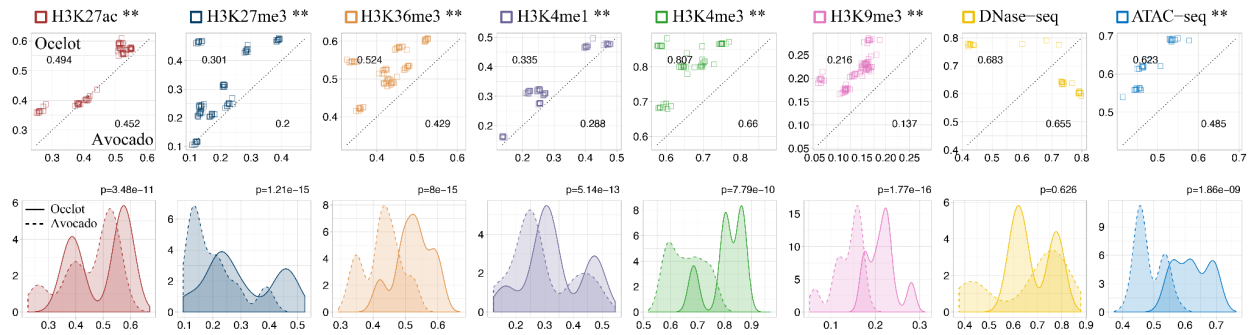**Asymmetric Predictive Relationships Across Histone Modifications**

Hongyang Li[1,*], Yuanfang Guan[1,*]

1. Department of Computational Medicine and Bioinformatics, University of Michigan, 100 Washtenaw Avenue, Ann Arbor, MI 48109, USA
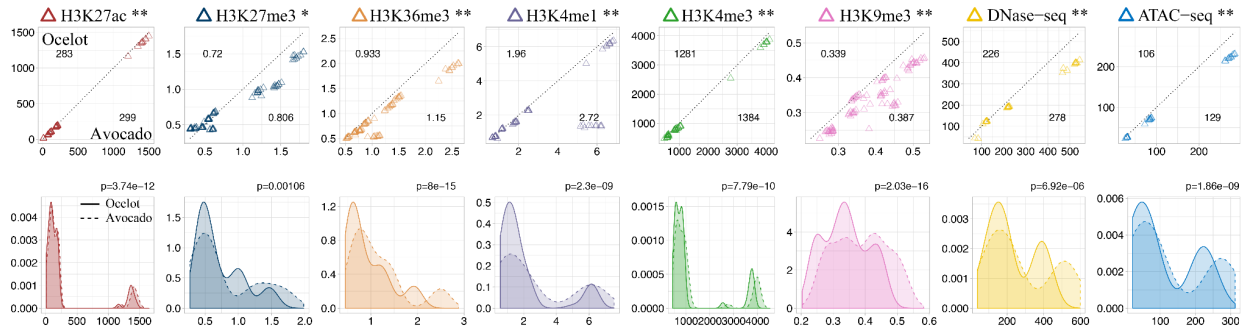
* Corresponding authors: hyangl@umich.edu or gyuanfan@umich.edu
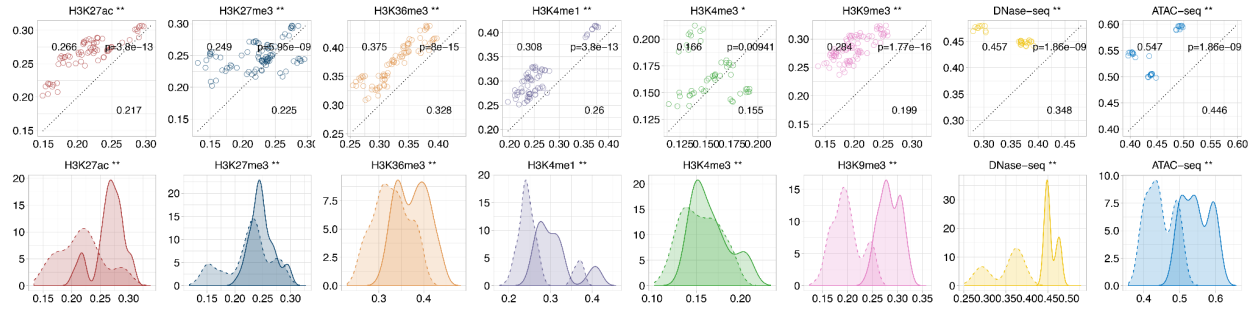
# Supplementary Figures



**Supplementary Figure 1: Benchmarking against Avocado using the global Pearson's correlation at the mark level.**

In the first row, we benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each symbol represents a testing sample. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001). In the second row, the overall distribution of this evaluation metric is shown as density plots, where the solid lines represent our method and the dashed lines represent Avocado.

**Supplementary Figure 2: Benchmarking against Avocado using the global mean squared error at the mark level.**

In the first row, we benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each symbol represents a testing sample. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001). In the second row, the overall distribution of this evaluation metric is shown as density plots, where the solid lines represent our method and the dashed lines represent Avocado.
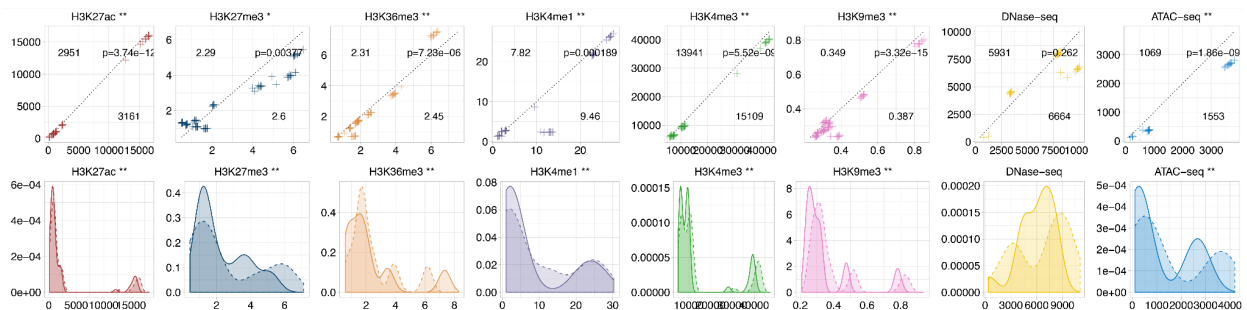
**Supplementary Figure 3: Benchmarking against Avocado using the global Spearman's correlation at the mark level.**

In the first row, we benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each symbol represents a testing sample. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001). In the second row, the overall distribution of this evaluation metric is shown as density plots, where the solid lines represent our method and the dashed lines represent Avocado.
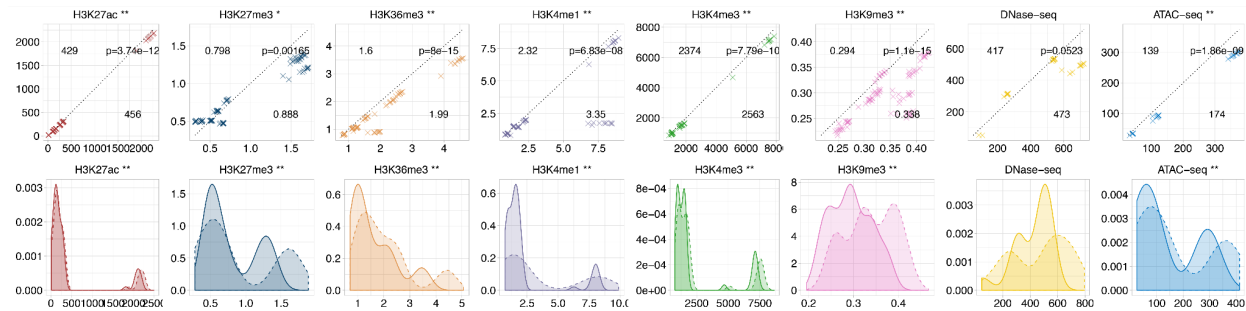
**Supplementary Figure 4: Benchmarking against Avocado using the local MSEs across promoter regions at the mark level.**
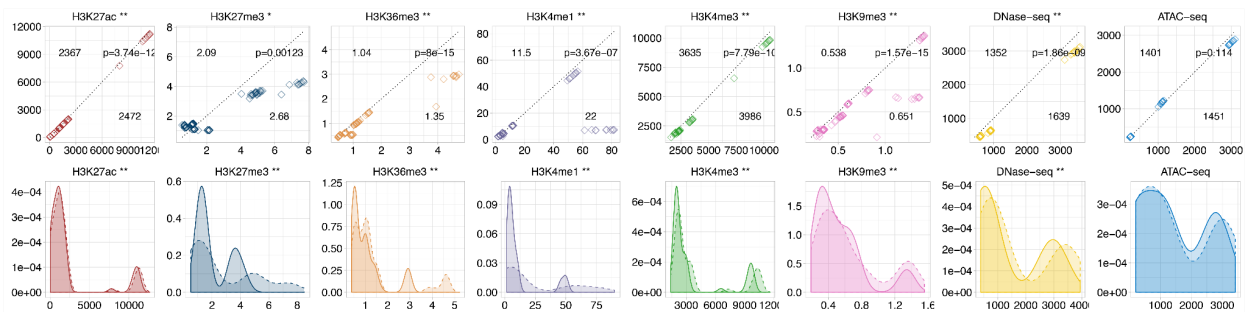
In the first row, we benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each symbol represents a testing sample. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001). In the second row, the overall distribution of this evaluation metric is shown as density plots, where the solid lines represent our method and the dashed lines represent Avocado.



**Supplementary Figure 5: Benchmarking against Avocado using the local MSEs across gene regions at the mark level.**
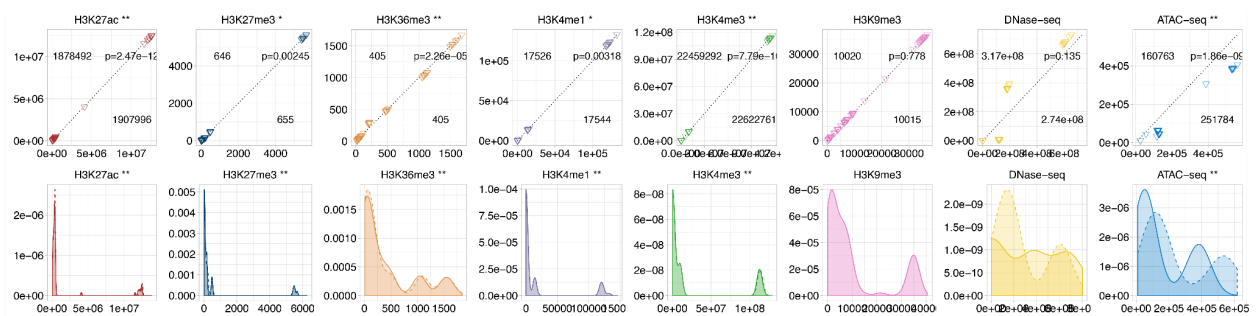
In the first row, we benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each symbol represents a testing sample. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly

different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001). In the second row, the overall distribution of this evaluation metric is shown as density plots, where the solid lines represent our method and the dashed lines represent Avocado.



**Supplementary Figure 6: Benchmarking against Avocado using the local MSEs across enhancer regions at the mark level.**

In the first row, we benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each symbol represents a testing sample. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001). In the second row, the overall distribution of this evaluation metric is shown as density plots, where the solid lines represent our method and the dashed lines represent Avocado.
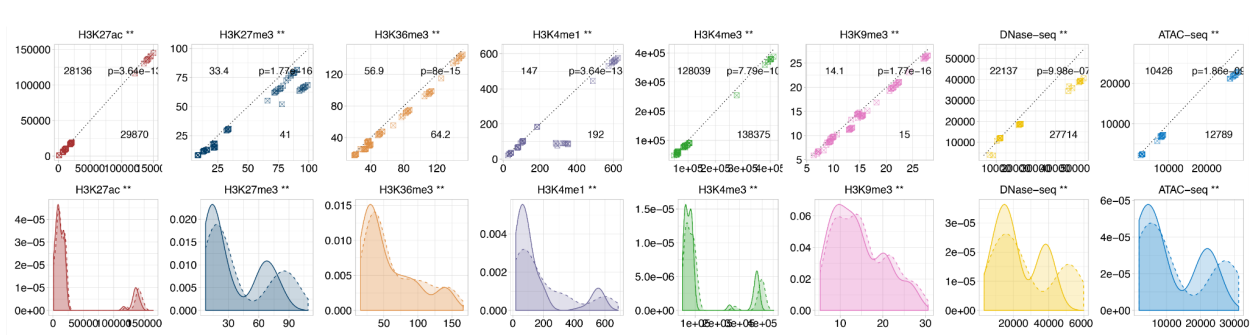
**Supplementary Figure 7: Benchmarking against Avocado using the global MSE weighted by the cross-cell-type variance at the mark level.**

In the first row, we benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each symbol represents a testing sample. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001). In the second row, the overall distribution of this evaluation metric is shown as density plots, where the solid lines represent our method and the dashed lines represent Avocado.
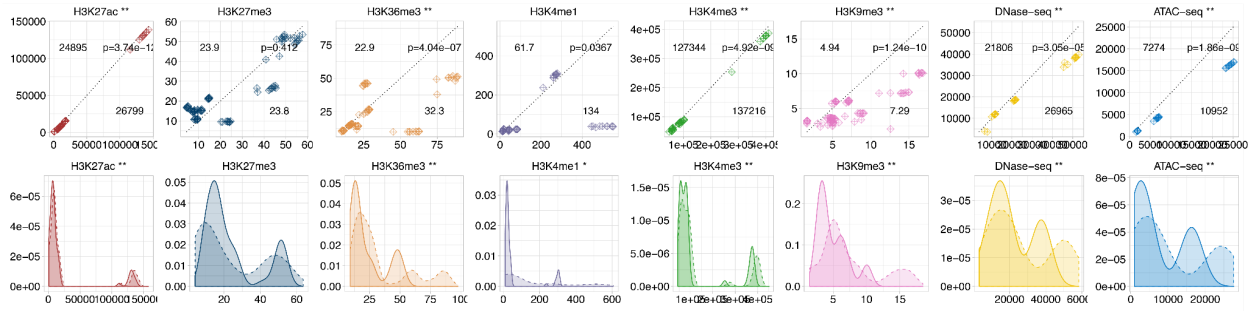
**Supplementary Figure 8: Benchmarking against Avocado using the local MSE across genomic regions with top 1% observed values.**

In the first row, we benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each symbol represents a testing sample. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001). In the second row, the overall distribution of this evaluation metric is shown as density plots, where the solid lines represent our method and the dashed lines represent Avocado.

