

Integrating multi-omics data to identify tissue-specific DNA methylation biomarkers for cancer risk

Yaohua Yang^{1†*}, Yaxin Chen^{2†}, Shuai Xu³, Xingyi Guo³, Guochong Jia³, Jie Ping³, Xiang Shu⁴, Tianying Zhao³, Fangcheng Yuan³, Gang Wang², Yufang Xie², Hang Ci², Hongmo Liu², Yawen Qi², Yongjun Liu⁵, Dan Liu², Weimin Li², Fei Ye⁶, Xiao-Ou Shu³, Wei Zheng³, Li Li⁷, Qiuyin Cai^{3*}, Jirong Long^{3*}

* Corresponding authors contact information:

Dr. Yaohua Yang, Center for Public Health Genomics, School of Medicine, University of Virginia. Address: 560 Ray C. Hunt Dr, Rm 4408, Charlottesville, VA 22903. Email:

vta8we@virginia.edu

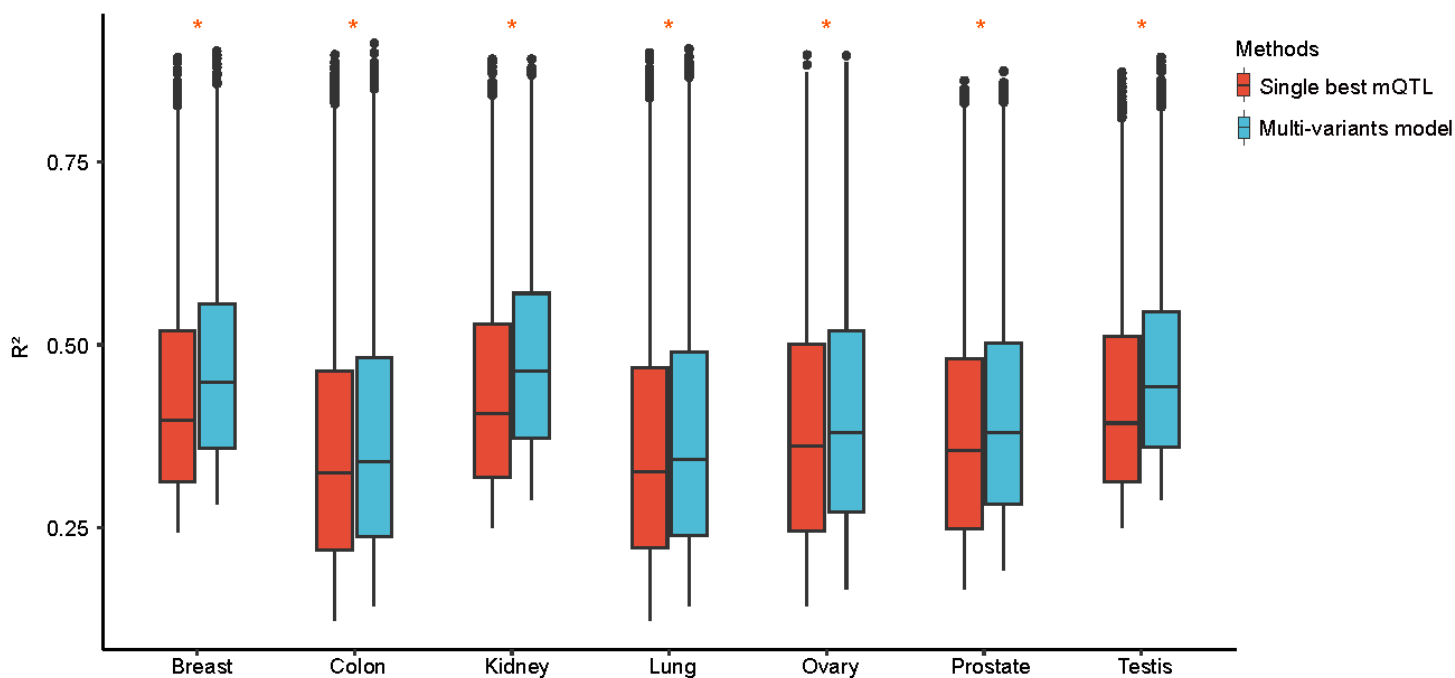
Dr. Qiuyin Cai, Division of Epidemiology, Department of Medicine, Vanderbilt University Medical Center. Address: 1161 21st Avenue South, MCN B-2104, Nashville, TN, 37232.

