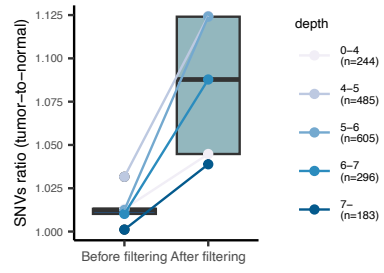
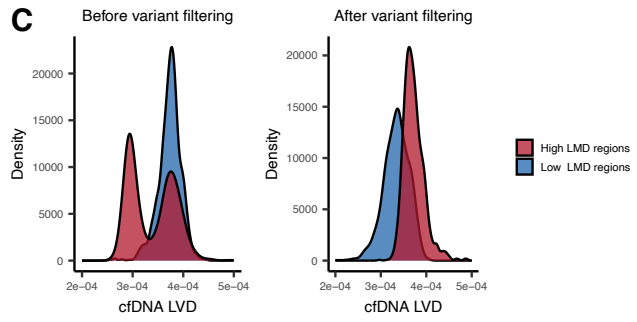
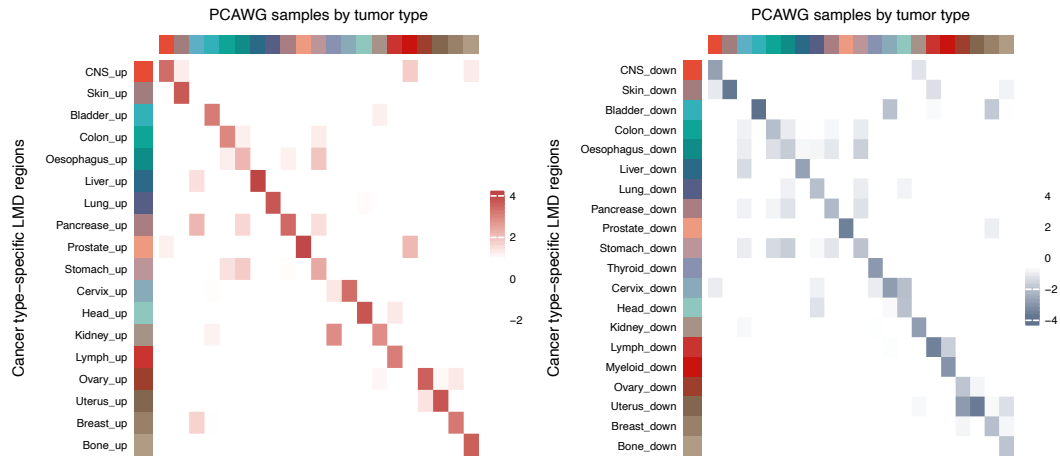
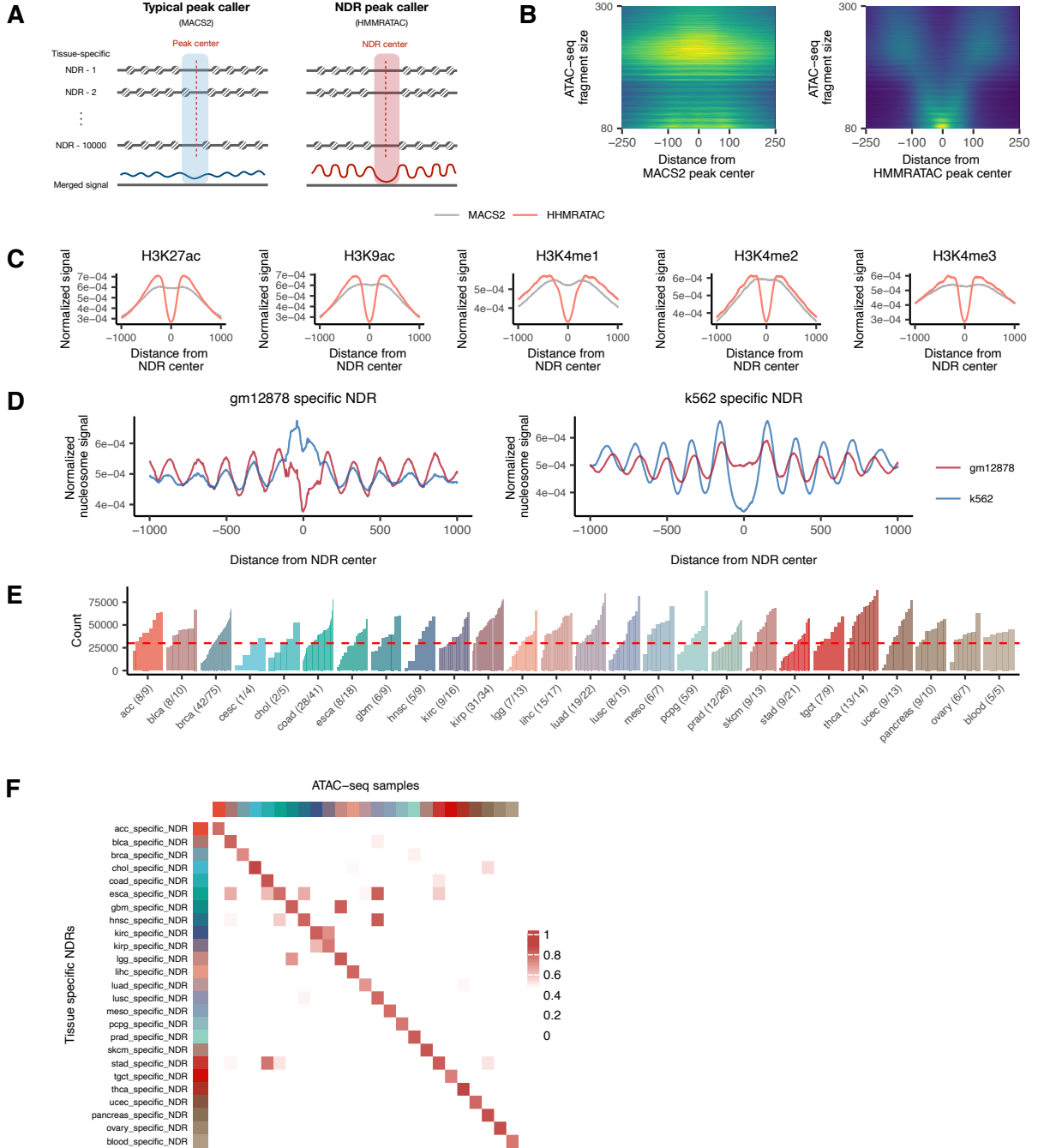