**Supplementary Figure 9: Benchmarking against Avocado using the local MSE across genomic regions with top 1% predicted values at the mark level.**

In the first row, we benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each symbol represents a testing sample. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001). In the second row, the overall distribution of this evaluation metric is shown as density plots, where the solid lines represent our method and the dashed lines represent Avocado.
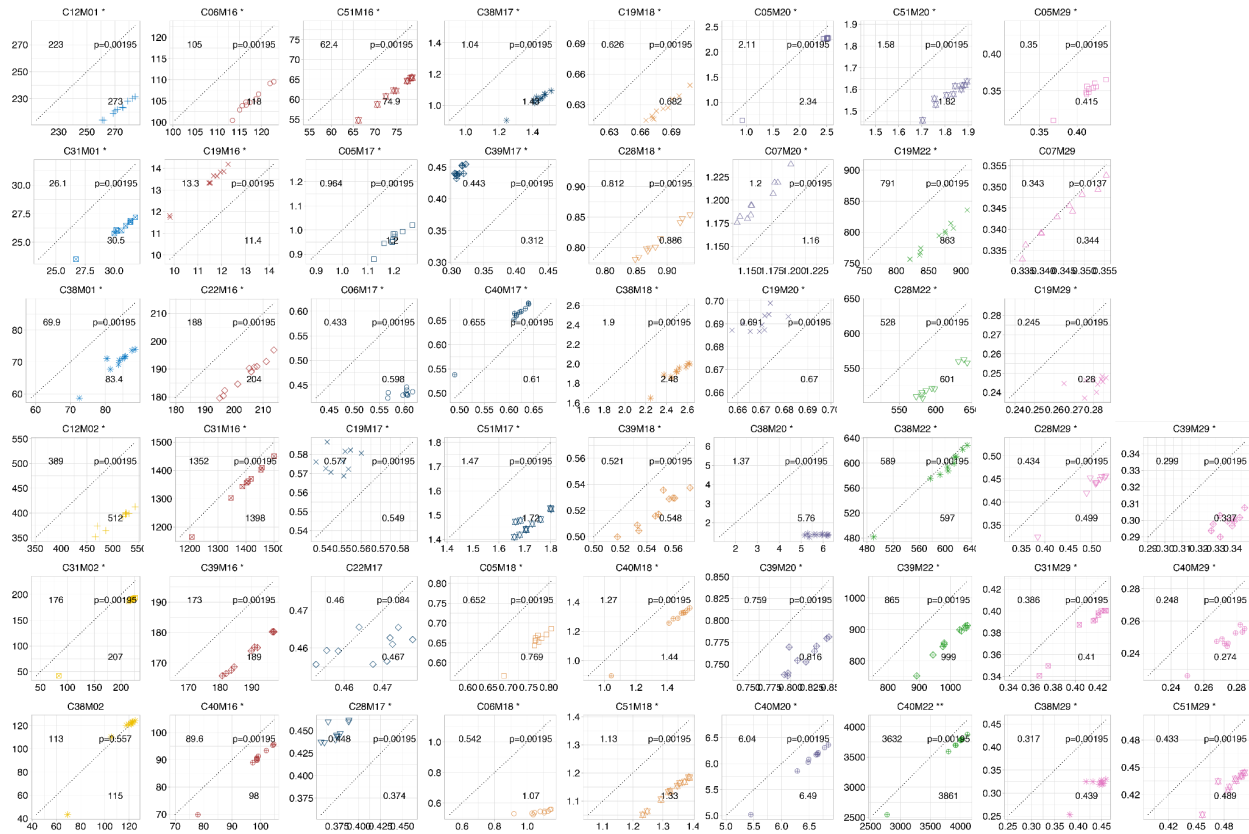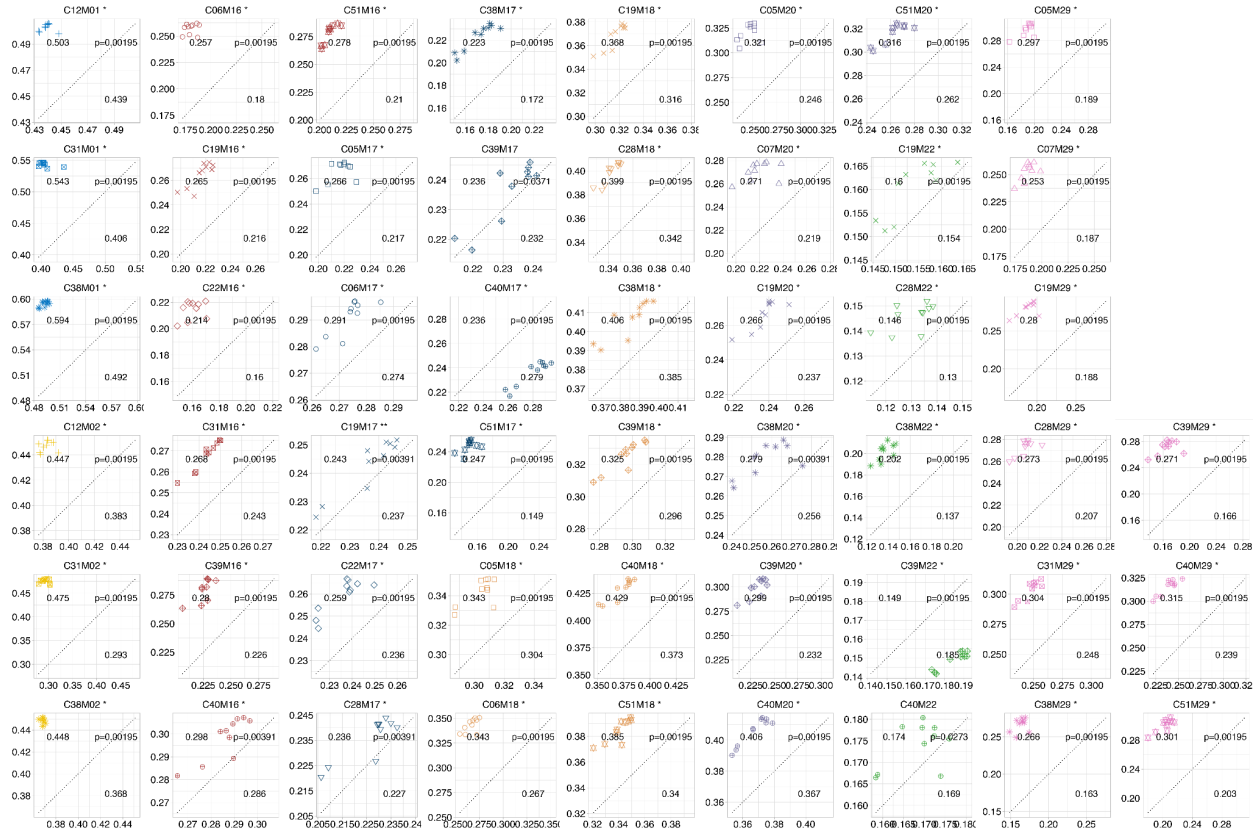
**Supplementary Figure 10: Benchmarking against Avocado using the global Pearson's correlation for each testing mark-cell type pair.**

We benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each figure contains 10 bootstrap sampling results from the human genome. The color represents the mark type and the symbol shape represents the cell type. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).

**Supplementary Figure 11: Benchmarking against Avocado using the global MSE for each testing mark-cell type pair.**

We benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each figure contains 10 bootstrap sampling results from the human genome. The color represents the mark type and the symbol shape represents the cell type. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
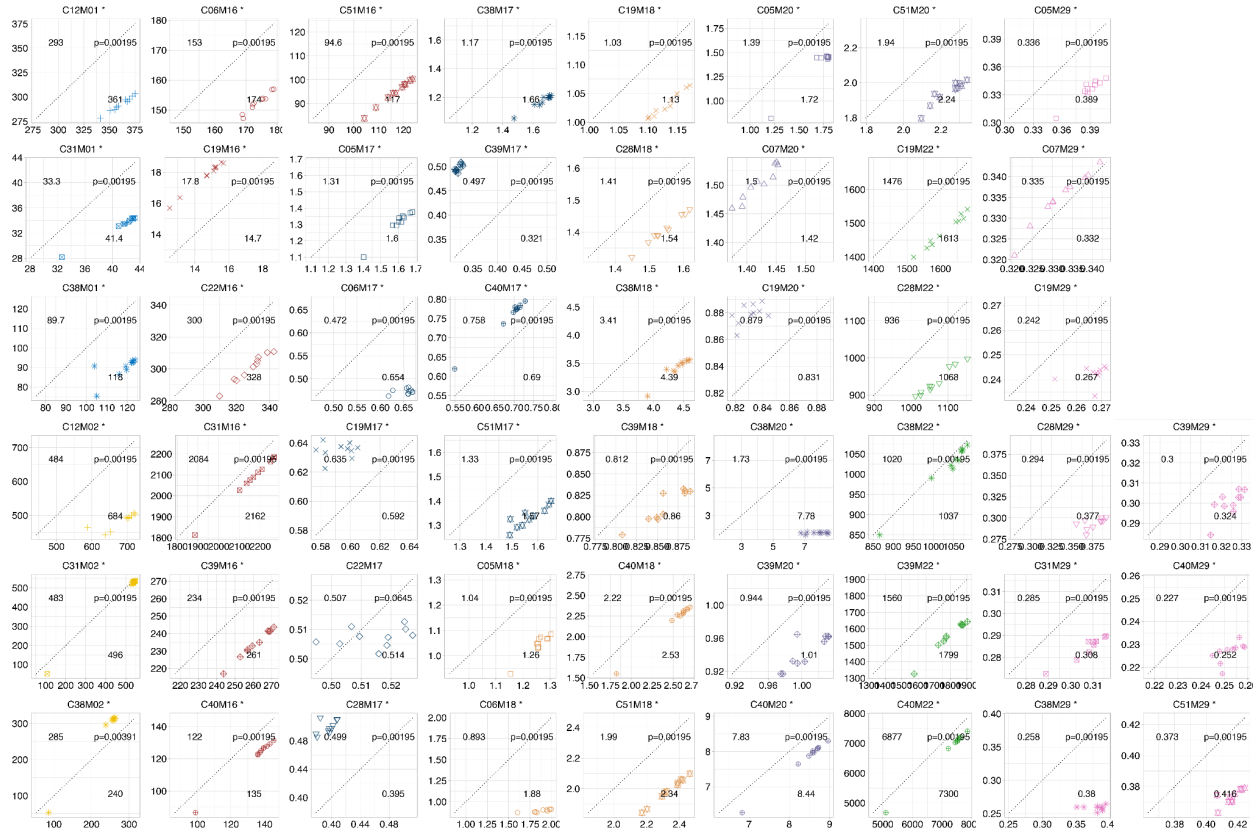
**Supplementary Figure 12: Benchmarking against Avocado using the global Spearman's correlation for each testing mark-cell type pair.**

We benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each figure contains 10 bootstrap sampling results from the human genome. The color represents the mark type and the symbol shape represents the cell type. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
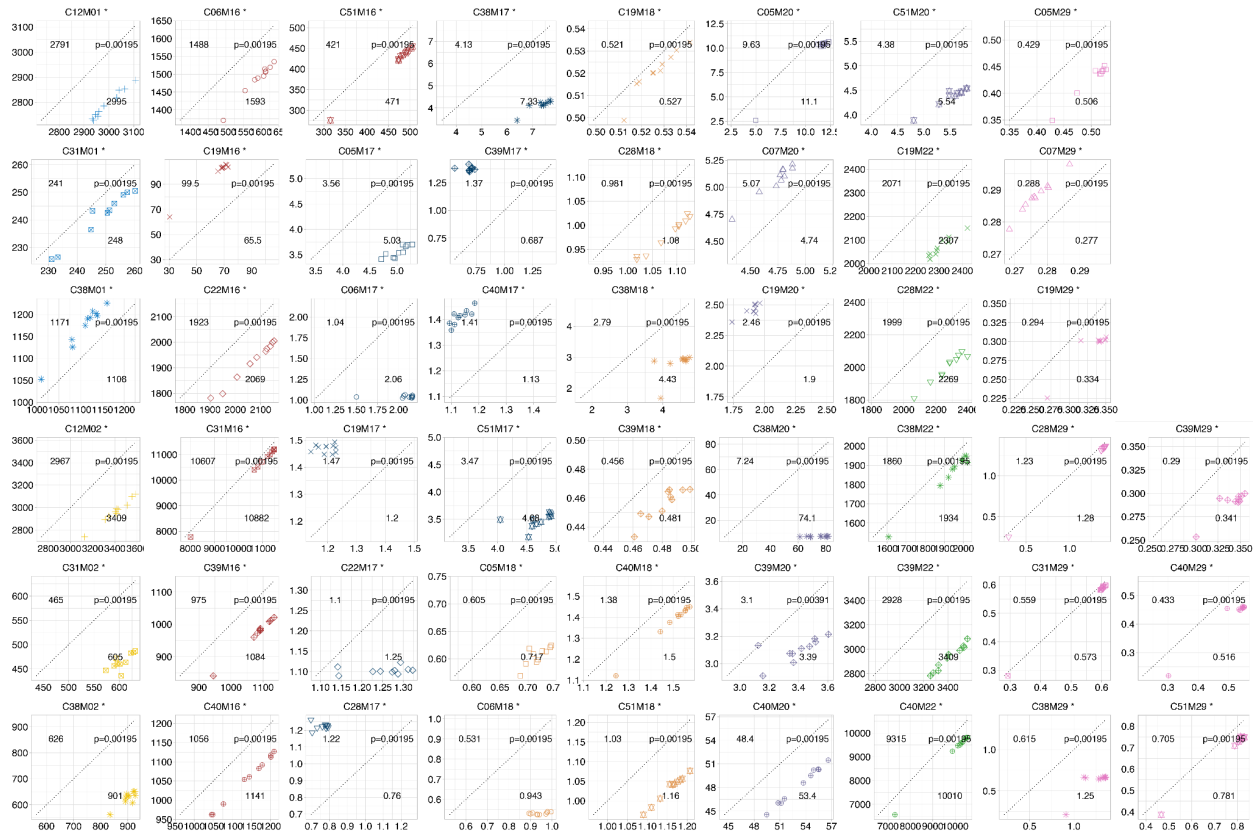
**Supplementary Figure 13: Benchmarking against Avocado using the local MSEs across promoter regions for each testing mark-cell type pair.**

We benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each figure contains 10 bootstrap sampling results from the human genome. The color represents the mark type and the symbol shape represents the cell type. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
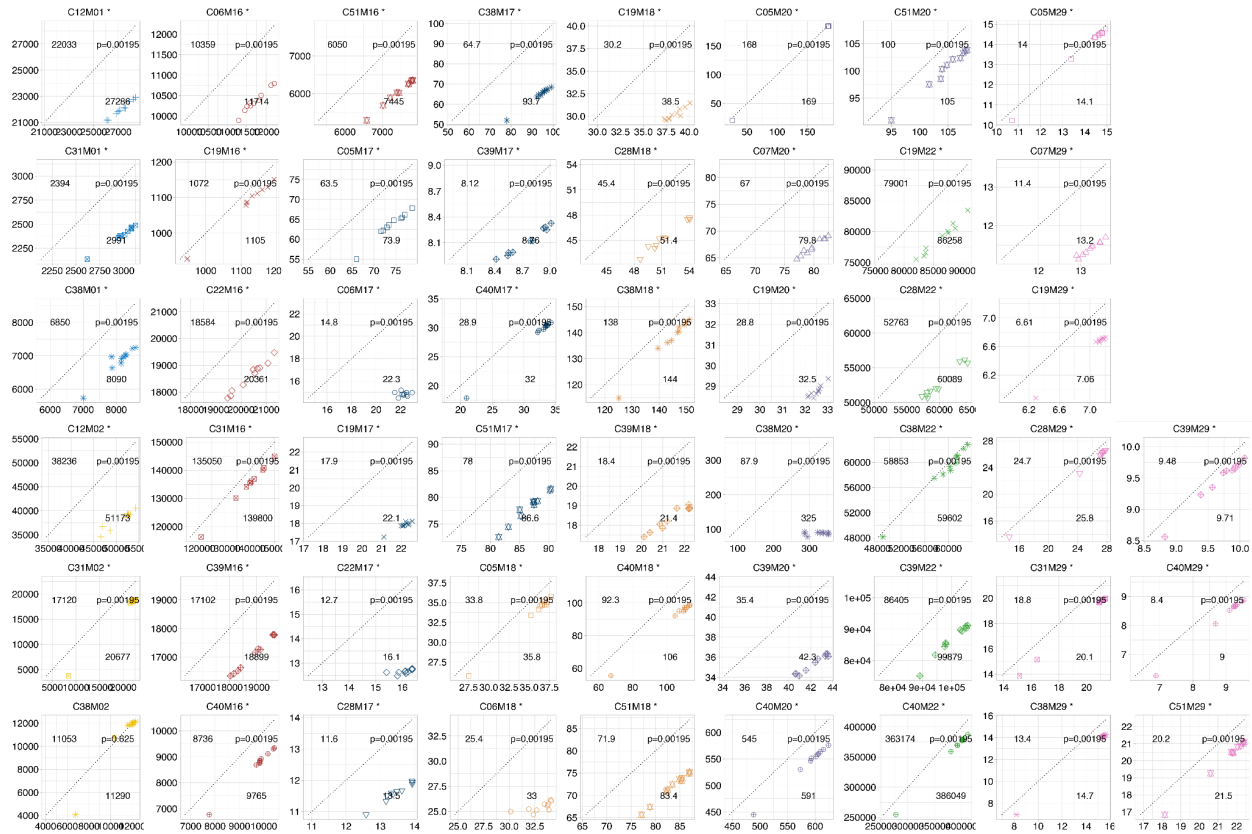
**Supplementary Figure 14: Benchmarking against Avocado using the local MSEs across gene regions for each testing mark-cell type pair.**

We benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each figure contains 10 bootstrap sampling results from the human genome. The color represents the mark type and the symbol shape represents the cell type. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
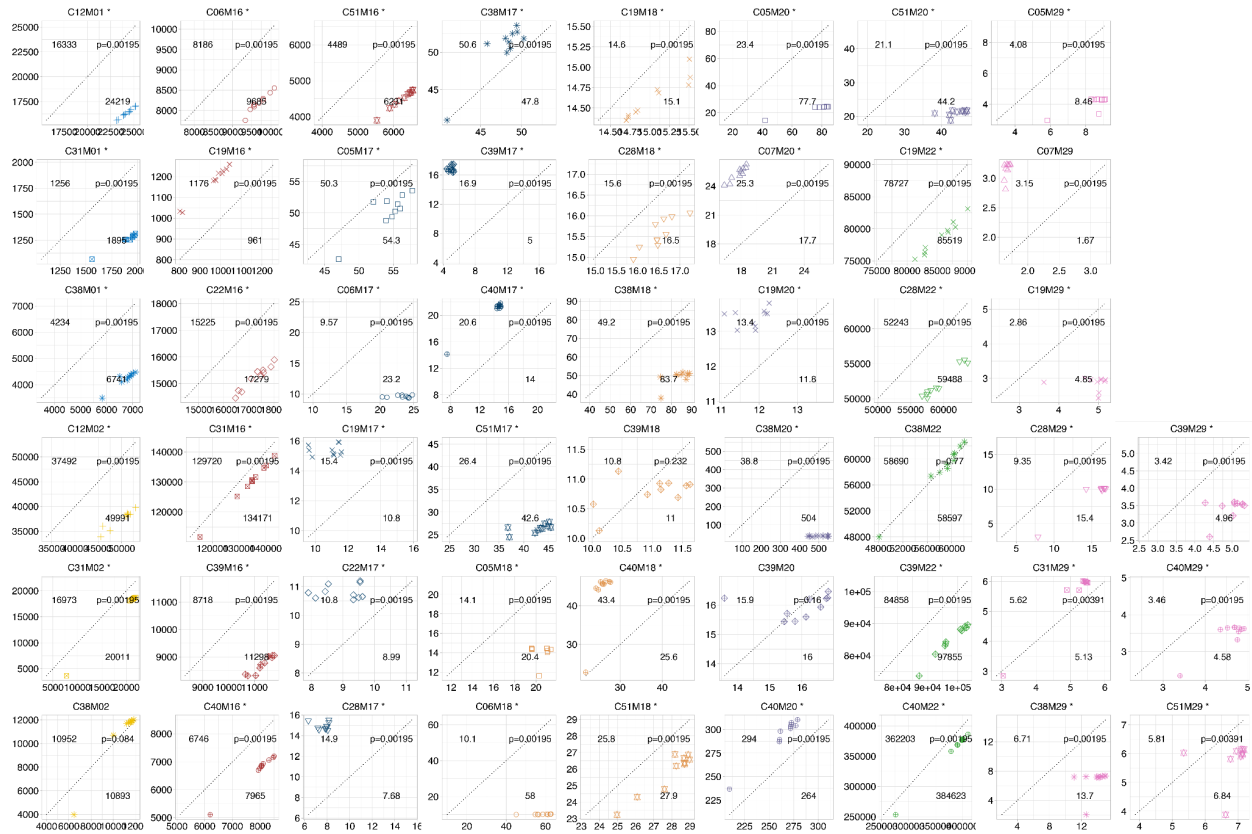
**Supplementary Figure 15: Benchmarking against Avocado using the local MSEs across enhancer regions for each testing mark-cell type pair.**

We benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each figure contains 10 bootstrap sampling results from the human genome. The color represents the mark type and the symbol shape represents the cell type. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
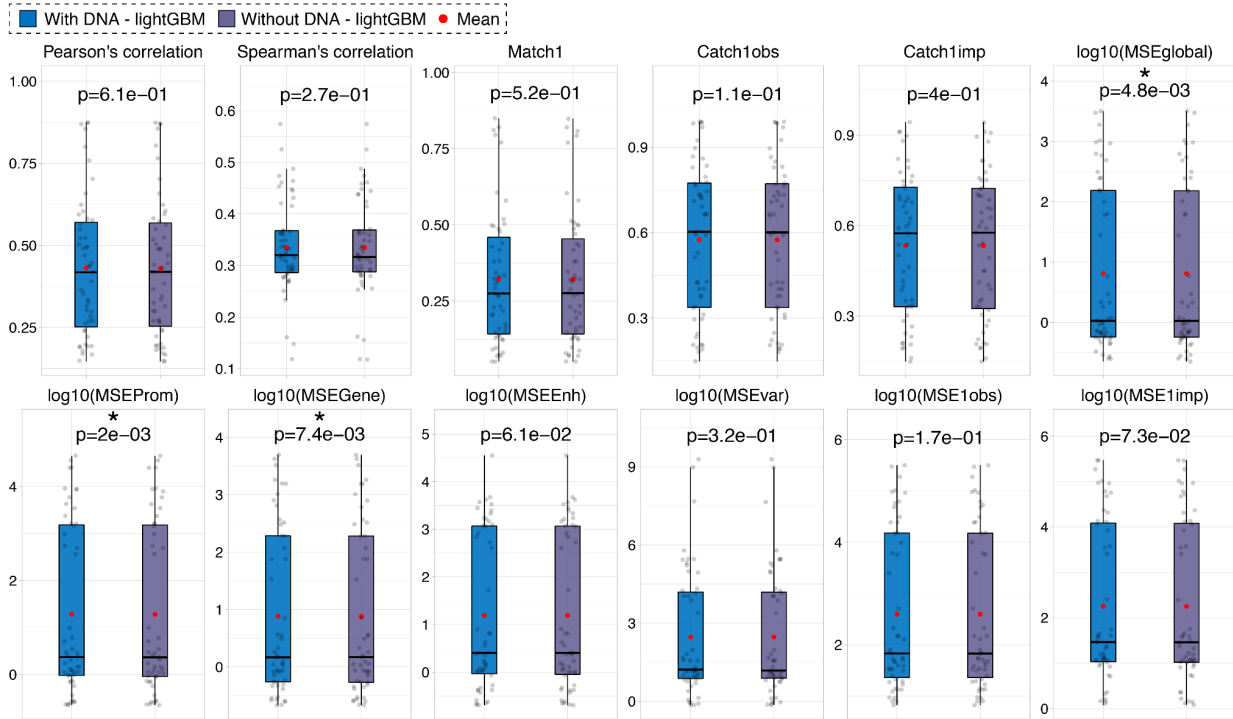
**Supplementary Figure 16: Benchmarking against Avocado using the global MSE weighted by the cross-cell-type variance for each testing mark-cell type pair.**

We benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each figure contains 10 bootstrap sampling results from the human genome. The color represents the mark type and the symbol shape represents the cell type. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
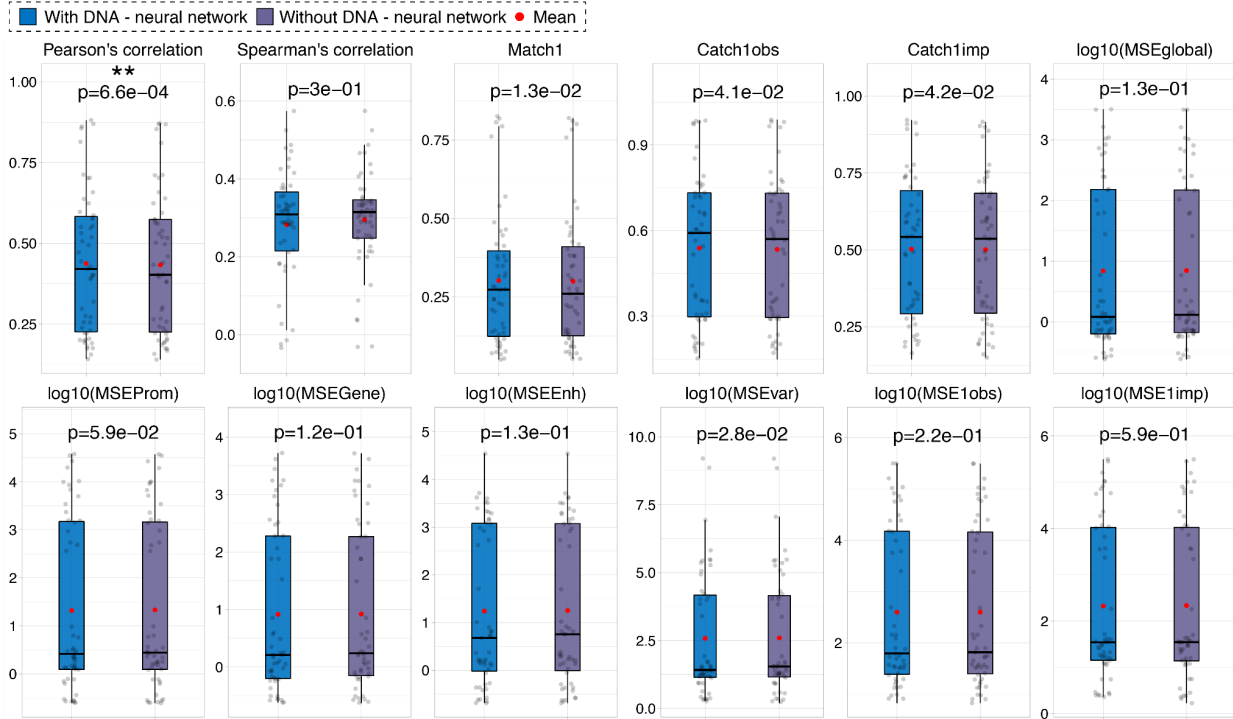
**Supplementary Figure 17: Benchmarking against Avocado using the local MSE across genomic regions with top 1% observed values for each testing mark-cell type pair.**

We benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each figure contains 10 bootstrap sampling results from the human genome. The color represents the mark type and the symbol shape represents the cell type. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
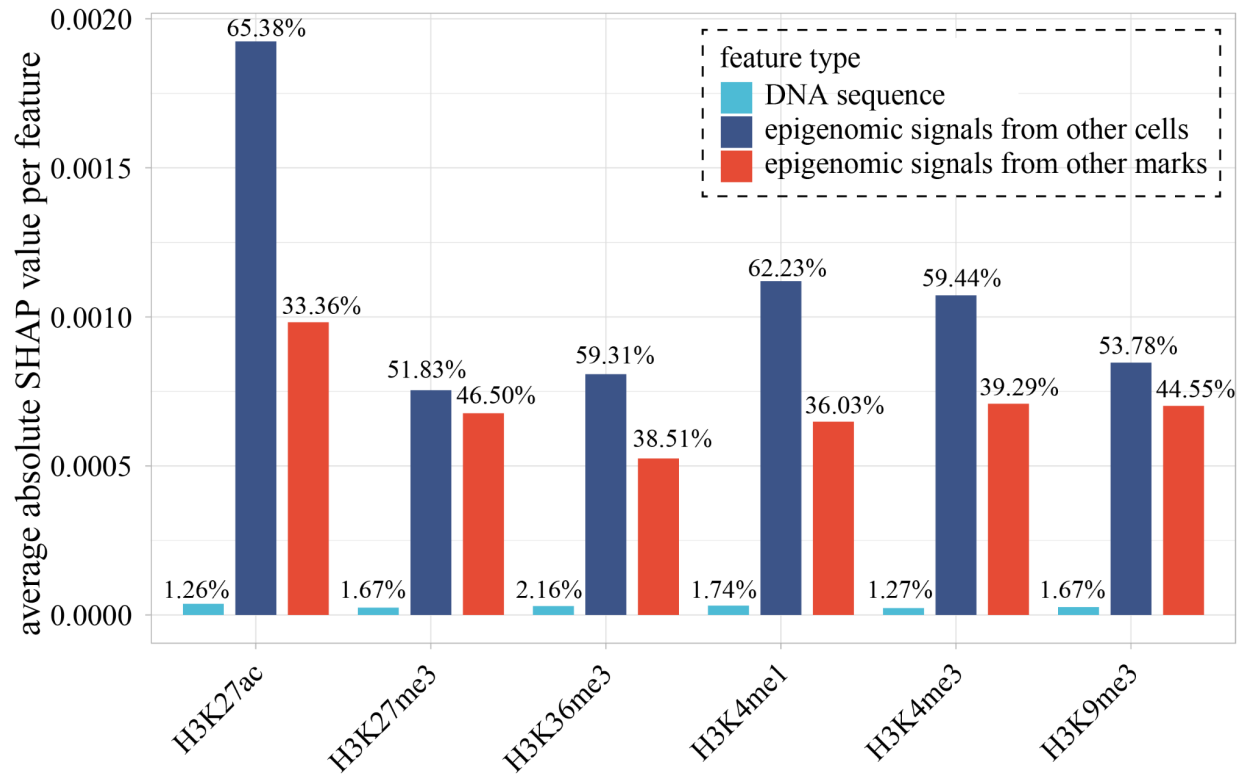
**Supplementary Figure 18: Benchmarking against Avocado using the local MSE across genomic regions with top 1% predicted values for each testing mark-cell type pair.**

We benchmarked our method (y-axis) against Avocado (x-axis) across the entire human genome on 51 held-out testing epigenomes covering six histone modifications, DNase-seq and ATAC-seq. Each figure contains 10 bootstrap sampling results from the human genome. The color represents the mark type and the symbol shape represents the cell type. If a symbol is above the diagonal dashed line, it means our method achieves a higher predictive performance than Avocado. The paired Wilcoxon signed-rank test was used to statistically compare two methods and the p-values are shown in the figures. Significantly different results are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).

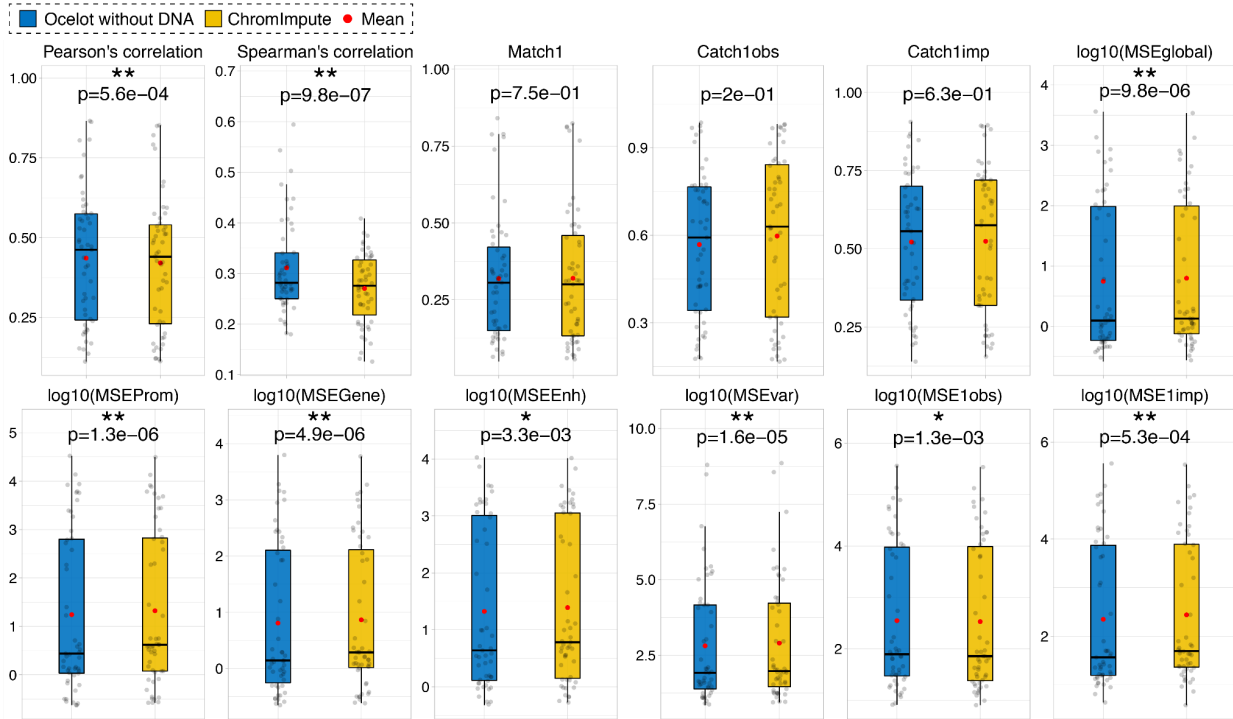**Supplementary Figure 19: Predictive performance comparison between lightGBM models with and without DNA sequence features**

Predictive performance was evaluated on the largest Chromosome 1 using 12 scoring metrics on the ENCODE Imputation Challenge testing set of 51 mark-cell type pairs. For each metric, the paired one-sided Wilcoxon signed-rank test was performed between two methods across 51 testing mark-cell type pairs. The significantly different ones are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
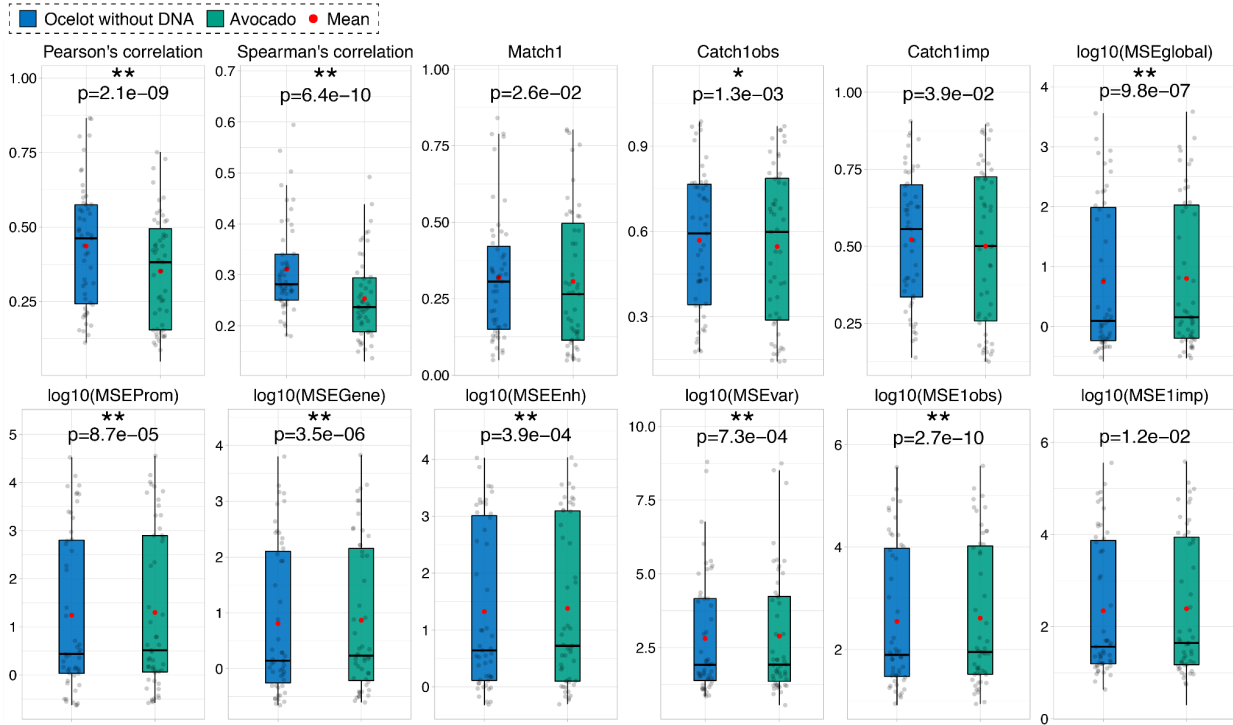
**Supplementary Figure 20: Predictive performance comparison between neural network models with and without DNA sequence features**

Predictive performance was evaluated on the largest Chromosome 1 using 12 scoring metrics on the ENCODE Imputation Challenge testing set of 51 mark-cell type pairs. For each metric, the paired one-sided Wilcoxon signed-rank test was performed between two methods across 51 testing mark-cell type pairs. The significantly different ones are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
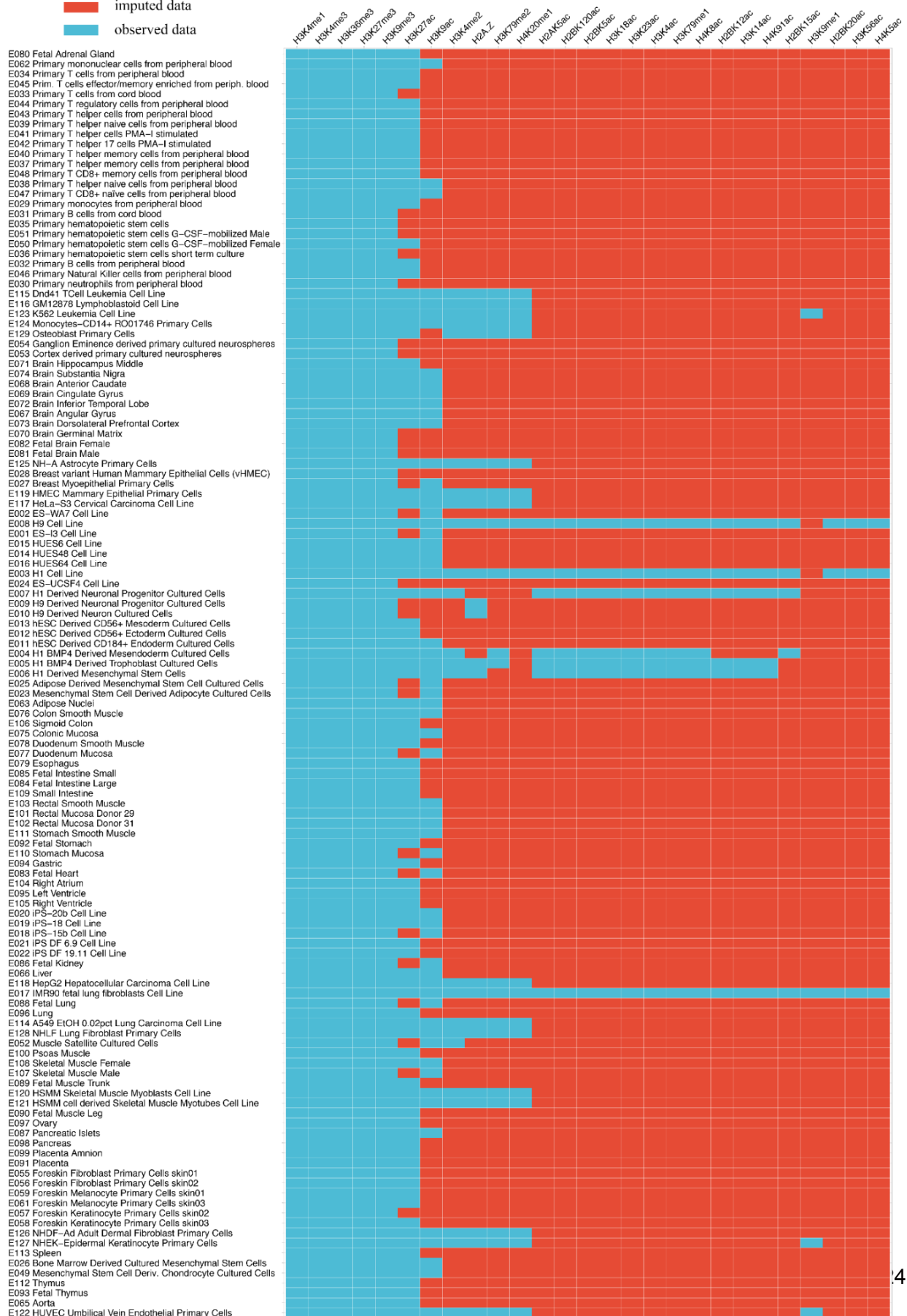
**Supplementary Figure 21: The contributions to predictions of different types of features.**

We used three types of input features in Ocelot: (1) DNA sequence, (2) epigenomic signals from other cell types, and (3) epigenomic signals from other marks. For each histone mark, the average absolute SHAP value per feature was calculated. The percentages are shown on top of each bar.

