

1 **Large-scale integration of omics and electronic health records to identify potential**
2 **risk protein biomarkers and therapeutic drugs for cancer prevention and intervention**

3
4 Qing Li^{1,2†}, Qingyuan Song^{3,4,†}, Zhishan Chen^{1,†}, Jungyoon Choi⁵, Victor Moreno^{6,7,8,9}, Jie
5 Ping¹, Wanqing Wen¹, Chao Li¹, Xiang Shu¹⁰, Jun Yan¹¹, Xiao-ou Shu¹, Qiuyin Cai¹, Jirong
6 Long¹, Jeroen R Huyghe¹², Rish Pai¹³, Stephen B Gruber¹⁴, Graham Casey¹⁵, Xusheng
7 Wang¹⁶, Adetunji T. Toriola¹⁷, Li Li¹⁸, Bhuminder Singh¹⁹, Ken S Lau¹⁹, Li Zhou²⁰, Chong
8 Wu²¹, Ulrike Peters^{12,22}, Wei Zheng¹, Quan Long^{2,23,24,25,26*}, Zhijun Yin^{3,4*}, Xingyi Guo^{1,3*}

9
10 **Affiliations:**

11 ¹ Division of Epidemiology, Department of Medicine, Vanderbilt Epidemiology Center,
12 Vanderbilt-Ingram Cancer Center, Vanderbilt University School of Medicine, Nashville, TN,
13 USA

14 ² Department of Biochemistry and Molecular Biology, University of Calgary, Calgary, Alberta,
15 Canada

16 ³ Department of Biomedical Informatics, Vanderbilt University School of Medicine, Nashville,
17 TN, USA

18 ⁴ Department of Department of Computer Science, Vanderbilt University, Nashville, TN, USA

19 ⁵ Division of Oncology and Hematology, Department of Internal Medicine, Korea University
20 Ansan Hospital, Korea University College of Medicine, Ansan 15355, Korea

21 ⁶ Oncology Data Analytics Program, Catalan Institute of Oncology (ICO), L'Hospitalet de
22 Llobregat, Barcelona, Spain

23 ⁷ Colorectal Cancer Group, ONCOBELL Program, Institut de Recerca Biomedica de
24 Bellvitge (IDIBELL), L'Hospitalet de Llobregat, Barcelona, Spain

25 ⁸ Department of Clinical Sciences, Faculty of Medicine and health Sciences and Universitat
26 de Barcelona Institute of Complex Systems (UBICS), University of Barcelona (UB),
27 L'Hospitalet de Llobregat, Barcelona, Spain

28 ⁹ Consortium for Biomedical Research in Epidemiology and Public Health (CIBERESP),
29 Madrid, Spain

30 ¹⁰ Department of Epidemiology and Biostatistics, Memorial Sloan Kettering Cancer Center,
31 New York, NY, USA

32 ¹¹ Physiology and Pharmacology, University of Calgary, Calgary, Alberta, Canada

33 ¹² Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, WA,
34 USA

35 ¹³ Mayo Clinic Arizona, Scottsdale, AZ, USA.

36 ¹⁴ Department of Preventive Medicine & USC Norris Comprehensive Cancer Center, Keck
37 School of Medicine, University of Southern California, Los Angeles, CA, USA

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

38 ¹⁵ Center for Public Health Genomics, University of Virginia, Charlottesville, VA, USA

39 ¹⁶ Department of Genetics, Genomics and Informatics, University of Tennessee Health
40 Science Center, Memphis, TN, USA

41 ¹⁷ Division of Public Health Sciences, Department of Surgery, Washington University School
42 of Medicine and Siteman Cancer Center, St. Louis, MO, USA

43 ¹⁸ Department of family medicine, School of Medicine, University of Virginia, Charlottesville,
44 VA, USA

45 ¹⁹ Epithelial Biology Center and Department of Cell and Developmental Biology, Vanderbilt
46 University School of Medicine, Nashville, TN, USA

47 ²⁰ Harvard Medical School, Boston, Massachusetts; Division of General Internal Medicine
48 and Primary Care, Department of Medicine, Brigham and Women's Hospital, Boston,
49 Massachusetts.

50 ²¹ Department of Biostatistics, the University of Texas MD Anderson, Houston, Texas, USA

51 ²² Department of Epidemiology, University of Washington School of Public Health, Seattle,
52 WA, USA

53 ²³ Department of Medical Genetics, Cumming School of Medicine, University of Calgary,
54 Calgary, Alberta, Canada

55 ²⁴ The Mathison Centre for Mental Health Research & Education, Hotchkiss Brain Institute,
56 Cumming School of Medicine, University of Calgary, Calgary, Alberta, Canada

57 ²⁵ Alberta Children's Hospital Research Institute, Cumming School of Medicine, University of
58 Calgary, Calgary, Alberta, Canada

59 ²⁶ Department of Mathematics and Statistics, Faculty of Science, University of Calgary,
60 Calgary, Alberta, Canada

61

62 † Author names shared co-first authorship

63 *Corresponding authors

64

65 Dr. Xingyi Guo

66 Division of Epidemiology, Department of Medicine, Vanderbilt Epidemiology Center, and
67 Vanderbilt-Ingram Cancer Center, Vanderbilt University School of Medicine

68 2525 West End Ave. Suite 330, Nashville, TN 37203

69 Tel: +1 615-936-3471; Fax: +1 615-332-0502

70 Email: xingyi.guo@vumc.org

71

72 Dr. Zhijun Yin

73 Department of Biomedical Informatics, Department of Computer Science

74 Vanderbilt University Medical Center

75 2525 West End Ave. Suite 1475, Nashville, TN 37203

76 Tel: +1 615-936-3690; Fax: +1 615-322-0102

77 Email: zhijun.yin.1@vumc.org

78

79 Dr. Quan Long

80 Department of Biochemistry and Molecular Biology,

81 University of Calgary,

82 3330 Hospital Drive NW

83 Room 1151 (1173), Health Sciences Centre

84 Calgary, AB T2N 4N1 Canada

85 Tel: 403-220-5580; Fax: 403-210-9538

86 Email: quan.long@ucalgary.ca

87 **Abstract**

88 Identifying risk protein targets and their therapeutic drugs is crucial for effective cancer
89 prevention. Here, we conduct integrative and fine-mapping analyses of large genome-wide
90 association studies data for breast, colorectal, lung, ovarian, pancreatic, and prostate
91 cancers, and characterize 710 lead variants independently associated with cancer risk.
92 Through mapping protein quantitative trait loci (pQTL) for these variants using plasma
93 proteomics data from over 75,000 participants, we identify 365 proteins associated with
94 cancer risk. Subsequent colocalization analysis identifies 101 proteins, including 74 not
95 reported in previous studies. We further characterize 36 potential druggable proteins for
96 cancers or other disease indications. Analyzing >3.5 million electronic health records, we
97 uncover five drugs (Haloperidol, Trazodone, Tranexamic Acid, Haloperidol, and Captopril)
98 associated with increased cancer risk and two drugs (Caffeine and Acetazolamide) linked to
99 reduced colorectal cancer risk. This study offers novel insights into therapeutic drugs
100 targeting risk proteins for cancer prevention and intervention.

101 Introduction

102
103 Human genetic research has not only advanced our understanding of disease
104 mechanisms but has also significantly contributed to drug discovery and development. Drugs
105 supported by genetic evidence exhibit enhanced therapeutic validity compared to those
106 lacking such support, highlighting the importance of incorporating genetic evidence in drug
107 development initiatives^{1,2}. Common risk variants implicated in diseases can dysregulate
108 nearby gene or protein expression, which can mimic the effects of therapeutic drugs on the
109 targetable proteins. These proteins could serve as potential targets for therapeutic
110 intervention³. Thus, concerted efforts for cancer prevention based on proteins influenced by
111 common polymorphisms that modulate cancer risk, are urgently needed⁴. To date, genome-
112 wide association studies (GWAS) have identified several hundred common genetic risk loci
113 for each of three prevalent cancer types: breast, colorectal, and prostate⁵⁻⁸, and several
114 dozen risk loci have been identified for other cancers, such as cancer of lung, pancreas, and
115 ovarian⁹⁻¹³. Previous research, including our work, has identified hundreds of putative
116 cancer susceptible genes potentially regulated by these risk variants, using methods such as
117 expression quantitative trait loci (eQTL) analysis^{8-12,14-20} and transcriptome-wide association
118 studies (TWAS)^{7,19,21-29}. However, most dysregulated gene expression has not been
119 thoroughly investigated at the protein level.

120
121 To deepen the understanding of causal mechanisms and enhance drug discovery
122 endeavors, it is imperative to explore data from transcriptomic to proteomic studies. Proteins,
123 the ultimate products of mRNA translation, play critical roles in cellular activities and
124 represent promising therapeutic targets, as evidenced by successful drug targeting of
125 enzymes, transporters, ion channels, and receptors³⁰. Recent studies include protein
126 quantitative trait loci (pQTL) mapping and Mendelian randomization (MR) analysis by
127 integrating cancer GWAS and blood proteomics data to identify potential risk proteins.
128 However, only a few dozen of cancer risk proteins have been reported, with a false
129 discovery rate < 0.05 ³¹⁻³⁶. Most reported proteins have not been directly linked to the GWAS-
130 identified risk variants in common cancer types. Furthermore, research is lacking in
131 integrating multiple population-scale proteomic studies like the recent emerging UK Biobank
132 Pharma Proteomics Project (UKB-PPP)³⁷, which offers an unprecedented opportunity to
133 establish extensive pQTL databases, accelerating therapeutic drug discovery for therapeutic
134 prevention and intervention in human cancers.

135
136 Traditional drug discovery faces numerous challenges, including escalating costs,
137 lengthy timelines, and high failure rates³⁸. Drug repurposing presents a promising strategy

138 by identifying new applications for existing drugs, leveraging their well-documented
139 characteristics³⁹. With the widespread adoption of modern electronic health record (EHR)
140 systems, vast amounts of real-world patient data are available to augment pre-clinical
141 outcomes and facilitate drug repurposing screening. Recently, drug repurposing using EHRs
142 has successfully discovered repurposing hypotheses for preventing Alzheimer's Disease⁴⁰,
143 reducing cancer mortality^{41,42}, treating COVID-19^{43,44}, and coronary artery disease⁴⁵.
144 However, for therapeutic drugs that have been used for a long term to treat disease
145 indications with evidence of affecting the expression of cancer risk proteins, their potential
146 association with the risk of human cancers remains largely unclear. Some of these drugs
147 may be linked to an increased cancer risk due to long-neglected side effects.

148
149 In this work, we integrate large GWAS data for breast, colorectal, lung, ovarian,
150 pancreatic, and prostate cancers and population-scale proteomics data from over 75,000
151 participants combined from Atherosclerosis Risk in Communities study (ARIC)⁴⁶, deCODE
152 genetics⁴⁷, and UKB-PPP to identify risk proteins associated with each cancer. We further
153 characterized therapeutic drugs based on druggable risk proteins targeted by approved
154 drugs or undergoing clinical trials for cancer treatment or other indications. We further
155 evaluate the effect of cancer risk for those drugs approved for the indications, using over 3.5
156 million EHR database at Vanderbilt University Medical Center (VUMC). Findings from this
157 study offer novel insights into therapeutic drugs targeting risk proteins for cancer prevention
158 and intervention.