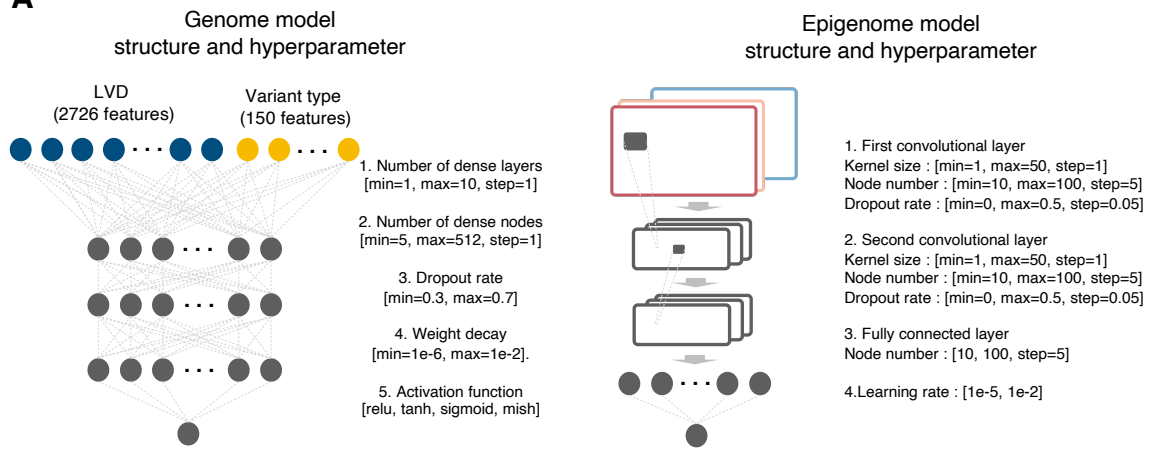
A**C****B**

Supplementary Fig. 1. Supporting data for genome model development. (A) The ratio of single nucleotide variants detected in cancer versus normal cfDNA samples before and after applying our variant filtering pipeline. Variant calling was performed at differing read depths. Each box indicates IQR and median, whiskers indicates 1.5 x IQR, black dots indicates outlier. (B) Heatmaps showing the average LMD values in high LMD regions (left) and low LMD regions (right), specifically identified for each cancer type. The color intensity is scaled to show the degree of relative mutation enrichment (left) or mutation depletion (left). A total of 2,754 PCAWG samples were used for this analysis. (C) Distribution of the LVD in cfDNA in the above identified high or low LMD regions in the matching cancer type before and after applying our variant filtering pipeline. Source data are provided as a Source Data file.

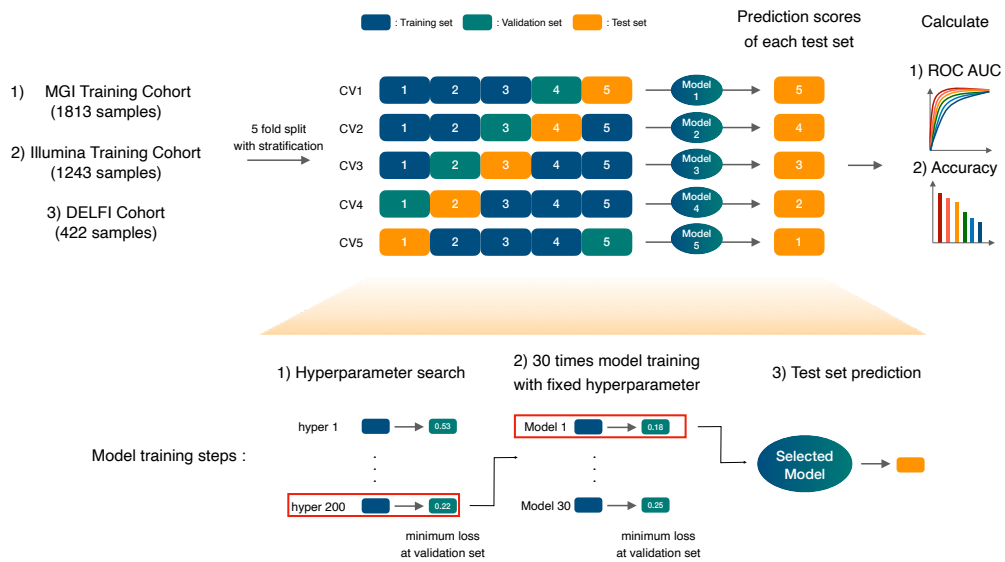


Supplementary Fig. 2. Supporting data for epigenome model development. (A) Illustration of peak call process for typical peak caller (MACS2) and NDR peak caller (HMMRATAC). (B) V-plot image constructed for the NDRs identified by ATAC-seq peak calling by MACS2 (left) and HMMRATAC (right). Data from the GM12878 cell line are shown. (C) ChIP-seq signals of histone modification (H3K27ac, H3K9ac, H3K4me1, H3K4me2, and H3K4me3) centered on the NDRs identified by HMMRATAC (red) and MACS2 (gray). (D) Nucleosome occupancy at GM12878-specific (left) and K562-specific (right) NDRs in GM12878 cell line (red) and in K562 cell line (blue). (E) The number of the NDRs identified by HMMRATAC in 431 samples grouped into 23 tumor tissues, 2 normal tissues, or PBMCs. The red horizontal dashed line indicates the peak count threshold for sample filtering. The filtered and total number of samples is denoted in parentheses. (F) Heatmap showing the average peak scores of the NDRs specific to each tissue group. The color intensity is proportional to the peak score. Source data are provided as a Source Data file.

