

1 Assessment of Polygenic Architecture and Risk Prediction based on Common Variants Across 2 Fourteen Cancers

3
4 Yan Dora Zhang^{1,2}, Amber N. Wilcox^{3,4}, Haoyu Zhang^{3,5}, Parichoy Pal Choudhury³, Douglas F.
5 Easton^{6,7}, Roger L. Milne^{8,9,10}, Jacques Simard¹¹, Per Hall^{12,13}, Kyriaki Michailidou^{7,14}, Joe Dennis⁷,
6 Marjanka K. Schmidt^{15,16}, Jenny Chang-Claude^{17,18}, Puya Gharahkhani¹⁹, David Whiteman²⁰,
7 Peter T. Campbell²¹, Michael Hoffmeister²², Mark Jenkins⁹, Ulrike Peters²³, Li Hsu²³, Stephen B.
8 Gruber²⁴, Graham Casey²⁵, Stephanie L. Schmit²⁶, Tracy A. O'Mara²⁷, Amanda B. Spurdle²⁷,
9 Deborah J. Thompson⁷, Ian Tomlinson^{28,29}, Immaculata De Vivo^{30,31}, Maria Teresa Landi³,
10 Matthew H. Law¹⁹, Mark M. Iles³², Florence Demenais³³, Rajiv Kumar³⁴, Stuart MacGregor¹⁹, D.
11 Timothy. Bishop³⁵, Sarah V. Ward³⁶, Melissa L. Bondy³⁷, Richard Houlston³⁸, John K. Wiencke³⁹,
12 Beatrice Melin⁴⁰, Jill Barnholtz-Sloan⁴¹, Ben Kinnersley³⁸, Margaret R. Wrensch³⁹, Christopher I.
13 Amos⁴², Rayjean J. Hung⁴³, Paul Brennan⁴⁴, James McKay⁴⁴, Neil E. Caporaso³, Sonja I. Berndt³,
14 Brenda M. Birmann³⁰, Nicola J. Camp⁴⁵, Peter Kraft⁴⁶, Nathaniel Rothman³, Susan L. Slager⁴⁷,
15 Andrew Berchuck⁴⁸, Paul DP. Pharoah^{6,7}, Thomas A. Sellers²⁶, Simon A. Gayther⁴⁹, Celeste L.
16 Pearce^{50,24}, Ellen L. Goode⁵¹, Joellen M. Schildkraut⁵², Kirsten B. Moysich⁵³, Laufey T.
17 Amundadottir⁵⁴, Eric J. Jacobs²¹, Alison P. Klein⁵⁵, Gloria M. Petersen⁵⁶, Harvey A. Risch⁵⁷, Rachel
18 Z. Stolzenberg-Solomon³, Brian M. Wolpin⁵⁸, Donghui Li⁵⁹, Rosalind A. Eeles⁶⁰, Christopher A.
19 Haiman²⁴, Zsofia Kote-Jarai⁶⁰, Fredrick R. Schumacher⁶¹, Ali Amin Al Olama^{62,63}, Mark P. Purdue³,
20 Ghislaine Scelo⁴⁴, Marlene D. Dalgaard^{64,65}, Mark H. Greene⁶⁶, Tom Grotmol⁶⁷, Peter A.
21 Kanetsky²⁶, Katherine A. McGlynn³, Katherine L. Nathanson⁶⁸, Clare Turnbull³⁸, Fredrik
22 Wiklund⁶⁹, BCAC, BEACON, CCFR, CORECT, ECAC, GECCO, GenoMEL, GICC, ILCCO, INTEGRAL,
23 InterLymph, OCAC, Oral Cancer GWAS, PanC4, PanScan, PRACTICAL, Renal Cancer GWAS, TECAC,
24 Stephen J. Chanock³, Nilanjan Chatterjee^{5,55*†}, Montserrat Garcia-Closas^{3†}

25
26 ¹Department of Statistics and Actuarial Science, Faculty of Science, The University of Hong Kong, Hong Kong SAR,
27 China, ²Centre for PanorOmic Sciences, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong Kong
28 SAR, China, ³Division of Cancer Epidemiology and Genetics, National Cancer Institute, Rockville, MD, USA,
29 ⁴Department of Epidemiology, Gillings School of Global Public Health, University of North Carolina at Chapel Hill,
30 Chapel Hill, NC, USA, ⁵Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins University,
31 Baltimore, MD, USA, ⁶Centre for Cancer Genetic Epidemiology, Department of Oncology, University of Cambridge,
32 Cambridge, UK, ⁷Centre for Cancer Genetic Epidemiology, Department of Public Health and Primary Care,
33 University of Cambridge, Cambridge, UK, ⁸Cancer Epidemiology Division, Cancer Council Victoria, Melbourne,
34 Victoria, Australia, ⁹Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health,
35 The University of Melbourne, Melbourne, Victoria, Australia, ¹⁰Precision Medicine, School of Clinical Sciences at
36 Monash Health, Monash University, Clayton, Victoria, Australia, ¹¹Centre Hospitalier Universitaire de Québec–
37 Université Laval Research Center, Québec City, QC, Canada, ¹²Department of Medical Epidemiology and
38 Biostatistics, Karolinska Institutet, Stockholm, Sweden, ¹³Department of Oncology, Södersjukhuset, Stockholm,
39 Sweden, ¹⁴Department of Electron Microscopy/Molecular Pathology and The Cyprus School of Molecular Medicine,
40 The Cyprus Institute of Neurology & Genetics, Nicosia, Cyprus, ¹⁵Division of Molecular Pathology, The Netherlands
41 Cancer Institute - Antoni van Leeuwenhoek Hospital, Amsterdam, The Netherlands, ¹⁶Division of Psychosocial
42 Research and Epidemiology, The Netherlands Cancer Institute - Antoni van Leeuwenhoek hospital, Amsterdam,
43 The Netherlands, ¹⁷Division of Cancer Epidemiology, German Cancer Research Center (DKFZ), Heidelberg, Germany,
44 ¹⁸Cancer Epidemiology Group, University Cancer Center Hamburg (UCCH), University Medical Center Hamburg-
45 Eppendorf, Hamburg, Germany, ¹⁹Statistical Genetics, QIMR Berghofer Medical Research Institute, Brisbane,
46 Australia, ²⁰Cancer Control, QIMR Berghofer Medical Research Institute, Brisbane, Australia, ²¹Behavioral and
47 Epidemiology Research Group, American Cancer Society, Atlanta, GA, USA, ²²Division of Clinical Epidemiology and
48 Aging Research, German Cancer Research Center (DKFZ), Heidelberg, Germany, ²³Public Health Sciences Division,

49 Fred Hutchinson Cancer Research Center, Seattle, WA, USA, ²⁴Department of Preventive Medicine, USC Norris
50 Comprehensive Cancer Center, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA,
51 ²⁵Department of Public Health Sciences, Center for Public Health Genomics, University of Virginia, Charlottesville,
52 VA, USA, ²⁶Department of Cancer Epidemiology, H. Lee Moffitt Cancer Center and Research Institution, Tampa, FL,
53 USA, ²⁷Genetics and Computational Biology Division, QIMR Berghofer Medical Research Institute, Brisbane,
54 Australia, ²⁸Institute of Cancer and Genomic Sciences, University of Birmingham, Birmingham, UK, ²⁹Wellcome
55 Trust Centre for Human Genetics and Oxford NIHR Biomedical Research Centre, University of Oxford, Oxford, UK,
56 ³⁰Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital and Harvard
57 Medical School, Boston, MA, USA, ³¹Department of Epidemiology, Harvard T.H. Chan School of Public Health,
58 Boston, MA, USA, ³²Section of Epidemiology and Biostatistics, Leeds Institute of Cancer and Pathology, University
59 of Leeds, Leeds, UK, ³³Université de Paris, UMRS-1124, Institut National de la Santé et de la Recherche Médicale
60 (INSERM), F-75006 Paris, France, ³⁴Division of Molecular Genetic Epidemiology, German Cancer Research Center
61 (DKFZ), Heidelberg, Germany, ³⁵Division of Haematology and Immunology, Leeds Institute of Medical Research,
62 University of Leeds, Leeds, UK, ³⁶Centre for Genetic Origins of Health and Disease, School of Biomedical Sciences,
63 The University of Western Australia, Perth, Australia, ³⁷Department of Medicine, Section of Epidemiology and
64 Population Sciences, Baylor College of Medicine, Houston, TX, USA, ³⁸Division of Genetics and Epidemiology, The
65 Institute of Cancer Research, London, UK, ³⁹Department of Neurological Surgery, School of Medicine, University of
66 California, San Francisco, San Francisco, CA, USA, ⁴⁰Department of Radiation Sciences Oncology, Umeå University,
67 Umeå, Sweden, ⁴¹Case Comprehensive Cancer Center, Case Western Reserve University School of Medicine,
68 Cleveland, OH, USA, ⁴²Institute for Clinical and Translational Research, Dan L. Duncan Comprehensive Cancer
69 Center, Baylor College of Medicine, Houston, TX, USA, ⁴³Lunenfeld-Tanenbaum Research Institute, Sinai Health
70 System, Toronto, ON, Canada, ⁴⁴International Agency for Research on Cancer, World Health Organization, Lyon,
71 France, ⁴⁵Division of Hematology and Hematological Malignancies, University of Utah School of Medicine, Salt Lake
72 City, UT, USA, ⁴⁶Program in Genetic Epidemiology and Statistical Genetics, Harvard T.H. Chan School of Public
73 Health, Boston, MA, USA, ⁴⁷Department of Health Sciences Research, Division of Biomedical Statistics & Informatics,
74 Mayo Clinic, Rochester, MN, USA, ⁴⁸Department of Gynecologic Oncology, Duke University Medical Center,
75 Durham, NC, USA, ⁴⁹Center for Bioinformatics and Functional Genomics and the Cedars Sinai Genomics Core,
76 Cedars-Sinai Medical Center, Los Angeles, CA, USA, ⁵⁰Department of Epidemiology, University of Michigan School
77 of Public Health, Ann Arbor, MI, USA, ⁵¹Department of Health Science Research, Division of Epidemiology, Mayo
78 Clinic, Rochester, MN, USA, ⁵²Rollins School of Public Health, Emory University, Atlanta, GA, USA, ⁵³Division of
79 Cancer Prevention and Control, Roswell Park Cancer Institute, Buffalo, NY, USA, ⁵⁴Laboratory of Translational
80 Genomics, Division of Cancer Epidemiology and Genetics, National Cancer Institute, National Institutes of Health,
81 Bethesda, MD, USA, ⁵⁵Department of Oncology, Sidney Kimmel Comprehensive Cancer Center, Johns Hopkins
82 School of Medicine, Baltimore, MD, USA, ⁵⁶Department of Health Sciences Research, Division of Epidemiology,
83 Mayo Clinic, Rochester, MN, USA, ⁵⁷Chronic Disease Epidemiology, Yale School of Medicine, New Haven, CT, USA,
84 ⁵⁸Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA, ⁵⁹GI Medical Oncology
85 Department, Division of Cancer Medicine, The University of Texas MD Anderson Cancer Center, Houston, TX, USA,
86 ⁶⁰Division of Genetics and Epidemiology, The Institute of Cancer Research, Sutton, Surrey, UK, ⁶¹Department of
87 Population and Quantitative Health Sciences, Case Western Reserve University School of Medicine, Cleveland, OH,
88 USA, ⁶²Centre for Cancer Genetic Epidemiology, Department of Public Health and Primary Care, Strangeways
89 Research Laboratory, University of Cambridge, Cambridge, UK, ⁶³Department of Clinical Neurosciences, University
90 of Cambridge, Cambridge, UK, ⁶⁴Department of Growth and Reproduction, Copenhagen University Hospital
91 (Rigshospitalet), Copenhagen, Denmark, ⁶⁵Department of Health Technology, Technical University of Denmark,
92 Lyngby, Denmark, ⁶⁶Clinical Genetics Branch, Division of Cancer Genetics and Epidemiology, National Cancer
93 Institute, Rockville, MD, USA, ⁶⁷Cancer Registry of Norway, Oslo, Norway, ⁶⁸Department of Medicine, Division of
94 Translational Health and Human Genetics, University of Pennsylvania, Philadelphia, PA, USA, ⁶⁹Department of
95 Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden

96

97 † These authors jointly supervised this work

98

99 *Corresponding author:

100 Nilanjan Chatterjee (nchatte2@jhu.edu)

101 **Abstract**

102 We analyzed summary-level data from genome-wide association studies (GWAS) of European
103 ancestry across fourteen cancer sites to estimate the number of common susceptibility variants
104 (polygenicity), as well as the distribution of their associated effect sizes. All cancers evaluated
105 showed high degree of polygenicity, involving at a minimum of thousands of independent
106 susceptibility variants. We project that sample sizes, required to explain 80% of GWAS
107 heritability varies from 60,000 cases for testicular to over 1,000,000 cases for lung cancer. The
108 maximum relative risk achievable for subjects at the 99th risk percentile of underlying polygenic
109 risk scores, compared to average risk, ranges from 12 for testicular to 2.5 for ovarian cancer.
110 We show that polygenic risk scores have potential for risk stratification for relatively common
111 cancers, such as breast, colon and prostate, but less so for others because of modest
112 heritability and lower disease incidence.

113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131

132 **Introduction**

133 Genome-wide association studies (GWAS) have led to the identification of hundreds of
134 independent cancer susceptibility loci containing common, low-risk variants^{1,2}. The number of
135 discoveries varies widely across cancers, largely driven by available sample size, which reflects,
136 in part, disease incidence in the general population. However, specific cancers, e.g., chronic
137 lymphoid leukemia (CLL)³ and testicular cancer⁴, are notable for unexpectedly high numbers of
138 genome-wide significant discoveries from GWAS of relatively small sample size. Previous
139 studies have also reported that these two cancers have high heritability⁵. Across cancer types,
140 polygenic risk scores (PRS) show varying levels of risk stratification depending on the heritability
141 explained by the identified variants and the disease incidence rates in the population⁶⁻¹². Their
142 potential clinical utility would depend not only on the level of risk stratification, but also on
143 other factors such as the availability of appropriate risk-reducing interventions for those
144 identified as at high risk.

145
146 Estimation of heritability due to additive effects of all single nucleotide polymorphisms (SNPs)
147 included in GWAS arrays¹³, referred to as GWAS heritability in this article, have shown that
148 common variants have substantial potential to identify individuals at different levels of risk for
149 many cancer types¹⁴. It remains, however, unclear how large the sample sizes of GWAS need to
150 be to reap the full potential of PRS-based risk prediction. Herein, we apply our recently
151 published method¹⁵ to estimate the degree of polygenicity and the effect-size distribution
152 associated with common variants (MAF>0.05) across fourteen different cancer types, based on
153 summary-level association statistics from available GWAS¹⁶⁻²⁸ from populations of European

154 ancestry (**Supplementary Table 1**). From these inferred parameters, we then provide
155 projections of the expected number of common variants to be discovered and predictive
156 performance of associated PRS as a function of increasing sample size for future GWAS. Finally,
157 by incorporating age-specific incidence²⁹ from population-based cancer registries, we explore
158 the magnitude of absolute risk stratification potentially achievable by PRS.

159

160 **Results**

161 **Cancer Polygenicity**

162 We found that cancers are highly polygenic, like other complex traits^{15,30,31}. Estimates of the
163 number of susceptibility variants with independent risk associations vary from ~1,000 to 7,500
164 between the fourteen cancer sites (**Table 1**). For comparability, effect-size distributions are
165 shown in groups of similarly-sized GWAS with similar power for detecting associations (**Figure**
166 **1**). For GWAS with <10,000 cancer cases (group 1), CLL and testicular cancer are each
167 associated with 2,000-2,500 variants and characterized by a much larger proportion of variants
168 with larger estimated effect sizes than for the other group 1 cancers, as reflected by wider
169 effect size distribution with heavier tails (**Figure 1, Table 1**). GWAS heritability estimates
170 indicate that, in aggregate, common variants explain a high degree of variation of risk for these
171 two cancers. In contrast, in group 1, esophageal and oropharyngeal cancers are associated with
172 a larger proportion of variants with substantially smaller effect sizes, compared with CLL and
173 testicular cancers in group 1.

174

175 For GWAS with 10,000-25,000 cases (group 2), melanoma is noteworthy because it is associated
176 with a wider effect size distribution than other group 2 cancers. The estimated number of
177 susceptibility variants in this group ranges from 1,000 to 2,000. GWAS heritability estimates
178 indicate that aggregated common variants make a relatively small contribution to ovarian and
179 endometrial cancer susceptibility. Finally, for the three GWAS with >25,000 cases each (group
180 3), prostate cancer is remarkable for having more variants with large effect sizes, namely, the
181 underlying effect-size distribution has a heavier tail, compared with cancers of the breast and
182 lung (**Figure 1**). In this group, all three cancer types tend to have large numbers of associated
183 variants (>4,500) compared with cancer sites in other groups, but this pattern could partially be
184 due to the very large sample sizes of group 3 GWAS¹⁵.

185

186 For a large majority of the fourteen cancer sites, a two-component normal-mixture model for
187 non-null effects provides a substantially better fit to observed summary-statistics than a single-
188 normal distribution; this indicates the presence of a fraction of variants with distinctly larger
189 effect sizes than the remaining (**Supplementary Figures 1-2**). In contrast, a single normal
190 distribution appears to be adequate for esophageal and oropharyngeal cancer, indicating the
191 presence of a large number of variants with a continuum of small effects, similar to our
192 previous findings for traits related to mental health and abilities¹⁵. Across all fourteen cancers,
193 the predicted number of discoveries and their associated genetic variance explained for current
194 GWAS sample sizes match well to those observed empirically (**Supplementary Table 2**),
195 indicating good fit of our model to the observed data.