159

160 **Results**

161 **Overall analysis workflow**

162 In Figure 1, we outlined several main steps of a comprehensive integrative analysis
163 of GWAS, pQTLs, druggable proteins, and EHR data. First, we examined previously
164 identified risk loci from six cancer types based on the most recent GWAS in breast (N =
165 247,173), ovary (N = 63,347), prostate (N = 140,306), colorectum (N = 254,791), lung (N =
166 85,716), and pancreas (N = 21,536). Through additional fine-mapping analysis using
167 SuSiE⁴⁸, we characterized the most significantly associated variants (the lead variants) with
168 independent association signals at each risk locus for each cancer (**Fig. 1a; Online
169 Methods**). Second, we analyzed *cis*-pQTL results for the lead variants using proteomics
170 data from individuals of European descent from ARIC⁴⁶, deCODE⁴⁷, and UKB-PPP³⁷. We
171 conducted fixed-effect meta-analyses of summary statistics *cis*-pQTLs from ARIC⁴⁶ and
172 deCODE⁴⁷ through the same SOMAscan[®] platform (covering > 4,500 proteins) . We

173 combined them with the pQTL results from UKB-PPP through the Olink platform to identify
174 potential risk proteins. For proteins that satisfied the significance threshold after multiple
175 testing corrections, we further performed colocalization analyses to determine cancer risk
176 proteins with high confidence through evaluating the likelihood of shared causal variants
177 between pQTLs and GWAS (**Fig. 1a**). Third, these risk proteins with evidence of
178 colocalization were further annotated based on drug-protein information from four
179 drugs/compounds databases (DrugBank⁴⁹, ChEMBL⁵⁰, the Therapeutic Target Database⁵¹
180 (TTD) and OpenTargets⁵² (**Fig. 1b**). We next identified druggable proteins that are
181 therapeutic targets of approved drugs or undergoing clinical trials for cancer treatment or
182 other indications. Finally, we focused on drugs approved for other indications. We built
183 emulation of treated-control drug trials under the Inverse Probability of Treatment weighting
184 (IPTW) framework⁵³ through the analysis of over 3.5 million EHRs at VUMC. In these
185 emulations, we used the Cox proportional hazard model for each trial to evaluate the hazard
186 ratio (HR) of the specific cancer risk between the treated focal drug and the control drug.
187 The significance of each focal drug (HR and *P* value) was derived from a random-effects
188 meta-analysis of results across its balanced trials (**Fig. 1c, Online Methods**).

189

190 **Characterizing lead variants for breast, ovarian, prostate, colorectal, lung, and** 191 **pancreas cancers**

192 To characterize lead variants at each locus for each cancer type, we collected the
193 reported risk variants from previous fine-mapping or GWAS. Using breast cancer as an
194 example, we included 196 lead variants with independent association signals at loci from a
195 previous fine-mapping study based on conditional association analysis⁵⁴ and additional 32
196 genetic variants identified from a recent GWAS⁶. We then performed additional fine-mapping
197 analysis using SuSiE⁴⁸ based on summary statistics of GWAS (N = 247,173) from the Breast
198 Cancer Association Consortium (BCAC, **Supplementary Table 1**). After integrating the
199 previous results with new fine-mapping efforts, we identified 227 lead variants with
200 independently associated with cancer risk at each locus through several processing steps,
201 which included lead variant selection ($P < 1 \times 10^{-6}$ in European populations) and evaluating a
202 linkage disequilibrium (LD) ($r^2 < 0.1$) among the identified risk variants (**Extended Data Fig.**
203 **1; Online Methods**). Similarly, we characterized lead variants from previous GWAS and our
204 fine-mapping studies for colorectal⁵⁵ and other cancers (**Extended Data Fig. 1;**
205 **Supplementary Table 1; Online Methods**). In our analysis, we identified 710 lead variants,
206 including 227 for breast cancer, 213 for colorectal cancer, 213 for prostate cancer, 26 for
207 lung cancer, 13 for ovarian cancer and 18 for pancreatic cancer (**Fig. 2 and Supplementary**
208 **Tables 2 and 3; Online Methods**).

209

210 **Identifying cancer risk proteins from pQTLs mapping and colocalization analyses**

211 We mapped the 710 lead variants to *cis*-pQTLs to identify cancer risk proteins. At a
212 Bonferroni-corrected $P < 0.05$, we identified a total of 459 pQTL association signals
213 (corresponding 365 proteins after combined proteins unique for each cancer) for 222 lead
214 variants across six cancer types, including 74 for breast, 127 for colorectal, 37 for lung, 5 for
215 ovarian, 9 for pancreatic, and 113 for prostate cancer (**Fig. 2; Supplementary Table 4**).
216 Notably, 312 of the identified proteins (85.4% of 365) among these cancer types have not
217 been reported in previous proteomics-based MR studies^{32-36,56,57} (**Supplementary Table 5**).
218 Furthermore, through analysis of the identified proteins commonly observed in multiple
219 cancers, we found that 60 proteins were commonly observed in at least two of these six
220 cancers. In particular, we observed that several well-known cancer-related proteins, such as
221 HLA-A and HLA-E, were linked to lead variants located in major histocompatibility complex
222 (MHC) in breast, colorectal and lung cancers, highlighting the potential role of these proteins
223 in cancer pleiotropy and shared cancer risk mechanisms (**Fig. 2**).

224

225 A further colocalization analysis identified 101 proteins after combined proteins
226 unique for each cancer that showed strong evidence supported by either colocalization or
227 SMR+HEIDI analysis (**Online Methods**). Specifically, we identified 23 proteins for breast, 38
228 proteins for colorectal cancer, 7 proteins for lung, 2 for ovarian, 2 for pancreatic and 29 for
229 prostate cancer, respectively (**Fig. 3a-b, Supplementary Table 6**). Of these, 74 proteins
230 (73.2% of 101) have not been previously linked to cancer risk (**Supplementary Table 7**). Of
231 note, 71 proteins were only assayed by either SOMAscan[®] (n=32) or Olink platform (n=39).
232 For the remaining 22 significant proteins commonly assayed, all showed a pQTL significance
233 signal with a minimal nominal $P < 1 \times 10^{-5}$ in both ARIC+deCODE and UKB-PPP ($r = 0.66$, P
234 $= 2 \times 10^{-4}$; **Fig. 3c**). In particular, seven proteins were highlighted as cancer-driver proteins^{58,59}
235 and Cancer Gene Census (CGC)⁶⁰, including ALDH2, HLA-A and SUB1 for breast cancer,
236 ALDH2 and HLA-A for colorectal cancer, NT5C2 for lung cancer, and NT5C2, RNF43,
237 TYRO3 and USP28 for prostate cancer.

238

239 **Cancer risk proteins supported by functional genomics analyses**

240 Of the identified 101 proteins among the six cancers, we next examined whether they
241 are supported by functional genomics analyses. Specifically, we first evaluated xQTL (i.e.,
242 eQTLs, alternative splicing - sQTLs, and alternative polyadenylation - apaQTLs) results in
243 their respective target tissues and whole blood samples (**Online Methods**). We found 63
244 proteins that were supported by at least one xQTLs at a nominal $P < 0.05$, including 12 for
245 breast (52% of 23), 22 for colorectal (57% of 38), 5 for lung (71% of 7), 2 for ovarian, 2 for
246 pancreatic, and 20 (68% of 29) for prostate cancer (**Supplementary Table 7**). Second, we

247 used functional genomic data generated in their cancer-related tissues/cells (i.e., promoter
248 and enhancer) to characterize putative functional variants that are in strong LD ($r^2 > 0.8$ in
249 the European population) with the lead variants (**Online Methods**) Our results showed that
250 17 genes were likely regulated by the closest putative regulatory variants with either
251 promoter and/or enhancer activities (**Supplementary Table 8**). We further investigated the
252 potential distal regulatory effects of putative functional variants on these genes by analyzing
253 chromatin-chromatin interaction data (**Online Methods**). We found that 39 genes were
254 regulated distally by putative functional variants through long-term promoter-enhancer
255 interactions (**Supplementary Table 9**). Lastly, we examined differential protein expression
256 between normal and tumor tissues available for breast, colon, lung and pancreatic cancers
257 using data from Clinical Proteomic Tumor Analysis Consortium (CPTAC). We showed
258 evidence of the 18 identified proteins with consistent association directions supported by
259 significantly differential expression at a nominal $P < 0.05$, including 3 for breast cancer, 10
260 for colorectal cancer, 4 for lung cancer and 1 for pancreatic cancer. Similarly, we showed
261 evidence of the 40 identified proteins supported by significantly differential mRNA expression
262 using data from The Cancer Genome Atlas Program (TCGA), including 6 for breast cancer,
263 15 for colorectal cancer, 3 for lung cancer, 1 for pancreatic cancer and 15 for prostate
264 cancer. Taken together, our analysis provided additional evidence that most of the identified
265 proteins partially or wholly supported by functional genomics analyses (**Supplementary**
266 **Table 10**).

267

268 **Identifying druggable proteins**

269 Using data from DrugBank⁴⁹, ChEMBL⁵⁰, the Therapeutic Target Database⁵¹ (TTD)
270 and OpenTargets⁵², we comprehensively annotated our proteins as therapeutic targets of
271 approved or clinical-stage drugs (**Online Methods**). Of the 101 proteins among the six
272 cancers, we identified 36 druggable proteins potentially targeted by 404 approved drugs or
273 undergoing clinical trials for cancer treatment or other indications (**Fig. 4, Supplementary**
274 **Table 11**). Specifically, we found 19 proteins targeted by 133 drugs either approved or under
275 clinical trials to treat cancers (**Fig. 5, Supplementary Table 12**). Our results also provide
276 evidence that the remaining druggable proteins are targeted by 197 drugs used for treating
277 indications other than cancer (**Extended Data Fig. 3**).

278

279 **Evaluating associations of drugs approved for indications with cancer risk**

280 We next evaluated the effect on cancer risk of therapeutic drugs that have been used
281 long-term to treat indications based on real-world EHRs from the VUMC Synthetic Derivative
282 (SD) database. Given a focal drug, we first emulated its trials by building control patient
283 groups who were exposed to similar treated drugs under the same ATC-L2 category (**Online**

284 **Methods; Supplementary Table 13).** To mimic randomized controlled trials (RCT) to
285 evaluate the focal drug's effect, we applied the Inverse Probability of Treatment Weighting
286 (IPTW) framework⁶¹ to create a pseudo-population wherein confounding variables are evenly
287 distributed between the treated and control groups (**Online Methods**). After discarding the
288 trials with less than 500 eligible patients in either of the groups, we analyzed 14 treated
289 drugs with 335 balanced trials. Our analysis revealed that five drugs were linked to an
290 increased risk of cancer: Haloperidol (HR = 1.76; $P = 1.6 \times 10^{-23}$; targeting HLA-A protein)
291 and Trazodone (HR = 1.32; $P = 2.3 \times 10^{-12}$; targeting HLA-A protein) for breast cancer;
292 Tranexamic Acid (HR = 1.53; $P = 1.1 \times 10^{-3}$; targeting PLG protein) and Sirolimus (HR =
293 1.71; $P = 1.1 \times 10^{-28}$; targeting TYRO3 protein) for prostate cancer; and Haloperidol (HR =
294 2.62, $P = 6.6 \times 10^{-20}$, targeting HLA-A protein) and Captopril (HR = 1.65; $P = 2.2 \times 10^{-9}$;
295 targeting TF protein) for colorectal cancer (**Fig. 6**). In contrast, we also found that two drugs
296 associated with a decreased risk of colorectal cancer: Caffeine (HR = 0.74, $P = 9.3 \times 10^{-5}$,
297 targeting ALDH2 protein) and Acetazolamide (HR = 0.72; $P = 1.1 \times 10^{-20}$; targeting HLA-A
298 protein) (**Fig. 6**).

299

300 **Discussion**

301 In this study, we conducted a comprehensive investigation of cancer risk proteins by
302 integrating lead variants and pQTLs for six common cancer types using large-scale GWAS
303 and population-based proteomics data. Through pQTL mapping and subsequent
304 colocalization analysis, we identified 101 risk proteins across the six cancer types, with over
305 three-quarters of them not previously linked to cancer susceptibility. Moreover, most of the
306 proteins we identified are supported by functional genomics analyses. Our findings not only
307 significantly expand the pool of known cancer risk proteins but also offer new insights into
308 the biology and susceptibility of common cancers.

309

310 Through analysis of drug-protein interaction databases, we identified 36 druggable
311 proteins potentially targeted by 404 therapeutic drugs. Among these, 30 drugs have already
312 received approval for cancer treatment, while 73 are currently undergoing clinical trials for
313 cancer treatment. These findings offer genetic evidence supporting the effectiveness of
314 certain drugs and suggest potential opportunities for repurposing them to treat additional
315 cancers that share common risk proteins. However, it's crucial to acknowledge that while the
316 cancer risk proteins identified in our study hold promise as therapeutic targets for cancer
317 treatment, drugs may also have adverse effects, potentially exacerbating cancer
318 development through these targets (i.e., depending on their inhibitory or promotive effects)⁶².

319 Additionally, our analysis characterized 197 drugs used for indications other than cancer,
320 which may influence cancer risk due to their interactions with cancer-risk proteins. Overall,
321 our findings have the potential to accelerate therapeutic drug discovery for the prevention
322 and intervention of human cancers.