**Supplementary Figure 22: Predictive performance comparison between Ocelot without DNA sequence features and ChromImpute**

Predictive performance was evaluated on the whole human genome using 12 scoring metrics on the ENCODE Imputation Challenge testing set of 51 mark-cell type pairs. For each testing pair, we calculated genome-wide evaluation scores through concatenating signals of 23 chromosomes. For each metric, the paired one-sided Wilcoxon signed-rank test was performed between two methods across 51 testing mark-cell type pairs. The significantly different ones are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
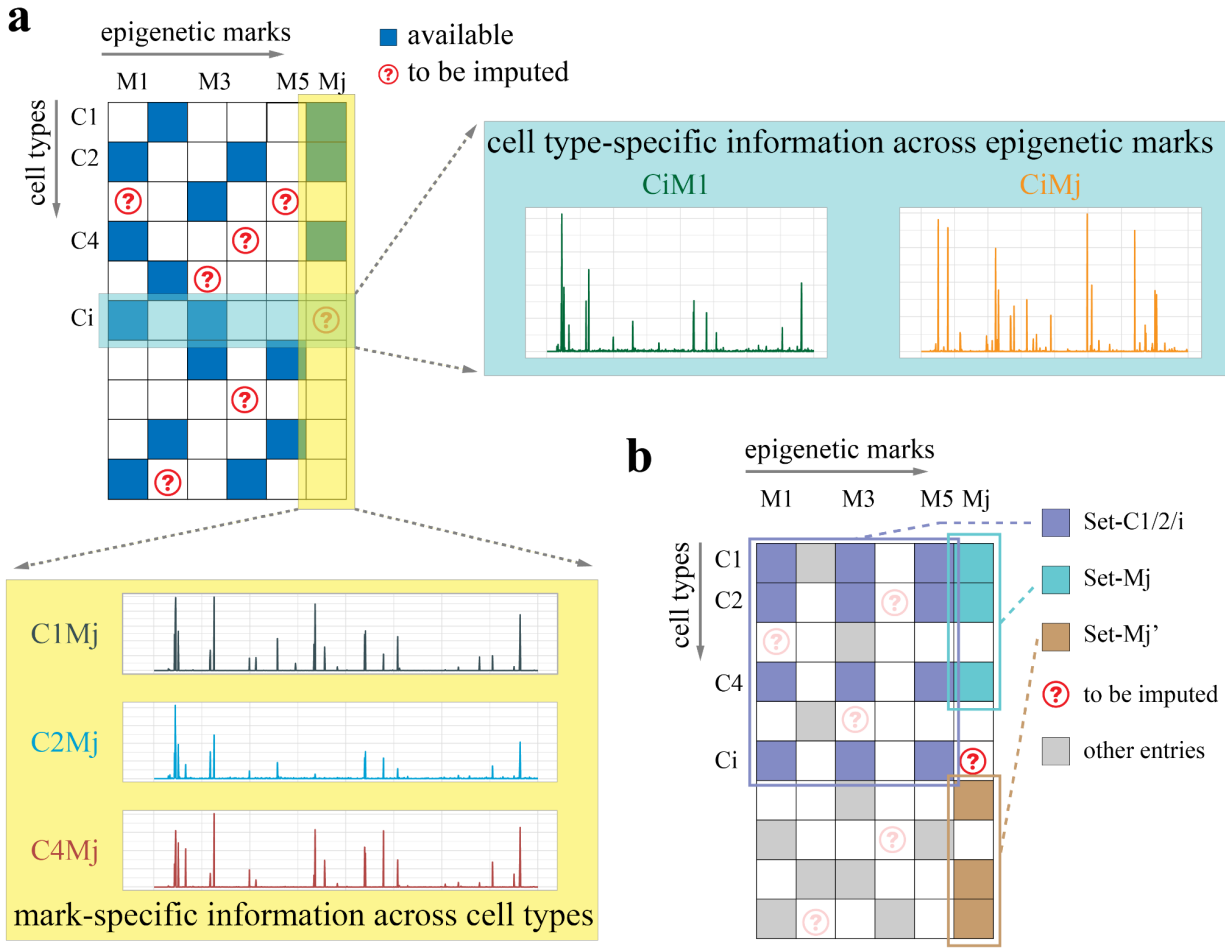
**Supplementary Figure 23: Predictive performance comparison between Ocelot without DNA sequence features and Avocado**

Predictive performance was evaluated on the whole human genome using 12 scoring metrics on the ENCODE Imputation Challenge testing set of 51 mark-cell type pairs. For each testing pair, we calculated genome-wide evaluation scores through concatenating signals of 23 chromosomes. For each metric, the paired one-sided Wilcoxon signed-rank test was performed between two methods across 51 testing mark-cell type pairs. The significantly different ones are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
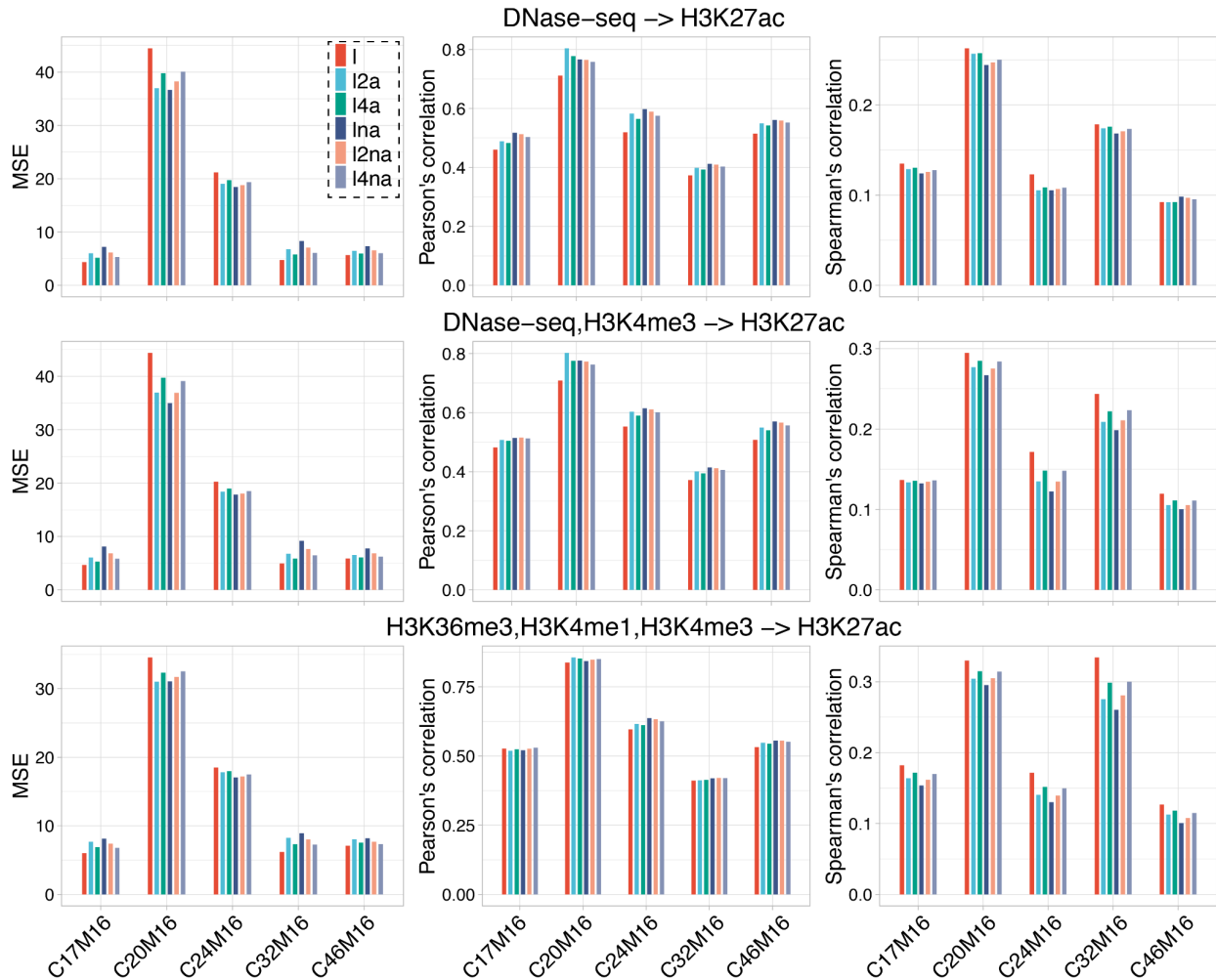
Roadmap Dataset

**Supplementary Figure 24: Application of Ocelot to impute missing entries and complete the Roadmap Epigenomics dataset.**

The Roadmap Epigenomics dataset covers 27 histone marks in 127 cell and tissue conditions after excluding four histone marks (H3K23me2, H2AK9ac, H3T11ph, H4K12ac) that only have one or two observed signal tracks. This is because with too few observed whole-genome tracks, we could not train a solid lightGBM or neural network model. A total of 974 (28.40%) whole-genome profiles were observed (blue blocks) and used to build machine learning models to predict the remaining 2,455 (71.60%) missing profiles (red blocks).
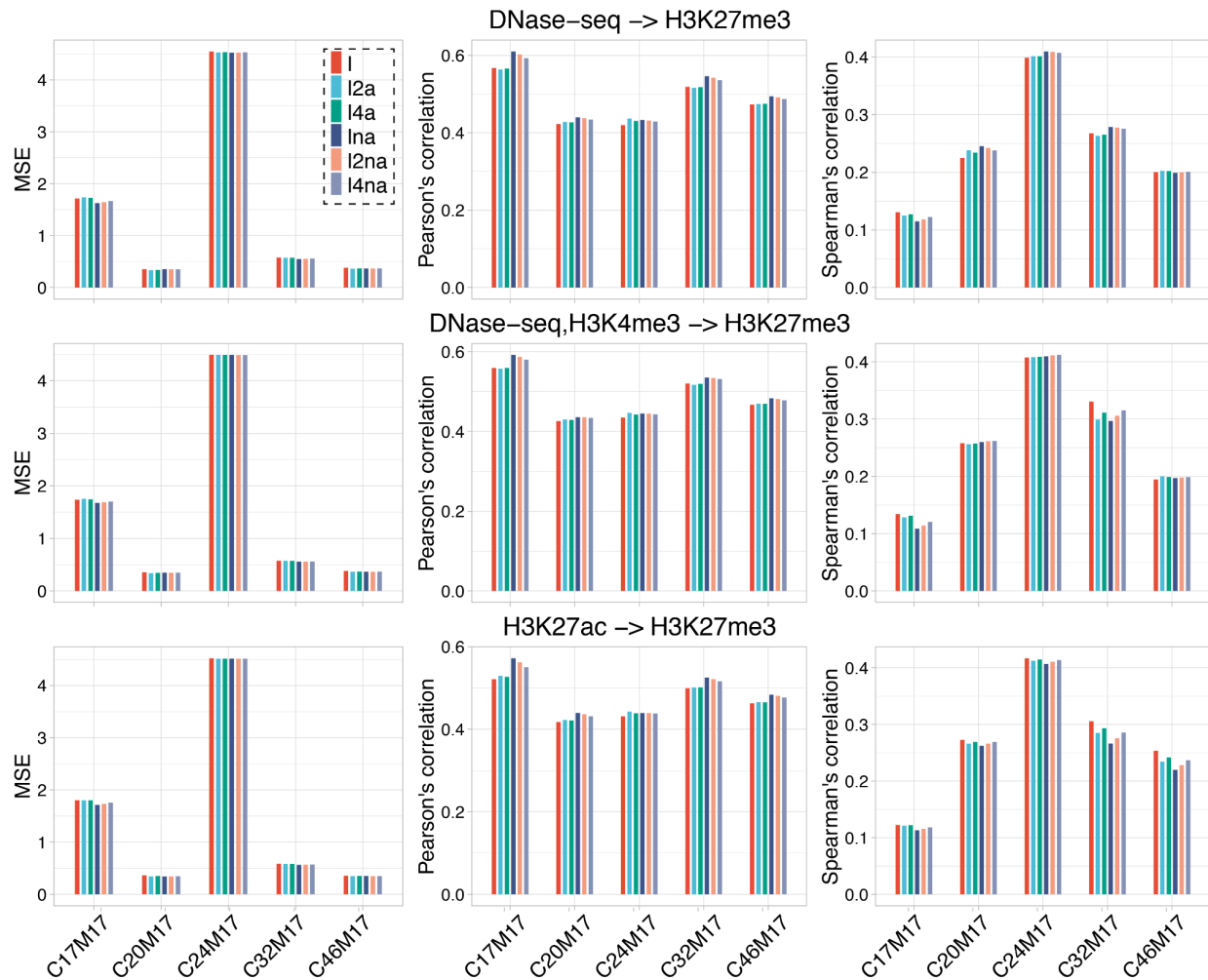
**Supplementary Figure 25: Data partition for machine learning models.**