196

197 **Future GWAS Projections**

198 GWAS heritability estimates indicate that the potential of PRS for risk discrimination in the
199 population varies widely among cancer types (**Table 1**). The area under the curve (AUC)
200 statistics associated with the best achievable PRS varies from 64% (endometrial and ovarian
201 cancer) to 88% (testicular cancer), and in the range of 70 to 80% for most cancers. The
202 percentage of GWAS heritability explained by known variants varies widely, depending on study
203 sample size and the underlying trait genetic architecture (**Figure 2**). Known variants explain
204 more than a quarter of heritability for cancer sites based on very large sample sizes (e.g., breast
205 and prostate cancer) or for cancer sites that have susceptibility variants with relatively large
206 effect sizes (e.g., CLL, melanoma and testicular cancer). Oropharyngeal cancer, in contrast, has
207 both a small sample size and small effect sizes; its percentage heritability currently explained is
208 almost zero.

209
210 The sample size needed to identify common variants that could explain approximately 80% of
211 the total GWAS heritability for the cancers evaluated is generally very large, requiring 200,000
212 to 1,000,000 cancer cases, with a comparable number of controls (**Figure 2**). However, for three
213 sites, namely, testicular cancer, CLL, and melanoma, the required sample size is smaller, 60,000,
214 80,000 and 110,000 cases, respectively, due to the large effect sizes of their associated variants.
215 By quadrupling the sample sizes of currently published GWAS, the percentage of GWAS
216 heritability explained would rise to more than 40% across all cancers, except for oropharyngeal
217 cancer. Such sample size increases would also lead to appreciable improvements in PRS
218 discriminatory power across all these sites (**Figures 3-4**). For cancers that were found to be the

219 most polygenic and that had small effect sizes (e.g., cancers of breast, lung and oropharynx),
220 improvement would occur at a slower rates as sample sizes increase, and these sites would
221 require the largest sample sizes to generate PRSs with discriminatory power close to theoretical
222 limits. Of note, for a number of cancers, the achievable relative risks for subjects at the 99th
223 percentile of PRS distribution compared with those at average risk, are comparable to those for
224 monogenic disorders³² (e.g., relative-risk more than 3-4 fold) (**Figure 4**). Across all fourteen
225 cancer types, inclusion of SNPs using more liberal but optimized *p*-value thresholds (**see**
226 **Methods**) would improve performance of PRS-based risk prediction versus using the stringent
227 genome-wide significance level, but the anticipated gains would be generally modest
228 (**Supplementary Figures 3-4**).

229
230 Projections of residual lifetime cancer risks for the US non-Hispanic white population show that
231 the discriminatory power of PRS built from current or foreseeable studies will depend heavily
232 on the underlying cancer incidence in the population (**Figure 5, Supplementary Figures 5-7**).

233 The potential clinical utility of PRS depends on the degree of risk stratification and specific
234 prevention or early detection strategies for a given cancer, should they exist. For common
235 cancers, such as breast, colorectal and prostate, a PRS with even modest discriminatory power
236 (maximum AUC of approximately 70%, **Figure 3**) can provide substantial stratification of
237 absolute risk in the population. In contrast, for CLL and testicular cancer, even though its PRS
238 could achieve a higher AUC (e.g. in the range 80-90%, **Figure 3**), the degree of absolute risk
239 stratification will be modest because of the infrequency of these cancers. Thus, a PRS by itself
240 has the least impact on risk stratification for cancer sites that are infrequent or/and that have

241 low heritability. However, it is possible that PRS could have clinical utility for some of these
242 cancers in the presence or in combination with other risk factors and biomarkers. For example,
243 a PRS for lung cancer may provide larger stratification for absolute risk among smokers than
244 never smokers because of the higher baseline risk in smokers.

245

246 **Discussion**

247 Our study is subject to several limitations. We may have underestimated the number of
248 underlying common susceptibility loci, especially for those cancers for which current GWAS
249 have small sample sizes¹⁵. Thus, the interpretation of comparisons of the underlying genetic
250 architecture across cancer types with very different sample sizes requires caution. Nevertheless,
251 the major patterns are unlikely to be due to differences in sample size. For example, we
252 estimated oropharyngeal and esophageal cancers to be two of the most polygenic sites, though
253 the GWAS sample sizes for these two sites were relatively small. Further, Q-Q plots of observed
254 and expected p -values indicate that the inferred models for effect-size distributions explain
255 observed GWAS summary-statistics well, regardless of GWAS sample size. Another important
256 limitation is that we only included data from subjects of European ancestry, since GWAS data
257 for other ancestries are currently too small to permit reliable projections for most cancer sites.
258 In addition, several cancers (e.g., lung, ovary, glioma, and breast) consist of etiologically
259 heterogeneous subtypes that were not considered in our analyses due to lack of adequate
260 sample sizes for appropriate subtypes for most of these cancer sites. Further studies of
261 ancestry- and subtype-specific genetic architectures are needed to address these limitations.

262

263 In our projections, we assume standard agnostic association analysis of SNPs without
264 incorporating any external information on population genetics or functional characteristics of
265 SNPs. It is, however, possible to incorporate various types of external information to improve
266 power for discovery of associations³³⁻³⁶ and genetic risk prediction³⁷. We have evaluated the
267 merit of future GWAS only in terms of their ability to explain heritability and improve risk
268 prediction. However, current and future discoveries have other major implications, including
269 provident insights to biological pathways and mechanisms, potential gene-environment
270 interactions and understanding causal relationships through Mendelian Randomization
271 analyses³⁸. A number of these cancers are known to have rare high-penetrant risk variants, but
272 for this study we have focused on estimating effect-size distribution associated with common
273 variants. Furthermore, heritability analysis indicate that uncommon and rare variants could
274 explain a substantial fraction of the variation of complex traits³⁹, and thus, it is likely that there
275 are many unknown uncommon and rare variants associated with these cancers as well. In the
276 future, characterization of heritability and effect-size distribution associated with the full
277 spectrum of allele frequencies will require individual level sequencing data on a substantially
278 larger number of cases and controls.

279

280 The observed differences in the underlying genetic architecture of susceptibility across cancers
281 could be due to various factors, including the effect of negative selection^{30,40}, tissue-specific
282 genetic regulation of gene-expression⁴¹, cell of origin⁴², the number of biological steps needed
283 to transition from normal to malignant tissue⁴³, mediation of genetic effects by underlying
284 environmental exposures⁴⁴, and the presence of heterogeneous cancer-specific

285 subtypes^{21,25,27,28}. A number of cancer types, including those of lung, oropharynx and esophagus,
286 which were associated with large numbers of SNPs with small average effect sizes, have known
287 strong environmental risk factors and distinct etiologic subtypes. It is also noteworthy that
288 testicular cancer also stands out for a large number of discoveries in cross-tissue expression
289 quantitative trait loci analyses, likely indicating a stronger association of SNPs on gene
290 expression levels for this tissue compared to others⁴¹.

291

292 In conclusion, our comprehensive analysis of fourteen cancer sites in adults of European
293 ancestry reveals that while all sites have polygenic influences, there is substantial diversity
294 observed in their underlying genetic architectures, which reflects important biology and also
295 influences the utility of polygenic risk prediction for individual cancers. Our projections for
296 future yields of GWAS across these cancers provide a roadmap for important returns from
297 future investment in research, including the potential clinical utility of polygenic risk prediction
298 for stratification of absolute risks in the population.

299

300 **Methods**

301

302 **Description of GWAS studies.** We analyzed summary data from GWAS studies across fourteen
303 cancer types. For select cancer sites^{26,28}, we downloaded publicly available genome-wide
304 summary-level statistics from the latest consortium-based analyses. For others, we obtained
305 access to data through collaborative efforts with individual consortia. Details about individual
306 studies, including the number of cases and controls, are provided in **Supplementary Table 1**.

307

308 **LD Reference Panel Selection.** We consider a reference panel with ~1.07 million SNPs included
309 in the HapMap3 and which had minor allele frequency > 0.05 in the 1000 Genome European
310 Ancestry sample. Based on known LD among common variants, we expect these set of variants
311 to provide high coverage for all common variants for European ancestry population and thus
312 loss of information due to imperfect tagging of causal variants to be fairly minimal.

313

314 **Quality control for summary GWAS data.** Across all cancers, we applied several filtering steps
315 analogous to those used earlier for estimation of heritability^{45,46} and effect-size distribution
316 using summary-level data¹⁵. First, we restricted analysis to SNPs within a set of reference ~1.07
317 million SNPs included in the HapMap3 and which had minor allele frequency > 0.05 in the 1000
318 Genome European Ancestry sample. Second, we excluded SNPs having substantial amounts of
319 missing genotype data: sample sizes less than 0.67 times the 90th percentile of the distribution
320 of sample sizes across all SNPs. Third, we excluded SNPs within the major histocompatibility
321 complex (MHC) region (i.e., SNPs between 26,000,000 and 34,000,000 base pairs on
322 chromosome six) which is known to have very complex allelic architecture and can have
323 uncharacteristically large effects on some traits. Fourth, we removed regions that have SNPs
324 with extremely large effect sizes to reduce possible undue influence of them on estimation of
325 parameters associated with overall effect-size distributions. Using PLINK --clump, we identify all
326 top SNPs which have associated chi-square statistics greater than 80 (i.e., odds-ratio (in
327 standardized scale) >2.19) and removed all SNPs which were within 1MB distance of or had an
328 estimated squared LD larger than 0.1 with those top SNPs. We added back the contribution of

329 these top independent SNPs in the final reporting of the total number of susceptibility SNPs,
330 estimates of total heritability, and various projections we made as a function of sample size of
331 the GWAS.

332

333 **Statistical model.** We inferred common variant genetic architecture of the different cancers
334 using GENESIS¹⁵, a method we recently developed to characterize underlying effect-size
335 distributions in terms of the total number of susceptibility SNPs (polygenicity) and a normal
336 mixture model for the distribution of their effects. Specifically, it is assumed that standardized
337 effects of common SNPs in an underlying logistic regression model on the risk of a cancer can
338 be specified in the mixture distribution in the form $\beta_m \sim (1 - \pi_c)\delta_0 + \pi_c N(0, \sigma^2)$ (two-
339 component model), or $\beta_m \sim (1 - \pi_c)\delta_0 + \pi_c [p_1 N(0, \sigma_1^2) + p_2 N(0, \sigma_2^2)]$ (three-component
340 model) where δ_0 is the Dirac delta function indicating that a fraction, $1 - \pi_c$, of the SNPs have
341 null effects, and remaining π_c fraction of SNPs have non-null effects. Under the three-
342 component model, $p_2 = 1 - p_1$ denotes the proportion of SNPs allocated to mixture
343 component with larger variance component (assuming $\sigma_2^2 > \sigma_1^2$) models. Under these models,
344 $M\pi_c$ characterizes the degree of polygenicity, i.e., the number of susceptibility SNPs with
345 independent effects on disease risk. Under both models, we defined “GWAS heritability” of a
346 disease as $h^2 = M\pi_c E(\beta^2)$, where $E(\beta^2)$ denotes the average variance size of the non-null
347 SNPs. We observed that under the above model, h^2 is also the population variance of the
348 underlying “true” polygenic risk score, defined as $PRS = \sum_{m=1}^M \beta_m G_m$, where G_m denotes the
349 standardized genotype associated with the m -th SNP. Under the two-component model, which
350 assumes a single normal distribution for the effect of all susceptibility SNPs, $E(\beta^2) = \sigma^2$. Under

351 the three-component model, which allows mixture of two-normal distributions with distinct
352 variance components and thus can better accommodate the presence of a group of
353 susceptibility SNPs with much larger effects than others, we have $p_1\sigma_1^2 + p_2\sigma_2^2$. Under the
354 three-component model, we use the fraction $v = p_1\sigma_1^2 / (p_1\sigma_1^2 + p_2\sigma_2^2)$ to characterize the
355 proportion of heritability explained by SNPs associated with the larger variance component
356 parameter. As we removed SNPs with extremely large effects ($\chi_i^2 > 80$) and the associated
357 regions from the analysis, in reporting the final heritability estimates, we added back the
358 contribution of the independent top SNPs from these excluded regions as $\sum_i (\hat{\beta}_i^2 - \tau_i^2)$ where $\hat{\beta}$
359 is the estimate of log-odds-ratio (in standardized scale) and τ_i is the corresponding standard
360 error for the i -th SNP.

361

362 **Genetic variance projection.** Given the estimated effect-size distribution, we calculated
363 expected discoveries and genetic variance explained using

364 $ED = M\hat{\pi}_c \int_{\beta} \text{pow}_{\alpha,n}(\beta) \sum_{h=1}^H \hat{p}_h N(0, \hat{\sigma}_h^2) d\beta$ and

365 $EV = M\hat{\pi}_c \int_{\beta} \beta^2 \text{pow}_{\alpha,n}(\beta) \sum_{h=1}^H \hat{p}_h N(0, \hat{\sigma}_h^2) d\beta$, respectively, at $\alpha = 5 \times 10^{-8}$ for a GWAS of

366 sample size n , where $\text{pow}_{\alpha,n}(\beta) = 1 - \Phi\left(\frac{c_{\alpha}}{2} - \sqrt{n}\beta\right) + \Phi\left(-\frac{c_{\alpha}}{2} - \sqrt{n}\beta\right)$ with $\Phi(\cdot)$ the

367 standard normal cumulative density function and $c_{\alpha} = \Phi^{-1}(1 - \alpha)$ the α -th quantile for the

368 standard normal distribution. Similar to heritability calculations, we added back the

369 contributions of **independent** top SNPs with very large effects to the number of expected

370 discoveries and associated variances explained by the quantities $\sum_i \text{pow}_{\alpha,n}(\hat{\beta}_i)$ and

371 $h^{-2} \sum_i (\hat{\beta}_i^2 - \tau_i^2) \text{pow}_{\alpha,n}(\hat{\beta}_i)$. We observed that for projections involving sample sizes bigger
372 than the current study $\text{pow}_{\alpha,n}(\hat{\beta}_i)$ for the large effect SNPs will all be very close to 1.0.

373

374 **Projection for AUC and relative risk at top 1%.** As we quantify heritability in terms of the
375 variability of the underlying “true” polygenic risk-score, we used the formula^{12,47,48}, $\text{AUC} =$

376 $\Phi\left(\sqrt{\frac{h^2}{2}}\right)$ to characterize the best discriminatory power achievable in limiting using common

377 variant PRS. We used the same formula to calculate the AUC associated with PRSs that could

378 be built using SNPs either reaching genome-wide significance (p -value $< 5 \times 10^{-8}$) or a weaker

379 but optimized threshold, for a GWAS of given sample size based on the projected variance of

380 the respective PRS. **Given sample size of GWAS and an effect-size distribution for the underlying**

381 **cancer, an optimal threshold for SNP selection that will maximize the expected predictive**

382 **performance of PRS is calculated using analytic formula we have derive earlier⁴⁸.** The relative

383 risk for those estimated to be at the 99th percentile or higher of the distribution of a PRS

384 (compared to the average risk of the population) was calculated using the formula¹² $\exp\left(-\frac{h^2}{2} +$

385 $\Phi^{-1}(0.99)\sqrt{h^2}\right)$ where h^2 is the population variance of the PRS.

386

387 **Absolute risk projection.** For each cancer site, we projected the distribution of residual lifetime

388 risk (up to age 80 years) for Non-Hispanic White individuals in the general US population

389 according to PRSs which could be built from GWAS of different sample sizes. For any given age,

390 we first obtain the distribution of residual lifetime risks based on a model for absolute risks

391 developed using the iCARE tool that we have described earlier^{12,29}. The iCARE tool uses

392 projected standard deviations of PRS at different GWAS sample sizes and age-specific cancer
393 incidence rates available from the US National Cancer Institute-Surveillance, Epidemiology, and
394 End Results Program (NCI-SEER) (2015) to obtain absolute risk distributions. In deriving absolute
395 risks, we adjusted for competing risk of mortality due to other causes using the age-specific
396 mortality rates from the Center for Disease Control (CDC) WONDER database (2016). We then
397 weighted the projected residual lifetime risk distribution at different baseline ages (in five-year
398 categories) based on the US population distribution of ages within 30 to 75 years, as observed
399 in the estimated 2016 US Census. For cancers of the reproductive system, weights were based
400 on the age distributions among males or females, as appropriate.

401

402 **Data Availability**

403 The data that support the findings of this study are available by application from the
404 participating consortia: BCAC (bcac@medschl.cam.ac.uk), BEACON (P Gharahkhani), ColonCFR
405 (M Jenkins), GECCO/CORECT (U Peters), ECAC (TA O'Mara), GenoMEL (M Iles), GICC (R
406 Houlston), ILLCO/INTEGRAL (C Amos), InterLymph (S Berndt), OCAC (PDP Pharoah), Oral Cancer
407 GWAS (P Brennan), PanC4/PanScan (LT Amundadottir), PRACTICAL (Data Access
408 Committee/<http://practical.icr.ac.uk/>), Renal Cancer GWAS (MP Purdue, P Brennan), TECAC (KA
409 McGlynn). For breast and prostate cancers, summary GWAS data can also be downloaded from
410 [http://bcac.ccge.medschl.cam.ac.uk/bcacdata/oncoarray/gwas-icogs-and-oncoarray-summary-](http://bcac.ccge.medschl.cam.ac.uk/bcacdata/oncoarray/gwas-icogs-and-oncoarray-summary-results/)
411 [results/](http://bcac.ccge.medschl.cam.ac.uk/bcacdata/oncoarray/gwas-icogs-and-oncoarray-summary-results/) and http://practical.icr.ac.uk/blog/?page_id=8164.

412

413 **Code Availability**

414 The code for running the analysis in the paper is freely available from the CancerEffectSize
415 GitHub repository (<https://github.com/yandorazhang/CancerEffectSize>).

416

417 **Acknowledgement**

418 The research was supported by a RO1 grant from NHGRI (1 R01 HG010480-01) and the
419 intramural program of the National Cancer Institute.