323

324 We uncovered five non-cancer drugs associated with increased cancer risk:
325 Tranexamic Acid (PLG), Sirolimus (TYRO3), Haloperidol (HLA-A), Trazodone (HLA-A), and
326 Captopril (TF). Our genetic evidence or previous pharmacological studies further support
327 these findings. Specifically, Tranexamic Acid, an antifibrinolytic agent used to block the
328 breakdown of blood clots and prevent bleeding, was associated with an increased risk of
329 prostate cancer. Our findings suggest its potential inhibition of the protein expression of
330 human plasminogen (PLG), based on data from the ChEMBL database. Our GWAS and
331 pQTL results indicate that PLG may serve as a potential tumor suppressor, supported by
332 evidence of the risk allele C of rs9347480 being associated with increased prostate cancer
333 risk ($P = 1.41 \times 10^{-07}$) and decreased protein expression ($P = 4.57 \times 10^{-33}$). Additionally, this
334 protein also shows notable evidence of decreased gene expression ($P = 0.038$) in prostate
335 tumor samples compared to normal samples, as observed in data from TCGA. The drug
336 Sirolimus, primarily used to treat immune system and eye diseases, can potentially affect
337 receptor tyrosine kinases, TYRO3, a known protein important for prostate cancer
338 development^{63,64}. In breast and colorectal cancers, we identified two candidate drugs
339 (Haloperidol and Trazodone) targeting the major histocompatibility complex, HLA-A, an
340 essential protein for the immune system's defense against cancer development.
341 Interestingly, our analysis suggests that Haloperidol, a type of antipsychotic treatment, is
342 highly likely to increase the risk of both breast and colon cancer. Haloperidol, the first-
343 generation antipsychotics, has been reported to be a carcinogenic compound⁶⁵ and its
344 exposure of five years or more was associated with an increased risk of breast cancer in a
345 Finland nationwide study⁶². Consistently, another prior study showed its notably increased
346 risks of colorectal cancer in patients with schizophrenia who take antipsychotic
347 medications⁶⁶. We also found that Captopril, originally for cardiovascular diseases and
348 promisingly repurposed for cancer treatment in clinical trials and several studies⁶⁷⁻⁷⁰, has the
349 potential to increase the risk of colorectal cancer, aligning with a previous study⁷¹.

350

351 Conversely, our study also identified two non-cancer drugs (Caffeine and
352 Acetazolamide) associated with a reduced risk of colorectal cancer. Acetazolamide,
353 prioritized by the risk protein named TF, exhibited a notable effect in preventing colorectal
354 cancer development ($HR = 0.72$, $P = 1.1 \times 10^{-20}$). In line with our findings, prior studies
355 demonstrated its role in inhibiting cell viability, migration, and colony formation ability of

356 colorectal cancer cells⁷², as well as its ability to suppress the development of intestinal
357 polyps in Min Mice⁷³. In addition, Caffeine, a drug prioritized by the risk protein ALDH2 in
358 colorectal cancer, has been shown to exert a protective effect on colorectal cancer by prior
359 studies^{74,75}. However, such low-risk association may vary by colon subsites⁷⁶ and specific
360 populations⁷⁷. Of note, a clinical trial (NCT05692024) is undergoing the recruitment phase to
361 evaluate the effects of instant coffee on the gut microbiome, metabolome, liver fat, and
362 fibrosis in colorectal cancer patients.

363

364 Although the larger sample size of the European-ancestry study available for both
365 GWAS and proteomics enabled us to identify a larger number of association signals for risk
366 protein discovery based on colocalization analysis, our study was primarily limited to
367 individuals of European ancestry and further investigations are needed to assess the
368 relevance of these proteins in non-European populations. Millions of EHRs provide an
369 unprecedented opportunity to systematically evaluate non-cancer drugs' effect on risk of
370 cancer development. Especially these therapeutic drugs have been used for a long time to
371 treat diseases other than cancers, which can provide appropriate statistical power for the
372 analysis. While this approach is limited in only examining the cancer risk of common
373 approved drugs, it serves as an efficient complementary method to the pre-clinical data
374 analysis for cancer prevention and treatment. Despite the supportive evidence of
375 Acetazolamide *in vitro* and *in vivo*, it remains necessary to evaluate the effects of our
376 reported candidate drugs through both *in vitro* and *in vivo* assays in future investigations.
377

378 **Online Methods**

379 **Data resources**

380 The GWAS summary statistics data of European descendants for breast, prostate,
381 ovarian, and lung cancers were downloaded and compiled from their corresponding
382 consortia, including the Breast Cancer Association Consortium (BCAC)⁶ (N = 247,173,
383 133,384 cases and 113,789 controls), the Transdisciplinary Research of Cancer in Lung of
384 the International Lung Cancer Consortium (TRICL-ILCCO) and the Lung Cancer Cohort
385 Consortium (LC3)¹³ (N = 85,716, 29,266 cases and 56,450 controls), the Ovary Cancer
386 Association Consortium (OCAC)¹¹ (N = 63,347, 22,406 cases and 40,941 controls), and the
387 Pancreatic Cancer Case-Control Consortium (PanC4)¹⁰ (N = 21,536, 9,040 cases and
388 12,496 controls), and the Prostate Cancer Association Group Investigate Cancer Associated
389 Alterations in the Genome (PRACTICAL)⁷⁸ (N = 140,306, 79,194 cases and 61,112
390 controls). For colorectal cancer, we included GWAS data of 125,487 subjects from the
391 European population.^{19,79,80} In addition, the GWAS data (N = 254,791) consisting of 100,204

392 colorectal cancer cases and 154,587 controls from European and Asian populations⁷ were
393 also used in our analysis.

394 The large-scale *cis*-protein quantitative trait loci (*cis*-pQTLs) among European-
395 ancestry populations were analyzed based on three proteomics datasets: UKB-PPP³⁷ (N =
396 34,557, 2,922 plasma proteins), ARIC⁴⁴ (N = 7,213, 4,657 plasma proteins) and deCODE
397 genetics⁴⁵ (N = 35,559, 4,907 plasma proteins). Detailed descriptions of sample collection
398 and processes of the *cis*-pQTL analyses from the above proteomics datasets have been
399 described in previous studies^{37,46,47}.

400

401 We utilized the synthetic derivative (SD) database at Vanderbilt University Medical
402 Center (VUMC)⁸¹. This VUMC SD database contains de-identified clinical information
403 derived from Vanderbilt's electronic medical record. The SD has longitudinal clinical data for
404 over 3.5 million individuals, including patient demographics, medical history, laboratory
405 results, and medication history.

406

407 **Characterization of lead variants in six types of cancer**

408 In breast cancer, we included 196 strong independent association signals at $P < 1 \times$
409 10^{-6} from a fine-mapping study⁵⁴ and 32 risk variants from a GWAS⁶. We first combined the
410 reported lead variants from these two studies after removing those variants in LD ($r^2 < 0.1$ in
411 European populations). We further included additional lead variants from SuSiE fine-
412 mapping analysis on GWAS (N = 247,173)⁴⁸, with fine-mapping windows of 500 kilobases
413 (kb) and allowed a maximum of five causal variants. LD reference was based on the British-
414 ancestry UK Biobank samples (N = 337,000)⁸². We identified a credible set of causal
415 variants with a 95% posterior inclusion probability (95% PIP) for each independent risk
416 signal and a lead variant was represented by the variant with the minimum P . We included
417 additional lead variants from our SuSiE analysis with LD $r^2 < 0.1$ in European populations
418 with the above set of lead variants for those with independent risk-associated signals at
419 GWAS $P < 5 \times 10^{-8}$ and located in GWAS loci with independent risk-associated signals at P
420 $< 1 \times 10^{-6}$ in European populations.

421

422 For colorectal cancer, we analyzed 238 lead variants from our recent fine-mapping
423 study⁵⁵ based on the GWAS data from 254,791 participants in both European and Asian
424 populations. We characterized 233 lead variants with independent risk-associated signals at
425 minimal $P < 1 \times 10^{-6}$ in European populations, from the analysis based on GWAS from trans-
426 ancestry and European populations, respectively. For prostate cancer, we first identified lead
427 variants with independent risk-associated signals at GWAS $P < 5 \times 10^{-8}$ from our SuSiE fine-

428 mapping analysis on GWAS summary statistics (N = 140,306). We next included additional
429 GWAS-identified risk variants with $P < 1 \times 10^{-6}$ in European populations from the previous
430 trans-ancestry GWAS⁸ and $r^2 < 0.1$ with any lead variants from the above set in the fine-
431 mapping analysis. Similarly, we used the above strategy to characterize lead variants from
432 fine-mapping analysis for ovarian cancer (N = 63,347) and pancreatic cancer (N = 21,536).
433 We next included additional risk variants that are missed in the above set of lead variants
434 from previous GWAS for ovarian¹¹ and pancreatic cancer¹⁰. For lung cancer, we included 26
435 risk variants at $P < 1 \times 10^{-6}$ in European populations from the trans-ancestry GWAS⁹.

436

437 **Identification of putative target proteins for lead variants**

438 To identify potential cancer risk proteins, we mapped GWAS lead variants to *cis*-
439 pQTLs (+/- 500Kb region of a gene) results from three studies among European populations:
440 UK Biobank Pharma Proteomics Project³⁷, Atherosclerosis Risk in Communities study⁴⁴ and
441 deCODE genetics⁴⁵. To increase the power of pQTLs, we combined *cis*-pQTLs from the
442 ARIC and the deCODE (both assayed through SOMAscan[®] platform) via a fixed-effects
443 meta-analysis using META⁸³. *Cis*-pQTLs from the UKB-PPP (assayed through Olink
444 platform) were independently analyzed. In few cases where the lead variant did not overlap
445 with any *cis*-pQTLs, we substituted it with the correlated variant exhibiting the strongest
446 association signal. Putative cancer risk protein was defined based on pQTL significance at a
447 Bonferroni threshold of $P < 0.05$ (nominal $P = 2.3 \times 10^{-5}$, corresponding to 2,164 variant-
448 protein tests for UKB-PPP; nominal $P = 3.8 \times 10^{-5}$, corresponding to 1,322 variant-protein
449 tests for ARIC+deCODE).

450

451 **Colocalization analyses between pQTL and GWAS signals**

452 To identify cancer risk proteins, we conducted colocalization analysis using two
453 approaches: Bayesian method *coloc*⁸⁴ and summary data-based Mendelian Randomization
454 (SMR)⁸⁵. For the SMR approach, a followed HEIDI test is performed on significant SMR
455 results to determine if the colocalized signals can be explained by one single causal variant
456 or by multiple causal variants in the locus. For each protein, SNPs with $P < 0.5$ from GWAS,
457 MAF > 0.01, and within 50 kb of the lead variant were included. To estimate the posterior
458 probability (PP) of colocalization, we utilized the default priors and *coloc.abf* function. In our
459 study, we particularly focused on the assumption that one genetic variant is simultaneously
460 associated with both two traits, which was quantified by PP.H4. We considered a protein to
461 host one shared causal variant from GWAS and pQTLs if its *coloc* PP.H4 > 0.5. Additionally,
462 we also performed SMR+HEIDI analysis for significant *cis*-pQTL with default parameter
463 settings. Specifically, significant SMR+HEIDI results were defined as a tested locus with
464 Bonferroni-adjusted SMR $P < 0.05$ and HEIDI $P \geq 0.05$ (no obvious evidence of

465 heterogeneity of estimated effects or linkage). The above analyses were only conducted in
466 European populations available for both GWAS and proteomics in European populations.

467

468 **Functional genomic analyses**

469 For our identified cancer risk proteins, we examined their xQTLs, including eQTLs,
470 sQTLs, and apaQTL using the resource from the GTEx (version 8). We collected eQTLs and
471 sQTLs from six normal tissues and whole blood from GTEx studies, and we collected
472 apaQTLs from Li's work⁸⁶. A nominal P value < 0.05 for at least one xQTL in either tissue or
473 blood samples was considered supportive of the pQTL results.