**a,** The mark-specific and cell type-specific information is integrated into machine learning models to improve predictive performance. **b,** Based on the target entry to be imputed and available entries, the cell type - mark matrix can be partitioned into subset, serving as training targets, mark-specific features or cell type-specific features.
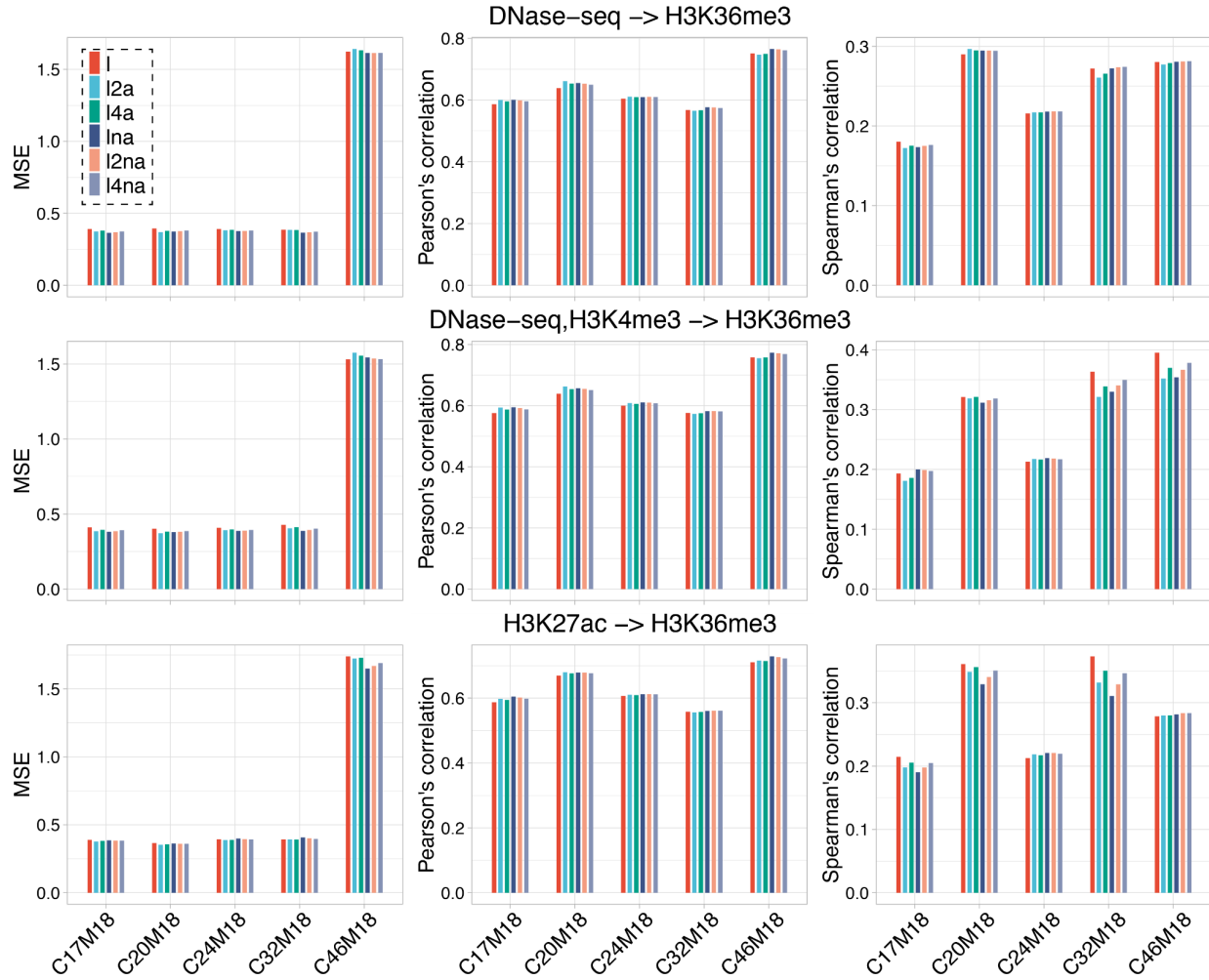
**Supplementary Figure 26: Cross-validation results of predicting H3K27ac based on the challenge training data**

Within the challenge training data, five cell types (C17, C20, C24, C32, C46) were held out as the testing set to evaluate predictive performance of different model design and ensemble weights. Three columns represent three evaluation metrics: (1) MSEglobal, (2) Pearson's correlation, and (3) Spearman's correlation. Multiple rows represent different model designs, where one or multiple marks are used as feature marks in machine learning models. Different colors represent different ensemble weights between lightGBM ("l"), neural network ("n") and average signal track ("a"). For example, "l4na" represents the ensemble weights of l:n:a = 4:1:1.
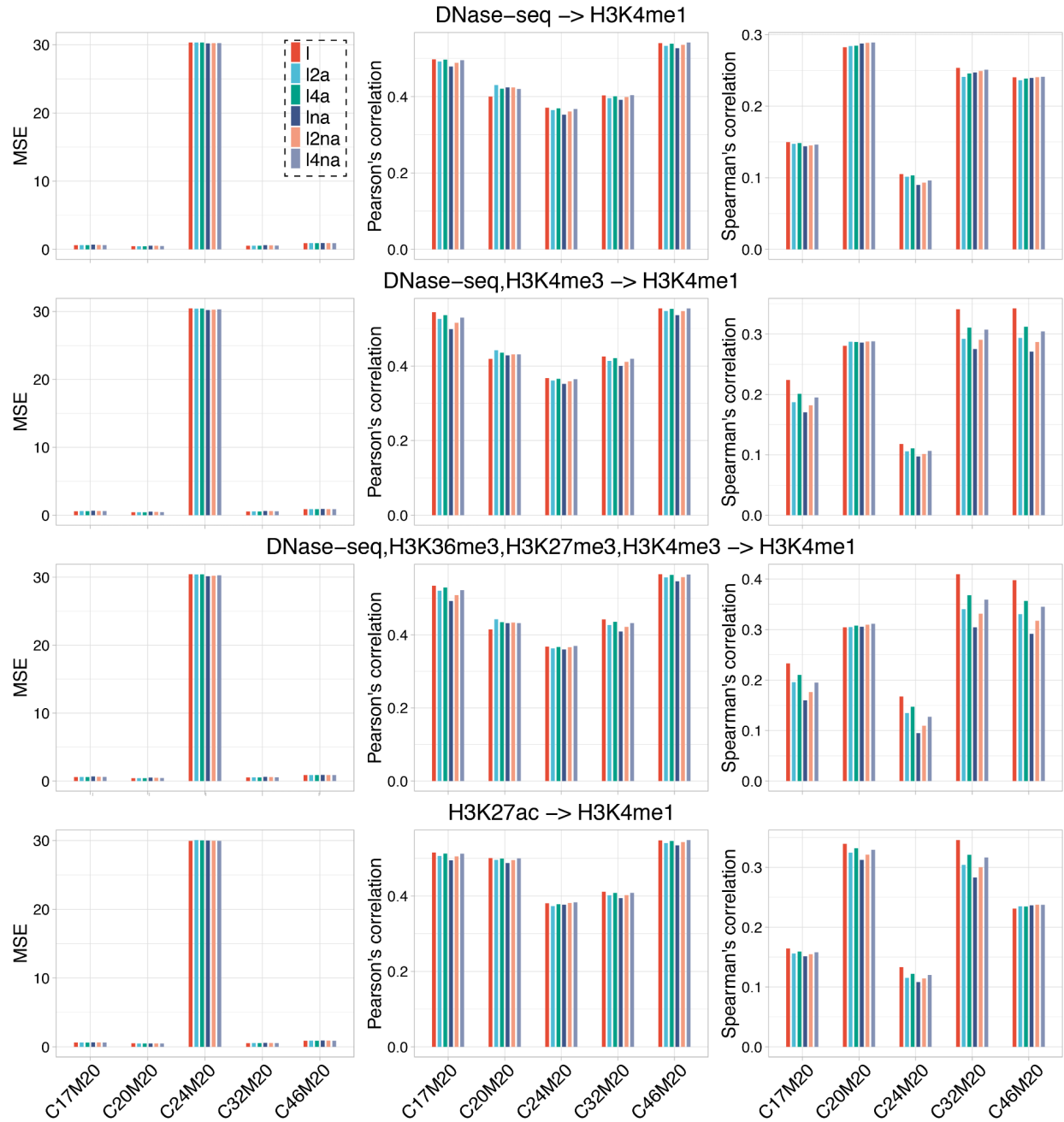
**Supplementary Figure 27: Cross-validation results of predicting H3K27me3 based on the challenge training data**

Within the challenge training data, five cell types (C17, C20, C24, C32, C46) were held out as the testing set to evaluate predictive performance of different model design and ensemble weights. Three columns represent three evaluation metrics: (1) MSEglobal, (2) Pearson's correlation, and (3) Spearman's correlation. Multiple rows represent different model designs, where one or multiple marks are used as feature marks in machine learning models. Different colors represent different ensemble weights between lightGBM ("l"), neural network ("n") and average signal track ("a"). For example, "l4na" represents the ensemble weights of l:n:a = 4:1:1.

**Supplementary Figure 28: Cross-validation results of predicting H3K36me3 based on the challenge training data**
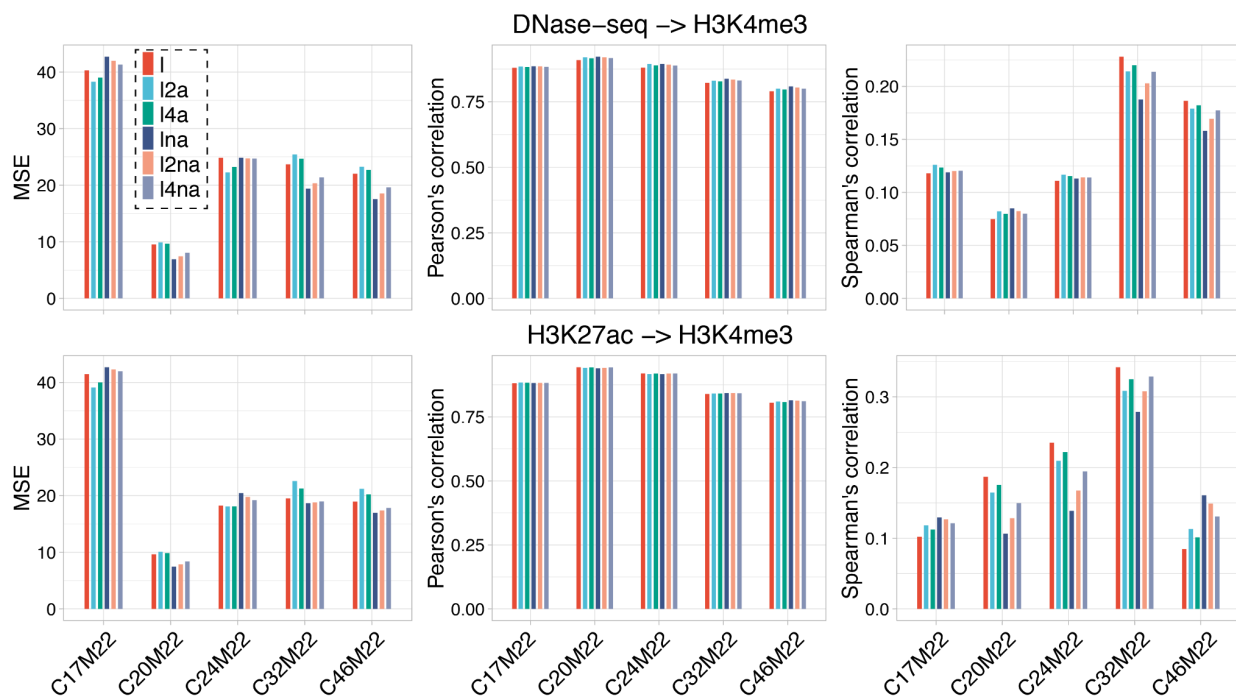
Within the challenge training data, five cell types (C17, C20, C24, C32, C46) were held out as the testing set to evaluate predictive performance of different model design and ensemble weights. Three columns represent three evaluation metrics: (1) MSEglobal, (2) Pearson's correlation, and (3) Spearman's correlation. Multiple rows represent different model designs, where one or multiple marks are used as feature marks in machine learning models. Different colors represent different ensemble weights between lightGBM ("l"), neural network ("n") and average signal track ("a"). For example, "l4na" represents the ensemble weights of l:n:a = 4:1:1.

**Supplementary Figure 29: Cross-validation results of predicting H3K4me1 based on the challenge training data**
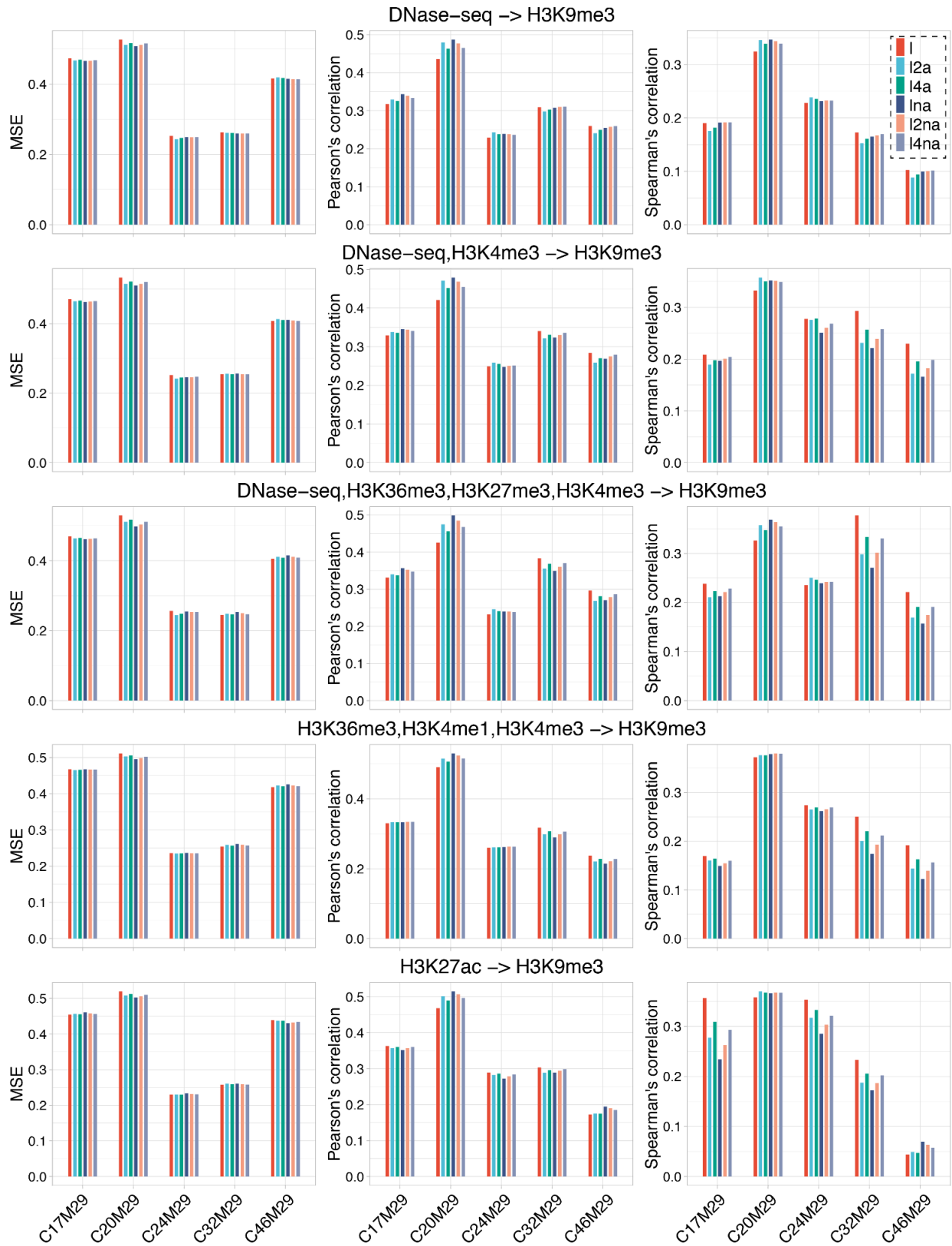
Within the challenge training data, five cell types (C17, C20, C24, C32, C46) were held out as the testing set to evaluate predictive performance of different model design and ensemble weights. Three columns represent three evaluation metrics: (1) MSEglobal, (2) Pearson's correlation, and (3) Spearman's correlation. Multiple rows represent different model designs, where one or multiple marks are used as feature marks in machine learning models. Different colors represent different ensemble weights between