420

421 **Author Contributions**

422 N.C. and M.G.C. conceived the project. Y.Z. and A.W. performed main analyses. Y.Z., N.C. and
423 M.G.C. wrote the first draft of the manuscript. BCAC, BEACON, CCFR, CORECT, ECAC, GECCO,
424 GenoMEL, GICC, ILCCO, Integral, InterLymph, OCAC, Oral Cancer GWAS, PANC4, PanScan,
425 PRACTICAL, Renal Cancer GWAS, and TECAC contributed data. P.P.C., R.L.M., M.K.S., M.J., U.P.,
426 L.H., S.L.S., T.A.O., A.B.S., D.J.T., M.H.L., M.M.I., F.D., S.M., S.V.W., M.R.W., C.I.A., S.I.B., B.M.B.,
427 N.J.C., P.D.P.P., T.A.S., L.T.A., E.J.J., H.A.R., R.Z.S.S., M.P.P., M.H.G., K.A.M., and S.J.C.
428 commented on earlier drafts of manuscript. H.Z., D.F.E., J.S., P.H., K.M., J.D., J.C.C., P.G., D.W.,
429 P.T.C., M.H., S.B.G., G.C., I.T., I.D.V., M.T.L., R.K., D.T.B., M.L.B., R.H., J.K.W., B.M., J.B.S., B.K.,
430 R.J.H., P.B., J.M., N.E.C., P.K., N.R., S.L.S., A.B., S.A.G., C.L.P., E.L.G., J.M.S., K.B.M., A.P.K., G.M.P.,
431 B.M.W., D.L., R.A.E., C.A.H., Z.K.J., F.R.S., A.A.A.O., G.S., M.D.D., T.G., P.A.K., K.L.N., C.T., and F.W.
432 reviewed the manuscript. All authors reviewed and approved drafts of manuscripts.

433

434 **Competing Interests**

435 We declare that none of the authors have competing financial or non-financial interests as
436 defined by Nature Research.

437

438 **Consortia**

439 **Breast Cancer Association Consortium (BCAC)**

440 Douglas F. Easton^{6,7}, Roger L. Milne^{8,9,10}, Jacques Simard¹¹, Per Hall^{12,13}, Kyriaki Michailidou^{7,14},
441 Joe Dennis⁷, Marjanka K. Schmidt^{15,16}, Jenny Chang-Claude^{17,18}

442

443 **Barrett's and Esophageal Adenocarcinoma Consortium (BEACON)**

444 Puya Gharahkhani¹⁹, David Whiteman²⁰

445

446 **Colon Cancer Family Registry (ColonCFR), Transdisciplinary Studies of Genetic Variation in**

447 **Colorectal Cancer (CORECT), Genetics and Epidemiology of Colorectal Cancer Consortium**

448 **(GECCO)**

449 Peter T. Campbell²¹, Michael Hoffmeister²², Mark Jenkins⁹, Ulrike Peters²³, Li Hsu²³, Stephen B.
450 Gruber²⁴, Graham Casey²⁵, Stephanie L. Schmit²⁶

451

452 **Endometrial Cancer Association Consortium (ECAC)**

453 Tracy A. O'Mara²⁷, Amanda B. Spurdle²⁷, Deborah J. Thompson⁷, Ian Tomlinson^{28,29}, Immaculata
454 De Vivo^{30,31}

455

456 **Melanoma Genetics Consortium (GenoMEL)**

457 Maria Teresa Landi³, Matthew H. Law¹⁹, Mark M. Iles³², Florence Demenais³³, Rajiv Kumar³⁴,
458 Stuart MacGregor¹⁹, D. Timothy. Bishop³⁵, Sarah V. Ward³⁶

459

460 **Glioma International Case-Control Study (GICC)**

461 Melissa L. Bondy³⁷, Richard Houlston³⁸, John K. Wiencke³⁹, Beatrice Melin⁴⁰, Jill Barnholtz-
462 Sloan⁴¹, Ben Kinnersley³⁸, Margaret R. Wrensch³⁹

463

464 **International Lung Cancer Consortium (ILCCO), Integrative Analysis of Lung Cancer Etiology
465 and Risk (INTEGRAL) Consortium**

466 Christopher I. Amos⁴², Rayjean J. Hung⁴³, Paul Brennan⁴⁴, James McKay⁴⁴, Neil E. Caporaso³

467

468 **International Consortium of Investigators Working on Non-Hodgkin's Lymphoma
469 Epidemiologic Studies (InterLymph)**

470 Sonja I. Berndt³, Brenda M. Birmann³⁰, Nicola J. Camp⁴⁵, Peter Kraft⁴⁶, Nathaniel Rothman³,
471 Susan L. Slager⁴⁷

472

473 **Ovarian Cancer Association Consortium (OCAC)**

474 Andrew Berchuck⁴⁸, Paul DP. Pharoah^{6,7}, Thomas A. Sellers²⁶, Simon A. Gayther⁴⁹, Celeste L.
475 Pearce^{50,24}, Ellen L. Goode⁵¹, Joellen M. Schildkraut⁵², Kirsten B. Moysich⁵³

476

477 **Pancreatic Cancer Cohort Consortium (PanScan), Pancreatic Cancer Case-Control Consortium
478 (PanC4)**

479 Laufey T. Amundadottir⁵⁴, Eric J. Jacobs²¹, Alison P. Klein⁵⁵, Gloria M. Petersen⁵⁶, Harvey A.
480 Risch⁵⁷, Rachel Z. Stolzenberg-Solomon³, Brian M. Wolpin⁵⁸, Donghui Li⁵⁹

481

482 **Prostate Cancer Association Group to Investigate Cancer Associated Alterations in the**
483 **Genome (PRACTICAL)**

484 Rosalind A. Eeles⁶⁰, Christopher A. Haiman²⁴, Zsofia Kote-Jarai⁶⁰, Fredrick R. Schumacher⁶¹, Ali
485 Amin Al Olama^{62,63}

486

487 **Renal Cancer GWAS**

488 Mark P. Purdue³, Ghislaine Scelo⁴⁴

489

490 **Testicular Cancer Consortium (TECAC)**

491 Marlene D. Dalgaard^{64,65}, Mark H. Greene⁶⁶, Tom Grotmol⁶⁷, Peter A. Kanetsky²⁶, Katherine A.
492 McGlynn³, Katherine L. Nathanson⁶⁸, Clare Turnbull³⁸, Fredrik Wiklund⁶⁹

493

494

495 **References**

496 1. Sud, A., Kinnersley, B. & Houlston, R. S. Genome-wide association studies of cancer:
497 current insights and future perspectives. *Nat Rev Cancer* **17**, 692-704 (2017).

498 2. Tam, V. et al. Benefits and limitations of genome-wide association studies. *Nat Rev Genet*
499 (2019).

500 3. Law, P. J. et al. Genome-wide association analysis implicates dysregulation of immunity
501 genes in chronic lymphocytic leukaemia. *Nat Commun* **8**, 14175 (2017).

- 502 4. Litchfield, K. et al. Identification of 19 new risk loci and potential regulatory mechanisms
503 influencing susceptibility to testicular germ cell tumor. *Nat Genet* **49**, 1133-1140 (2017).
- 504 5. Mucci, L. A. et al. Familial Risk and Heritability of Cancer Among Twins in Nordic
505 Countries. *JAMA* **315**, 68-76 (2016).
- 506 6. Maas, P. et al. Breast Cancer Risk From Modifiable and Nonmodifiable Risk Factors
507 Among White Women in the United States. *JAMA Oncol* **2**, 1295-1302 (2016).
- 508 7. Mavaddat, N. et al. Polygenic Risk Scores for Prediction of Breast Cancer and Breast
509 Cancer Subtypes. *Am J Hum Genet* **104**, 21-34 (2019).
- 510 8. Jeon, J. et al. Determining Risk of Colorectal Cancer and Starting Age of Screening Based
511 on Lifestyle, Environmental, and Genetic Factors. *Gastroenterology* **154**, 2152-2164.e19
512 (2018).
- 513 9. Seibert, T. M. et al. Polygenic hazard score to guide screening for aggressive prostate
514 cancer: development and validation in large scale cohorts. *BMJ* **360**, j5757 (2018).
- 515 10. Garcia-Closas, M. et al. Common genetic polymorphisms modify the effect of smoking on
516 absolute risk of bladder cancer. *Cancer Res* **73**, 2211-2220 (2013).
- 517 11. Turnbull, C., Sud, A. & Houlston, R. S. Cancer genetics, precision prevention and a call to
518 action. *Nat Genet* **50**, 1212-1218 (2018).
- 519 12. Chatterjee, N., Shi, J. & García-Closas, M. Developing and evaluating polygenic risk
520 prediction models for stratified disease prevention. *Nat Rev Genet* **17**, 392-406 (2016).
- 521 13. Yang, J. et al. Common SNPs explain a large proportion of the heritability for human
522 height. *Nat Genet* **42**, 565-569 (2010).
- 523 14. Sampson, J. N. et al. Analysis of Heritability and Shared Heritability Based on Genome-
524 Wide Association Studies for Thirteen Cancer Types. *J Natl Cancer Inst* **107**, djv279

- 525 (2015).
- 526 15. Zhang, Y., Qi, G., Park, J. H. & Chatterjee, N. Estimation of complex effect-size
527 distributions using summary-level statistics from genome-wide association studies across
528 32 complex traits. *Nat Genet* **50**, 1318-1326 (2018).
- 529 16. Berndt, S. I. et al. Meta-analysis of genome-wide association studies discovers multiple
530 loci for chronic lymphocytic leukemia. *Nat Commun* **7**, 10933 (2016).
- 531 17. Wang, Z. et al. Meta-analysis of five genome-wide association studies identifies multiple
532 new loci associated with testicular germ cell tumor. *Nat Genet* **49**, 1141-1147 (2017).
- 533 18. Lesueur, C. et al. Genome-wide association analyses identify new susceptibility loci for
534 oral cavity and pharyngeal cancer. *Nat Genet* **48**, 1544-1550 (2016).
- 535 19. Klein, A. P. et al. Genome-wide meta-analysis identifies five new susceptibility loci for
536 pancreatic cancer. *Nat Commun* **9**, 556 (2018).
- 537 20. Scelo, G. et al. Genome-wide association study identifies multiple risk loci for renal cell
538 carcinoma. *Nat Commun* **8**, 15724 (2017).
- 539 21. Melin, B. S. et al. Genome-wide association study of glioma subtypes identifies specific
540 differences in genetic susceptibility to glioblastoma and non-glioblastoma tumors. *Nat*
541 *Genet* **49**, 789-794 (2017).
- 542 22. Law, M. H. et al. Genome-wide meta-analysis identifies five new susceptibility loci for
543 cutaneous malignant melanoma. *Nat Genet* **47**, 987-995 (2015).
- 544 23. O'Mara, T. A. et al. Identification of nine new susceptibility loci for endometrial cancer.
545 *Nat Commun* **9**, 3166 (2018).
- 546 24. Schumacher, F. R. et al. Genome-wide association study of colorectal cancer identifies six
547 new susceptibility loci. *Nat Commun* **6**, 7138 (2015).

- 548 25. Phelan, C. M. et al. Identification of 12 new susceptibility loci for different histotypes of
549 epithelial ovarian cancer. *Nat Genet* **49**, 680-691 (2017).
- 550 26. Schumacher, F. R. et al. Association analyses of more than 140,000 men identify 63 new
551 prostate cancer susceptibility loci. *Nat Genet* **50**, 928-936 (2018).
- 552 27. McKay, J. D. et al. Large-scale association analysis identifies new lung cancer
553 susceptibility loci and heterogeneity in genetic susceptibility across histological subtypes.
554 *Nat Genet* **49**, 1126-1132 (2017).
- 555 28. Michailidou, K. et al. Association analysis identifies 65 new breast cancer risk loci. *Nature*
556 **551**, 92-94 (2017).
- 557 29. Choudhury, P. P. et al. iCARE: An R package to build, validate and apply absolute risk
558 models. *PloS one* 15.2 (2020): e0228198.
- 559 30. Zeng, J. et al. Signatures of negative selection in the genetic architecture of human complex
560 traits. *Nat Genet* **50**, 746-753 (2018).
- 561 31. Stahl, E. A. et al. Bayesian inference analyses of the polygenic architecture of rheumatoid
562 arthritis. *Nat Genet* **44**, 483-489 (2012).
- 563 32. Khera, A. V. et al. Genome-wide polygenic scores for common diseases identify
564 individuals with risk equivalent to monogenic mutations. *Nat Genet* **50**, 1219-1224 (2018).
- 565 33. Schork, A. J. et al. All SNPs are not created equal: genome-wide association studies reveal
566 a consistent pattern of enrichment among functionally annotated SNPs. *PLoS Genet* **9**,
567 (2013).
- 568 34. Pickrell, J. K. Joint analysis of functional genomic data and genome-wide association
569 studies of 18 human traits. *Am J Hum Genet* **94**, 559-573 (2014).
- 570 35. Andreassen, O. A. et al. Improved detection of common variants associated with

- 571 schizophrenia by leveraging pleiotropy with cardiovascular-disease risk factors. *Am J Hum*
572 *Genet* **92**, 197-209 (2013).
- 573 36. Andreassen, O. A. et al. Improved detection of common variants associated with
574 schizophrenia and bipolar disorder using pleiotropy-informed conditional false discovery
575 rate. *PLoS Genet* **9**, e1003455 (2013).
- 576 37. Hu, Y. et al. Leveraging functional annotations in genetic risk prediction for human
577 complex diseases. *PLoS Comput Biol* **13**, e1005589 (2017).
- 578 38. Davey Smith, G. & Hemani, G. Mendelian randomization: genetic anchors for causal
579 inference in epidemiological studies. *Hum Mol Genet* **23**, R89-98 (2014).
- 580 39. Wainschein, P. et al. Recovery of trait heritability from whole genome sequence data.
581 Preprint at <https://www.biorxiv.org/content/10.1101/588020v1> (2019).
582
- 583 40. O'Connor, L. J. et al. Extreme polygenicity of complex traits is explained by negative
584 selection. *The American Journal of Human Genetics* 105.3 (2019): 456-476.
- 585 41. GTEx, C. The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**, 580-585 (2013).
- 586 42. Visvader, J. E. Cells of origin in cancer. *Nature* **469**, 314-322 (2011).
- 587 43. Rizzo, A. A., Strickland, D. & Bouchard, S. The challenge of using virtual reality in
588 telerehabilitation. *Telemed J E Health* **10**, 184-195 (2004).
- 589 44. Hutter, C. M. et al. Gene-environment interactions in cancer epidemiology: a National
590 Cancer Institute Think Tank report. *Genet Epidemiol* **37**, 643-657 (2013).
- 591 45. Bulik-Sullivan, B. K. et al. LD Score regression distinguishes confounding from
592 polygenicity in genome-wide association studies. *Nat Genet* **47**, 291-295 (2015).
- 593 46. Zheng, J. et al. LD Hub: a centralized database and web interface to perform LD score
594 regression that maximizes the potential of summary level GWAS data for SNP heritability
595 and genetic correlation analysis. *Bioinformatics* **33**, 272-279 (2017).

- 596 47. Pharoah, P. D. et al. Polygenic susceptibility to breast cancer and implications for
597 prevention. *Nat Genet* **31**, 33-36 (2002).
- 598 48. Chatterjee, N. et al. Projecting the performance of risk prediction based on polygenic
599 analyses of genome-wide association studies. *Nat Genet* **45**, 400-5, 405e1 (2013).
- 600
- 601

602 **Figure 1. Estimated effect-size distributions for susceptibility SNPs across 14 cancer sites.** Effect-size
603 distribution of susceptibility SNPs is modelled using a two-component normal mixture model for all sites,
604 except esophageal and oropharyngeal cancers. For these sites, effect-sizes are modelled using a single
605 normal distribution that provided similar fit as the two-component normal mixture model (see
606 **Supplementary Figure 1 & 2**). SNPs with extremely large effects are excluded for effect-size distribution
607 estimation (see **Methods**). Plots are stratified by sample size of the GWAS for comparability.
608 Distributions with fatter tails imply the underlying traits have relatively greater number of susceptibility
609 SNPs with larger effects. Note here the effect-size distribution is plotted on the log scale of odds ratio (x-
610 axis). CLL = chronic lymphocytic leukemia.

611
612 **Figure 2. Projections of percentage of GWAS heritability explained by SNPs as sample size for GWAS**
613 **increases.** Results are shown for projections including SNPs at the optimized p -value threshold (solid
614 curve) and at genome-wide significance ($p < 5 \times 10^{-8}$) level (dashed curve). Colored dots correspond to
615 sample size for largest published GWAS and those for doubled and quadrupled sizes. For oropharyngeal
616 cancer, the projections at the “current sample size” are based on a sample size of 25K cases and 25K
617 controls. For breast and esophageal cancer, the projections at the “current sample size” are based on
618 the current largest GWAS sample sizes: 123K cases and 106K controls, and 10K cases and 17K controls,
619 respectively. For all other cancer sites, the projections at the “current sample size” are based on the
620 GWAS sample sizes in **Supplementary Table 1**. CLL = chronic lymphocytic leukemia.

621
622 **Figure 3. Projections of area under the curve (AUC) characterizing predictive performance of PRS as**
623 **sample size for GWAS increases.** Results are shown for PRS including SNPs at the optimized p -value
624 threshold. The dotted horizontal red line indicates the maximum AUC achievable according to the
625 estimate of GWAS heritability. Colored dots correspond to sample size for largest published GWAS and
626 those for doubled and quadrupled sizes. For oropharyngeal cancer, the projections at the “current
627 sample size” are based on a sample size of 25K cases and 25K controls. For breast and esophageal
628 cancer, the projections at the “current sample size” are based on the current largest GWAS sample sizes:
629 123K cases and 106K controls, and 10K cases and 17K controls, respectively. For all other cancer sites,
630 the projections at the “current sample size” are based on the GWAS sample sizes in **Supplementary**
631 **Table 1**. CLL = chronic lymphocytic leukemia.

632
633 **Figure 4. Projections of relative risks for individuals at or higher than 99th percentile of PRS as sample**
634 **size for GWAS increases.** Results are shown where PRS is built based on SNPs at optimized p -value
635 threshold. The dotted horizontal red line indicates the maximum relative risk achievable according to
636 estimate of GWAS heritability. Colored dots correspond to sample size for largest published GWAS and
637 those for doubled and quadrupled sizes. Y-axis is presented in log10 scale. For oropharyngeal cancer,
638 the projections at the “current sample size” are based on a sample size of 25K cases and 25K controls.
639 For breast and esophageal cancer, the projections at the “current sample size” are based on the current
640 largest GWAS sample sizes: 123K cases and 106K controls, and 10K cases and 17K controls, respectively.
641 For all other cancer sites, the projections at the “current sample size” are based on the GWAS sample
642 sizes in **Supplementary Table 1**. CLL = chronic lymphocytic leukemia.

