

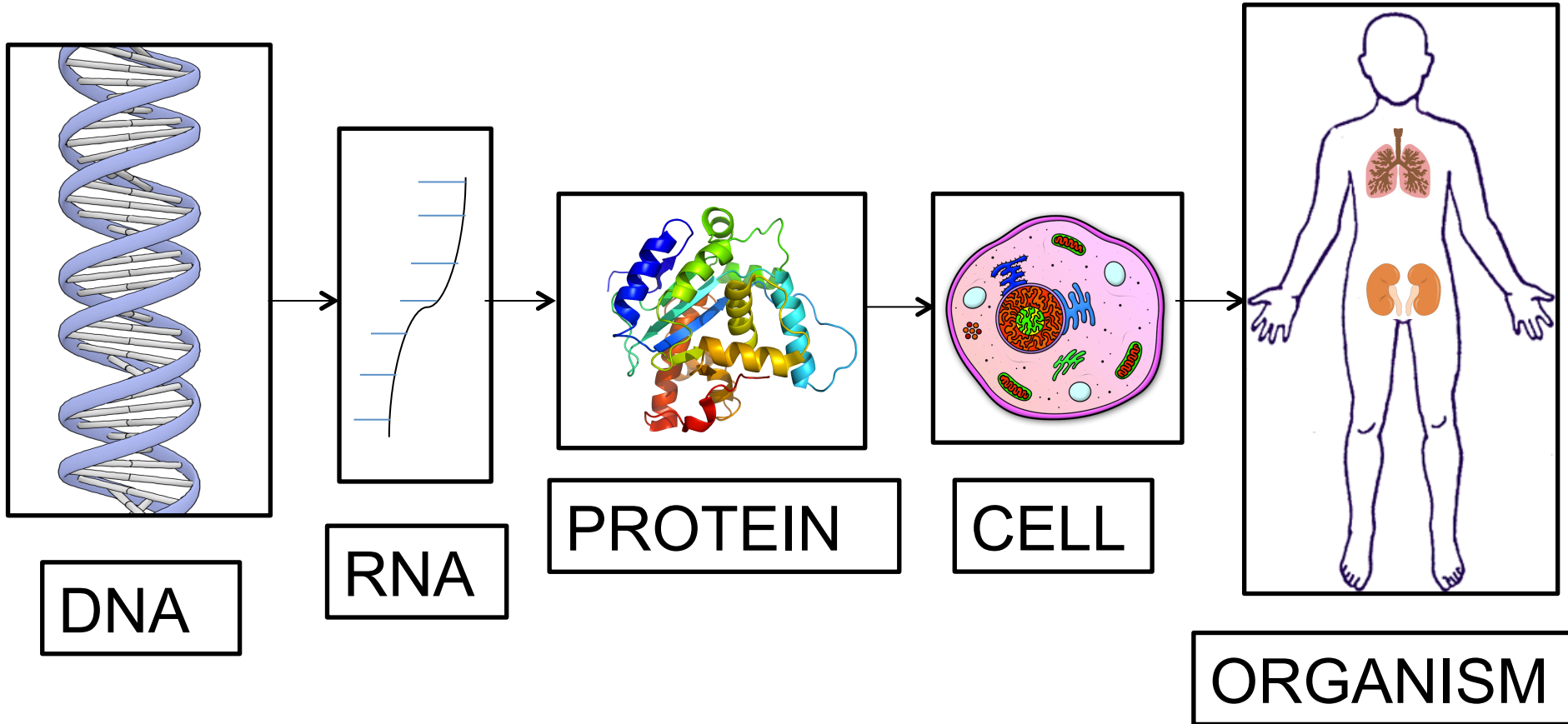


DeepChrome

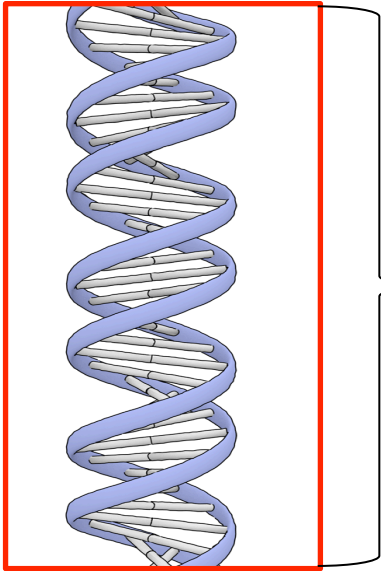
Deep-learning for predicting gene expression from histone modification