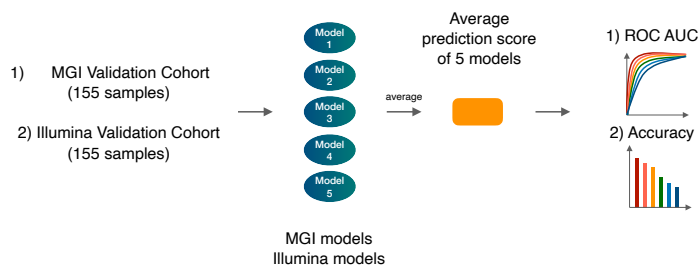
A



B

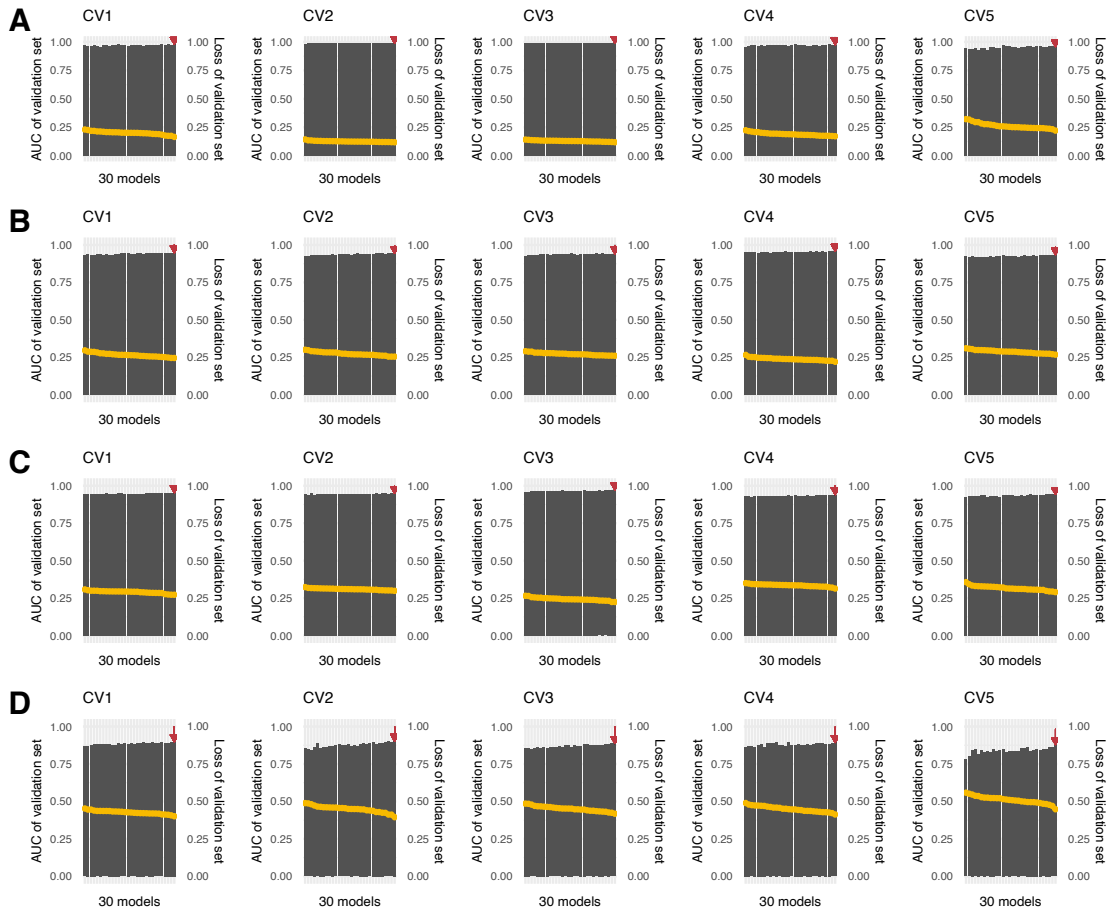


C



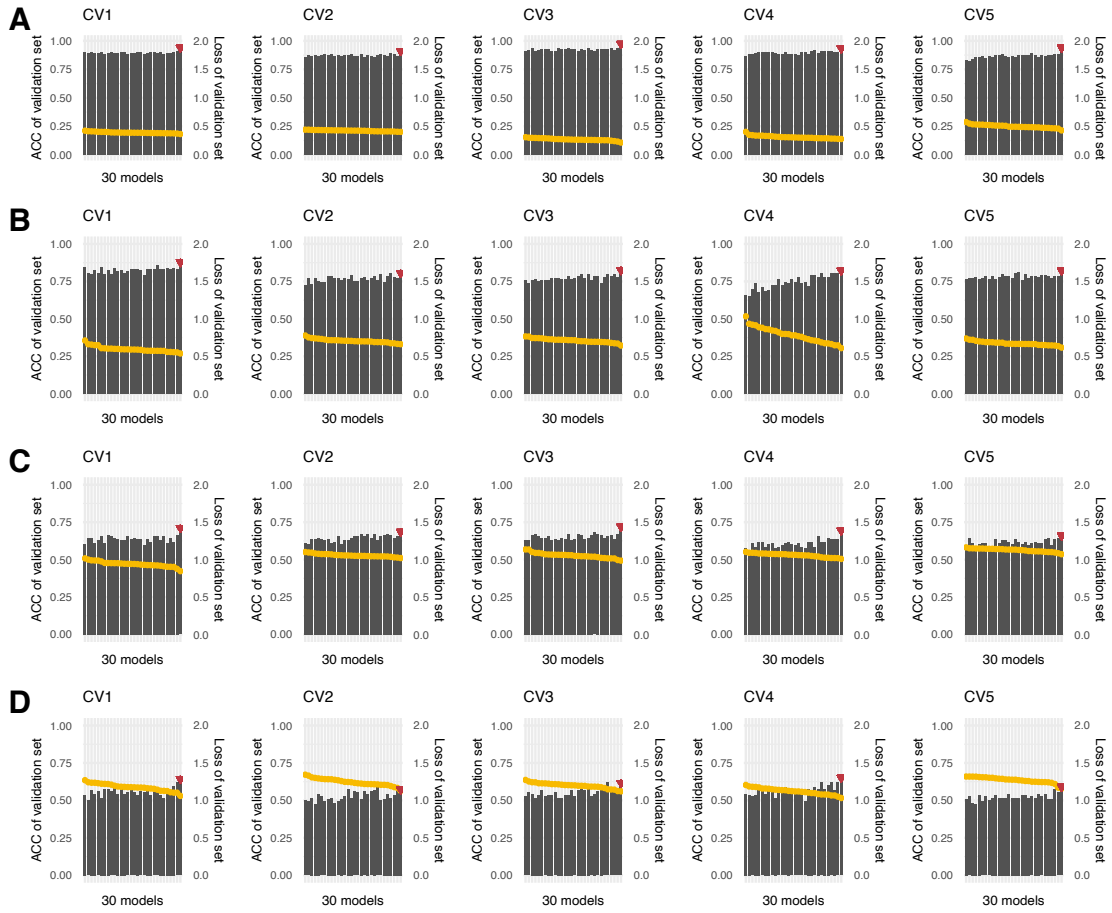
Supplementary Fig. 3. Schematic of the training and prediction of genome and epigenome models. (A) Illustration of the genome and epigenome model structure and hyperparameter space. (B) Schematic of the model training process using the training cohorts. Stratified five-fold cross-validation method was used to train and evaluate the models for each training cohort. At each split, the best hyperparameter that minimizes the validation loss was selected. Using the fixed hyperparameter, the model was trained 30 times and the model with lowest validation loss was chosen as the final model. (C) Schematic of the external prediction using the validation cohorts. The validation cohorts were predicted using the five models trained with the training cohort. The average of five prediction values was used as the final prediction score.

● Loss of validation set ■ AUC of validation set

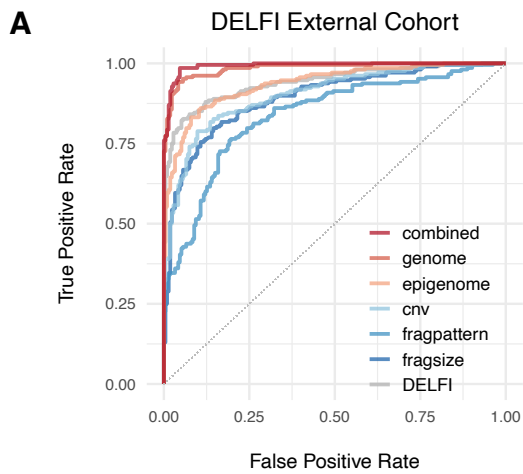


Supplementary Fig. 4. Validation performance of the cancer detection models trained with fixed hyperparameter in 30 different random states. (A-D) Validation loss and ROC-AUC of (A) the genome model of the MGI training cohort, (B) the epigenome model of the MGI training cohort, (C) the genome model of the Illumina training cohort, (D) the epigenome model of the Illumina training cohort. Bar plots (gray) represent ROC AUC of the validation set. Line plot (yellow) represent loss of the validation set. The red arrows indicate the final model with the minimum loss at the validation set. AUC, area under the curve; CV, cross validation. Source data are provided as a Source Data file.