643
644 **Figure 5. Projected distribution of average residual lifetime risk in the US population of Non-Hispanic**
645 **Whites aged 30 to 75 years.** The risk is obtained according to variation of polygenic risk scores. The
646 projections are shown for PRS built based on GWAS with current, doubled and quadrupled sample sizes
647 and the best PRS that corresponds to limits defined by heritability. The projections are obtained by
648 combining information on projected population variance of PRS, age-specific population incidence rate,
649 competing risk of mortality and current distribution of age according to US 2016 census. For

650 oropharyngeal cancer, the projections at the “current sample size” are based on a sample size of 25K
651 cases and 25K controls. For breast and esophageal cancer, the projections at the “current sample size”
652 are based on the current largest GWAS sample sizes: 123K cases and 106K controls, and 10K cases and
653 17K controls, respectively. For all other cancer sites, the projections at the “current sample size” are
654 based on the GWAS sample sizes in **Supplementary Table 1**. CLL = chronic lymphocytic leukemia.
655

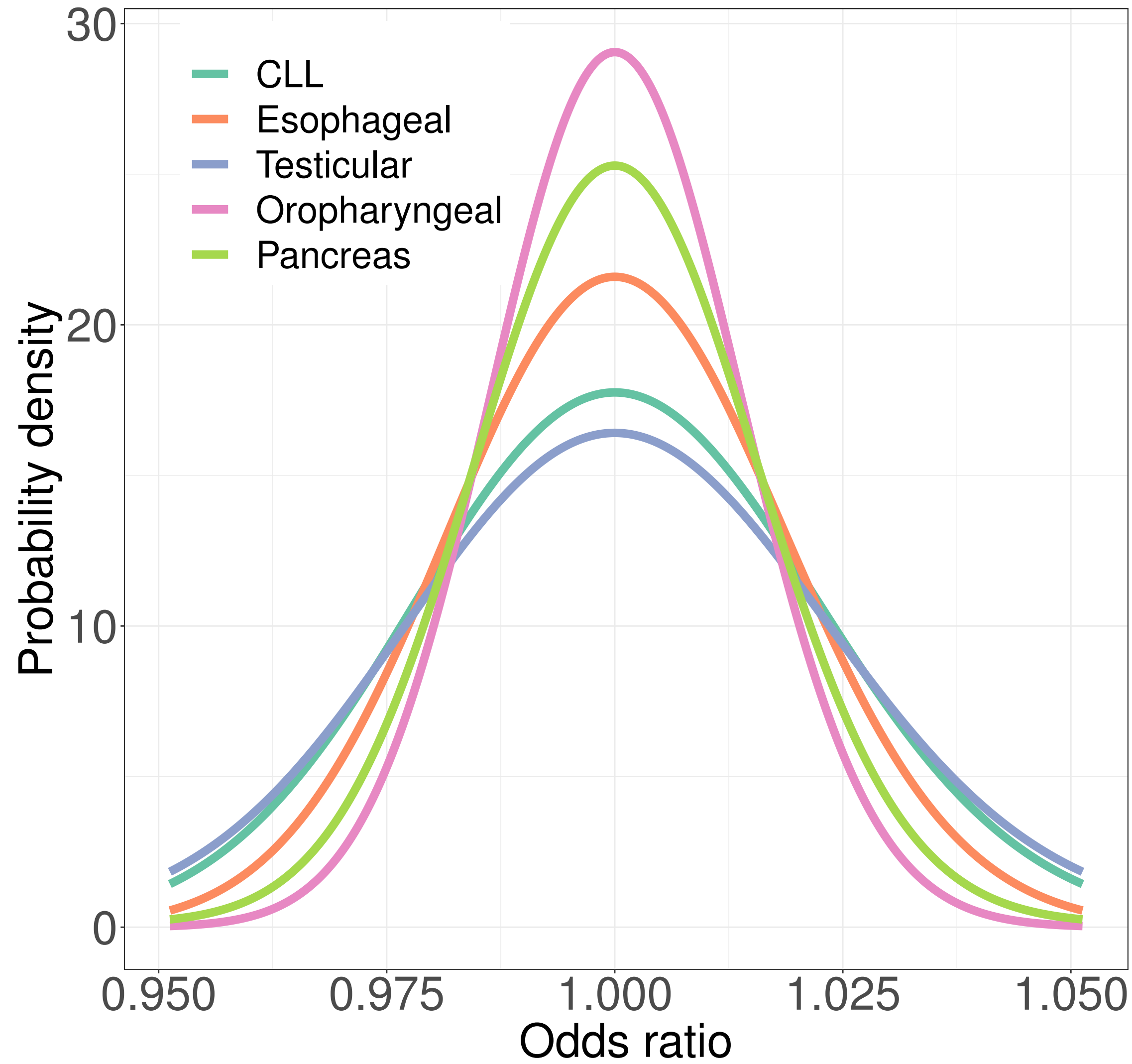
656 **Table 1: Estimated number of independent common susceptibility variants and heritability across 14**
 657 **cancer sites.**
 658

Number of cases in the analysis	Cancer site ^a	Total Number of susceptibility SNPs ^b (SE ^c)	Total heritability, in log-OR scale ^d (SE)	Average heritability explained per susceptibility SNP ^e (SE), in 10 ⁻⁴	Number of SNPs associated with larger variance component (SE)	% of heritability explained by SNPs with larger variance component	AUC associated with the best PRS ^f (SE)
<10,000	CLL ^g	2025 (1501)	1.62 (0.37)	7.2 (4.4)	52 (15)	41	0.82 (0.03)
<10,000	Esophageal	3641 (2515)	1.24 (0.36)	3.4 (1.9)	NA ^h	NA	0.78 (0.03)
<10,000	Testicular	2598 (2088)	2.81 (0.40)	9.2 (6.6)	196 (75)	54	0.88 (0.02)
<10,000	Oropharyngeal	3623 (2060)	0.68 (0.27)	1.9 (0.5)	NA	NA	0.72 (0.04)
<10,000	Pancreas	1757 (1490)	0.60 (0.16)	3.2 (2.2)	47 (27)	31	0.71 (0.03)
10,000 - 25,000	Renal	2220 (1555)	0.57 (0.12)	2.4 (1.4)	46 (36)	24	0.70 (0.02)
10,000 - 25,000	Glioma	2364 (1593)	0.87 (0.11)	2.2 (1.2)	61 (25)	55	0.75 (0.01)
10,000 - 25,000	Melanoma	1098 (533)	0.65 (0.09)	4.4 (1.6)	106 (58)	52	0.72 (0.01)
10,000 - 25,000	Colorectal	1484 (696)	0.43 (0.10)	2.9 (0.8)	14 (11)	7	0.68 (0.02)
10,000 - 25,000	Endometrial	1052 (772)	0.27 (0.07)	2.5 (1.3)	46 (34)	26	0.64 (0.02)
10,000 - 25,000	Ovarian	1015 (715)	0.24 (0.06)	2.2 (1.1)	49 (31)	36	0.64 (0.02)
>25,000	Lung	6096 (2750)	0.39 (0.06)	0.6 (0.2)	15 (7)	15	0.67 (0.01)
>25,000	Prostate	4530 (1052)	0.77 (0.04)	1.1 (0.2)	276 (99)	51	0.73 (0.01)
>25,000	Breast	7599 (1615)	0.60 (0.03)	0.6 (0.1)	587 (133)	56	0.71 (0.00)

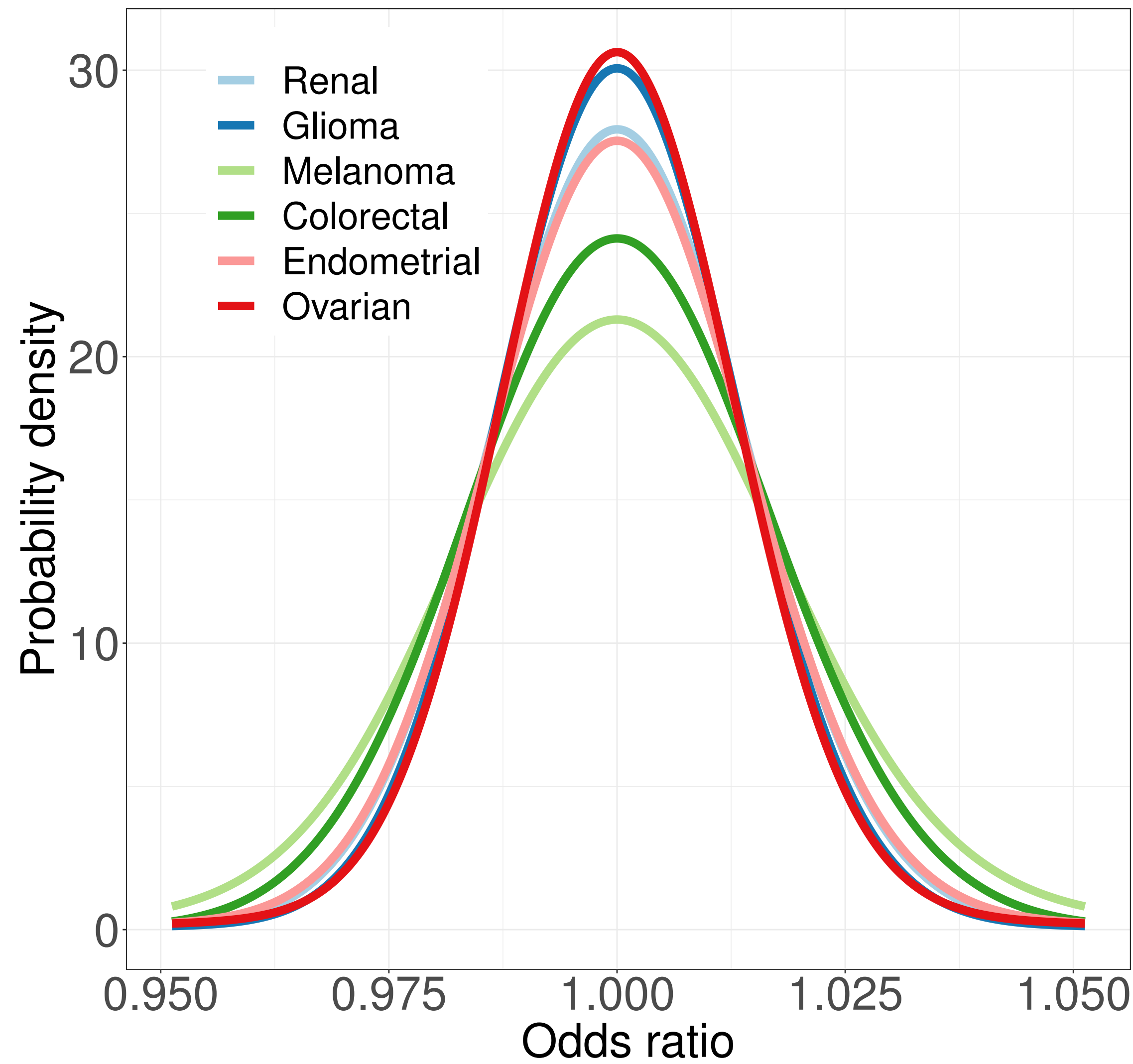
659 ^aAll results are reported using the best fitted (two- or three-component) normal mixture model for effect-size distributions, with respect to a
 660 reference panel of 1.07 million common SNPs included in the Hapmap3 panel after removal of MHC region. ^bSNP: single nucleotide
 661 polymorphism. ^cStandard errors. ^dTotal heritability is characterized by population variance of the underlying true PRS as
 662 $h^2 = \text{Var}(\sum_{m=1}^M \beta_m G_m) = M\pi_c E(\beta^2)$, where $E(\beta^2)$ denotes per-SNP effect-size of the non-null SNPs in the log-odds-ratio scale. ^eAverage
 663 heritability explained per susceptibility SNP excludes SNPs with extremely large effects (see **Methods**). ^fArea under the curve (AUC) associated
 664 with best PRS is calculated using the formula $\text{AUC} = \Phi(\sqrt{h^2/2})$ where $\Phi(\cdot)$ is the cumulative density function of standard normal distribution.
 665 ^gCLL = chronic lymphocytic leukemia. ^hNA indicates a two-component model is favourable compared to three-component model.