Email: qiuyin.cai@vanderbilt.edu

Dr. Jirong Long, Division of Epidemiology, Department of Medicine, Vanderbilt University Medical Center. Address: 2525 West End Avenue, Suite 800, Nashville, TN, 37203. Email:

jirong.long@vumc.org

† These authors contribute equally to this work



Supplementary Fig. 1. Prediction performance of models established by the present

study and mQTLs. CpGs that could be reliably predicted ($R > 0.10$ and $P < 0.05$) by both our

prediction models and mQTLs were used for comparison. The number of such CpGs is

$n = 24,604$ for breast tissue, $n = 85,285$ for colon tissue, $n = 31,296$ for kidney tissue, $n = 87,987$

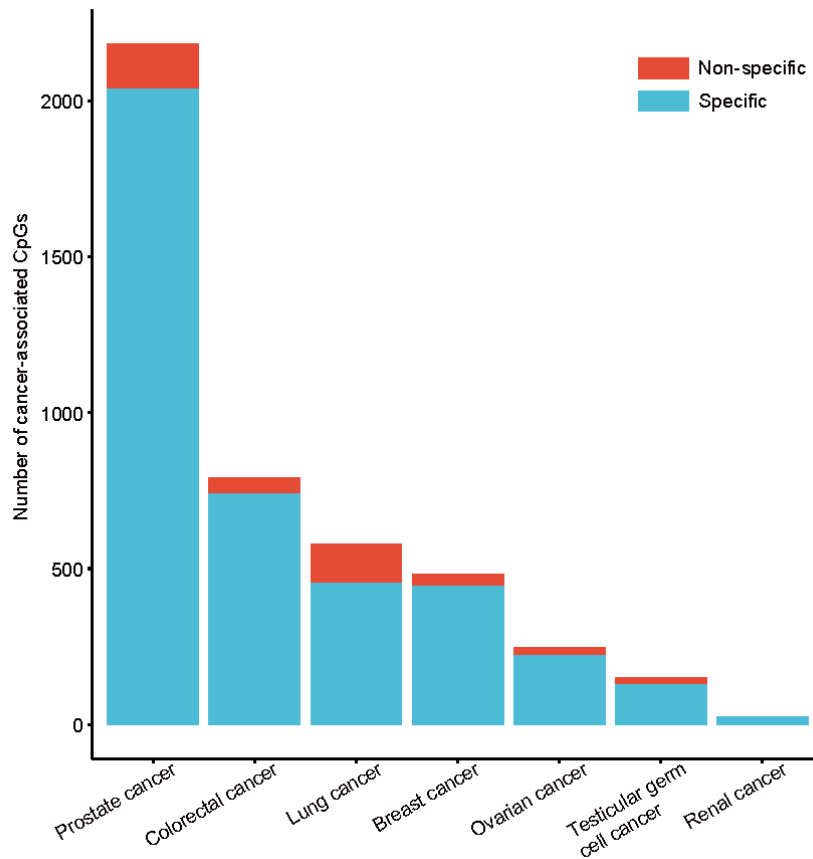
for lung tissue, $n = 76,161$ for ovary tissue, $n = 51,193$ for prostate tissue, and $n = 24,297$ for

testis tissue. Boxes represent the interquartile range, black bars are medians, whiskers extend

