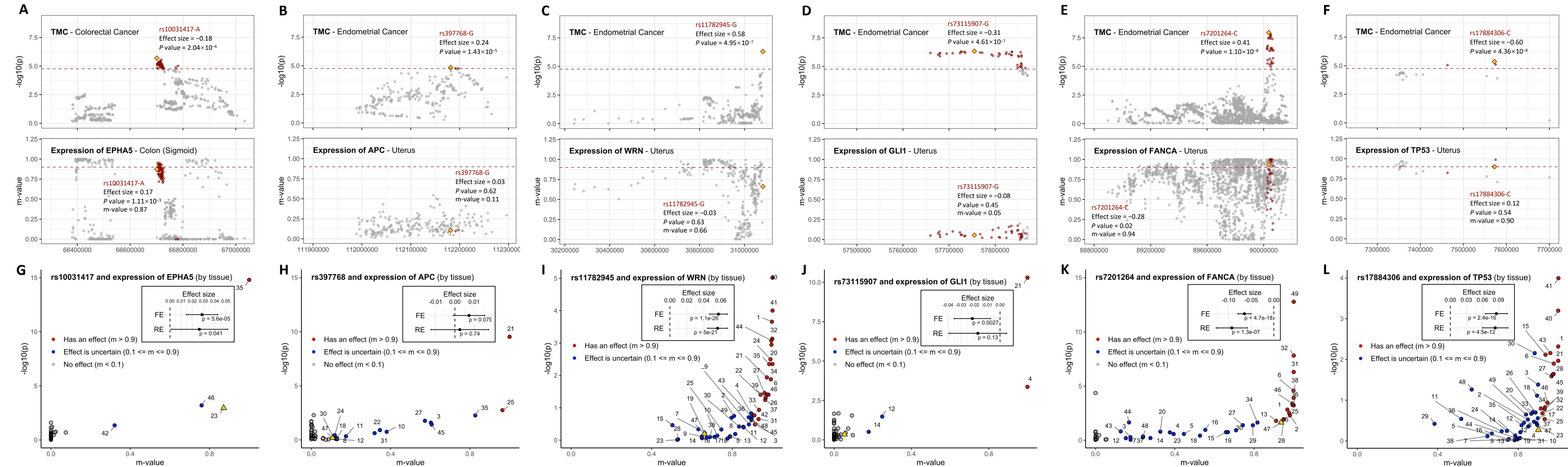
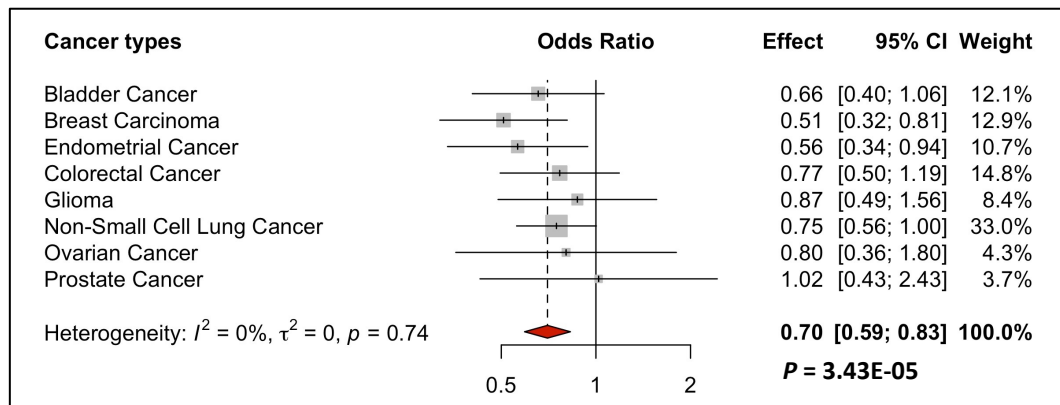
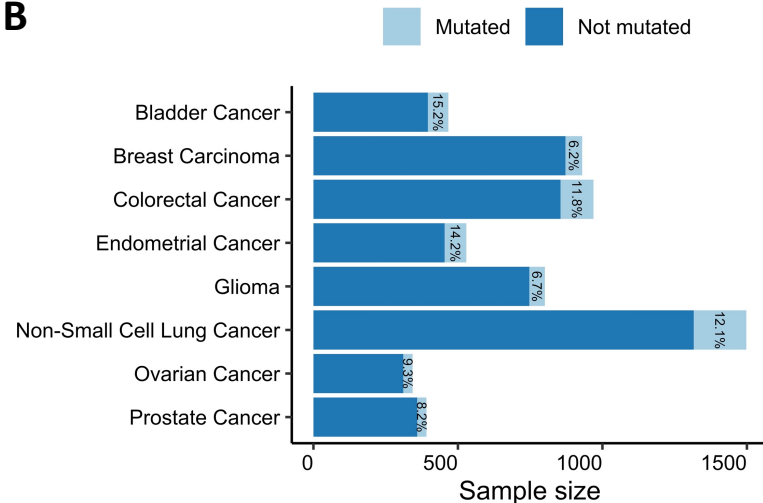
Figure 4