474

475 We identified putative regulatory variants in strong linkage disequilibrium (LD) ($r^2 >$
476 0.8 in European population) for lead variants with significant colocalization between GWAS
477 and *cis*-pQTL signals. Using the HaploReg tool⁸⁷, we annotated these variants with a variety
478 of epigenetic annotations, including regulatory chromatin states based on DNase and
479 histone ChIP-Seq from Roadmap Epigenomics Project, histone marks for promoter and
480 enhancer, binding sites of transcription factors, and gene annotation from the GENCODE
481 and RefSeq. We denoted variants as "Proximal" if they overlapped with these functional
482 annotations near the closest target gene. We analyzed a variety of chromatin-chromatin
483 interaction data, from 4D genome⁸⁸, FANTOM5⁸⁹, EnhancerAtlas⁹⁰, and super-enhancer⁹¹.
484 We examined the overlap between putative regulatory variants and enhancer elements in
485 corresponding cell lines or tissues of these six cancer types. We further determined
486 enhancer-promoter loops after combining these data with ChIP-seq data of the histone
487 modification H3K27ac (an active enhancer mark). We focused on interacted loops in which a
488 fragment overlapped an H3K27ac peak (enhancer-like elements). In contrast, the other
489 fragment overlapped the promoter of a gene (defined as a region of upstream 2kb and
490 downstream 100bp around transcript start site). We denoted variants as "Distal" if they
491 overlapped with these chromatin-chromatin variants.

492

493 For our identified cancer risk proteins, we assessed the statistical significance of their
494 differential protein expression between tumor and normal tissue in breast, colorectal, lung,
495 ovarian, and pancreatic cancer samples using data from CPTAC, accessed through the
496 UALCAN website^{92,93}. Similarly, we analyzed their differential gene expression between
497 tumor and normal tissue using data from TCGA, also through the UALCAN website.

498

499 **Inclusion of patients for a focal drug and its control drugs**

500 To evaluate the impact of a focal drug on cancer development, we conducted
501 comparisons between its effects and those of its control drugs. To minimize potential

502 confounding factors associated with a drug prescription, we selected control drugs that
503 belong to the same second-level Anatomical Therapeutic Chemical classification category
504 (ATC-L2) as the focal drug. We formulated emulation trials, each containing one treated
505 patient group (taking the focal drug) and one control patient group (taking the control drug).
506 One focal drug may have multiple trials, depending on the number of potential control drugs
507 belonging to the same ATC-L2 category. Next, we enrolled patients for the treated group and
508 the control group from the VUMC SD based on the following criteria: 1) patients aged ≥ 40 at
509 the time of the latest EHRs or the initial diagnosis of cancer; 2) availability of at least one
510 year of EHRs before the first prescription of the treated/control drug (index date); 3) for
511 cancer patients, a minimum of two exposures to the treated or control drug during the follow-
512 up period (from the index date to the three months before cancer diagnosis); 4) for non-
513 cancer individuals, a minimum of two exposures to control drugs during the follow-up period
514 (from the index date to the date of the latest EHRs). Finally, we precluded patients who were
515 prescribed both treated and control drugs and discarded the trials with less than 500 eligible
516 patients in either patient group⁴⁰.

517

518 **Emulation of treated-control drugs balanced trials**

519 In the IPTW framework⁶¹, individuals are assigned weights based on the inverse of
520 their propensity scores (PS), which represent their probability of being exposed to risk
521 factors or a specific intervention, such as a treated drug, based on their baseline
522 characteristics. In this study, we followed Zang's work⁴⁰ and trained a logistic regression
523 propensity score (LR-PS) model with L1 or L2 regularization on patients' treatment
524 assignments Z and covariates, including age, gender, comorbidities, etc. (**Supplementary**
525 **Table 14**). We trained and selected the logistic model (Eq.1) with the highest area under
526 curve (AUC) using a 10-folder cross-validation. We used the selected model to calculate all
527 patient's stabilized weights (Eq. 2). These weights are used to calculate the standardized
528 mean difference (SMD, Eq.3) of the covariate's prevalence in treated and control groups. A
529 covariate d is defined as unbalanced if $SMD(d) > 0.1$ in IPTW framework (Eqs. 3, 4). A trial
530 is balanced if it contains $\leq 10\%$ unbalanced covariates (Eq. 5).

531

532 The logistic regression is defined as follows:

$$\log\left(\frac{P(\mathbf{Z} = 1)}{1 - P(\mathbf{Z} = 1)}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \quad 1$$

$$P(\mathbf{Z} = 1) = \frac{1}{1 + e^{-(\sum \beta_j X_j + \beta_0)}}$$

533 where \mathbf{Z} refers to treatment assignment (1 for treated patient group and 0 for control patient
 534 group) and $\mathbf{X}(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$ for baseline covariates. The propensity score is defined as
 535 $P(\mathbf{Z} = 1 | \mathbf{X})$ and the stabilized IPTW of each individual is calculated as follows:

$$\mathbf{w} = \frac{\mathbf{Z} \times P(\mathbf{Z} = 1)}{P(\mathbf{Z} = 1 | \mathbf{X})} + \frac{(1 - \mathbf{Z}) \times (1 - P(\mathbf{Z} = 1))}{1 - P(\mathbf{Z} = 1 | \mathbf{X})} \quad 2$$

536

537 Standardized mean difference is calculated as following:

$$SMD(\mathbf{x}_{treat} - \mathbf{x}_{control}) = \frac{|\boldsymbol{\mu}_{treat} - \boldsymbol{\mu}_{control}|}{\sqrt{(\mathbf{S}_{treat}^2 + \mathbf{S}_{control}^2)/2}} \quad 3$$

538 $\mathbf{x}_{treat}, \mathbf{x}_{control} \in R^D$, representing vectors of D number of covariates of treated group and
 539 control group respectively; $\boldsymbol{\mu}_{treat}, \boldsymbol{\mu}_{control}$ are their sample means, and $\mathbf{S}_{treat}^2, \mathbf{S}_{control}^2$ are
 540 their sample variances. In IPTW framework, the weighted sample mean $\boldsymbol{\mu}_w$ and sample
 541 variance \mathbf{S}_w^2 are calculated as following:

$$\boldsymbol{\mu}_w = \frac{\sum \mathbf{w}_i \mathbf{x}_i}{\sum \mathbf{w}_i} \quad 4$$

$$\mathbf{S}_w^2 = \frac{\sum \mathbf{w}_i}{(\sum \mathbf{w}_i)^2 - \sum \mathbf{w}_i^2} \sum \mathbf{w}_i (\mathbf{x}_i - \boldsymbol{\mu}_w)^2$$

542

543 Number of unbalanced covariates are calculated as following:

$$n = \sum_{d=1}^D \mathbb{1}[SMD(d) > 0.1] \quad 5$$

544 where D is the total number of covariates in the model, d is one covariate

545

546 **Logistic regression propensity score (LR-PS) hyperparameter selection and model** 547 **training**

548 To select the optimal regulation penalty weight (λ), we applied 10-fold cross-
 549 validation on a list of lambda elements $\lambda \in [0.005, 0.01, 0.05, 0.1, 0.5]$. Specifically, the
 550 logistic model was trained on 9 training folders, and learnable parameters (β) were
 551 estimated through minimizing the binary cross-entropy loss with L1 (Eq. 6) or L2 penalty (Eq.
 552 7). On the left-one out validation folder (k), we calculated SMD_k (Eq. 3) values for
 553 D covariates based on individuals' weights (Eq. 2) as well as the number unbalanced
 554 covariates n_k (Eq. 5). In addition, we evaluated trained model's prediction performance
 555 using area under curve (AUC_k) on the validation dataset. The same processes were
 556 repeated 10 times. We defined the optimal hyperparameter value is the value generates the
 557 smallest averaged n_k . For two hyperparameter values generate approximate n_k , the one
 558 with larger averaged AUC_k is the optimal. Finally, we trained LR-PS model's learnable

559 parameters (β) on all subjects with the optimal hyperparameter and leveraged this trained
560 model to compute weights and the proportion of imbalanced covariates to pinpoint balanced
561 trials.

562

563 Binary cross-entropy loss function with LASSO (L1) penalization:

$$\operatorname{argmin}_{\beta}(\text{Loss}) = -\frac{1}{n} \left(\sum_{i=1}^n z_i \log(P(z_i = 1)) + (1 - z_i) \log(1 - P(z_i = 1)) \right) + \lambda_1 \sum_{j=1}^b |\beta_j| \quad 6$$

564 Binary cross-entropy loss function with ridge (L2) penalization:

$$\operatorname{argmin}_{\beta}(\text{Loss}) = -\frac{1}{n} \left(\sum_{i=1}^n z_i \log(P(z_i = 1)) + (1 - z_i) \log(1 - P(z_i = 1)) \right) + \lambda_2 \sum_{j=1}^b \beta_j^2 \quad 7$$

565 **Calculation of overall hazard ratio for cancer risk drug on cancer development risk**

566 We evaluate subjects' hazard of developing/preventing cancer in balanced treated-
567 control trials through survival analyses. We applied weighted Cox proportional hazard
568 model⁹⁴ using the lifelines 0.28.0 Python package to systematically evaluate the hazard ratio
569 of developing cancer for patients taking treated drug vs patients taking control drugs (time-
570 to-event). The time windows utilized in this study start from the earliest date in EHRs for
571 prescription of the treated/control drug to patients and end at the date of the diagnosis of
572 cancer (event) or the end of EHRs records (censored). We included unbalanced covariates
573 (if exist) into Cox models. For a treated drug, its overall hazard ratio and p-value were
574 obtained by applying a random effect meta-analysis on the hazard ratios from its eligible
575 trials using the meta 7.0 R package (**Supplementary Table 15**). We reported that a treated
576 drug has a significantly increasing or decreasing risk of cancer development, contrasting to
577 its control drugs, if the overall hazard ratio has a $P < 0.05$ after Bonferroni correction
578 (nominal $P = 3.5 \times 10^{-3}$ corresponding to 14 tests).

579

580 **Data availability**

581 Supplementary Table 1 provides the download information for the summary statistics of
582 GWAS data for the six common cancers, including breast, ovary, prostate, colorectum, lung,
583 and pancreas. Metadata and pQTL summary statistics from UKB-PPP can be downloaded
584 from Synapse: Project SynID: syn51364943; pQTL from ARIC⁴⁶ and deCODE genetics⁴⁷ can
585 be accessed through previous publications (PMID: 34857953 and PMID: 35501419).
586 Functional genomic data includes: TCGA and CPTAC differential expression results
587 accessible through <https://ualcan.path.uab.edu/index.html>; 4DGenome:
588 <https://4dgenome.research.chop.edu/>; Depmap : <https://depmap.org/portal/>; FANTOM5:
589 <http://fantom.gsc.riken.jp/5/>. HaploReg: <https://pubs.broadinstitute.org/mammals/haploreg/>.

590 GTEx: <https://gtexportal.org/home/>. GENCODE (v26.GRCh38) was downloaded from
591 https://www.genecodegenes.org/human/release_26.html. National Cancer Institute can be
592 accessed through <https://www.cancer.gov/about-cancer/treatment/drugs>; CGC can be
593 accessed accessed via COSMIC website: <https://cancer.sanger.ac.uk/census>. Drugs and
594 compounds data can be downloaded from the following URLs: ChEMBL:
595 <https://www.ebi.ac.uk/chembl/>; Therapeutic Target Database: <https://db.idrblab.net/ttd/>;
596 Open Targets: <https://www.opentargets.org/>; DrugBank: <https://go.drugbank.com/>. The EHR
597 data, containing de-identified clinical information, can be accessed through the VUMC SD
598 database. Data is available through restricted access for approved studies and researchers
599 who agree to specific conditions of use.