● Loss of validation set ■ ACC of validation set

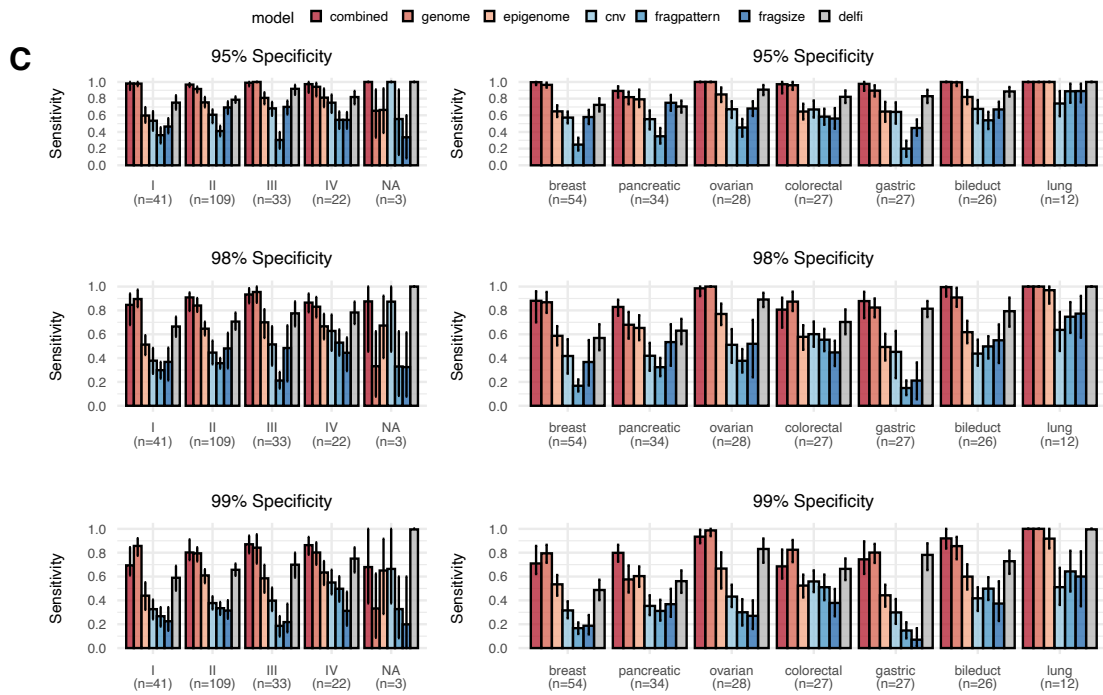


Supplementary Fig. 5. Validation performance of the cancer localization models trained with fixed hyperparameter in 30 different random states. (A-D) Validation accuracy and loss of (A) the genome model of the MGI training cohort, (B) the epigenome model of the MGI training cohort, (C) the genome model of the Illumina training cohort and (D) the epigenome model of the Illumina training cohort. Bar plots (gray) represent accuracy of the validation set. Line plot (yellow) represent loss of the validation set. The red arrows indicate the final model with the minimum loss at the validation set. ACC, accuracy; CV, cross-validation. Source data are provided as a Source Data file.