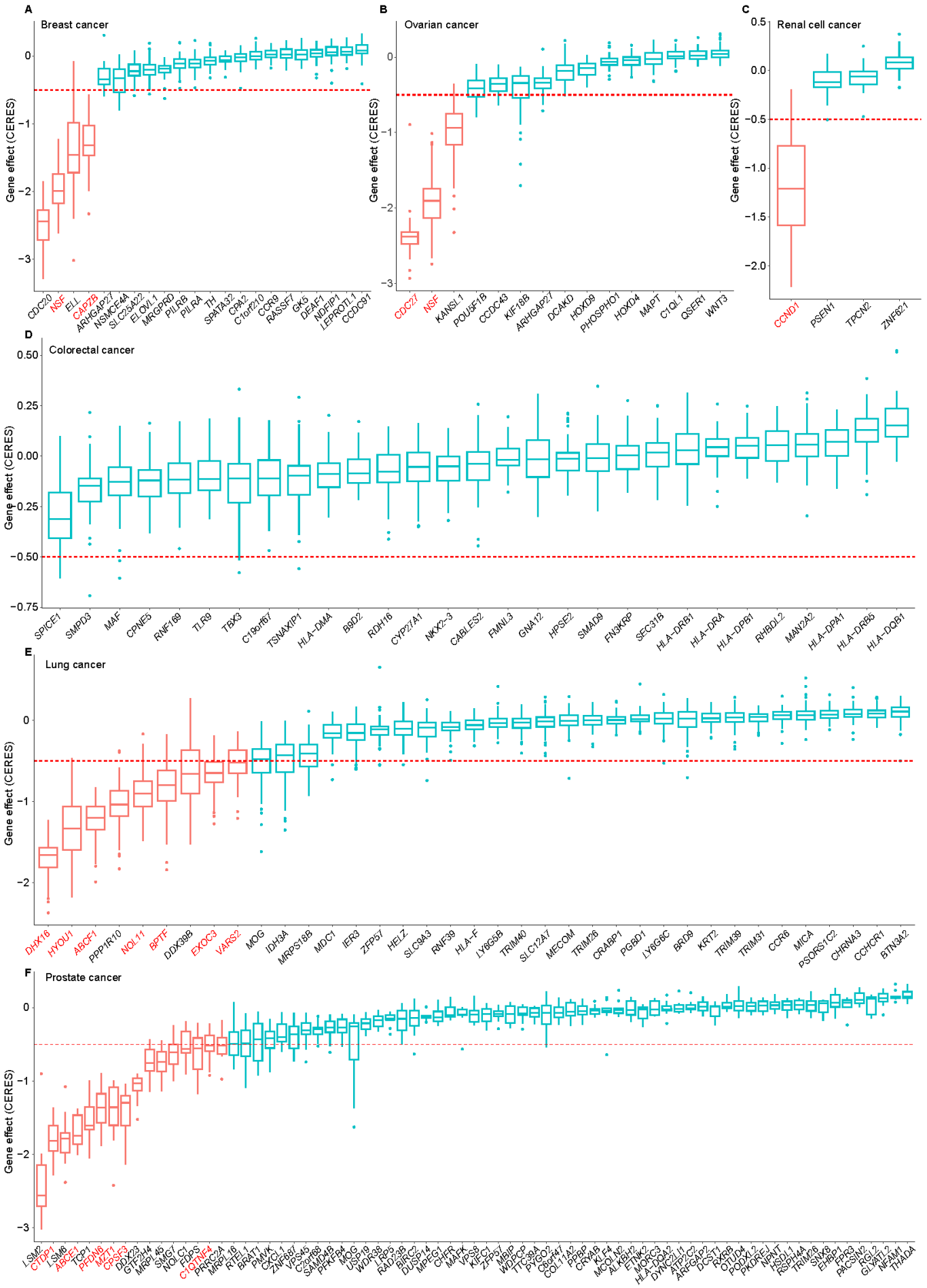
at most 1.5 times the interquartile range, and outliers are shown individually. Red asterisk

denotes $P \approx 0$ in two-sided paired t-test. mQTL, methylation quantitative trait loci. Source data