Figure 5

A



B



C

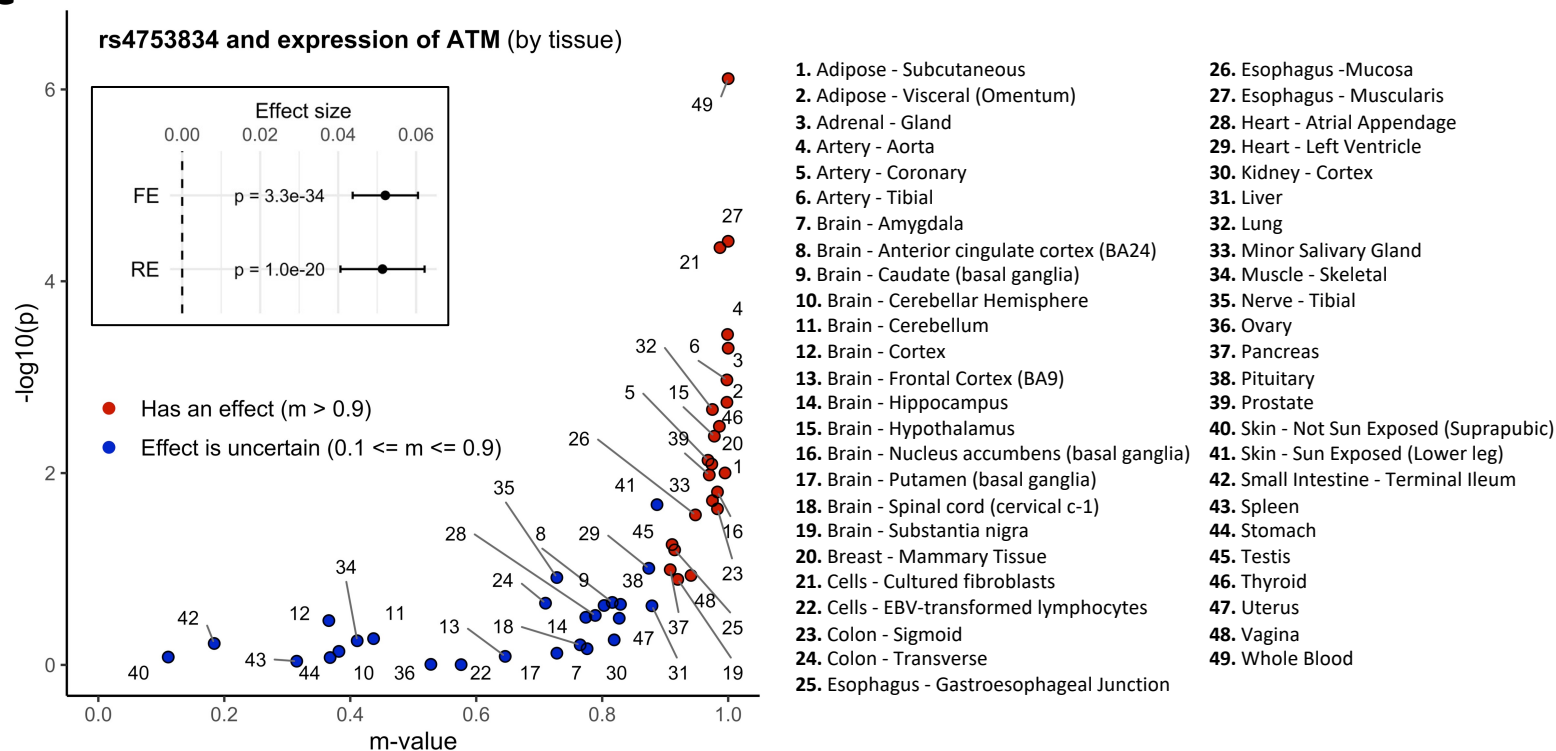
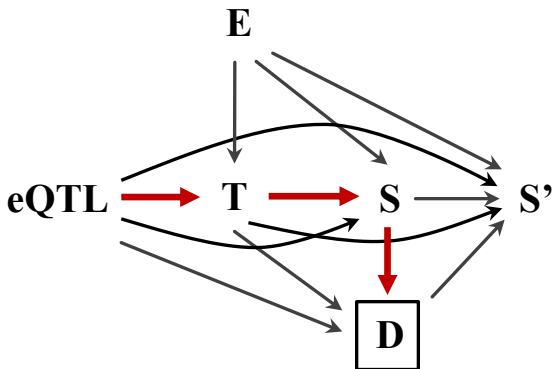
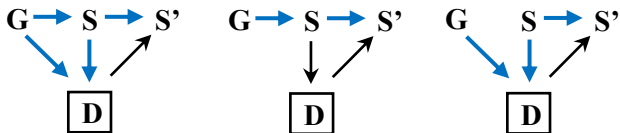


Figure 6

A



B



G induces S
and cancer

G induces S

No effect of
G on S