600

601 **Code availability**

602 The developed pipeline and main source R codes that are used in this work are available
603 from the GitHub website of Xingyi Guo's lab: https://github.com/XingyiGuo/PQTL_EHR/

604 **Declaration of Interests**

605 The authors declare no competing interests.

606

607 **Acknowledgments**

608 This work was supported by the US National Institutes of Health grant 1R37CA227130-01A1

609 and R01CA269589-01A1 to X.G. The data analyses were conducted using the Advanced

610 Computing Center for Research and Education (ACCRE) at Vanderbilt University. New

611 Frontiers in Research Fund (NFRFE-2018-00748) and NSERC Discovery Grant (RGPIN-

612 2024-04679) to Q.L. The computational infrastructure was partly supported by a Canada

613 Foundation for Innovation JELF grant (36605) to Q.L.

614 References

- 615 1. Nelson, M.R., *et al.* The support of human genetic evidence for approved drug indications.
616 *Nature genetics* **47**, 856-860 (2015).
- 617 2. Diogo, D., *et al.* Phenome-wide association studies across large population cohorts support
618 drug target validation. *Nature communications* **9**, 4285 (2018).
- 619 3. Finan, C., *et al.* The druggable genome and support for target identification and validation in
620 drug development. *Sci Transl Med* **9**(2017).
- 621 4. Peters, U. & Tomlinson, I. Utilizing Human Genetics to Develop Chemoprevention for Cancer-
622 Too Good an Opportunity to be Missed. *Cancer Prev Res (Phila)* **17**, 7-12 (2024).
- 623 5. Jia, G., *et al.* Genome- and transcriptome-wide association studies of 386,000 Asian and
624 European-ancestry women provide new insights into breast cancer genetics. *Am J Hum*
625 *Genet* **109**, 2185-2195 (2022).
- 626 6. Zhang, H., *et al.* Genome-wide association study identifies 32 novel breast cancer
627 susceptibility loci from overall and subtype-specific analyses. *Nature genetics* **52**, 572-581
628 (2020).
- 629 7. Fernandez-Rozadilla, C., *et al.* Deciphering colorectal cancer genetics through multi-omic
630 analysis of 100,204 cases and 154,587 controls of European and east Asian ancestries. *Nat*
631 *Genet* **55**, 89-99 (2023).
- 632 8. Conti, D.V., *et al.* Trans-ancestry genome-wide association meta-analysis of prostate cancer
633 identifies new susceptibility loci and informs genetic risk prediction. *Nature genetics* **53**, 65-
634 75 (2021).
- 635 9. Byun, J., *et al.* Cross-ancestry genome-wide meta-analysis of 61,047 cases and 947,237
636 controls identifies new susceptibility loci contributing to lung cancer. *Nature genetics* **54**,
637 1167-1177 (2022).
- 638 10. Klein, A.P., *et al.* Genome-wide meta-analysis identifies five new susceptibility loci for
639 pancreatic cancer. *Nat Commun* **9**, 556 (2018).
- 640 11. Phelan, C.M., *et al.* Identification of 12 new susceptibility loci for different histotypes of
641 epithelial ovarian cancer. *Nature genetics* **49**, 680-691 (2017).
- 642 12. Lawrenson, K., *et al.* Genome-wide association studies identify susceptibility loci for
643 epithelial ovarian cancer in east Asian women. *Gynecol Oncol* **153**, 343-355 (2019).
- 644 13. McKay, J.D., *et al.* Large-scale association analysis identifies new lung cancer susceptibility
645 loci and heterogeneity in genetic susceptibility across histological subtypes. *Nature genetics*
646 **49**, 1126-1132 (2017).
- 647 14. Wen, W., *et al.* Genetic variations of DNA bindings of FOXA1 and co-factors in breast cancer
648 susceptibility. *Nature communications* **12**, 5318 (2021).
- 649 15. Moreno, V., *et al.* Colon-specific eQTL analysis to inform on functional SNPs. *Br J Cancer* **119**,
650 971-977 (2018).
- 651 16. Chen, Z., *et al.* Identifying Putative Susceptibility Genes and Evaluating Their Associations
652 with Somatic Mutations in Human Cancers. *American journal of human genetics* **105**, 477-
653 492 (2019).
- 654 17. Guo, X., *et al.* A Comprehensive cis-eQTL Analysis Revealed Target Genes in Breast Cancer
655 Susceptibility Loci Identified in Genome-wide Association Studies. *American journal of*
656 *human genetics* **102**, 890-903 (2018).
- 657 18. He, J., *et al.* Integrating transcription factor occupancy with transcriptome-wide association
658 analysis identifies susceptibility genes in human cancers. *Nature communications* **13**, 7118
659 (2022).
- 660 19. Guo, X., *et al.* Identifying Novel Susceptibility Genes for Colorectal Cancer Risk From a
661 Transcriptome-Wide Association Study of 125,478 Subjects. *Gastroenterology* **160**, 1164-
662 1178 e1166 (2021).
- 663 20. Yuan, Y., *et al.* Multi-omics analysis to identify susceptibility genes for colorectal cancer.
664 *Human molecular genetics* **30**, 321-330 (2021).

- 665 21. Bien, S.A., *et al.* Genetic variant predictors of gene expression provide new insight into risk
666 of colorectal cancer. *Hum Genet* **138**, 307-326 (2019).
- 667 22. Wu, L., *et al.* A transcriptome-wide association study of 229,000 women identifies new
668 candidate susceptibility genes for breast cancer. *Nature genetics* **50**, 968-978 (2018).
- 669 23. Gao, G., *et al.* A joint transcriptome-wide association study across multiple tissues identifies
670 candidate breast cancer susceptibility genes. *Am J Hum Genet* **110**, 950-962 (2023).
- 671 24. Mancuso, N., *et al.* Large-scale transcriptome-wide association study identifies new prostate
672 cancer risk regions. *Nature communications* **9**, 4079 (2018).
- 673 25. Liu, D., *et al.* A transcriptome-wide association study identifies novel candidate susceptibility
674 genes for prostate cancer risk. *Int J Cancer* **150**, 80-90 (2022).
- 675 26. Bosse, Y., *et al.* Transcriptome-wide association study reveals candidate causal genes for
676 lung cancer. *International journal of cancer* **146**, 1862-1878 (2020).
- 677 27. Lu, Y., *et al.* A Transcriptome-Wide Association Study Among 97,898 Women to Identify
678 Candidate Susceptibility Genes for Epithelial Ovarian Cancer Risk. *Cancer research* **78**, 5419-
679 5430 (2018).
- 680 28. Gusev, A., *et al.* A transcriptome-wide association study of high-grade serous epithelial
681 ovarian cancer identifies new susceptibility genes and splice variants. *Nature genetics* **51**,
682 815-823 (2019).
- 683 29. Zhong, J., *et al.* A Transcriptome-Wide Association Study Identifies Novel Candidate
684 Susceptibility Genes for Pancreatic Cancer. *J Natl Cancer Inst* **112**, 1003-1012 (2020).
- 685 30. Zheng, C.J., Han, L.Y., Yap, C.W., Ji, Z.L., Cao, Z.W. & Chen, Y.Z. Therapeutic targets: Progress
686 of their exploration and investigation of their characteristics (vol 58, pg 259, 2006).
687 *Pharmacol Rev* **58**, 682-682 (2006).
- 688 31. Zhu, J., *et al.* Associations between Genetically Predicted Blood Protein Biomarkers and
689 Pancreatic Cancer Risk. *Cancer epidemiology, biomarkers & prevention : a publication of the*
690 *American Association for Cancer Research, cosponsored by the American Society of*
691 *Preventive Oncology* **29**, 1501-1508 (2020).
- 692 32. Shu, X., *et al.* Evaluation of associations between genetically predicted circulating protein
693 biomarkers and breast cancer risk. *International journal of cancer* **146**, 2130-2138 (2020).
- 694 33. Wu, L., *et al.* Analysis of Over 140,000 European Descendants Identifies Genetically
695 Predicted Blood Protein Biomarkers Associated with Prostate Cancer Risk. *Cancer research*
696 **79**, 4592-4598 (2019).
- 697 34. Gregga, I., *et al.* Predicted proteome association studies of breast, prostate, ovarian, and
698 endometrial cancers implicate plasma protein regulation in cancer susceptibility. *Cancer*
699 *Epidemiol Biomarkers Prev* (2023).
- 700 35. Jia, G., *et al.* Identification of target proteins for breast cancer genetic risk loci and blood risk
701 biomarkers in a large study by integrating genomic and proteomic data. *Int J Cancer* **152**,
702 2314-2320 (2023).
- 703 36. Considine, D.P.C., *et al.* Genetically predicted circulating protein biomarkers and ovarian
704 cancer risk. *Gynecol Oncol* **160**, 506-513 (2021).
- 705 37. Sun, B.B., *et al.* Plasma proteomic associations with genetics and health in the UK Biobank.
706 *Nature* (2023).
- 707 38. Tautermann, C.S. Current and Future Challenges in Modern Drug Discovery. *Methods Mol*
708 *Biol* **2114**, 1-17 (2020).
- 709 39. Pushpakom, S., *et al.* Drug repurposing: progress, challenges and recommendations. *Nat Rev*
710 *Drug Discov* **18**, 41-58 (2019).
- 711 40. Zang, C., *et al.* High-throughput target trial emulation for Alzheimer's disease drug
712 repurposing with real-world data. *Nature communications* **14**, 8180 (2023).
- 713 41. Wu, Y., *et al.* Discovery of Noncancer Drug Effects on Survival in Electronic Health Records of
714 Patients With Cancer: A New Paradigm for Drug Repurposing. *JCO Clin Cancer Inform* **3**, 1-9
715 (2019).

- 716 42. Xu, H., *et al.* Validating drug repurposing signals using electronic health records: a case study
717 of metformin associated with reduced cancer mortality. *J Am Med Inform Assoc* **22**, 179-191
718 (2015).
- 719 43. Bejan, C.A., Cahill, K.N., Staso, P.J., Choi, L., Peterson, J.F. & Phillips, E.J. DrugWAS: Drug-wide
720 Association Studies for COVID-19 Drug Repurposing. *Clin Pharmacol Ther* **110**, 1537-1546
721 (2021).
- 722 44. Reznikov, L.R., *et al.* Identification of antiviral antihistamines for COVID-19 repurposing.
723 *Biochem Biophys Res Commun* **538**, 173-179 (2021).
- 724 45. Liu, R., Wei, L. & Zhang, P. A deep learning framework for drug repurposing via emulating
725 clinical trials on real-world patient data. *Nat Mach Intell* **3**, 68-75 (2021).
- 726 46. Zhang, J., *et al.* Plasma proteome analyses in individuals of European and African ancestry
727 identify cis-pQTLs and models for proteome-wide association studies. *Nat Genet* **54**, 593-602
728 (2022).
- 729 47. Ferkingstad, E., *et al.* Large-scale integration of the plasma proteome with genetics and
730 disease. *Nat Genet* **53**, 1712-1721 (2021).
- 731 48. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable
732 selection in regression, with application to genetic fine mapping. *J R Stat Soc Series B Stat*
733 *Methodol* **82**, 1273-1300 (2020).
- 734 49. Law, V., *et al.* DrugBank 4.0: shedding new light on drug metabolism. *Nucleic Acids Res* **42**,
735 D1091-1097 (2014).
- 736 50. Gaulton, A., *et al.* ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic*
737 *acids research* **40**, D1100-1107 (2012).
- 738 51. Zhou, Y., *et al.* TTD: Therapeutic Target Database describing target druggability information.
739 *Nucleic Acids Res* (2023).
- 740 52. Ochoa, D., *et al.* The next-generation Open Targets Platform: reimaged, redesigned,
741 rebuilt. *Nucleic Acids Res* **51**, D1353-D1359 (2023).
- 742 53. Chesnaye, N.C., *et al.* An introduction to inverse probability of treatment weighting in
743 observational research. *Clin Kidney J* **15**, 14-20 (2022).
- 744 54. Fachal, L., *et al.* Fine-mapping of 150 breast cancer risk regions identifies 191 likely target
745 genes. *Nature genetics* **52**, 56-73 (2020).
- 746 55. Chen, Z., *et al.* Fine-mapping analysis including over 254,000 East Asian and European
747 descendants identifies 136 putative colorectal cancer susceptibility genes. *Nat Commun* **15**,
748 3557 (2024).
- 749 56. Shu, X., *et al.* Associations between circulating proteins and risk of breast cancer by intrinsic
750 subtypes: a Mendelian randomisation analysis. *Br J Cancer* **127**, 1507-1514 (2022).
- 751 57. Sun, J., *et al.* Identification of novel protein biomarkers and drug targets for colorectal
752 cancer by integrating human plasma proteome with genome. *Genome Med* **15**(2023).
- 753 58. Bailey, M.H., *et al.* Comprehensive Characterization of Cancer Driver Genes and Mutations.
754 *Cell* **174**, 1034-1035 (2018).
- 755 59. Dietlein, F., *et al.* Identification of cancer driver genes based on nucleotide context. *Nature*
756 *genetics* **52**, 208-218 (2020).
- 757 60. Sondka, Z., Bamford, S., Cole, C.G., Ward, S.A., Dunham, I. & Forbes, S.A. The COSMIC Cancer
758 Gene Census: describing genetic dysfunction across all human cancers. *Nat Rev Cancer* **18**,
759 696-705 (2018).
- 760 61. Austin, P.C. & Stuart, E.A. Moving towards best practice when using inverse probability of
761 treatment weighting (IPTW) using the propensity score to estimate causal treatment effects
762 in observational studies. *Stat Med* **34**, 3661-3679 (2015).
- 763 62. Taipale, H., Solmi, M., Lahtenvuo, M., Tanskanen, A., Correll, C.U. & Tiihonen, J.
764 Antipsychotic use and risk of breast cancer in women with schizophrenia: a nationwide
765 nested case-control study in Finland. *Lancet Psychiatry* **8**, 883-891 (2021).