666
 667
 668
 669

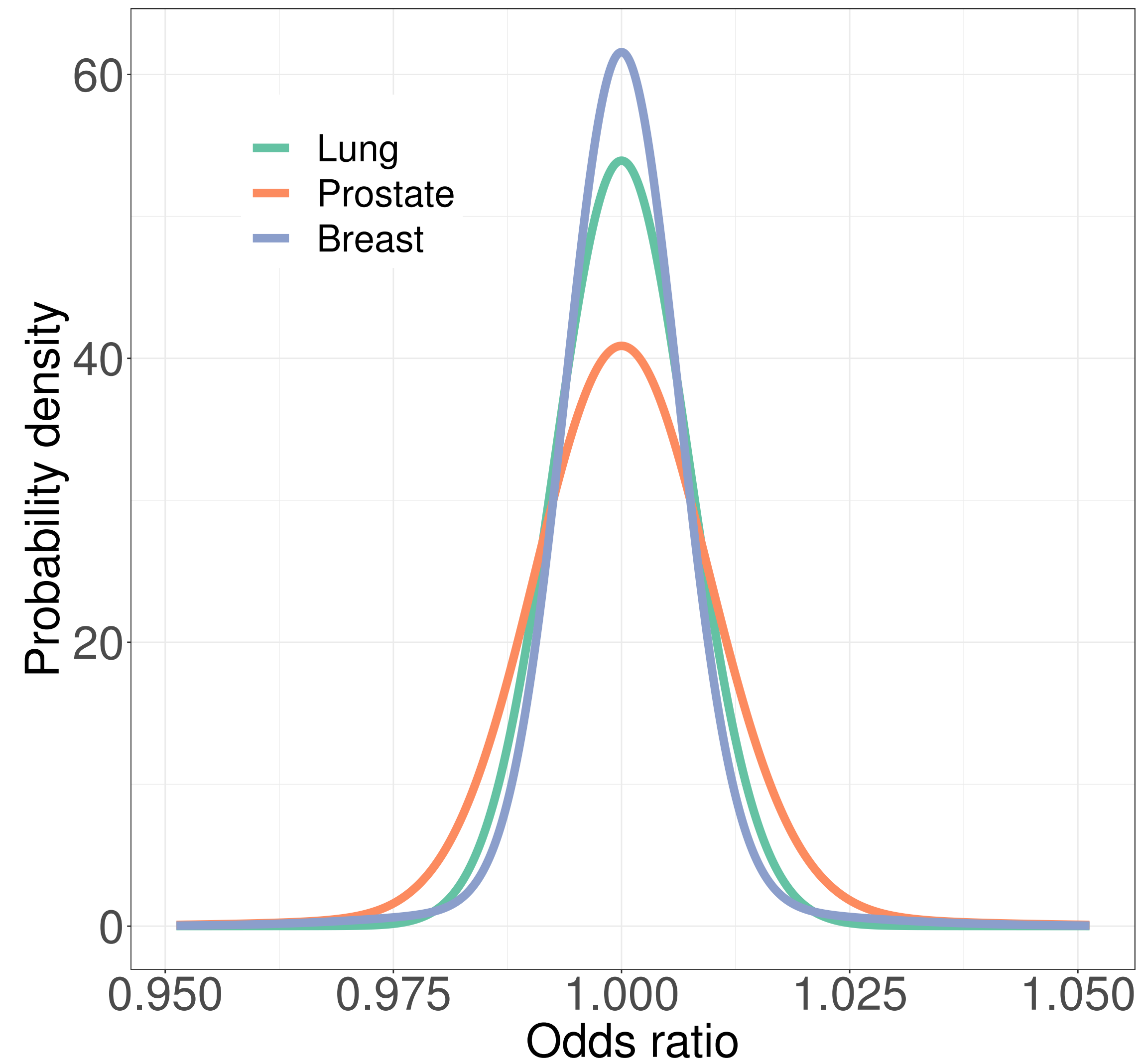
Cancer sites with <10,000 cases

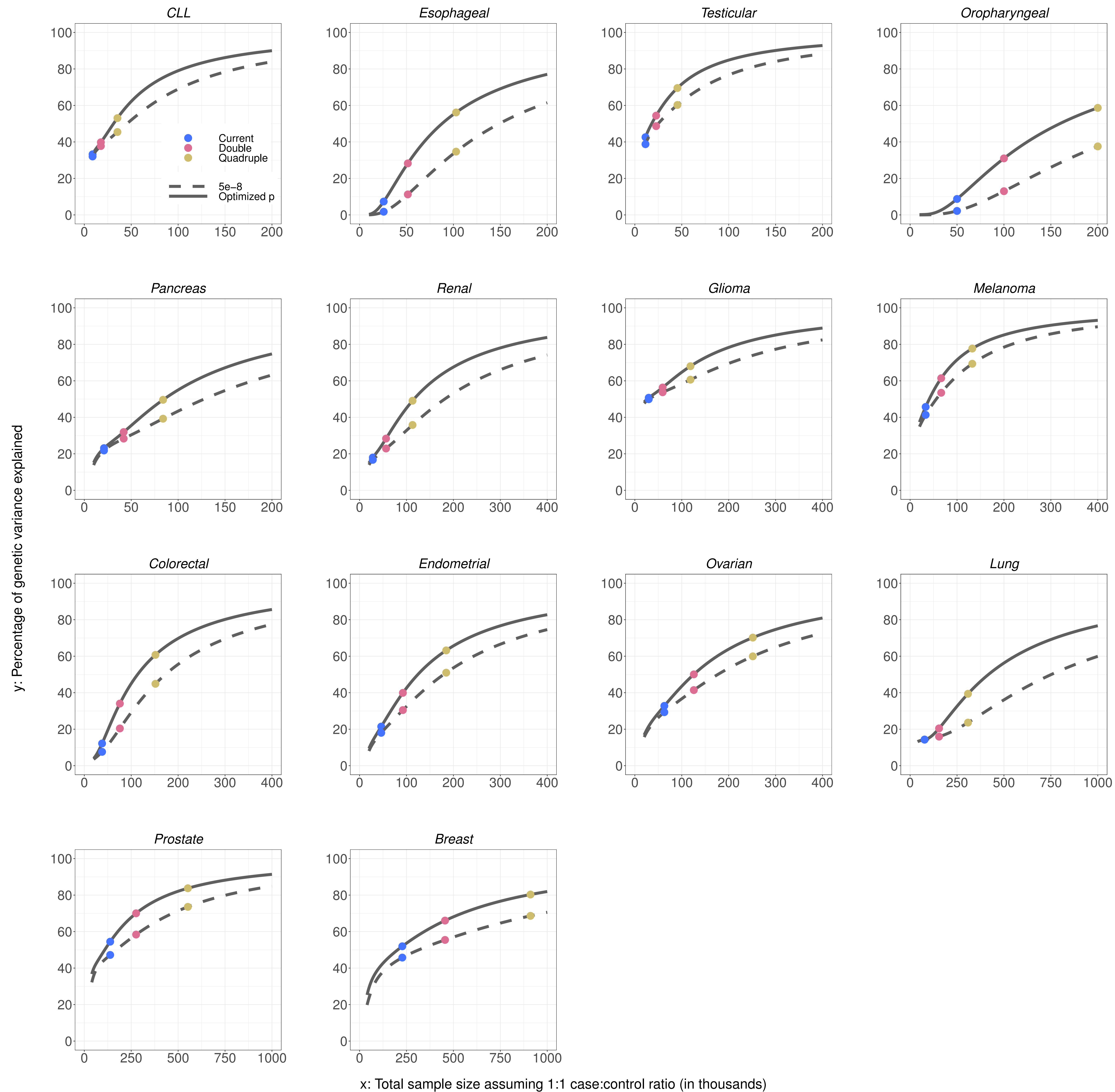


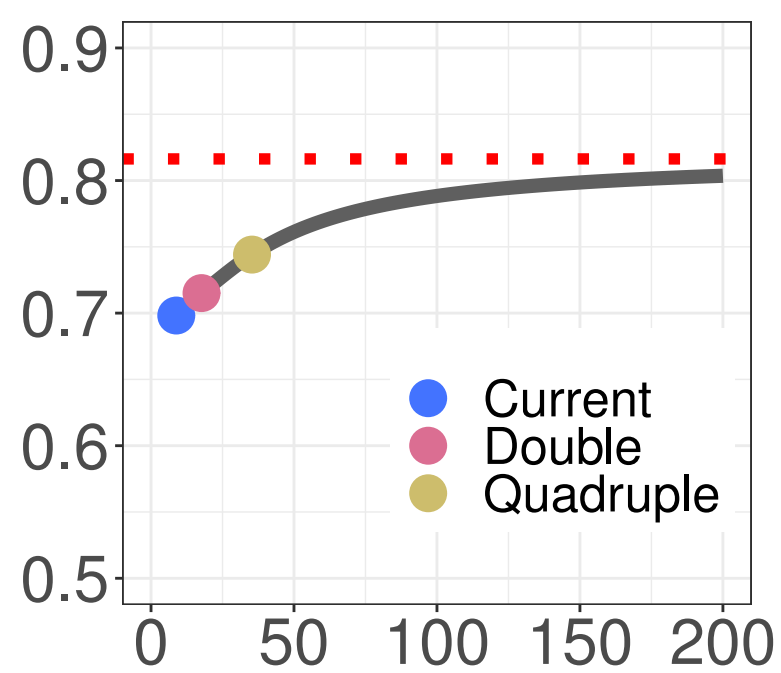
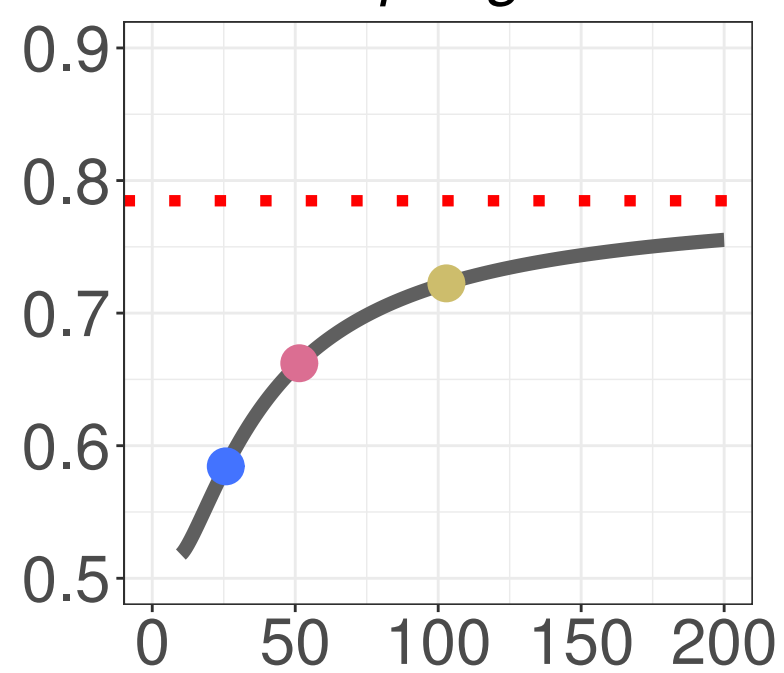
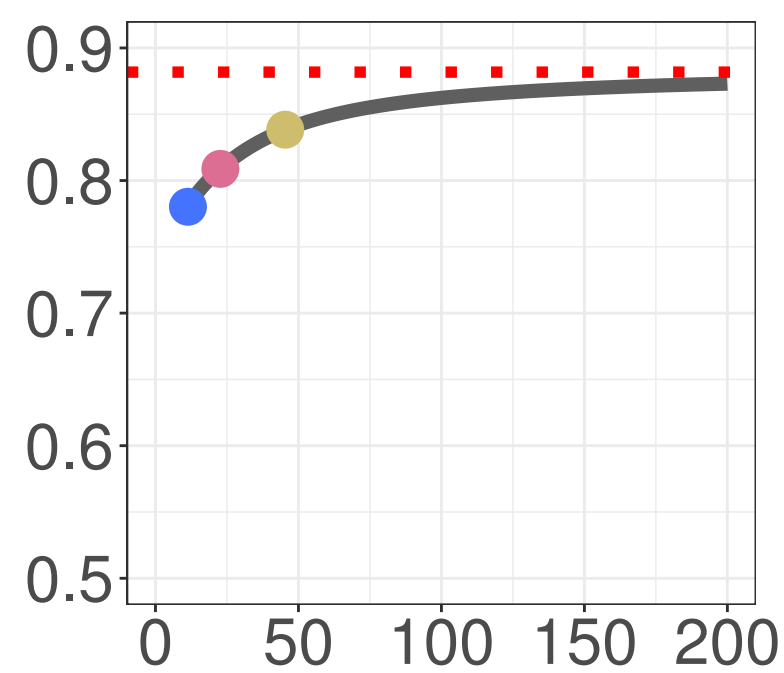
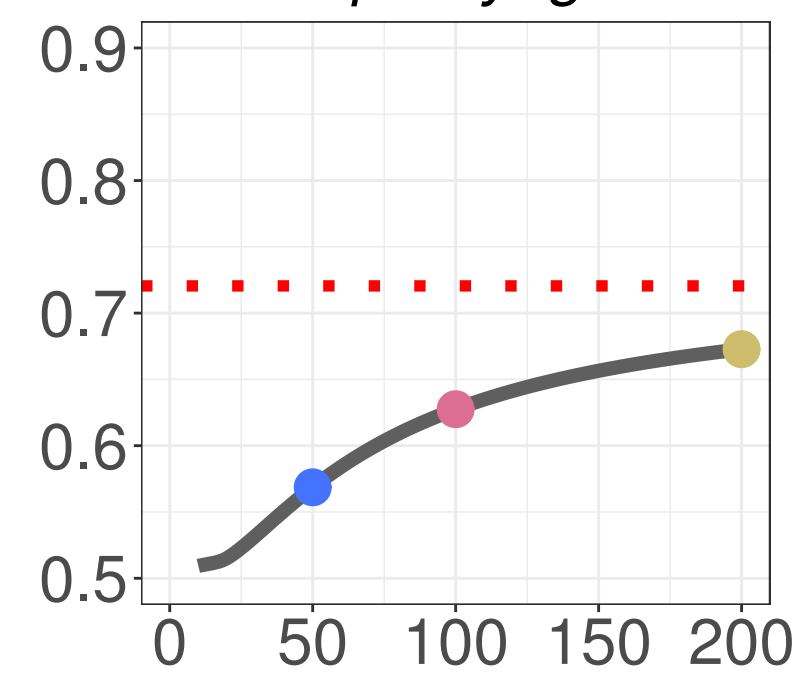
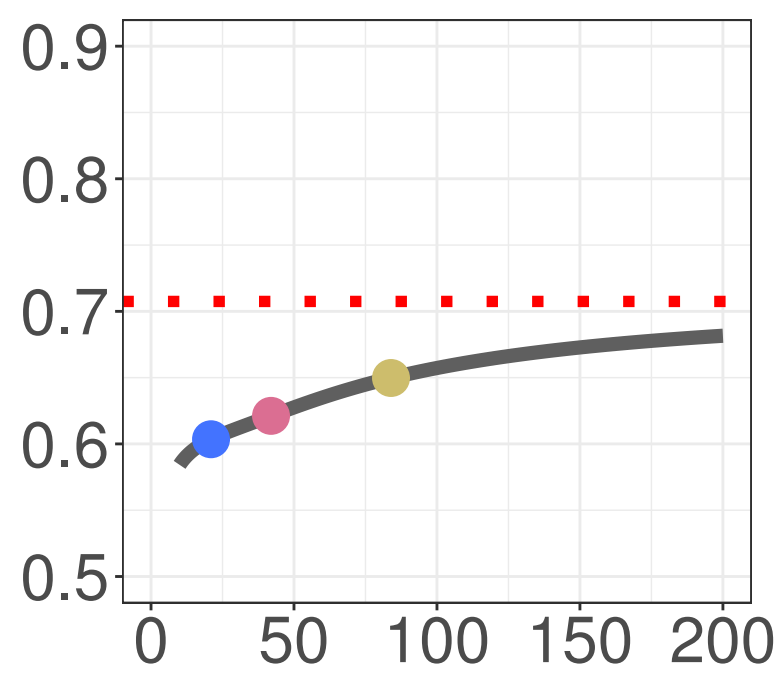
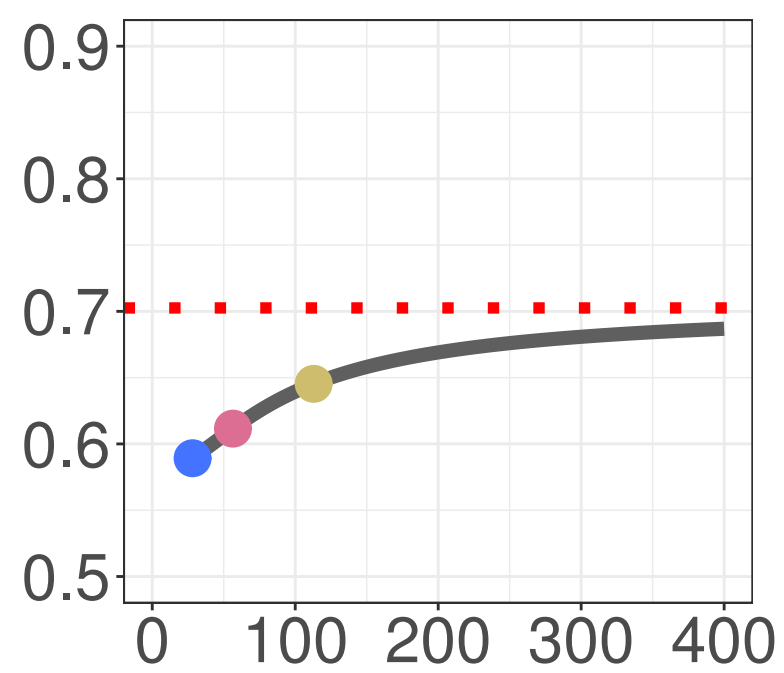