30

lightGBM ("l"), neural network ("n") and average signal track ("a"). For example, "l4na" represents the ensemble weights of l:n:a = 4:1:1.
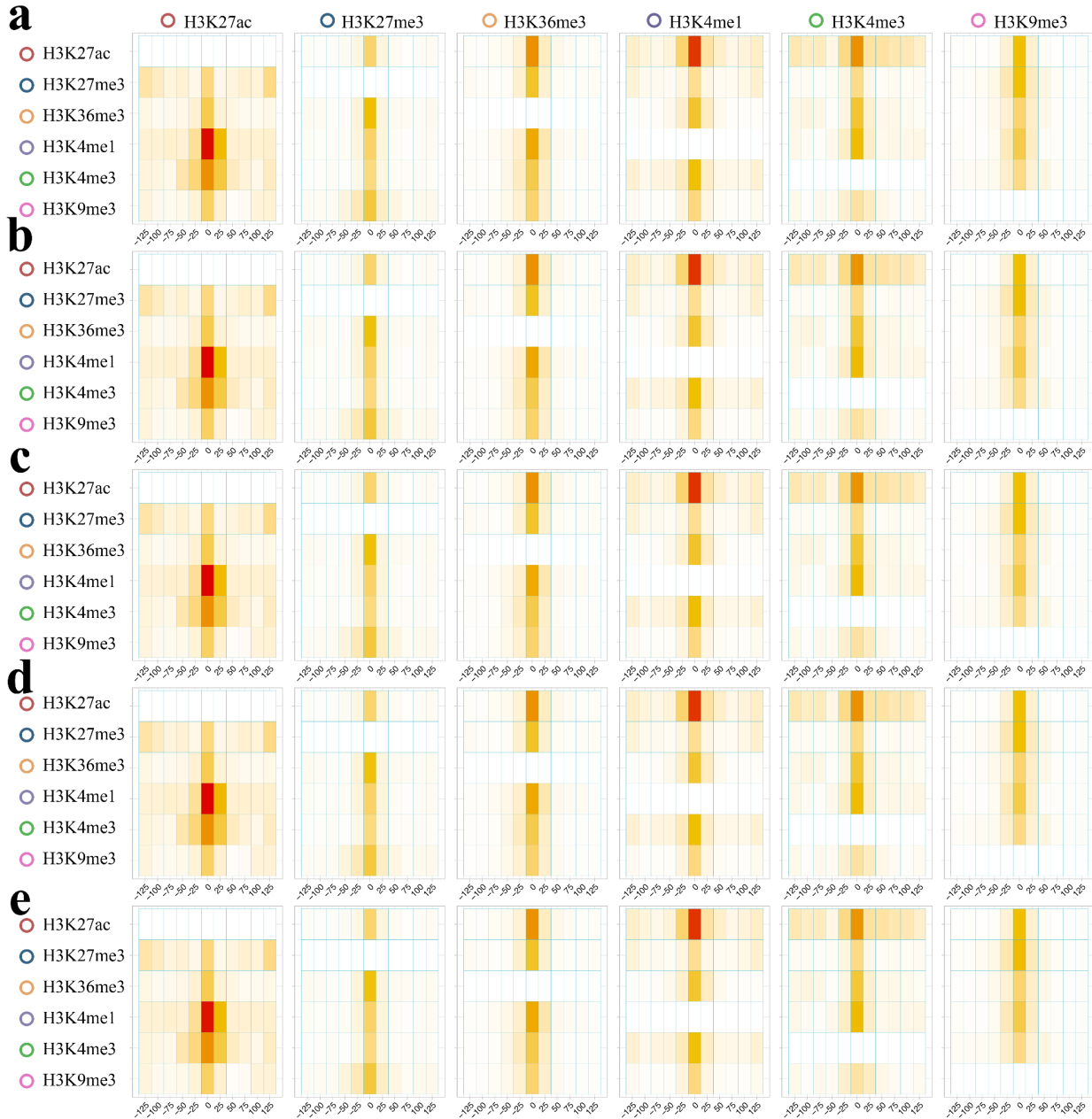


**Supplementary Figure 30: Cross-validation results of predicting H3K4me3 based on the challenge training data**

Within the challenge training data, five cell types (C17, C20, C24, C32, C46) were held out as the testing set to evaluate predictive performance of different model design and ensemble weights. Three columns represent three evaluation metrics: (1) MSEglobal, (2) Pearson's correlation, and (3) Spearman's correlation. Multiple rows represent different model designs, where one or multiple marks are used as feature marks in machine learning models. Different colors represent different ensemble weights between lightGBM ("l"), neural network ("n") and average signal track ("a"). For example, "l4na" represents the ensemble weights of l:n:a = 4:1:1.
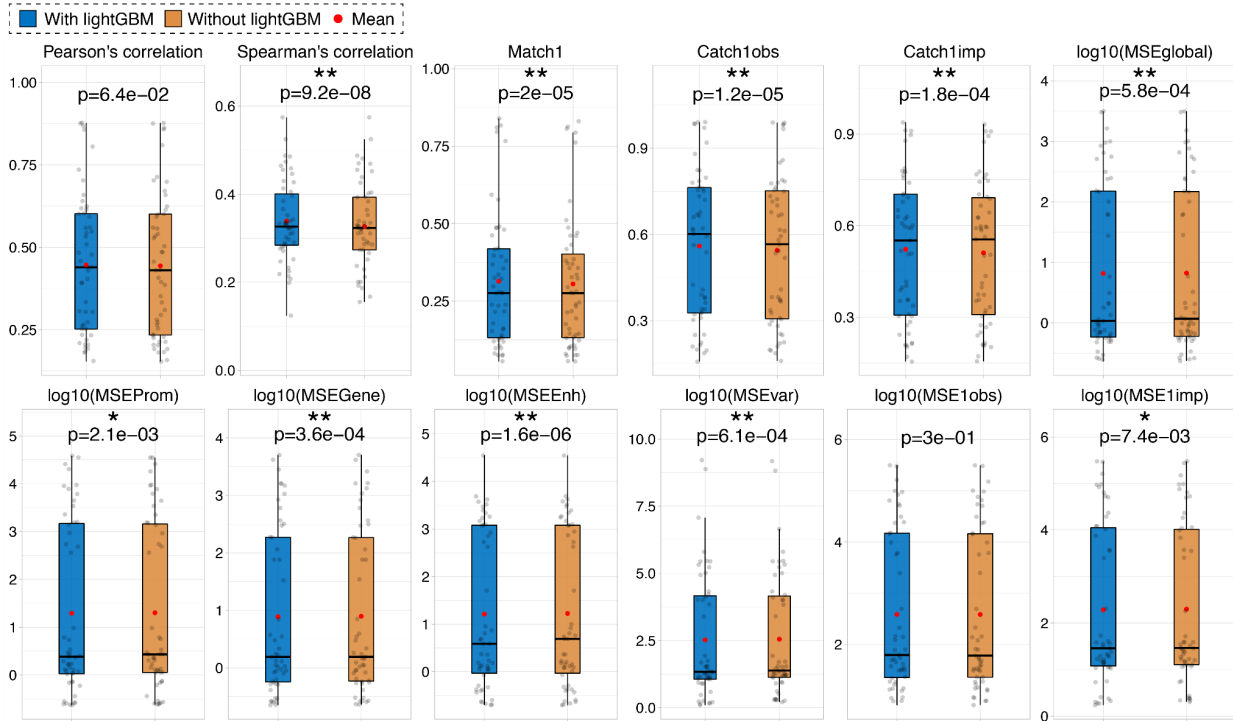
**Supplementary Figure 31: Cross-validation results of predicting H3K9me3 based on the challenge training data**

Within the challenge training data, five cell types (C17, C20, C24, C32, C46) were held out as the testing set to evaluate predictive performance of different model design and ensemble weights. Three columns represent three evaluation metrics: (1) MSEglobal, (2) Pearson's correlation, and (3) Spearman's correlation. Multiple rows represent different model designs, where one or multiple marks are used as feature marks in machine learning models. Different colors represent different ensemble weights between lightGBM ("l"), neural network ("n") and average signal track ("a"). For example, "l4na" represents the ensemble weights of l:n:a = 4:1:1.
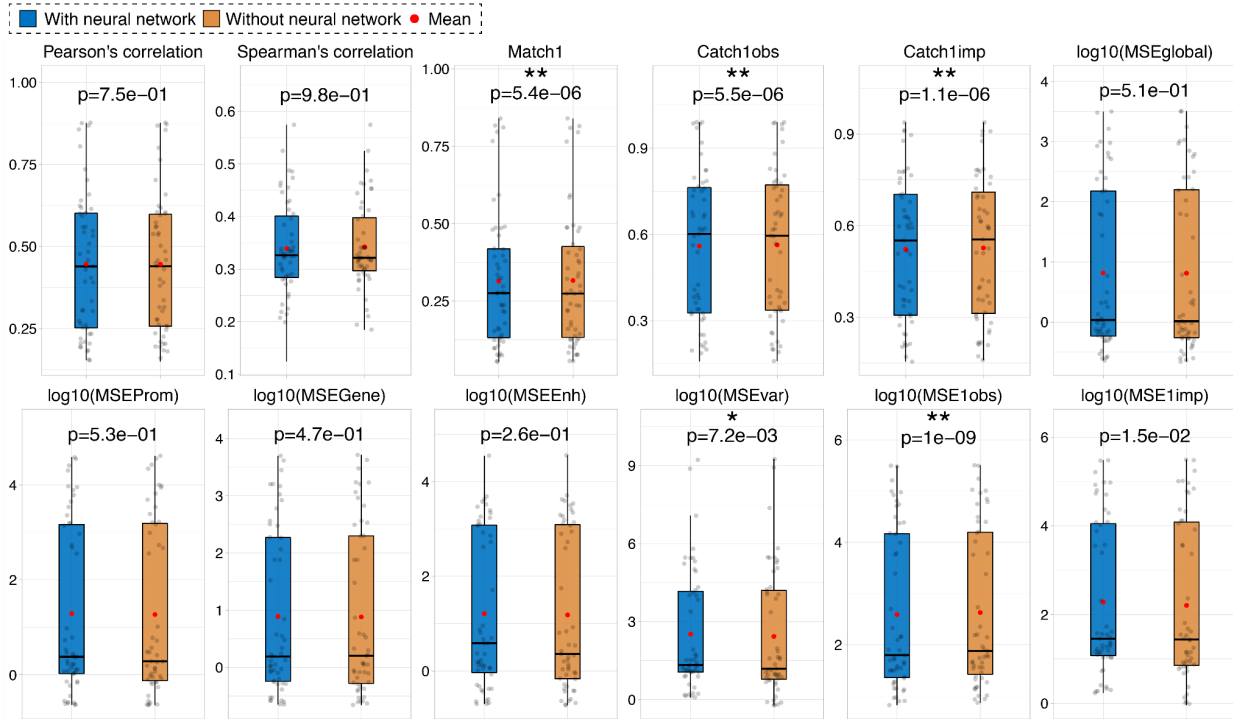
**Supplementary Figure 32: The SHAP heatmaps among six histone marks based on different subsets of testing cell types.**

Five cell types (C17, C20, C24, C32, C46) were held out as the testing set to perform SHAP analysis. In contrast to the SHAP heatmap in **Fig. 3a** that was derived from all five cell types, here only a subset of four cell types were used to calculate the SHAP heatmap. C17, C20, C24, C32 and C46 were excluded from the SHAP analysis in panels **a, b, c, d,** and **e,** respectively.
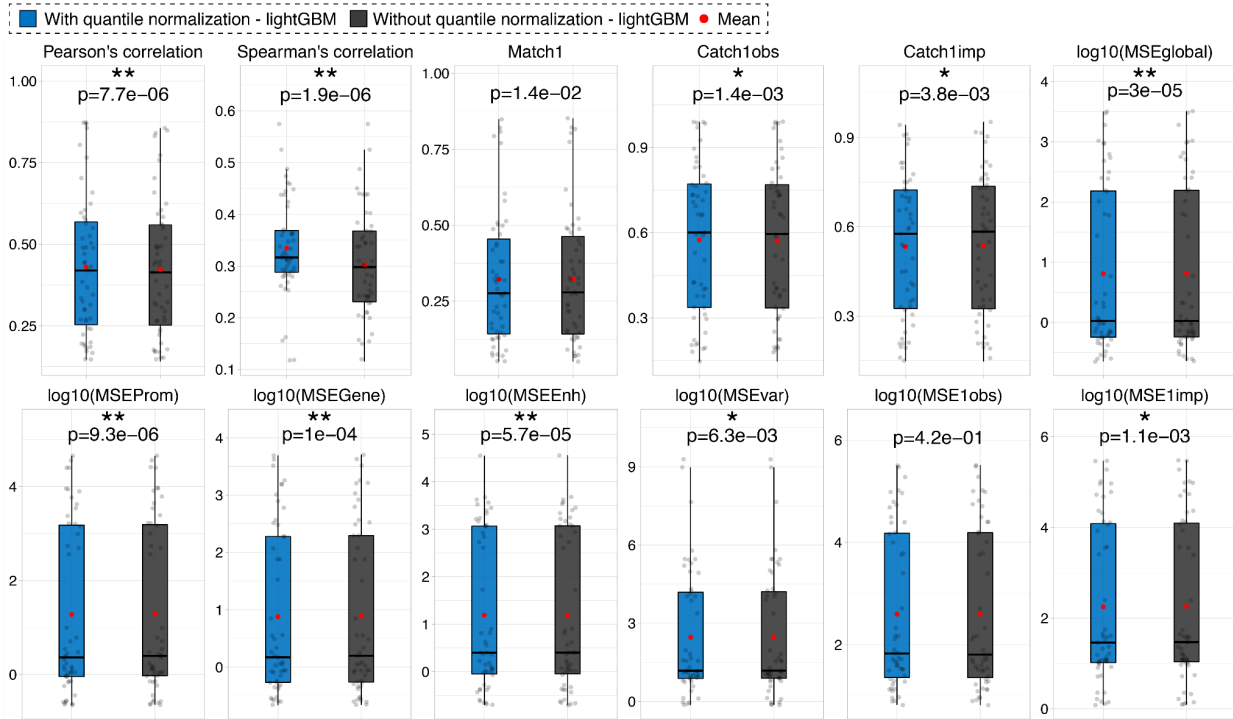
**Supplementary Figure 33: Predictive performance comparison between ensemble models with and without lightGBM models**

Predictive performance was evaluated on the largest Chromosome 1 using 12 scoring metrics on the ENCODE Imputation Challenge testing set of 51 mark-cell type pairs. For each metric, the paired one-sided Wilcoxon signed-rank test was performed between two methods across 51 testing mark-cell type pairs. The significantly different ones are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).