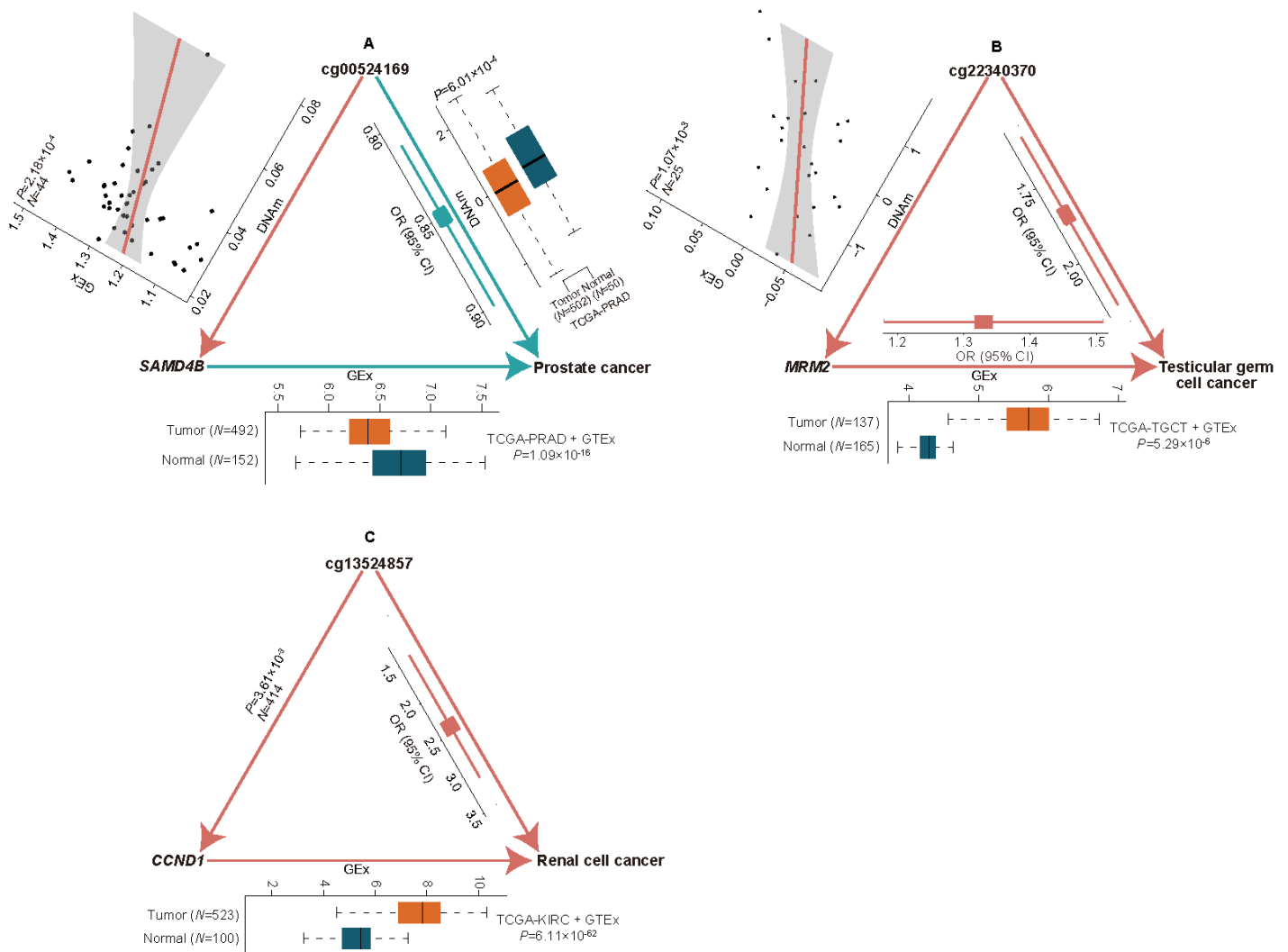
are provided as a Source Data file.