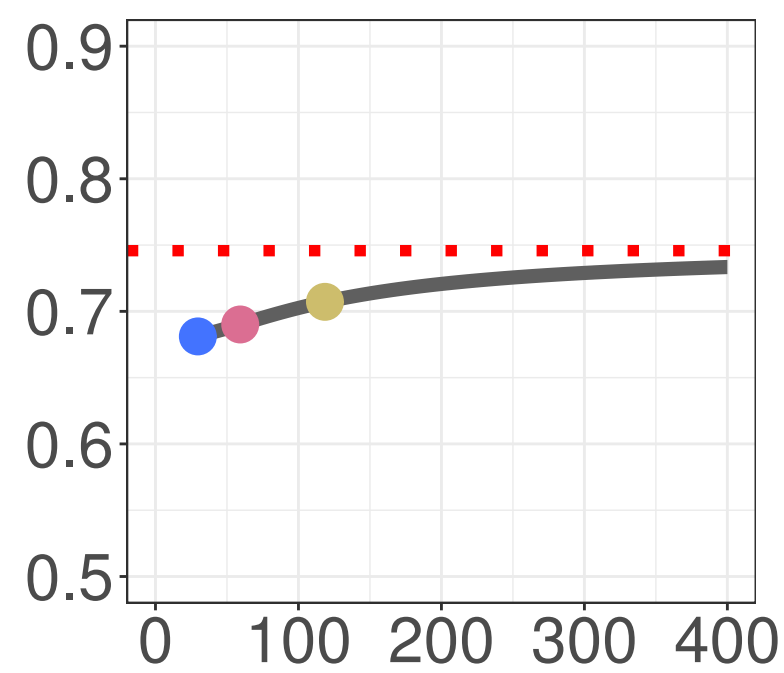
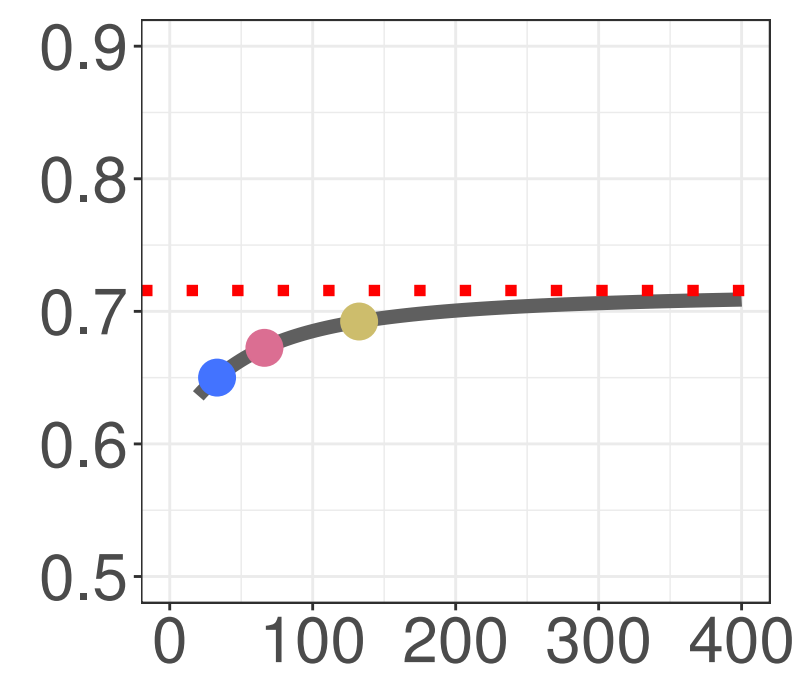
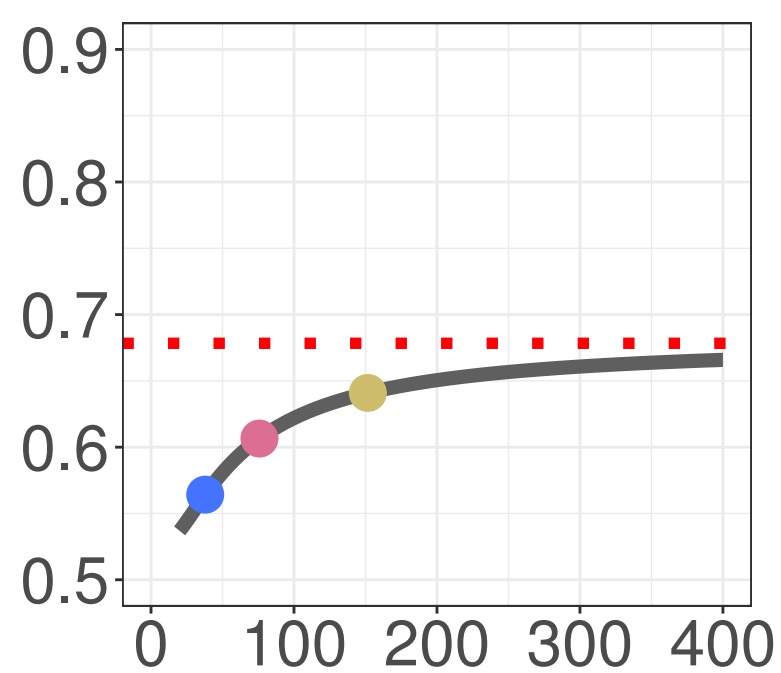
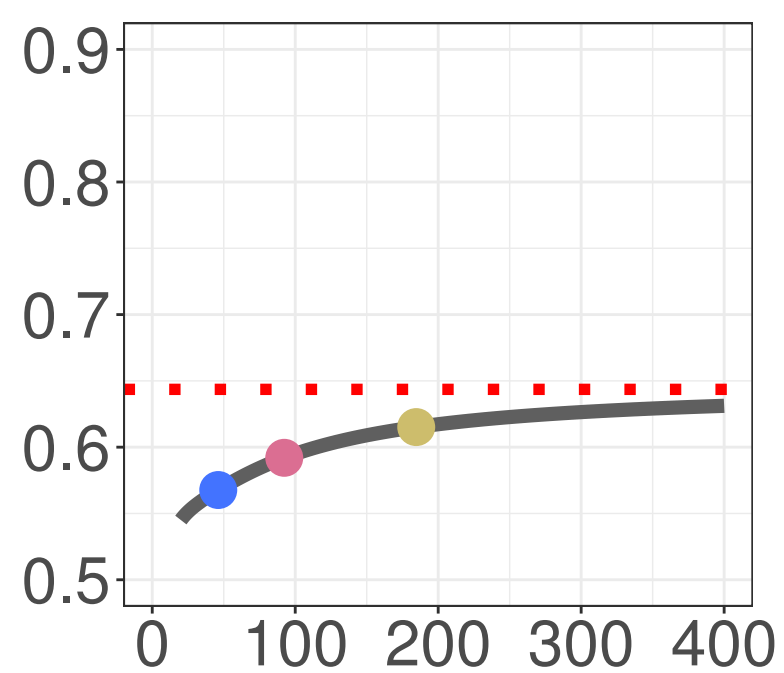
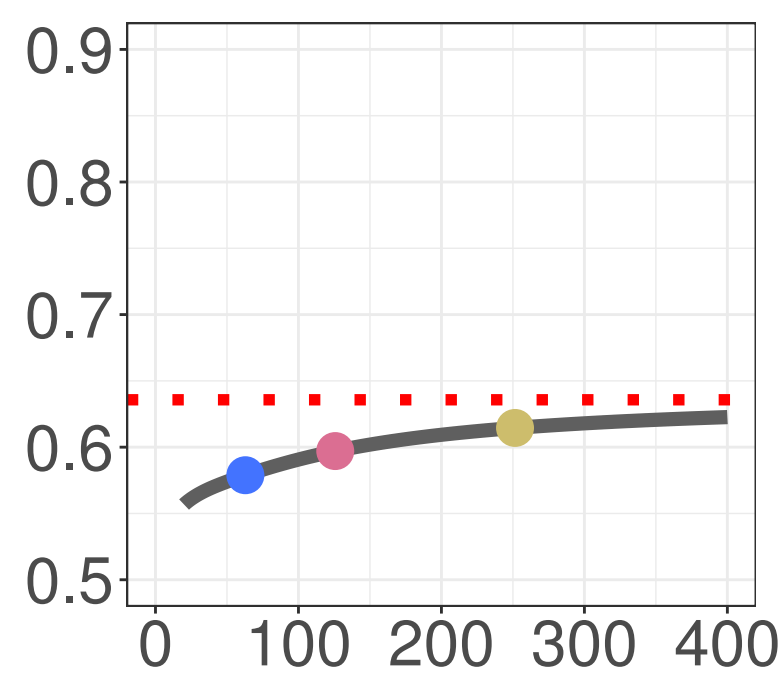
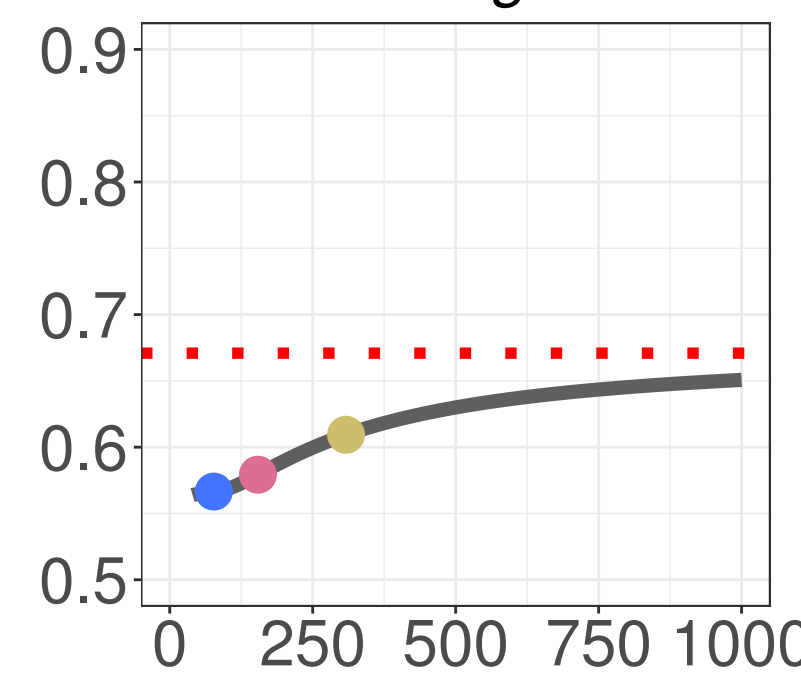
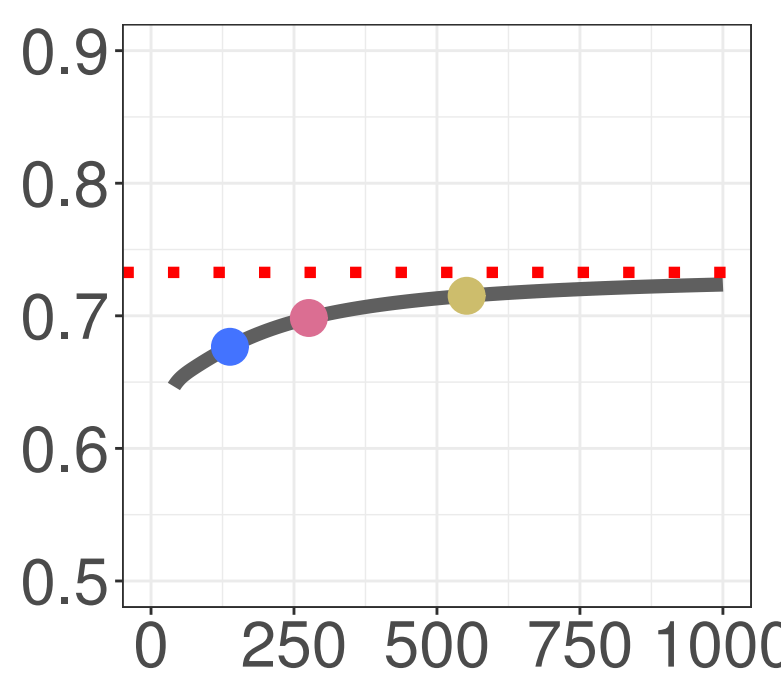
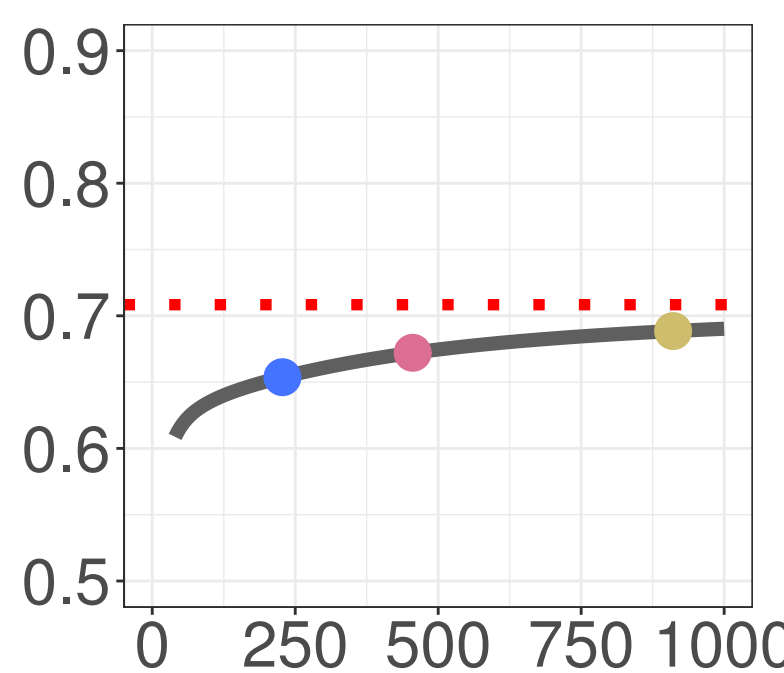
Cancer sites with 10,000–25,000 cases