**Supplementary Figure 34: Predictive performance comparison between ensemble models with and without neural network models**

Predictive performance was evaluated on the largest Chromosome 1 using 12 scoring metrics on the ENCODE Imputation Challenge testing set of 51 mark-cell type pairs. For each metric, the paired one-sided Wilcoxon signed-rank test was performed between two methods across 51 testing mark-cell type pairs. The significantly different ones are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
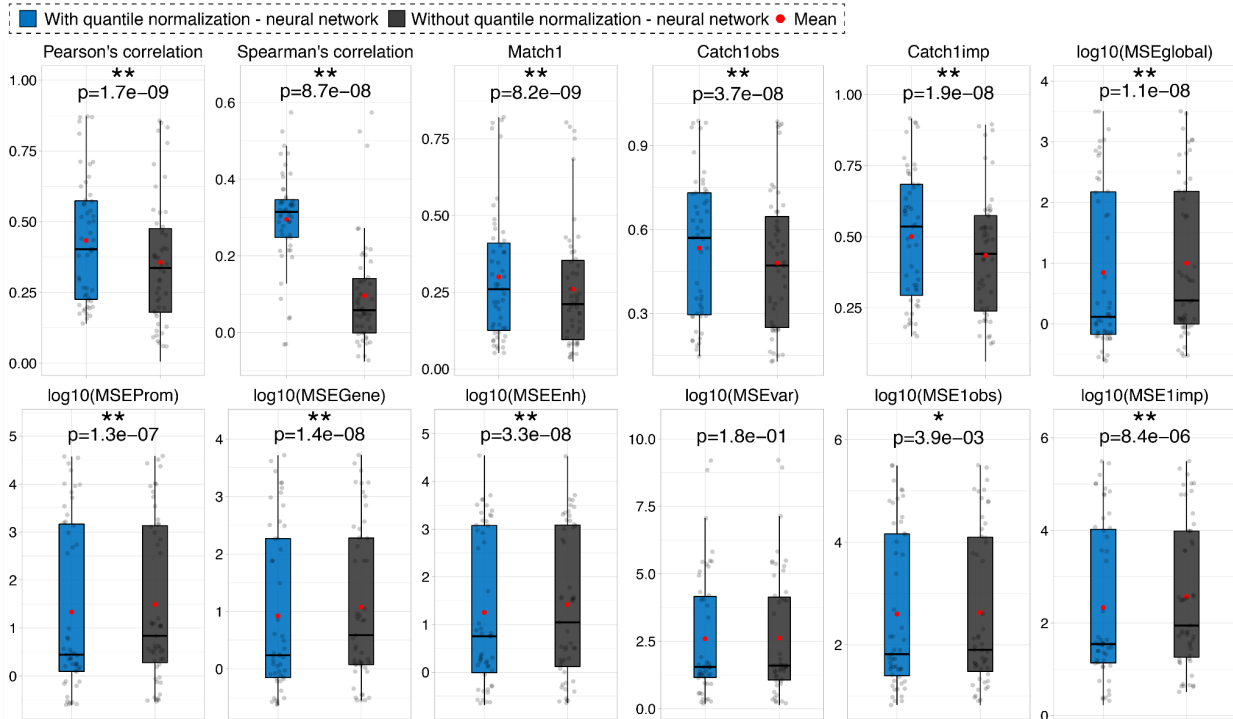
**Supplementary Figure 35: Predictive performance comparison between lightGBM models with and without quantile normalization**

Predictive performance was evaluated on the largest Chromosome 1 using 12 scoring metrics on the ENCODE Imputation Challenge testing set of 51 mark-cell type pairs. For each metric, the paired one-sided Wilcoxon signed-rank test was performed between two methods across 51 testing mark-cell type pairs. The significantly different ones are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
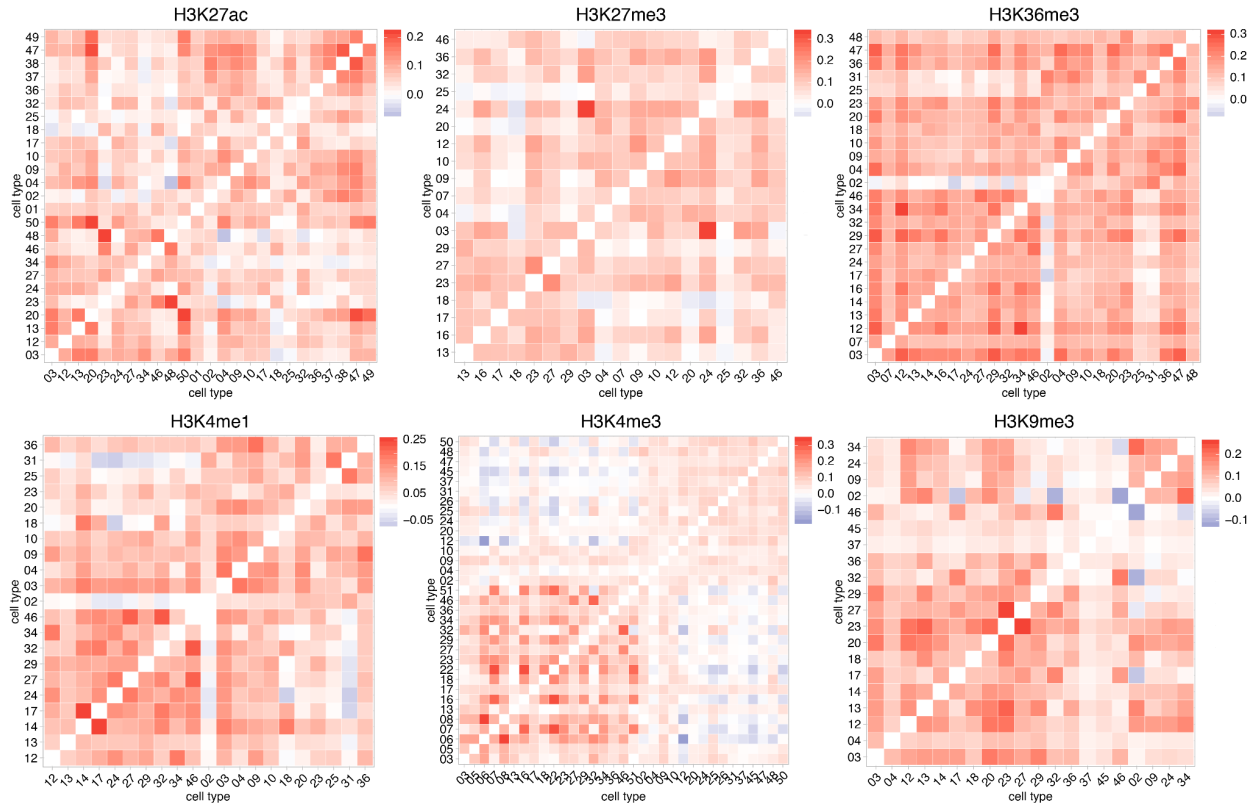
**Supplementary Figure 36: Predictive performance comparison between neural network models with and without quantile normalization**

Predictive performance was evaluated on the largest Chromosome 1 using 12 scoring metrics on the ENCODE Imputation Challenge testing set of 51 mark-cell type pairs. For each metric, the paired one-sided Wilcoxon signed-rank test was performed between two methods across 51 testing mark-cell type pairs. The significantly different ones are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).
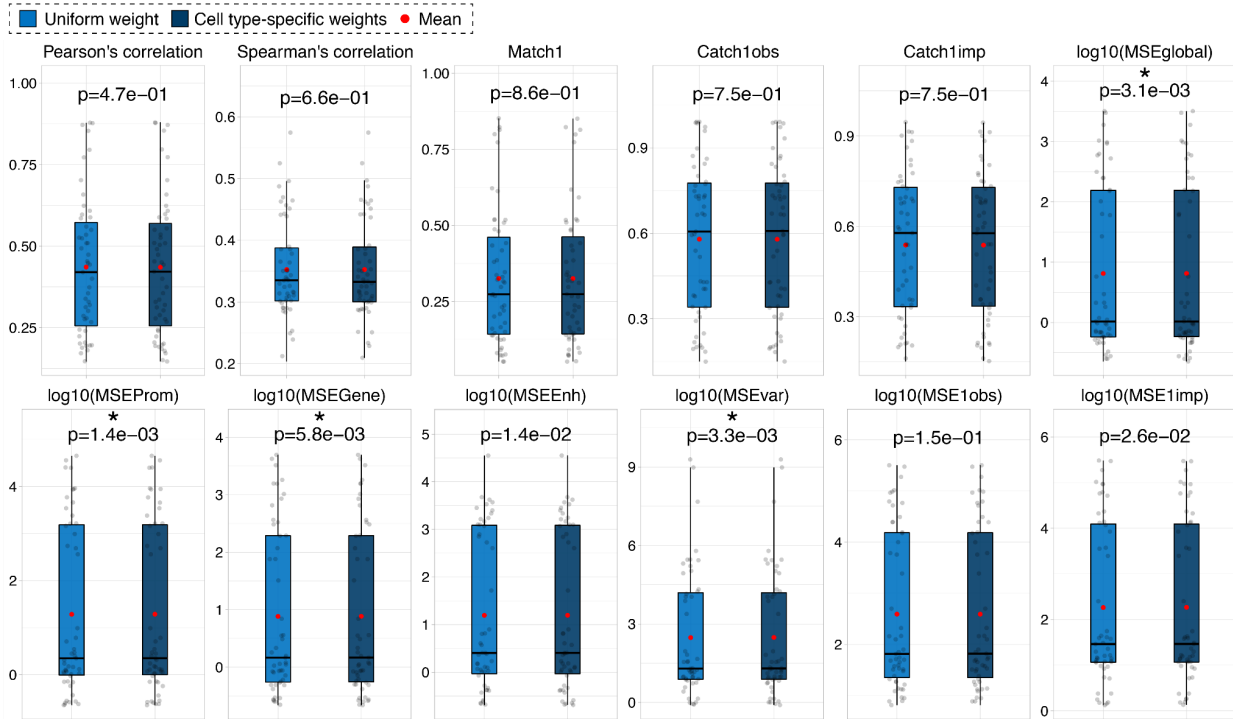
**Supplementary Figure 37: The pairwise Pearson's correlations of epigenomic signals among multiple training cell types.**

For each histone mark and each pair of available cell types, the pairwise Pearson correlation of epigenomics signals was calculated. Positive and negative correlations are shown in red and blue respectively.

**Supplementary Figure 38: Predictive performance comparison between ensemble predictions based on a uniform weight and ensemble predictions based on cell type-specific weights.**

The cell type-specific weights were defined by Pearson's correlations among cell types of marks that were used as features in machine learning models. Predictive performance was evaluated on the largest Chromosome 1 using 12 scoring metrics on the ENCODE Imputation Challenge testing set of 51 mark-cell type pairs. For each metric, the paired one-sided Wilcoxon signed-rank test was performed between two methods across 51 testing mark-cell type pairs. The significantly different ones are labeled by asterisks (* p-value < 0.01 and ** p-value < 0.001).

**Supplementary Table Legends**

**Supplementary Table 1: The SHAP values from 10 bootstrap samplings.**

The results from 10 bootstraps are shown in rows labeled by "bootstrap 0", "bootstrap 1", …, "bootstrap 9". The average values are shown in the last 10 rows labeled by "average". The histone marks in the first row represent that they are serving as the target being predicted and the histone marks along the column represent that they are serving as predictors to predict others. The distances are calculated from the target bin to be predicted.

**Supplementary Table 2:  Predictive performance comparison between Ocelot, ChromImpute and Avocado.**

**Supplementary Table 3:  Predictive performance comparison between lightGBM models with and without DNA sequence features.**

**Supplementary Table 4:  Predictive performance comparison between neural network (NN) models with and without DNA sequence features.**

**Supplementary Table 5:  Predictive performance comparison between Ocelot without DNA sequence features, ChromImpute and Avocado.**

**Supplementary Table 6: The data matrix of the ENCODE Imputation Challenge.**

The marks are shown in the column and the tissue and cell types are shown in the row. The green blocks with "T" are the training data and the red blocks with "B" are the blind testing data. The other empty blocks are unavailable mark-cell type combinations.

**Supplementary Table 7: The model design and ensemble weights of Ocelot final submission in the ENCODE Imputation Challenge.**

Based on the availability of training data in **Supplementary Table 6**, we design a specific model that includes as many feature marks as possible in predicting a target mark in a specific cell type. The ensemble weights of (1) lightGBM, (2) Neural Network (NN), and (3) the quantile normalized average signal of a mark in all training cell types were selected based on the cross-validation result on the training data during the challenge.

**Supplementary Table 8: The SHAP values of 5 subsets.**

**Supplementary Table 9: The PPI values of 5 subsets.**

**Supplementary Table 10: Predictive performance comparison between ensemble models, ensemble models without lightGBM models, and ensemble models without neural network (NN) models.**

**Supplementary Table 11: Predictive performance comparison between lightGBM models with and without quantile normalization (QN).**

**Supplementary Table 12: Predictive performance comparison between neural network (NN) models with and without quantile normalization (QN).**

**Supplementary Table 13: Predictive performance comparison between lightGBM predictions based on a uniform weight and lightGBM predictions based on cell type-specific weights.**

**Supplementary Table 14: The predictive performance of the top two solutions in the ENCODE Imputation Challenge.**

Nine columns represent nine evaluation metrics used in the challenge and the median scores across 51 testing mark-cell type pairs are listed. The percentage improvements (increase in correlation metrics and decrease in MSE metrics) are shown in the third row.

**Supplementary Table 15: The accession numbers of the ENCODE3 data used in this study.**