Supplementary Fig. 2. Tissue-specificity of cancer-associated CpGs. For each cancer, association analyses of genetically predicted DNA methylation (DNAm) with cancer risk were performed using SPrediXcan and significant associations were identified at Bonferroni-corrected two-side $P < 0.05$, corresponding to 4.93×10^{-7} for breast cancer (n=101,497 CpGs tested; n=539,198 subjects), 2.53×10^{-7} for colorectal cancer (n=197,947 CpGs tested; n=254,791 subjects), 3.98×10^{-7} for renal cell cancer (n=125,745 CpGs tested; n=31,190 subjects), 2.55×10^{-7} for lung cancer (n=195,764 CpGs tested; n=887,170 subjects), 2.66×10^{-7} for ovarian cancer (n=187,911 CpGs tested; n=70,668 subjects), 3.28×10^{-7} for prostate cancer (n=152,341 CpGs tested; n=944,762 subjects), and 4.22×10^{-7} for testicular germ cell cancer (n=118,568 CpGs tested; n=28,135 subjects). CpGs that were exclusively associated with a particular cancer at Bonferroni-corrected $P < 0.05$ are considered tissue specific. Source data are provided as a Source Data file.



Supplementary Fig. 3. Essential roles in cell proliferation of genes associated with cancer-associated CpGs. Boxplots show effects of genes on cell proliferation using experimental data from Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR) screens. Boxes represent the interquartile range, black bars are medians, whiskers extend at most 1.5 times the interquartile range, and outliers are shown individually. The dashed red line in each plot denotes the essentiality threshold of median CERES values < -0.50 . Genes shown in CpG-gene-cancer trios revealed by downstream analyses (see **Methods**) are highlighted in red. **A**, breast cancer (n=48 cell lines). **B**, ovarian cancer (n=57 cell lines). **C**, renal cell cancer (n=32 cell lines). **D**, colorectal cancer (n=57 cell lines). **E**, lung cancer (n=114 cell lines). **F**, prostate cancer (n=10 cell lines). Source data are provided as a Source Data file.



Supplementary Fig. 4. Additional examples of CpG-gene-cancer trios suggesting DNA methylation influencing cancer risk by modulating *cis*-gene expression. Association analyses of genetically predicted DNA methylation (DNAm) or gene expression (GEx) with cancer risk were performed using SPrediXcan. Differential DNAm or GEx analyses were conducted using linear mixed models. Association analyses between DNAm and GEx were performed using linear regression. Red arrows, lines, and blocks denote positive associations, while green ones denote negative associations. All statistical tests were two-sided and multiple comparisons were Bonferroni- or false discovery rate (FDR)-adjusted. The odds ratio (OR) and 95% confidence interval (CI) for cancer risk per standard deviation (SD) increase in genetically predicted DNAm or GEx are displayed as a block with error bands. In

boxplots, boxes represent the interquartile range, black bars are medians, and whiskers extend at most 1.5 times the interquartile range. In the scatter plot displaying the association between DNAm and GEx, directly measured DNAm and GEx values after quantile- and inverse-normalization are presented. n, number of samples, TCGA, The Cancer Genome Atlas; GTEx, Gene-Tissue Expression consortium; PRAD, prostate adenocarcinoma; TGCT, testicular germ cell tumors; KIRC, kidney renal clear cell carcinoma. **A**, DNAm at cg00524169 may decrease prostate cancer risk by promoting the expression of *SAMD4B*. The sample size was 552 for tumor-normal differential DNAm analyses, 44 for DNAm-GEx correlation analyses, 644 for tumor-normal differential GEx analyses, and 944,762 for association analyses of predicted DNAm and GEx with prostate cancer risk, respectively. **B**, DNAm at cg22340370 may elevate testicular germ cell cancer risk by promoting the expression of *MRM2*. The sample size was 25 for DNAm-GEx correlation analyses, 302 for tumor-normal differential GEx analyses, and 28,135 for association analyses of predicted DNAm and GEx with testicular germ cell cancer risk, respectively. **C**, DNAm at cg13524857 may increase renal cell cancer risk by promoting the expression of *CCND1*. The sample size was 414 for DNAm-GEx correlation analyses, 623 for tumor-normal differential GEx analyses, and 31,190 for association analyses of predicted DNAm and GEx with renal cell cancer risk, respectively. Source data are provided as a Source Data file.