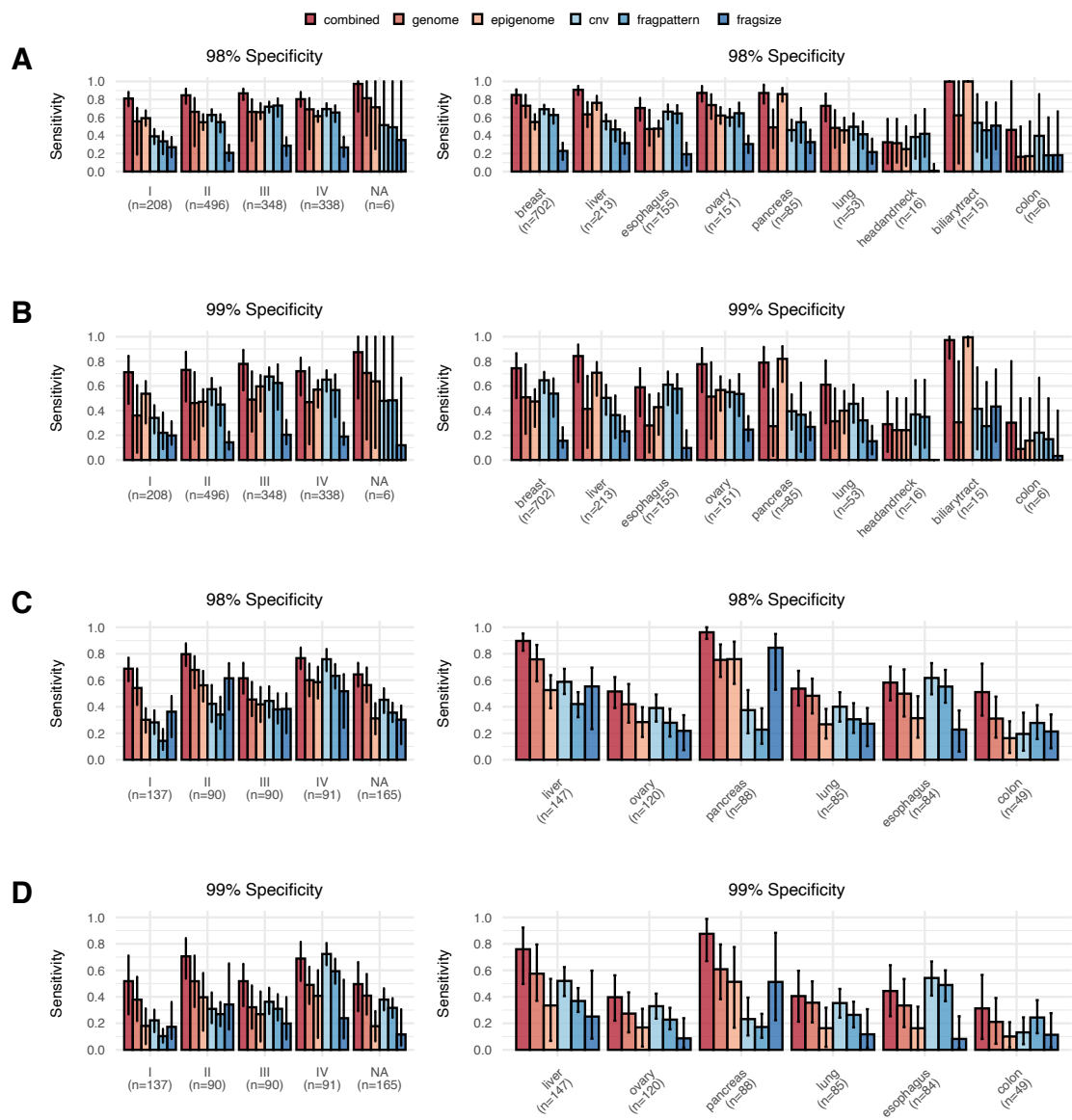


B

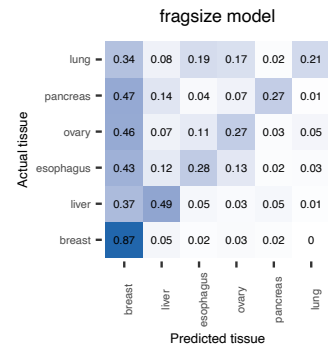
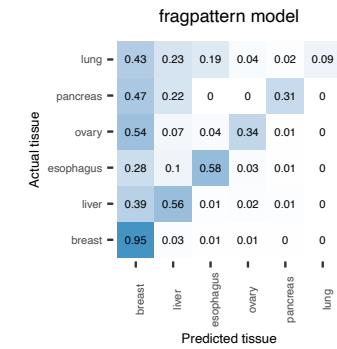
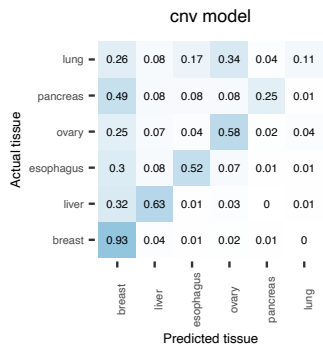
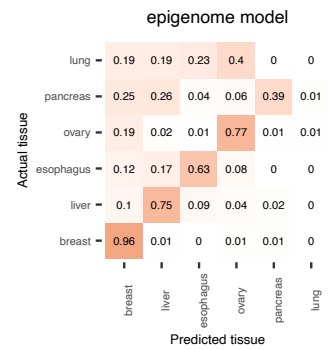
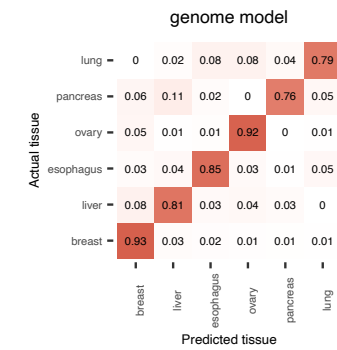
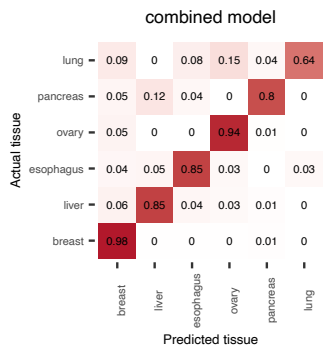
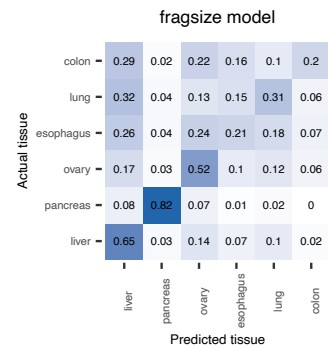
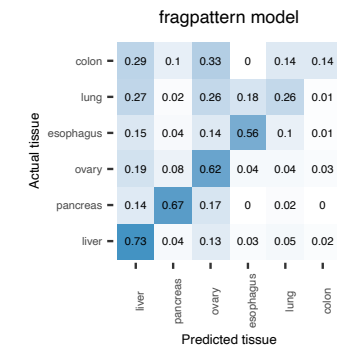
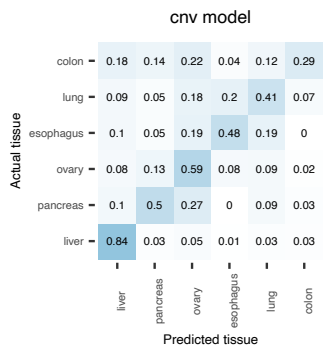
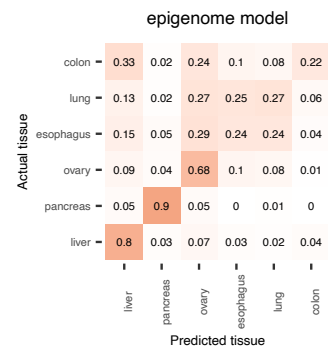
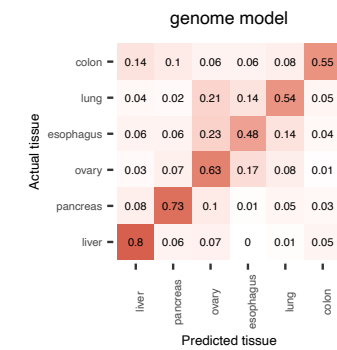
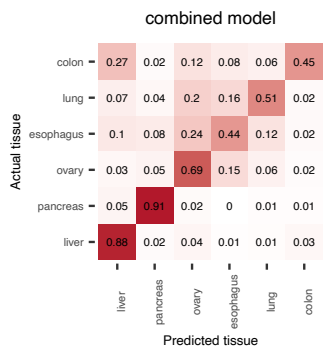
DELFI	External Cohort	
	auc	95% CI
combined	0.993	0.987-0.998
genome	0.986	0.976-0.993
epigenome	0.937	0.913-0.957
cnv	0.903	0.875-0.932
fragpattern	0.838	0.798-0.874
fragsize	0.893	0.863-0.923
DELFI	0.941	0.916-0.962