- 766 63. Wu, G., *et al.* Targeting Gas6/TAM in cancer cells and tumor microenvironment. *Mol Cancer*
767 **17**, 20 (2018).
- 768 64. Jansen, F.H., *et al.* Profiling of antibody production against xenograft-released proteins by
769 protein microarrays discovers prostate cancer markers. *J Proteome Res* **11**, 728-735 (2012).
- 770 65. Chiang, J.Y., *et al.* Haloperidol Instigates Endometrial Carcinogenesis and Cancer Progression
771 by the NF-kappaB/CSF-1 Signaling Cascade. *Cancers (Basel)* **14**(2022).
- 772 66. Hippisley-Cox, J., Vinogradova, Y., Coupland, C. & Parker, C. Risk of malignancy in patients
773 with schizophrenia or bipolar disorder: nested case-control study. *Arch Gen Psychiatry* **64**,
774 1368-1376 (2007).
- 775 67. Koh, S.L., Ager, E.I., Costa, P.L., Malcontenti-Wilson, C., Muralidharan, V. & Christophi, C.
776 Blockade of the renin-angiotensin system inhibits growth of colorectal cancer liver
777 metastases in the regenerating liver. *Clin Exp Metastasis* **31**, 395-405 (2014).
- 778 68. Childers, W.K. Interactions of the renin-angiotensin system in colorectal cancer and
779 metastasis. *Int J Colorectal Dis* **30**, 749-752 (2015).
- 780 69. Riddiough, G.E., *et al.* Captopril, a Renin-Angiotensin System Inhibitor, Attenuates Features
781 of Tumor Invasion and Down-Regulates C-Myc Expression in a Mouse Model of Colorectal
782 Cancer Liver Metastasis. *Cancers (Basel)* **13**(2021).
- 783 70. Kristensen, K.B., Hicks, B., Azoulay, L. & Pottgard, A. Use of ACE (Angiotensin-Converting
784 Enzyme) Inhibitors and Risk of Lung Cancer: A Nationwide Nested Case-Control Study. *Circ*
785 *Cardiovasc Qual Outcomes* **14**, e006687 (2021).
- 786 71. Yarmolinsky, J., *et al.* Genetically proxied therapeutic inhibition of antihypertensive drug
787 targets and risk of common cancers: A mendelian randomization analysis. *PLoS Med* **19**,
788 e1003897 (2022).
- 789 72. Karakus, F., Eyol, E., Yilmaz, K. & Unuvar, S. Inhibition of cell proliferation, migration and
790 colony formation of LS174T Cells by carbonic anhydrase inhibitor. *Afr Health Sci* **18**, 1303-
791 1310 (2018).
- 792 73. Noma, N., *et al.* Impact of Acetazolamide, a Carbonic Anhydrase Inhibitor, on the
793 Development of Intestinal Polyps in Min Mice. *Int J Mol Sci* **18**(2017).
- 794 74. Schmit, S.L., Rennert, H.S., Rennert, G., Gruber, S.B.J.C.E., Biomarkers & Prevention. Coffee
795 consumption and the risk of colorectal cancer. **25**, 634-639 (2016).
- 796 75. Sartini, M., *et al.* Coffee Consumption and Risk of Colorectal Cancer: A Systematic Review
797 and Meta-Analysis of Prospective Studies. *Nutrients* **11**(2019).
- 798 76. Um, C.Y., McCullough, M.L., Guinter, M.A., Campbell, P.T., Jacobs, E.J. & Gapstur, S.M.
799 Coffee consumption and risk of colorectal cancer in the Cancer Prevention Study-II Nutrition
800 Cohort. *Cancer Epidemiol* **67**, 101730 (2020).
- 801 77. Micek, A., Gniadek, A., Kawalec, P. & Brzostek, T. Coffee consumption and colorectal cancer
802 risk: a dose-response meta-analysis on prospective cohort studies. *Int J Food Sci Nutr* **70**,
803 986-1006 (2019).
- 804 78. Schumacher, F.R., *et al.* Association analyses of more than 140,000 men identify 63 new
805 prostate cancer susceptibility loci. *Nature genetics* **50**, 928-936 (2018).
- 806 79. Huyghe, J.R., *et al.* Discovery of common and rare genetic risk variants for colorectal cancer.
807 *Nature genetics* **51**, 76-87 (2019).
- 808 80. Chen, Z., *et al.* Novel insights into genetic susceptibility for colorectal cancer from
809 transcriptome-wide association and functional investigation. *J Natl Cancer Inst* **116**, 127-137
810 (2024).
- 811 81. Roden, D.M., *et al.* Development of a large-scale de-identified DNA biobank to enable
812 personalized medicine. *Clin Pharmacol Ther* **84**, 362-369 (2008).
- 813 82. Weissbrod, O., *et al.* Functionally informed fine-mapping and polygenic localization of
814 complex trait heritability. *Nat Genet* **52**, 1355-1363 (2020).
- 815 83. Schwarzer, G., Carpenter, J.R. & Rücker, G. *Meta-analysis with R*, (Springer, 2015).

- 816 84. Giambartolomei, C., *et al.* Bayesian test for colocalisation between pairs of genetic
817 association studies using summary statistics. *PLoS genetics* **10**, e1004383 (2014).
818 85. Zhu, Z., *et al.* Integration of summary data from GWAS and eQTL studies predicts complex
819 trait gene targets. *Nat Genet* **48**, 481-487 (2016).
820 86. Li, L., *et al.* An atlas of alternative polyadenylation quantitative trait loci contributing to
821 complex trait and disease heritability. *Nature genetics* **53**, 994-1005 (2021).
822 87. Ward, L.D. & Kellis, M. HaploReg: a resource for exploring chromatin states, conservation,
823 and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res*
824 **40**, D930-934 (2012).
825 88. He, B., Chen, C., Teng, L. & Tan, K. Global view of enhancer-promoter interactome in human
826 cells. *Proceedings of the National Academy of Sciences of the United States of America* **111**,
827 E2191-2199 (2014).
828 89. Lizio, M., *et al.* Gateways to the FANTOM5 promoter level mammalian expression atlas.
829 *Genome Biol* **16**, 22 (2015).
830 90. Gao, T. & Qian, J. EnhancerAtlas 2.0: an updated resource with enhancer annotation in 586
831 tissue/cell types across nine species. *Nucleic acids research* **48**, D58-D64 (2020).
832 91. Wang, Y., *et al.* SEdb 2.0: a comprehensive super-enhancer database of human and mouse.
833 *Nucleic Acids Res* **51**, D280-D290 (2023).
834 92. Chandrashekar, D.S., *et al.* UALCAN: A Portal for Facilitating Tumor Subgroup Gene
835 Expression and Survival Analyses. *Neoplasia* **19**, 649-658 (2017).
836 93. Chandrashekar, D.S., *et al.* UALCAN: An update to the integrated cancer data analysis
837 platform. *Neoplasia* **25**, 18-27 (2022).
838 94. Lin, D.Y. & Wei, L.-J. The robust inference for the Cox proportional hazards model. *J Am Stat*
839 *Assoc* **84**, 1074-1078 (1989).

840

841 **Figures**

842 **Fig. 1: Overview of the Analytical Framework.**

843 **a**, An illustration depicting the identification of proteins associated with the risk of the six
844 major cancers: breast, lung, colorectal, ovarian, pancreatic, and prostate. Population-based
845 proteomics data (for pQTLs) and GWAS data resources (for identifying lead variants) utilized
846 in this study are shown in the left panels. Meta-analyses of *cis*-pQTLs from ARIC and
847 deCODE, conducted through the SOMAscan® platform, were combined with pQTL results
848 from the UKB-PPP to identify potential risk proteins, as depicted in the middle panels.
849 Colocalization analyses between GWAS summary statistics and *cis*-pQTLs were performed
850 to identify cancer risk proteins with high confidence, as illustrated in the right panel. **b**, The
851 proteins with evidence of colocalization annotated based on drug-protein information from
852 four databases: DrugBank, ChEMBL, TTD, and OpenTargets. **c**, The framework for
853 evaluating the effects of drugs approved for indications on cancer risk. The Inverse
854 Probability of Treatment weighting (IPTW) framework was utilized to construct emulations of
855 treated-control drug trials based on millions of patients' Electronic Health Records stored at
856 VUMC SD (left pane). In these emulations, the Cox proportional hazard model was
857 conducted for each trial to assess the hazard ratio (HR) of cancer risk between the treated
858 focal drug and the control drug (right panels).

859

860 **Fig.2: Genome-wide distribution of lead variants and putative risk proteins among six**
861 **types of cancer.** Proteins identified for each cancer are represented by different colors.

862 Each circle represents a single lead variant-protein pair. Proteins marked with an asterisk (*)
863 denote multiple proteins associated with the lead variants. A dashed box highlights several
864 well-known cancer-related proteins, such as HLA-A and HLA-E, which are linked to lead
865 variants located in the major histocompatibility complex (MHC).

866

867 **Fig. 3: Identification of 101 cancer risk proteins through pQTL and colocalization**
868 **analysis**

869 **a**, Number of proteins showing evidence of colocalizations between pQTLs and GWAS
870 association signals for six cancer types. **b**, Percentage of proteins showing evidence of
871 colocalizations between pQTLs and GWAS summary statistics for six cancer types.

872 **c**, A plot illustrating the high consistency of pQTL p-values for 22 cancer risk proteins
873 between the ARIC+deCODE and the UKB-PPP (proteins commonly assayed from
874 SOMAscan® and Olink platforms).

875

876 **Fig. 4: A circular plot showing 36 druggable proteins potentially targeted by 404**

877 **approved drugs or undergoing clinical trials for cancer treatment or other indications**

878 Presented from inner to outer layers are cancer types, proteins, and drugs. Each drug-
879 protein interaction is annotated by DrugBank, ChEMBL, TTD, and OpenTargets, with lines in
880 different colors representing each database. Interactions where proteins are annotated by
881 two databases are linked to drugs with thick lines.

882

883 **Fig. 5: A circular plot showing 19 druggable proteins potentially targeted by 133**
884 **approved drugs or undergoing clinical trials for cancer treatment.** Presented from inner
885 to outer layers are cancer types, proteins, drugs and cancers. Drugs approved and
886 undergoing clinical trials for cancer treatment are highlighted on green and gray,
887 respectively. Drug approved indications are formatted in bold, while indications under clinical
888 trial are in regular font.

889

890 **Fig. 6: Drugs approved for treated indications showing significant effects on cancer**
891 **risk**

892 **a**, A table showing cancer risk alleles of lead variants, risk proteins, and drug name
893 approved for indications. Positive associations are indicated by upward arrows, while
894 negative associations are indicated by downward arrows. **b**, Boxplots showing differentially
895 expressed proteins between normal and tumor colon tissues using data from CPTAC.
896 **c**, An illustration of drugs linked to specific cancers based on the risk proteins targeted by
897 the drugs. **d**, Survival plots depict the statistically significant difference in the probability of
898 being cancer-free for patients in the treated group (taking a focal drug, shown in green)
899 compared to control groups (shown in purple). The shaded area represents the 95%
900 confidence interval. The overall hazard ratio and P-value for the focal drug, determined
901 through Cox proportional hazard models, are presented in the top right corner of each panel.
902

903 **Extended Data Fig. 1: A flowchart for characterizing lead variants with independent**
904 **risk signals in six cancer types.** The analysis for each of the six major cancers: breast,
905 lung, colorectal, ovarian, pancreatic, and prostate is separated by dashed lines. The detailed
906 protocols of new efforts from our additional fine-mapping analysis using SuSiE and a
907 collection of previously identified risk variants from GWAS or fine-mapping studies are
908 indicated in Box A, Box B and Box C, respectively.