Cancer sites with >25,000 cases



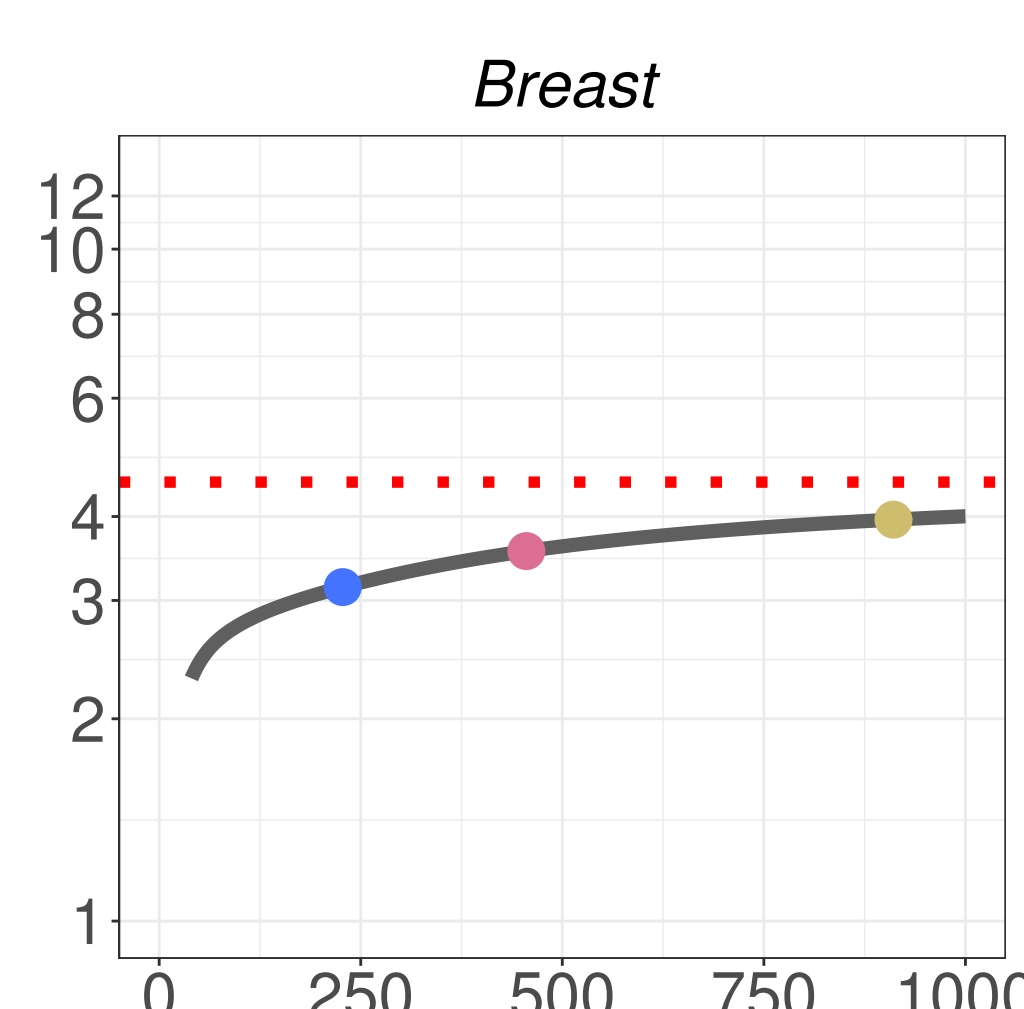
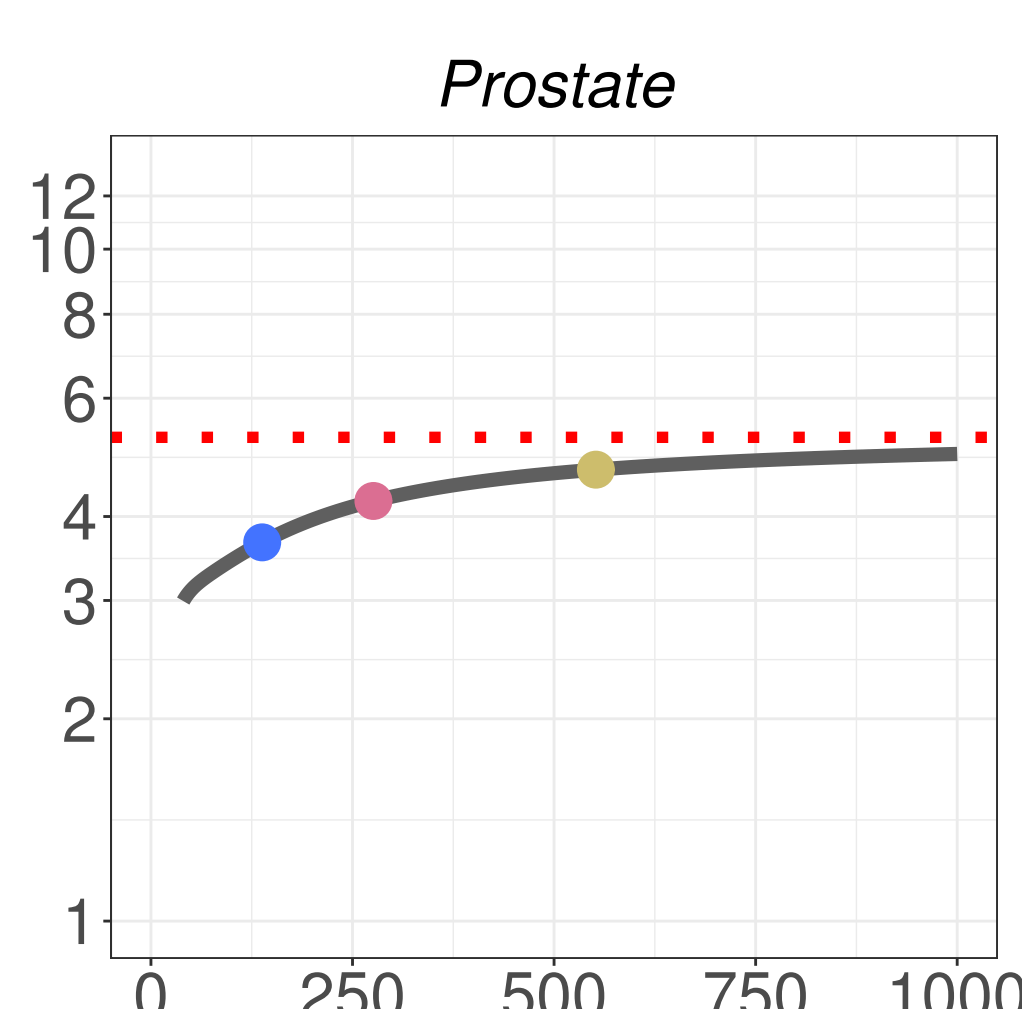
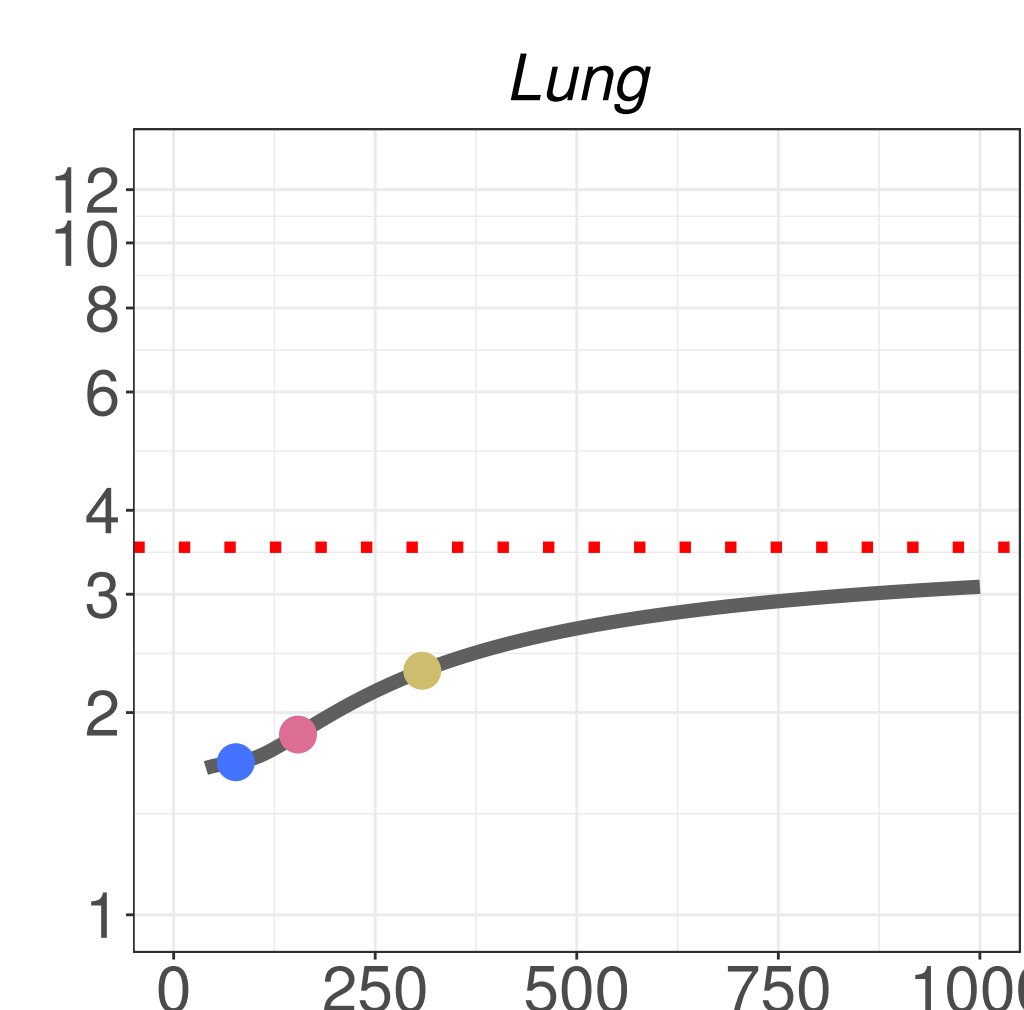
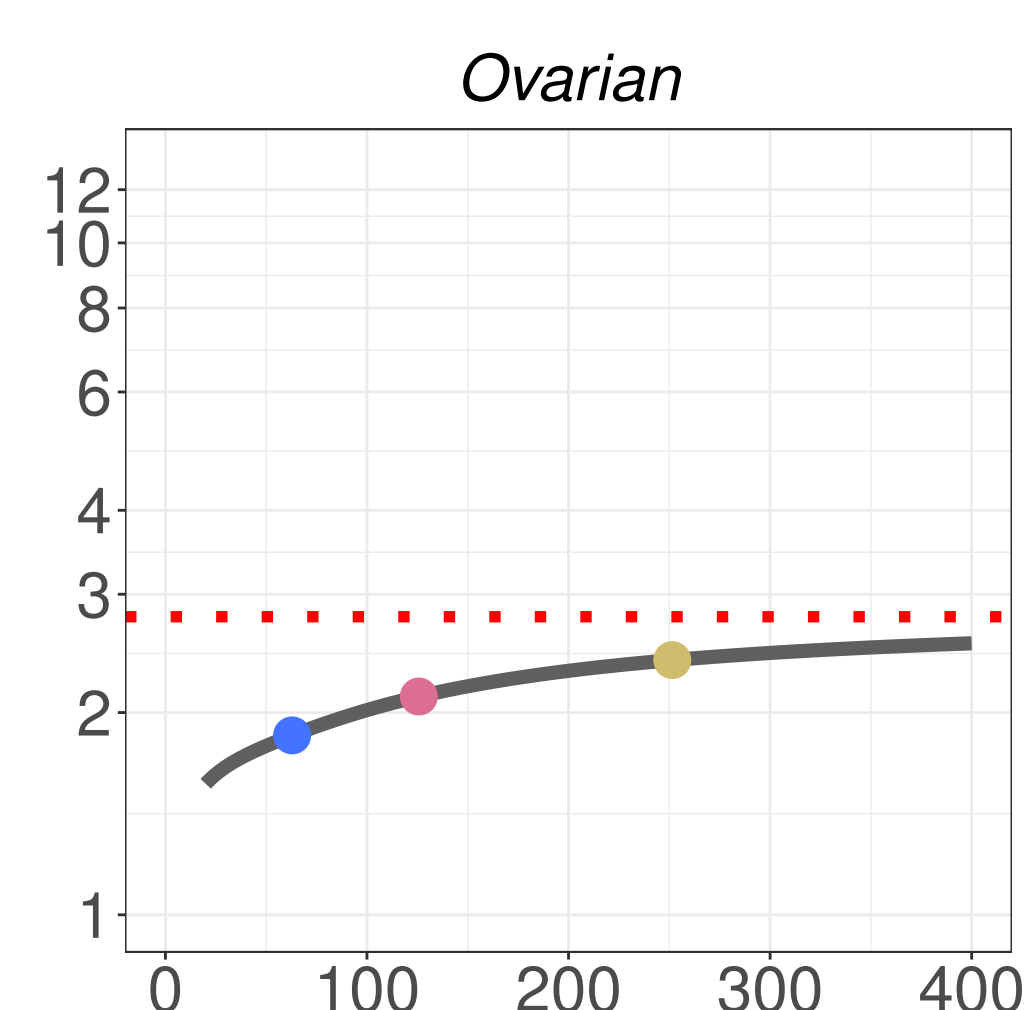
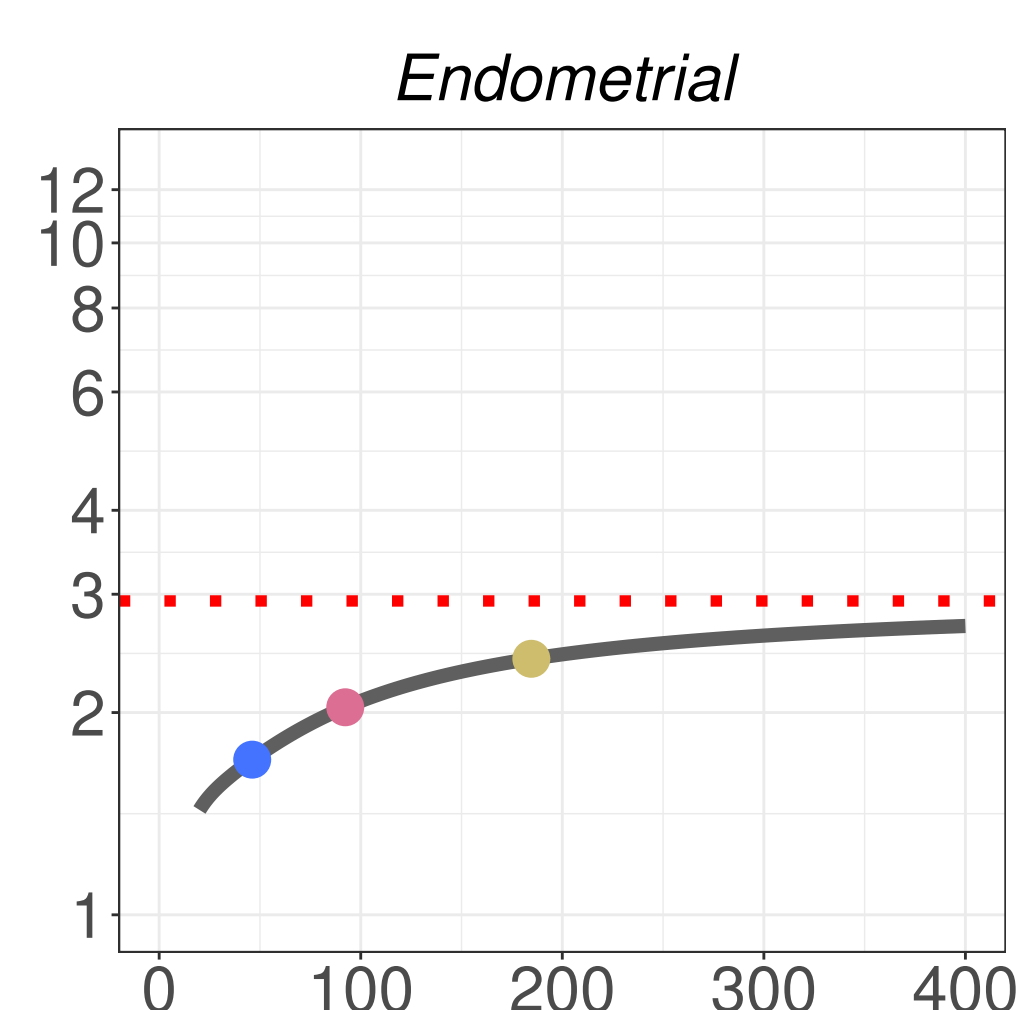
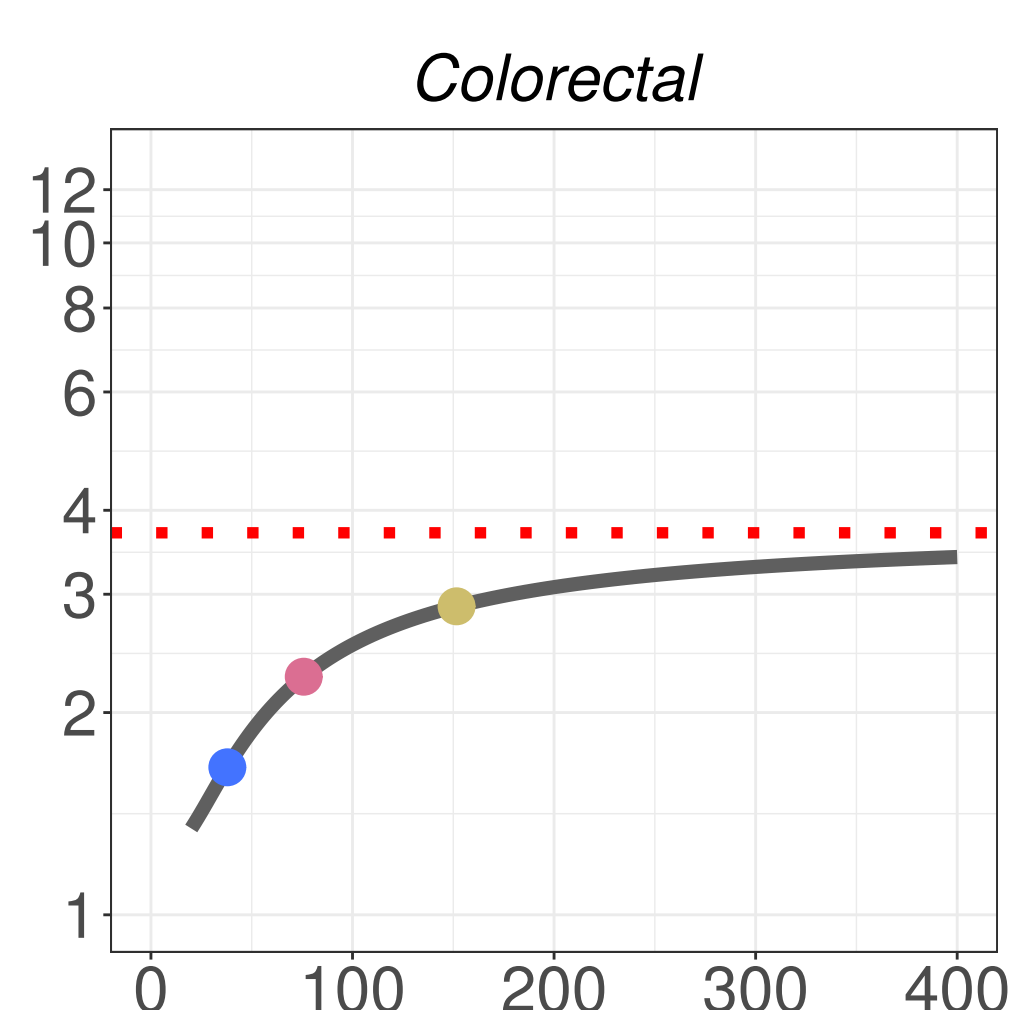
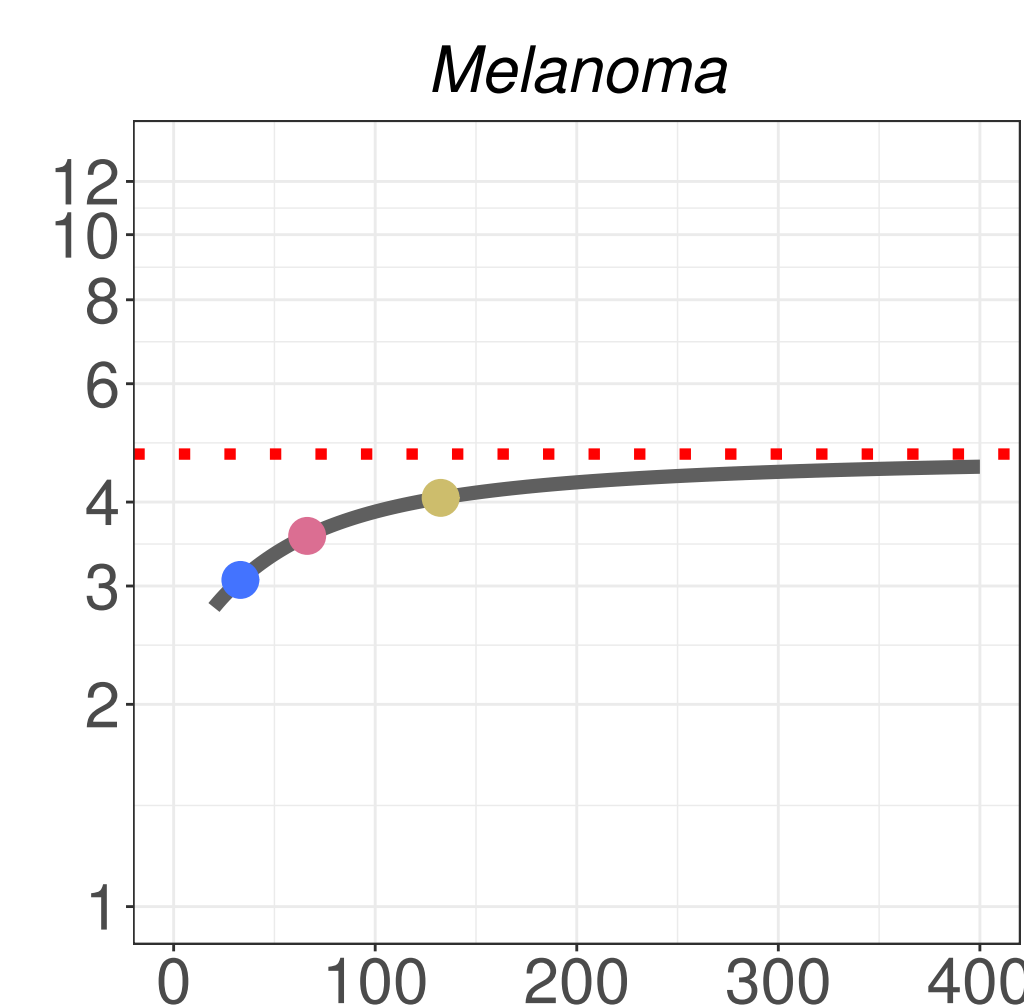
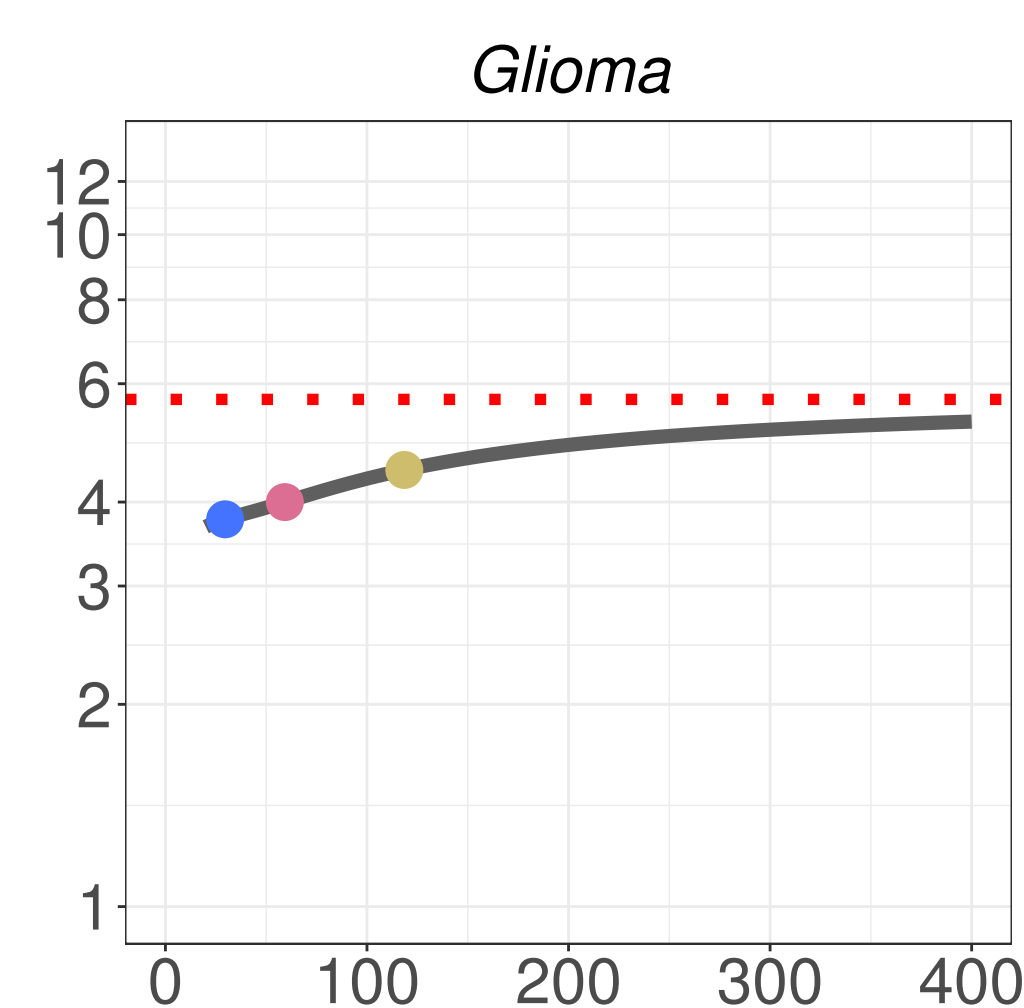
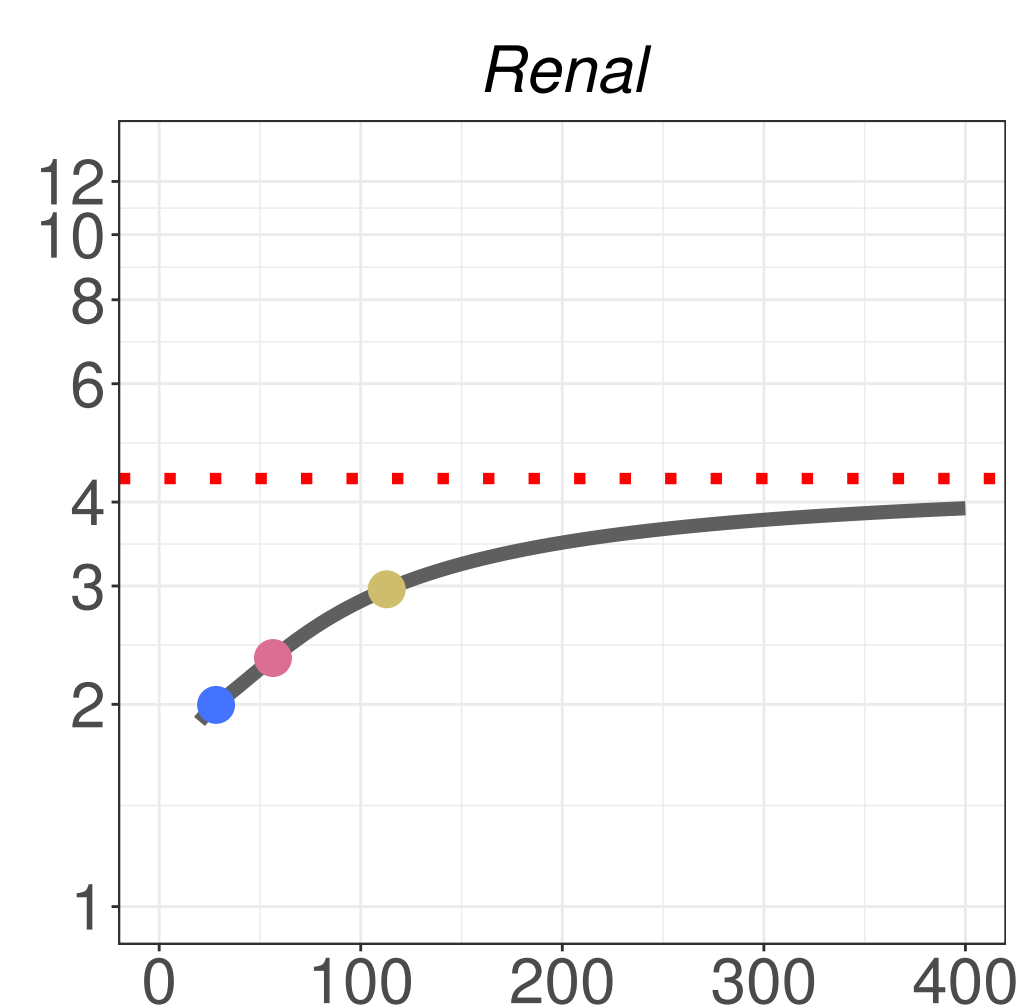
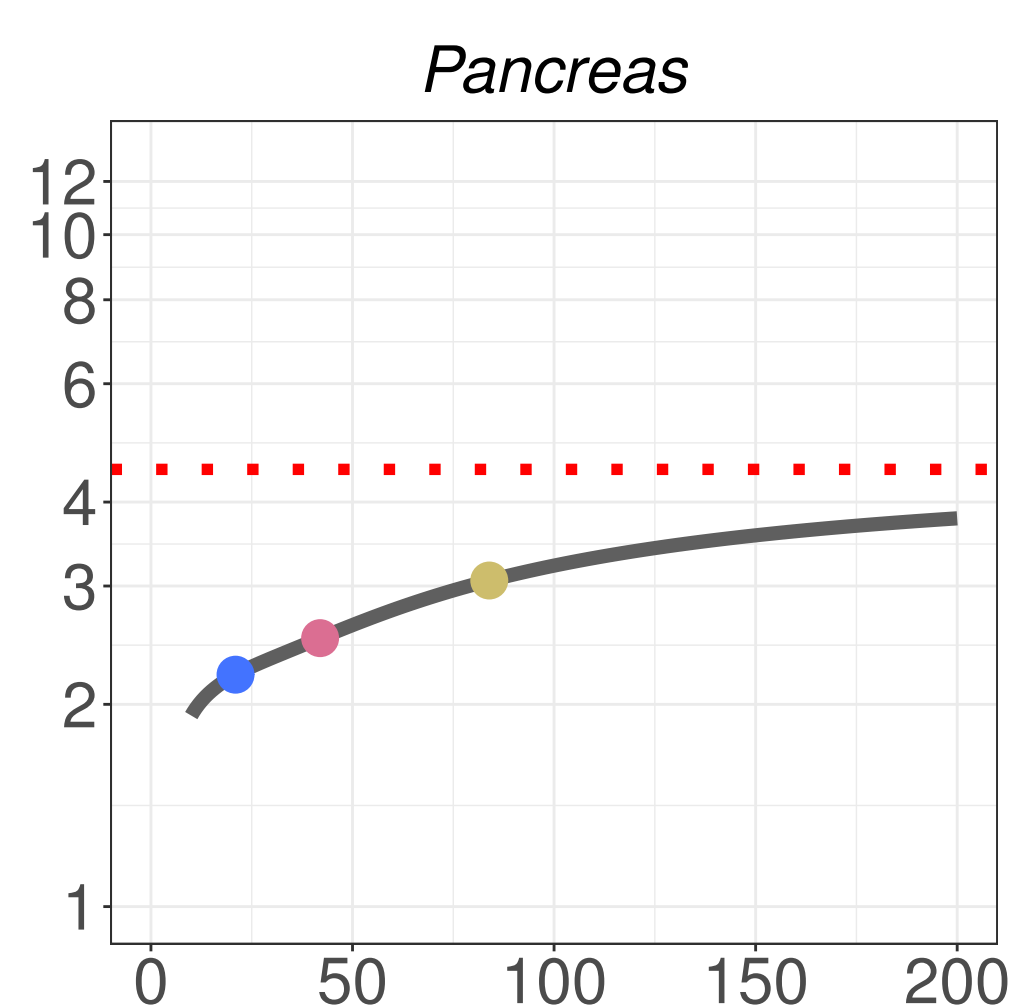
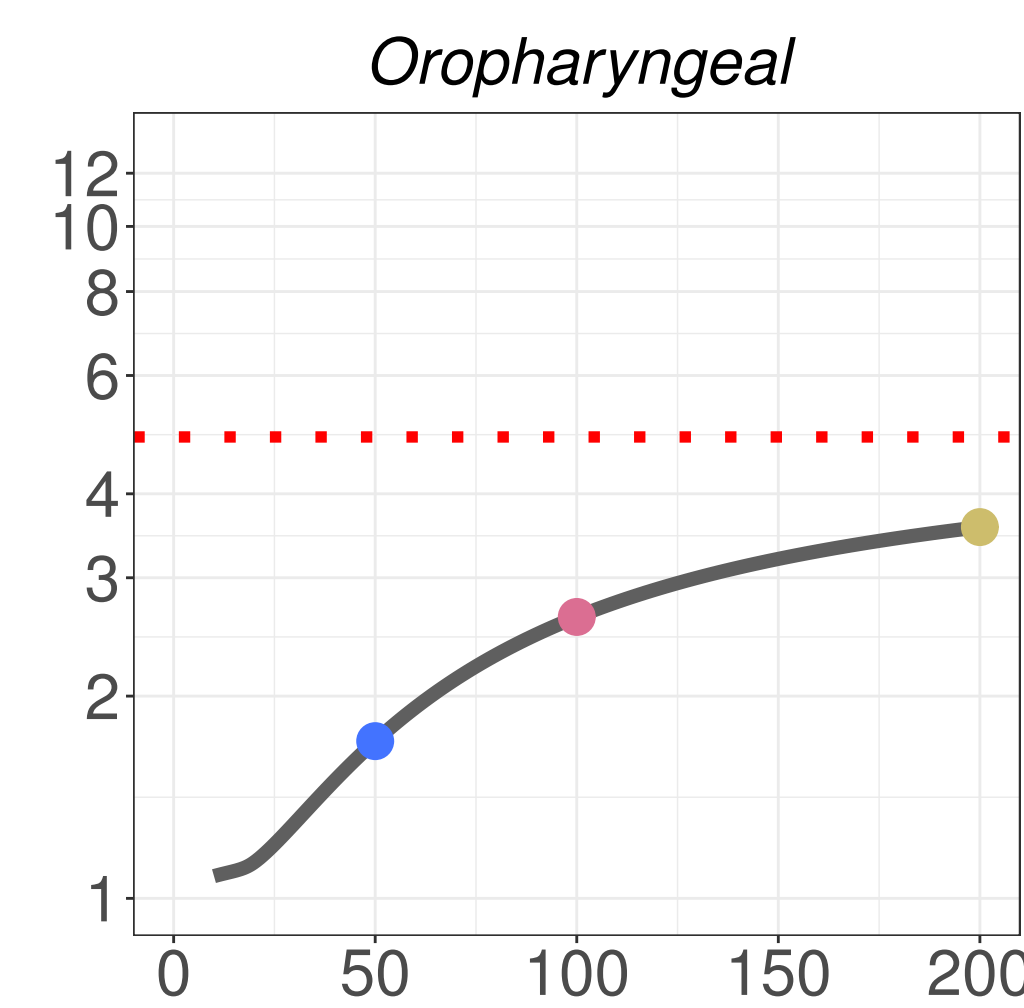