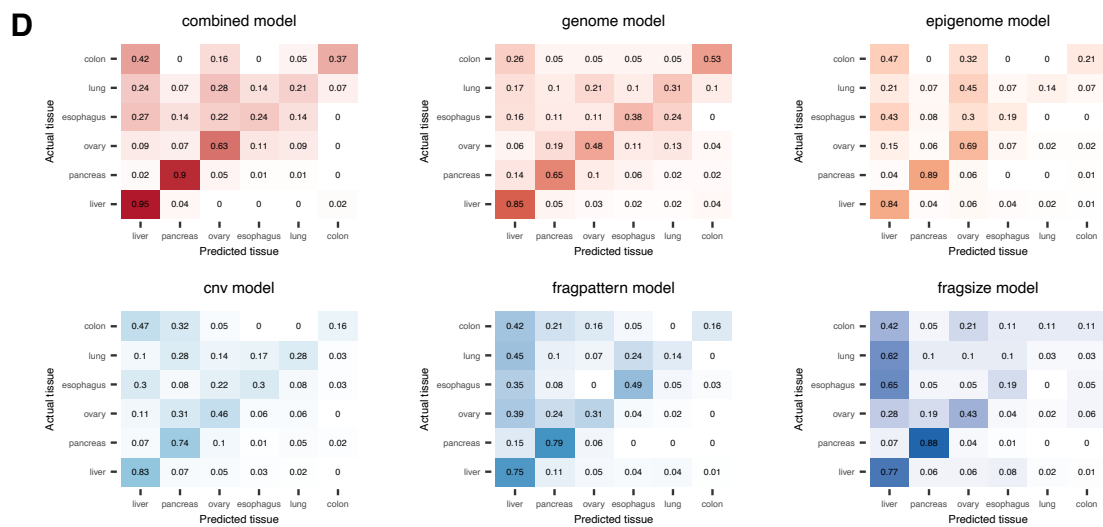
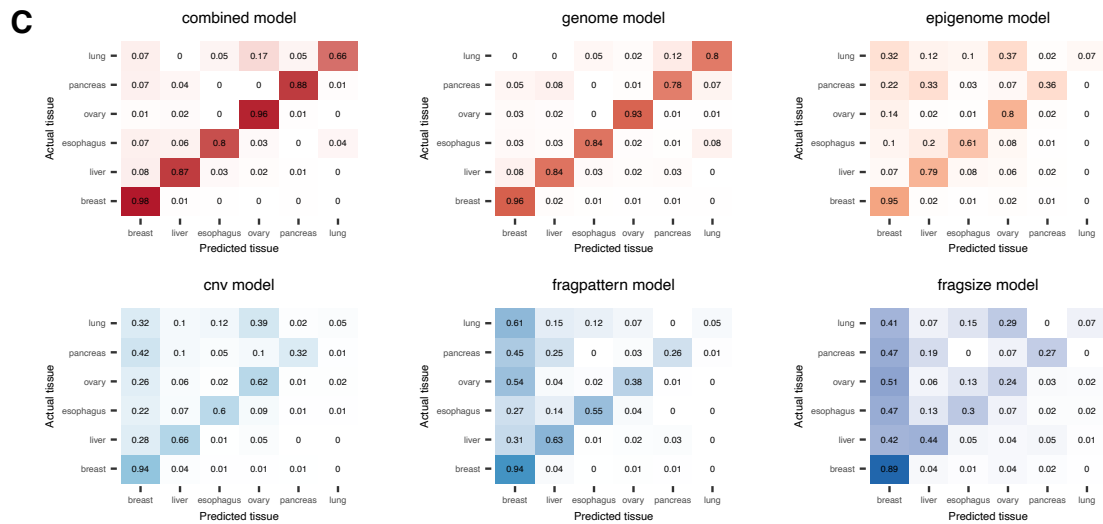
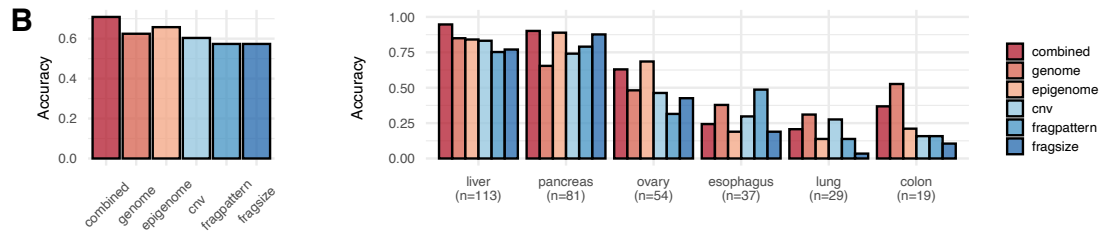
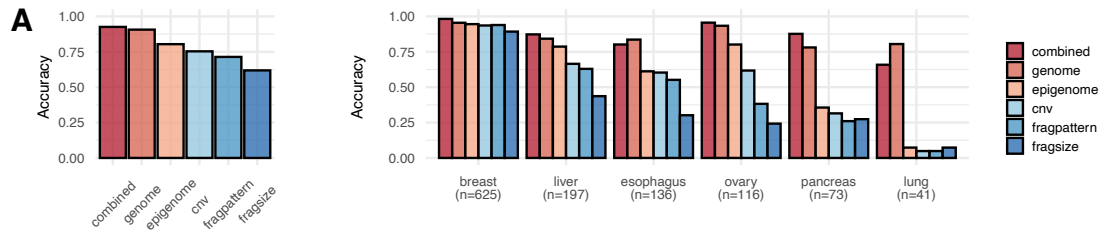
Supplementary Fig. 6. Performance of cancer detection on an external cohort. (A) ROC curve and (B) ROC-AUC table providing the 95% confidence interval for different models on the DELFI dataset⁹. A total of 208 cancer and 214 normal control samples were used. (C) Sensitivity values with the 95% confidence interval at 95%, 98%, and 99% specificity broken down by the tumor stage (left) and cancer type (right) for the DELFI dataset⁹. Confidence interval for sensitivity value was calculated from 1,000 bootstrapping samplings. (A-C) Our genome, epigenome, and combined models were compared with predictions based on fragmentation patterns⁹ (fragpattern), fragment size profiles⁸ (fragsize), and copy number variations⁵ (cnv). The gray ROC curve and bar graphs marked as DELFI correspond to the score provided by the authors⁹ combining fragmentation with other features. auc, area under the curve; NA, stage information not available; CI, Confidence interval. Source data are provided as a Source Data file.