909

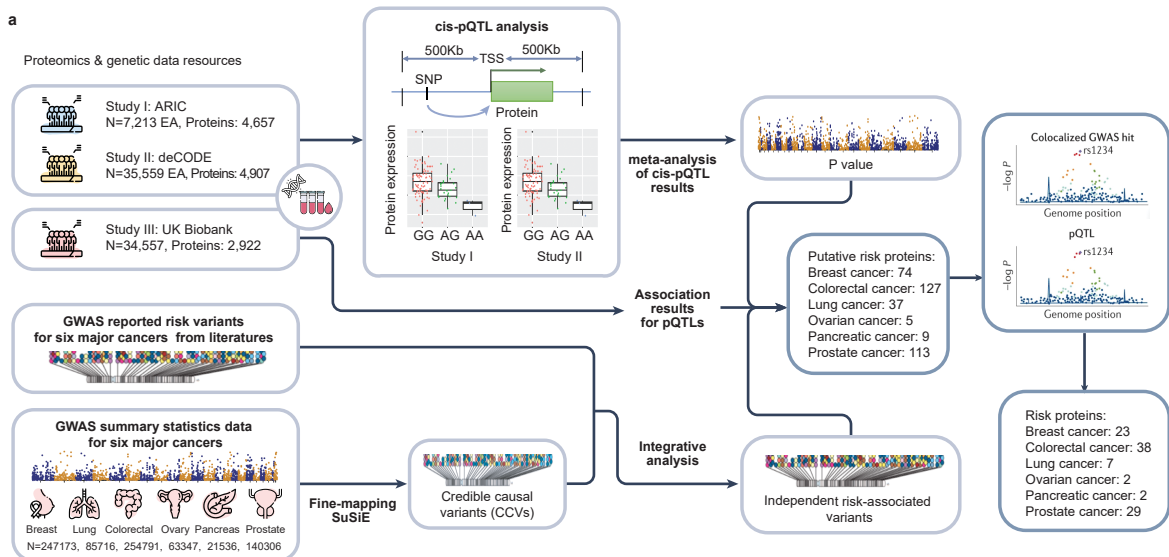
910 **Extended Data Fig. 2: Common cancer risk proteins identified across six cancer types**
911 **a,** A Venn plot showing common proteins in breast, colorectal, lung and prostate cancer. **b,**
912 A heatmap showing common proteins observed from at least two of six types of cancer.

913

914 **Extended Data Fig. 3: A circular plot showing 28 druggable proteins potentially**
915 **targeted by 197 approved drugs or undergoing clinical trials for treated indications**
916 **rather than cancers.** Presented from inner to outer layers are cancer types, proteins, drugs
917 and cancers. Drugs approved and undergoing clinical trials for cancer treatment are
918 highlighted on blue and gray, respectively. Drug approved indications are formatted in bold,
919 while indications under clinical trial are in regular font.

920

a



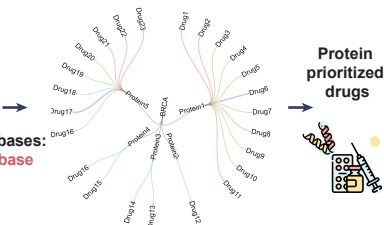
b

Proteins (drug targets):

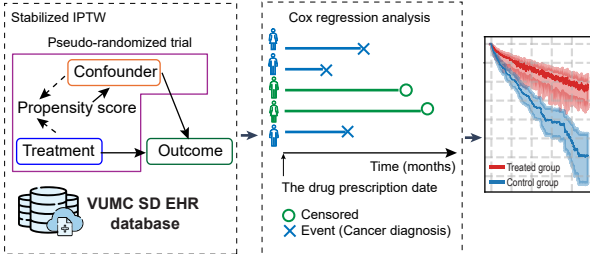
- Breast cancer: 9
- Colorectal cancer: 15
- Lung cancer: 3
- Ovarian cancer: 1
- Prostate cancer: 12

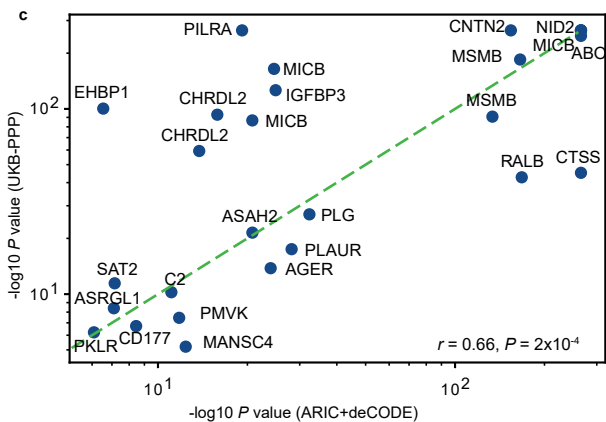
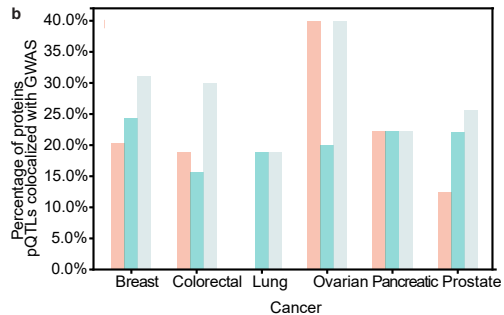
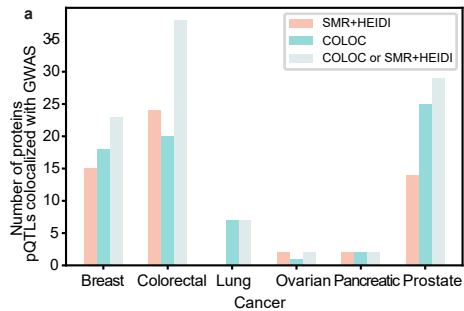
Drugs/compounds Databases:

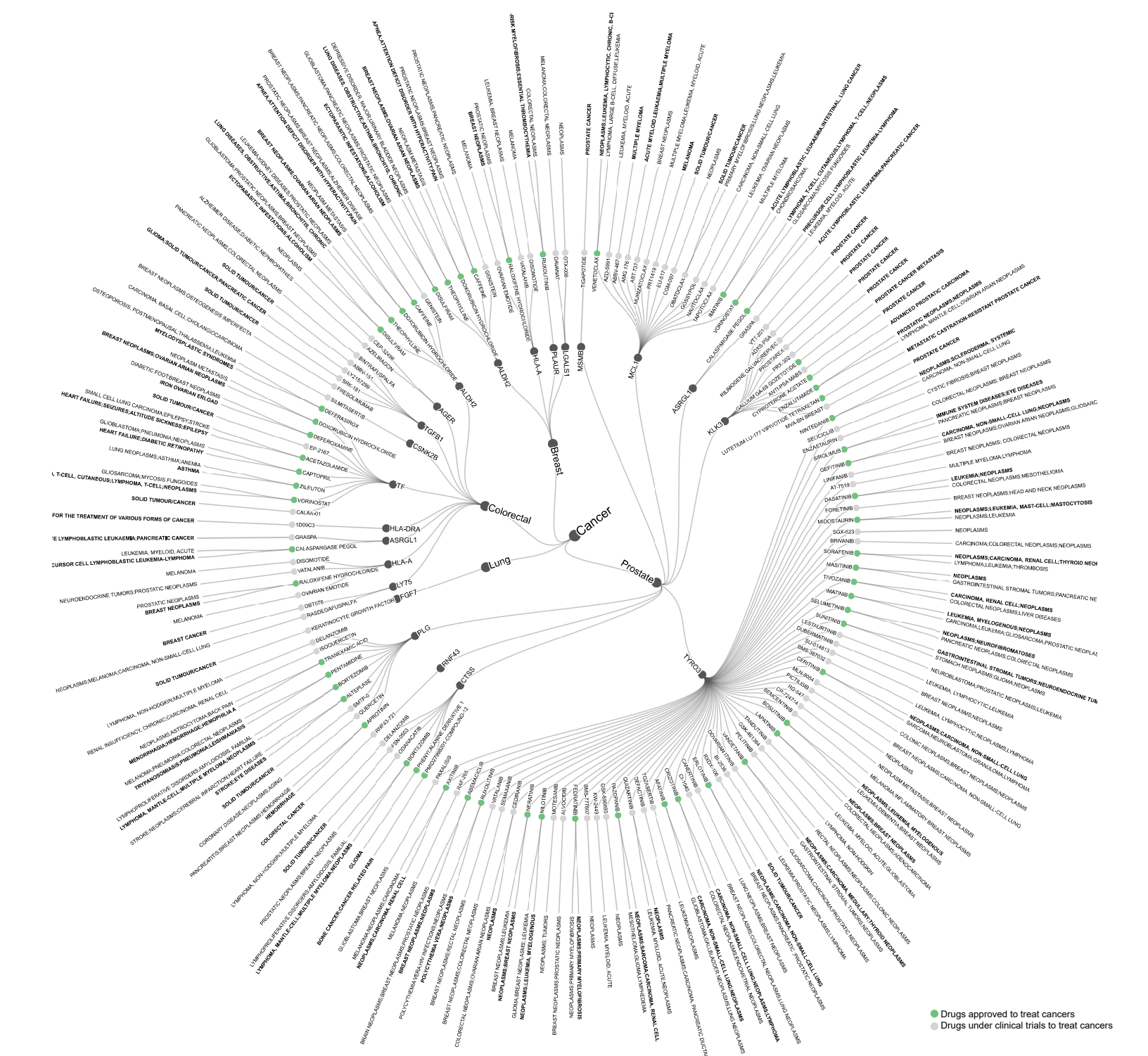
- Therapeutic Target Database
- Open Targets
- ChEMBL
- DrugBank



c



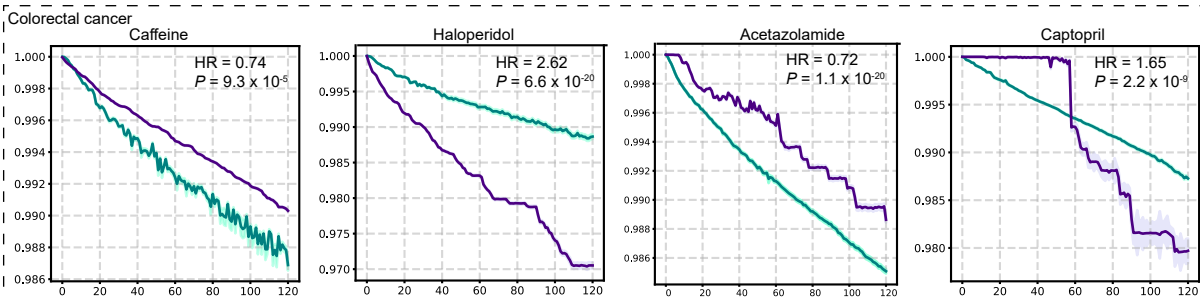
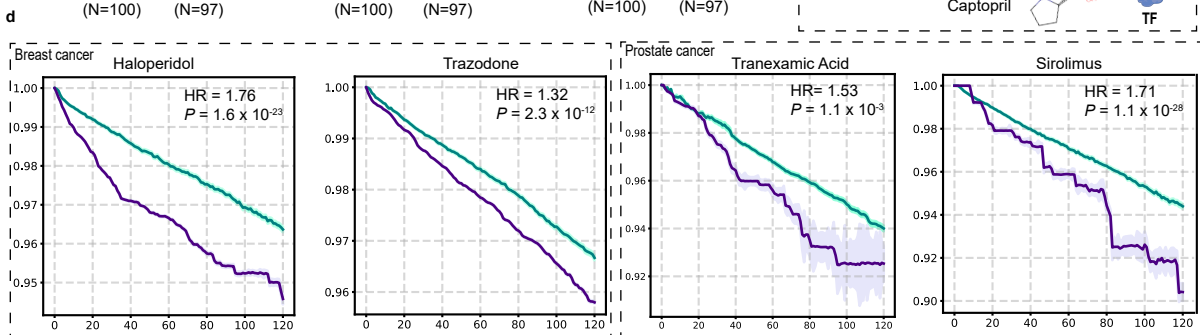
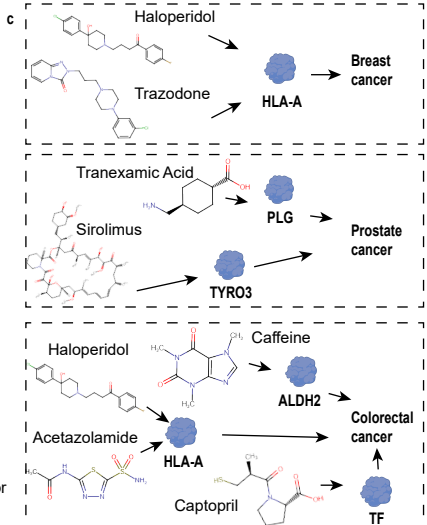
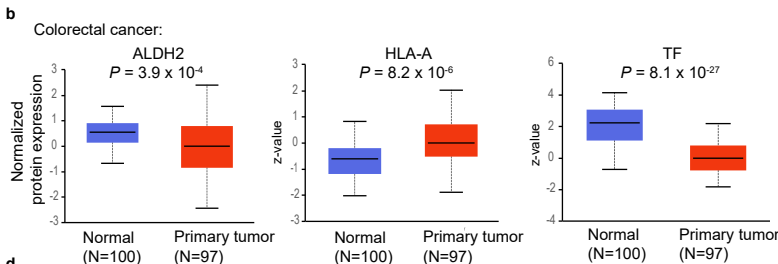