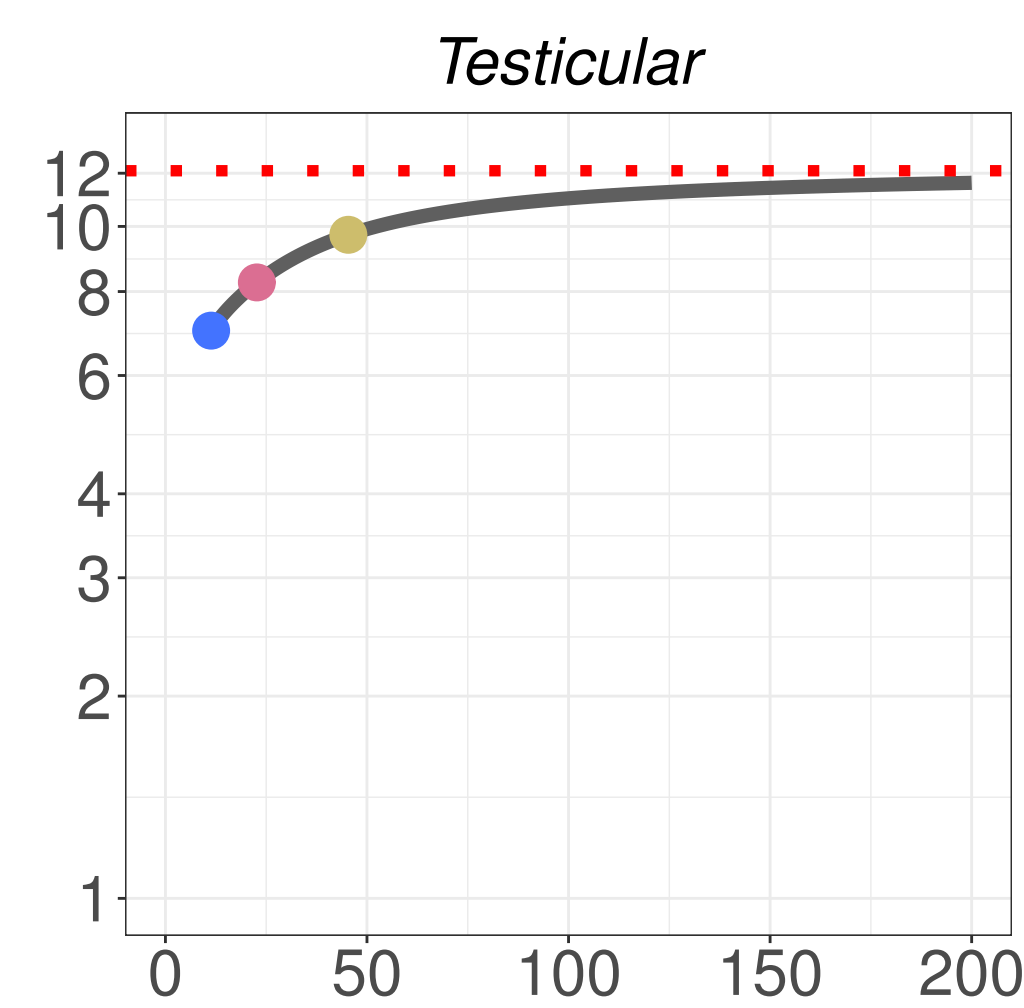
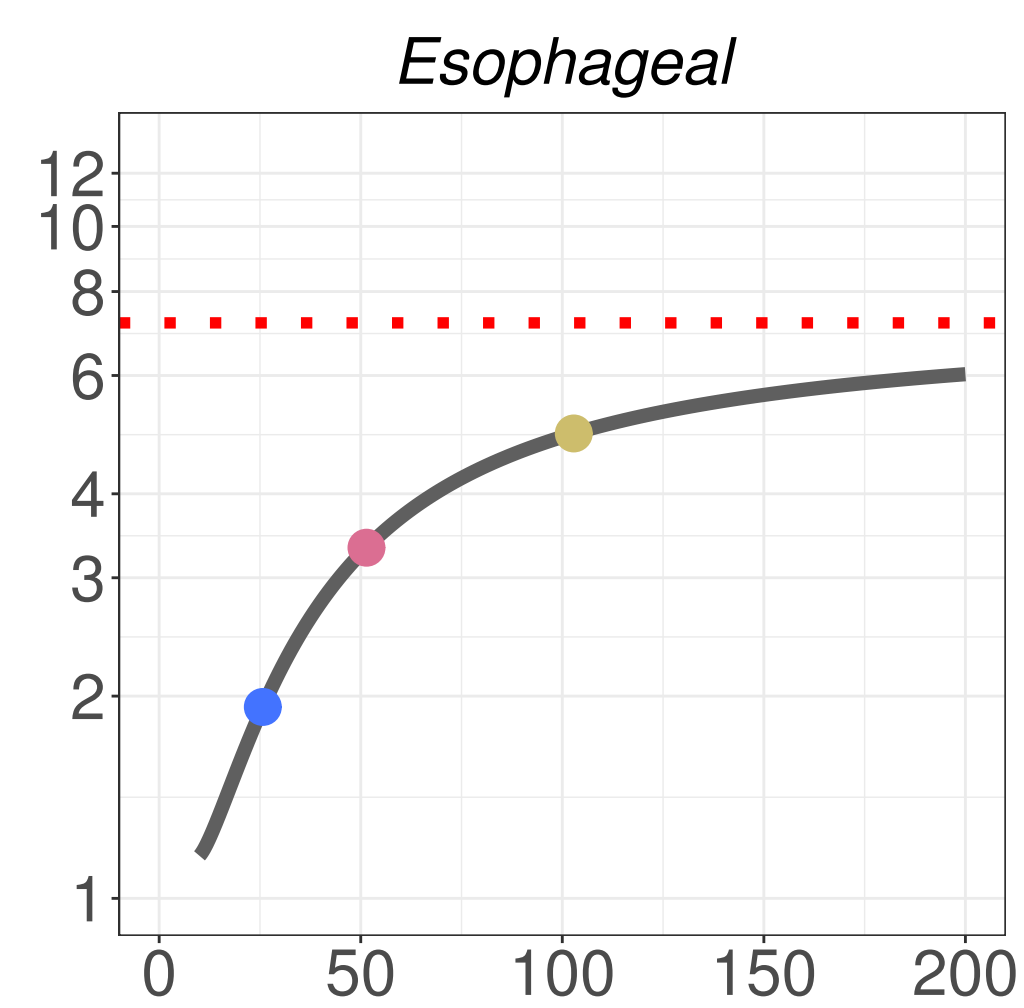
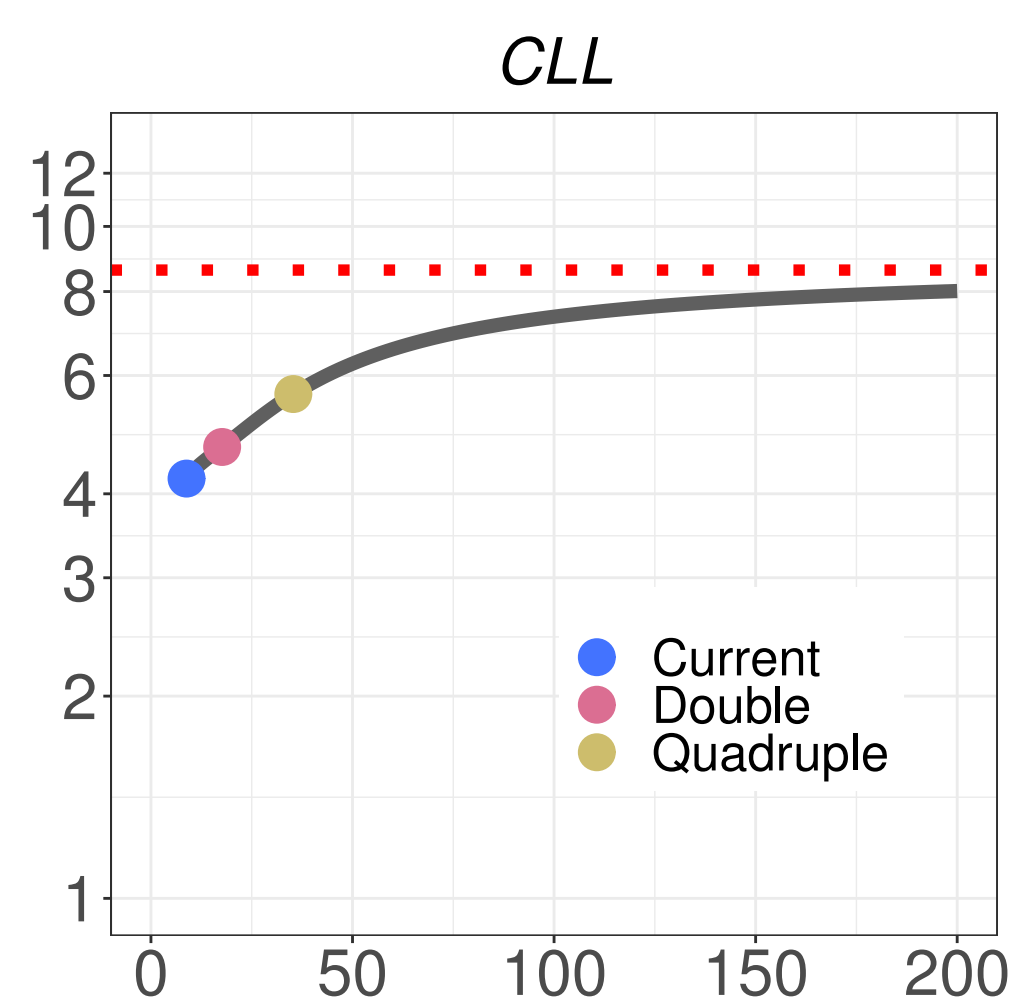


CLL*Esophageal**Testicular**Oropharyngeal**Pancreas**Renal**Glioma**Melanoma**Colorectal**Endometrial**Ovarian**Lung**Prostate**Breast*

y: AUC associated with the PRS

x: Total sample size assuming 1:1 case:control ratio (in thousands)

y: Relative risk for people at 99th centile compared to average risk of the population



x: Total sample size assuming 1:1 case:control ratio (in thousands)