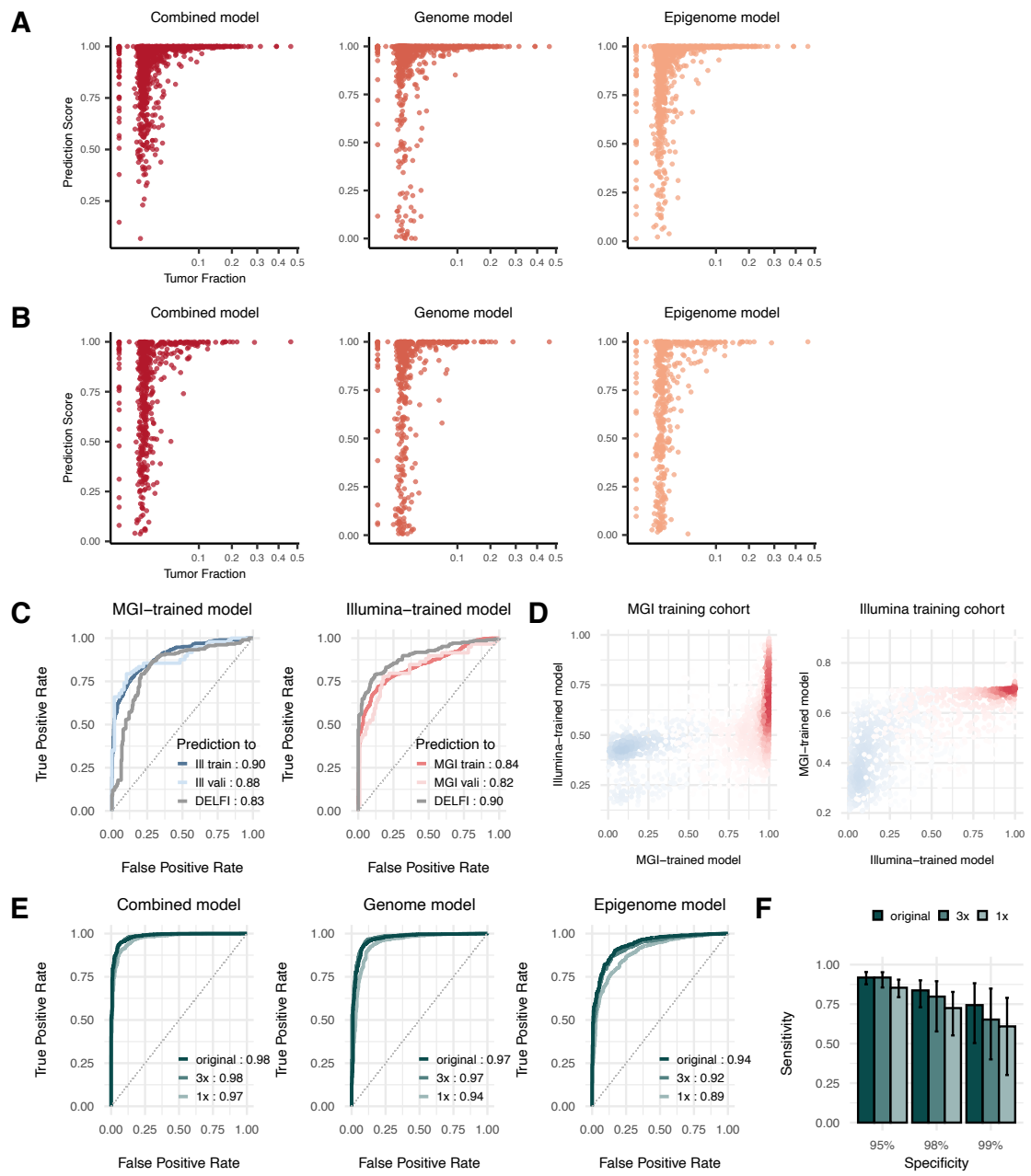
Supplementary Fig. 7. Sensitivity of cancer detection. (A-B) Sensitivity values with the 95% confidence interval at (A) 98% and (B) 99% specificity broken down by the tumor stage (left) and cancer type (right) on the MGI training cohort. (C-D) Sensitivity values with the 95% confidence interval at (C) 98% and (D) 99% specificity broken down by the tumor stage (left) and cancer type (right) on the Illumina training cohort. Confidence interval for sensitivity value was calculated from 1,000 bootstrapping samplings. Our genome, epigenome, and combined models were compared with predictions based on fragmentation patterns⁹ (fragpattern), fragment size profiles⁸ (fragsize), and copy number variations⁵ (cnv). NA, stage information not available. Source data are provided as a Source Data file.

A**B**

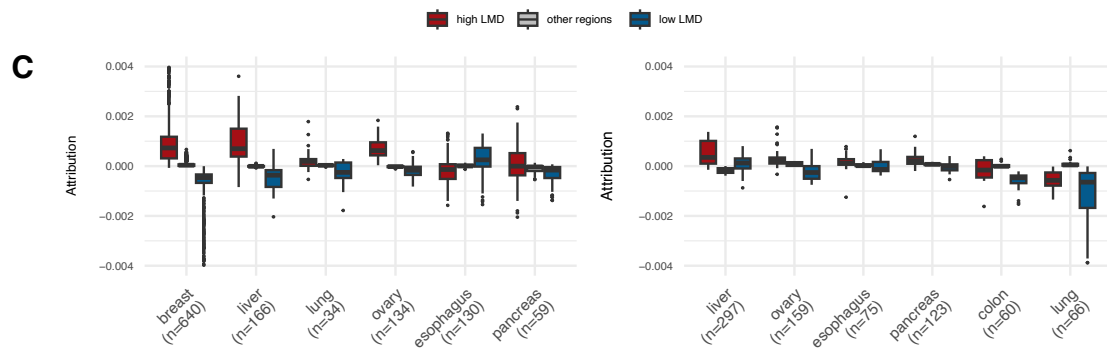
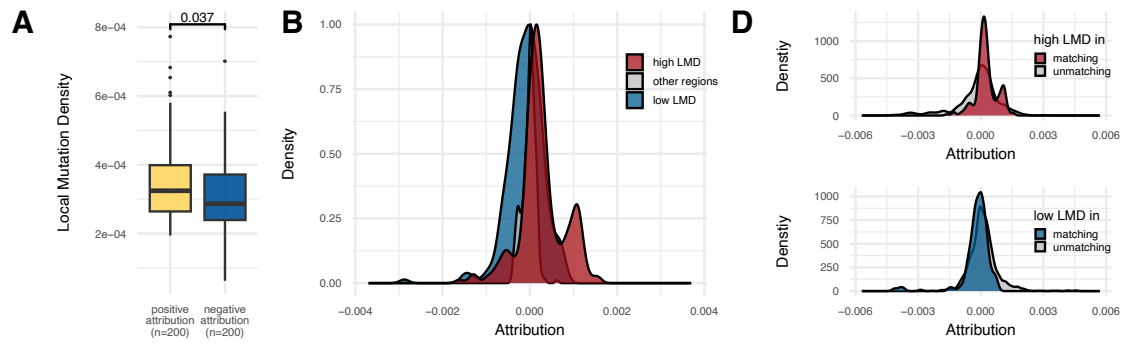
Supplementary Fig. 8. Confusion matrix for tissue-of-origin localization on the (A) MGI training cohort and (B) Illumina training cohort. The y-axis represents the actual site, and the x-axis represents the predicted site. The numbers in the cells of the matrix represent the proportion of samples of each cancer type localized to respective tumor sites. Our genome, epigenome, and combined models were compared with predictions based on fragmentation patterns⁹ (fragpattern), fragment size profiles⁸ (fragsize), and copy number variations³ (cnv). Source data are provided as a Source Data file.