● Drugs approved to treat cancers
 ● Drugs under clinical trials to treat cancers

a

Cancer	Protein	Lead variant	Risk allele	OR (GWAS)	Beta (pQTL)	Drug
Breast	HLA-A	rs79309050	T	1.13	↑ -0.47	↓ Haloperidol
	HLA-A	rs79309050	T	1.13	↑ -0.47	↓ Trazodone
Prostate	PLG	rs9347480	C	1.06	↑ -0.12	↓ Tranexamic Acid
	TYRO3	rs11561564	G	1.06	↑ 0.14	↑ Sirolimus
Colorectal	ALDH2	rs3858704	G	1.04	↑ -0.04	↓ Caffeine
	HLA-A	rs2517671	G	1.04	↑ 0.83	↑ Haloperidol
	HLA-A	rs2517671	G	1.04	↑ 0.83	↑ Acetazolamide
	TF	rs4854776	C	1.04	↑ -0.04	↓ Captopril



— Aggregated Cancer-free Probability (Control drugs) — 95% CI (Control drugs) — Aggregated Cancer-free Probability (Focal drug) — 95% CI (Focal drug)

Breast cancer

Study I (fine-mapping): 196 lead variants with strong independent association signals, at $P < 1 \times 10^{-6}$ (Fachal et al., Nature Genetics, 2020)

Study II (GWAS): 32 risk variants (Zhang et al., Nature Genetics, 2020)

Our fine-mapping analysis using SuSiE based on GWAS data in European populations (N=247,173)

1. Combining lead variants from study I and II, after removing those in linkage disequilibrium (LD) ($r^2 \geq 0.1$ in European populations).
2. Including additional lead variants from our SuSiE analysis with LD $r^2 < 0.1$ in European populations with the above set of lead variants
 - those with independent risk-associated signals at $P < 5 \times 10^{-8}$
 - those located in GWAS loci with independent risk-associated signals at $P < 1 \times 10^{-6}$

Identified 227 lead variants

Colorectal cancer

Our recent fine-mapping study: 238 lead variants (Chen et al., Nature Communications, 2024)

Box 1

Identified 213 lead variants

Lung cancer

Trans-ancestry GWAS: 36 risk variants (Byun et al., Nature Genetics, 2022)

Box 1

Identified 26 lead variants

Prostate cancer

GWAS data in European populations (N=140,306)

Box 2

Trans-ancestry GWAS (N=234,253) (Conti et al., Nature Genetics, 2021)

Box 3

Identified 213 lead variants

Ovarian cancer

GWAS data in European populations (N = 63,347)

Box 2

European population GWAS (Phelan et al., Nature Genetics, 2017)

Box 3

Identified 13 lead variants

Pancreatic cancer

GWAS data in European populations (N = 21,536)

Box 2

European population GWAS (Klein et al., Nature Communication, 2018)

Box 3

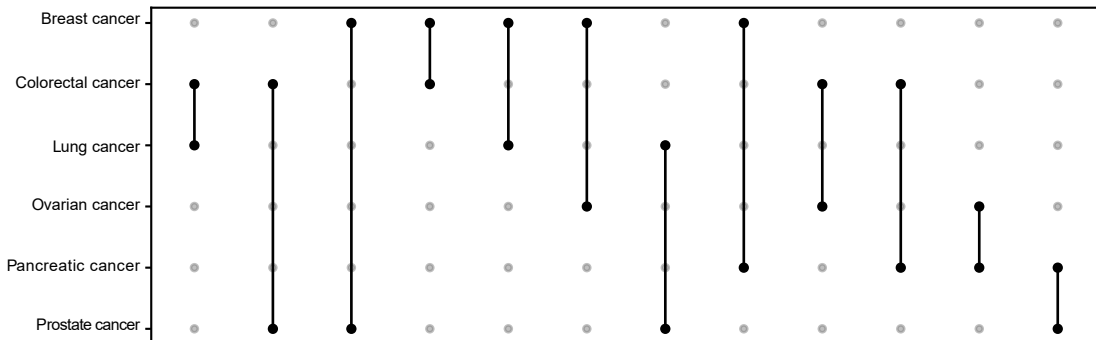
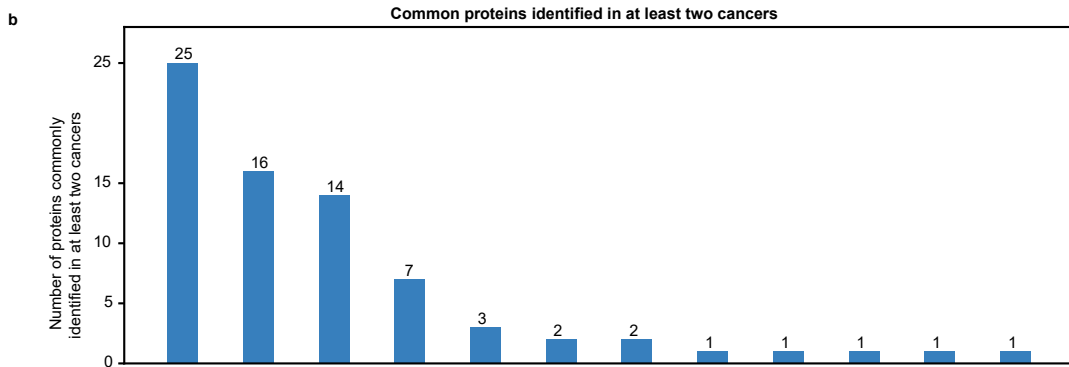
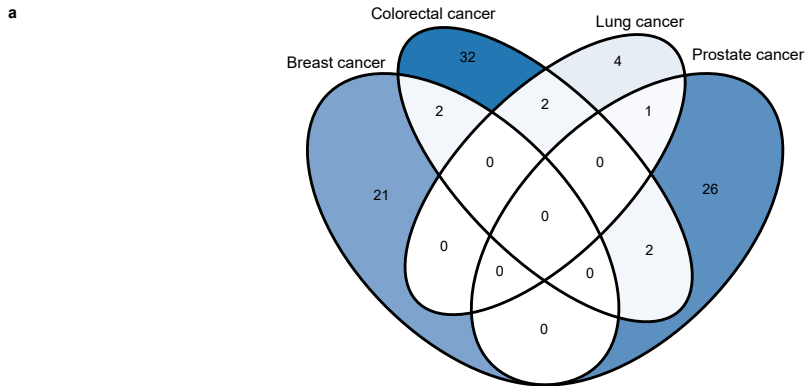
Identified 18 lead variants

Box 1
Only including lead variants with $P < 1 \times 10^{-6}$ from analysis in European populations or Trans-ancestry GWAS.

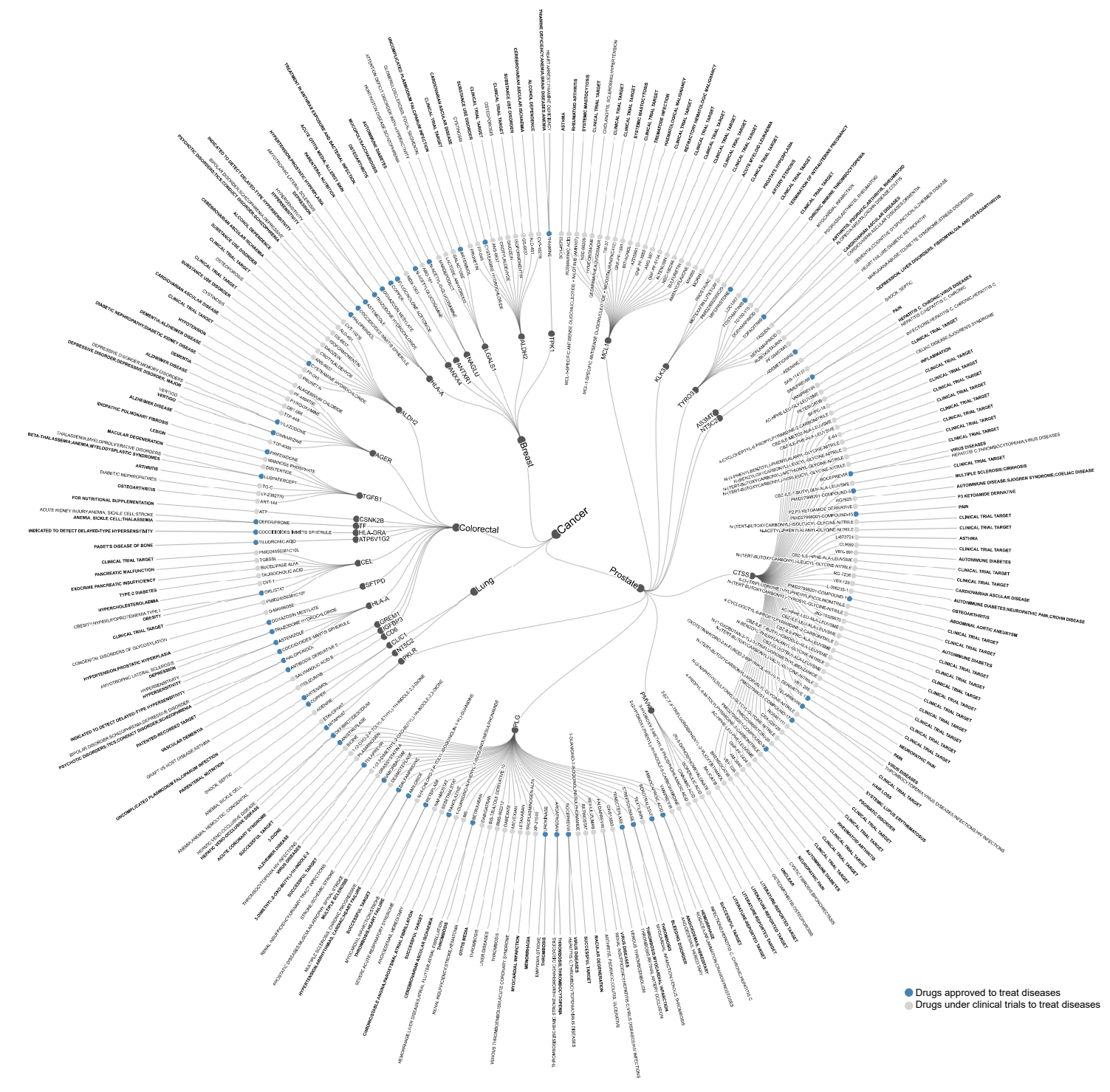
Box 2 – SuSiE analysis:
Including lead variants with independent risk-associated signals

- those with $P < 5 \times 10^{-8}$
- those with $P < 1 \times 10^{-6}$, located in GWAS loci (i.e., within a 1Mb region)

Box 3 – previous GWAS
Including additional GWAS-identified risk variants with $P < 1 \times 10^{-6}$ in European populations from the GWAS study, and LD $r^2 < 0.1$ in European populations with the above set of lead variants



Cancer



● Drugs approved to treat diseases
● Drugs under clinical trials to treat diseases