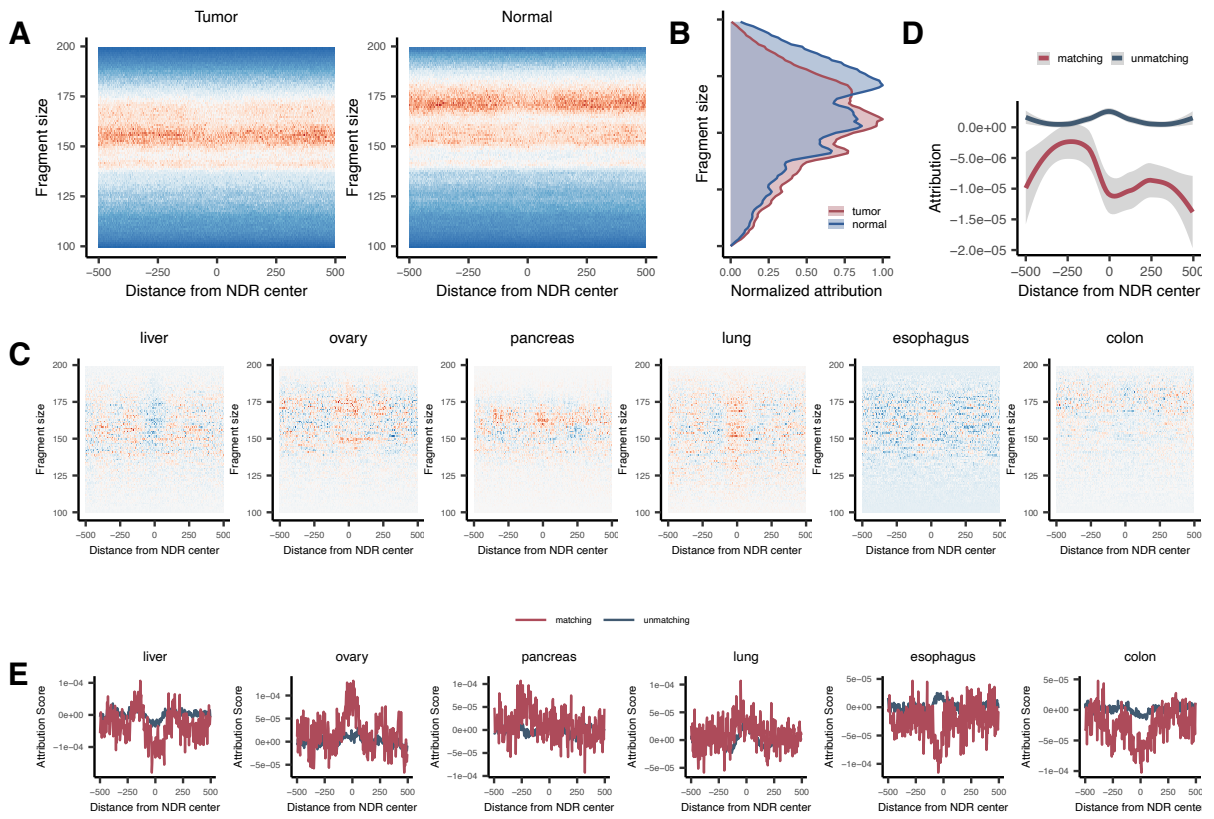
Supplementary Fig. 9. Performance of tissue-of-origin localization using samples predicted as cancer. (A) Average accuracy (left) and accuracy (right) on each cancer type for different models on the MGI training cohort. (B) Average accuracy (left) and accuracy (right) on each cancer type for different models on the Illumina training cohort. (C-D) Confusion matrix for localization using the combined model on the (C) MGI training cohort and (D) Illumina training cohort. The y-axis represents the actual site, and the x-axis represents the predicted site. The numbers in the cells of the matrix represent the proportion of samples of each cancer type localized to respective tumor sites. (A-D) Our genome, epigenome, and combined models were compared with predictions based on fragmentation patterns⁹ (fragpattern), fragment size profiles⁸ (fragsize), and copy number variations⁵ (cnv). Source data are provided as a Source Data file.



Supplementary Fig. 10. Additional analyses of our model. (A-B) Correlation between cfDNA tumor fraction and prediction score of combined, genome and epigenome model of the (A) MGI training cohort and (B) Illumina training cohort. (C-D) Combined models of the MGI training cohort and Illumina training cohort were used to predict the external cohort dataset. (C) ROC curve for the external cohort prediction. The combined model of the MGI training cohort was used to predict the Illumina training cohort, Illumina validation cohort and DELFI cohort (left). The combined model of the Illumina training cohort was used to predict the MGI training cohort, MGI validation cohort and DELFI cohort (right). (D) Density scatter plots which represent the relationship between the prediction score of the MGI trained model and the Illumina trained model on the MGI training cohort (left) and Illumina training cohort (right). The red dots indicate the density of cancer patient samples. The blue dots indicate the density of normal control samples. (E-F) Downsampled cfDNA WGS (3x, 1x) were compared with the original cfDNA WGS (5x) of the MGI training cohort. (E) ROC curve and (F) Sensitivity values with the 95% confidence Interval at 95%, 98%, 99% specificity for the different depth models of the MGI training cohort. Confidence interval for sensitivity value was calculated from 1,000 bootstrapping samplings. Illu train, Illumina training cohort; Illu vali, Illumina validation cohort; MGI train, MGI training cohort; MGI vali, MGI validation cohort. Source data are provided as a Source Data file.



Supplementary Fig. 11. Additional interpretation of the genome model. (A) LMD values obtained from tumor tissues in the LVD regions with positive or negative attribution for cancer samples in the genome model for cancer detection on the Illumina training cohort. *P* values from the Wilcoxon test are indicated (two-sided). Each box indicates IQR and median, whiskers indicates 1.5 x IQR, black dots indicates outlier. (B) Distribution of the attribution values assigned by the cancer localization genome model of the Illumina training cohort to high or low LMD regions in comparison to other regions of the PCAWG cancer type matching the given prediction label. (C) The same plot as Fig. 5B and Supplementary Fig. 11B, broken down by cancer types for the MGI training cohort (left) and Illumina training cohort (right). Box elements are same as (A). (D) Comparison of the attribution values assigned by the cancer localization genome model to high (upper) or low (lower) LMD regions of the PCAWG cancer types matching versus unmatching with the given prediction label. Source data are provided as a Source Data file.



Supplementary Fig. 12. Additional interpretation of the epigenome model. (A) Attribution values of the tissue-specific V-plots of the epigenome model for cancer detection on the Illumina training cohort. The average attribution values across the tissue-specific V-plots are compared for cancer samples (left) and normal samples (right). (B) Distribution of the normalized attribution values according to fragment size. (C) Attribution values mapped to the tissue-specific V-plots of the epigenome model for tissue-of origin localization on the Illumina training cohort. The average attribution values across the tissue-specific V-plots are shown for each cancer type. (D) Attribution values of the epigenome model for cancer localization on the Illumina training cohort, as averaged across the NDRs of the matching (red) versus unmatching (blue) cancer types, and plotted according to the distance from the NDR midpoints. Attribution values were smoothed using lowess regression. (E) The sample plot as (D) broken down by cancer types. Source data are provided as a Source Data file.