Ritambhara Singh, Jack Lanchantin,
Gabriel Robins, and Yanjun Qi

Biology in a Slide



DNA and Diseases



DNA

- Down Syndrome
- Parkinson's Disease
- Autism
- Muscular Atrophy
- Sickle Cell Disease

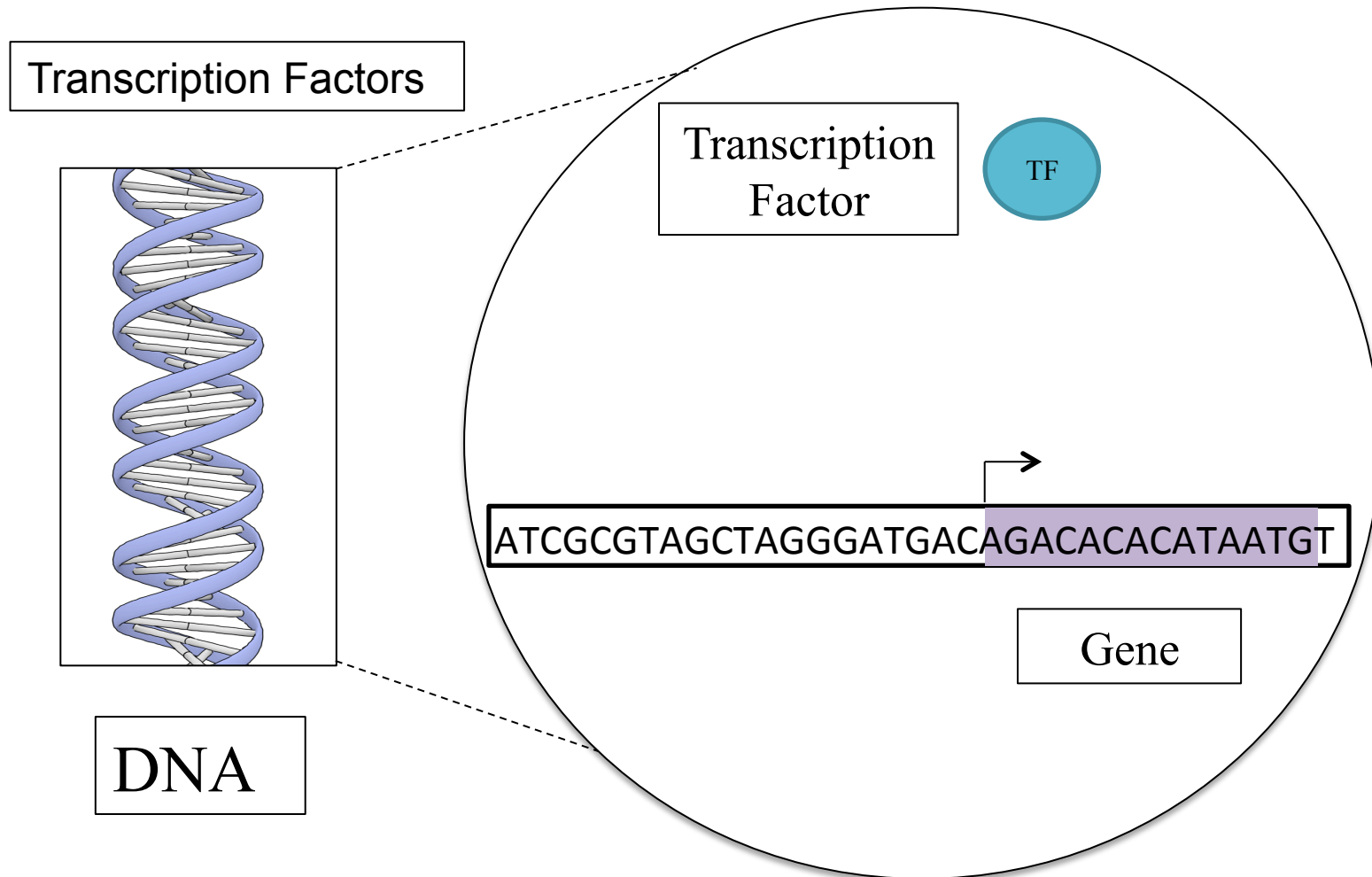
.....

.....



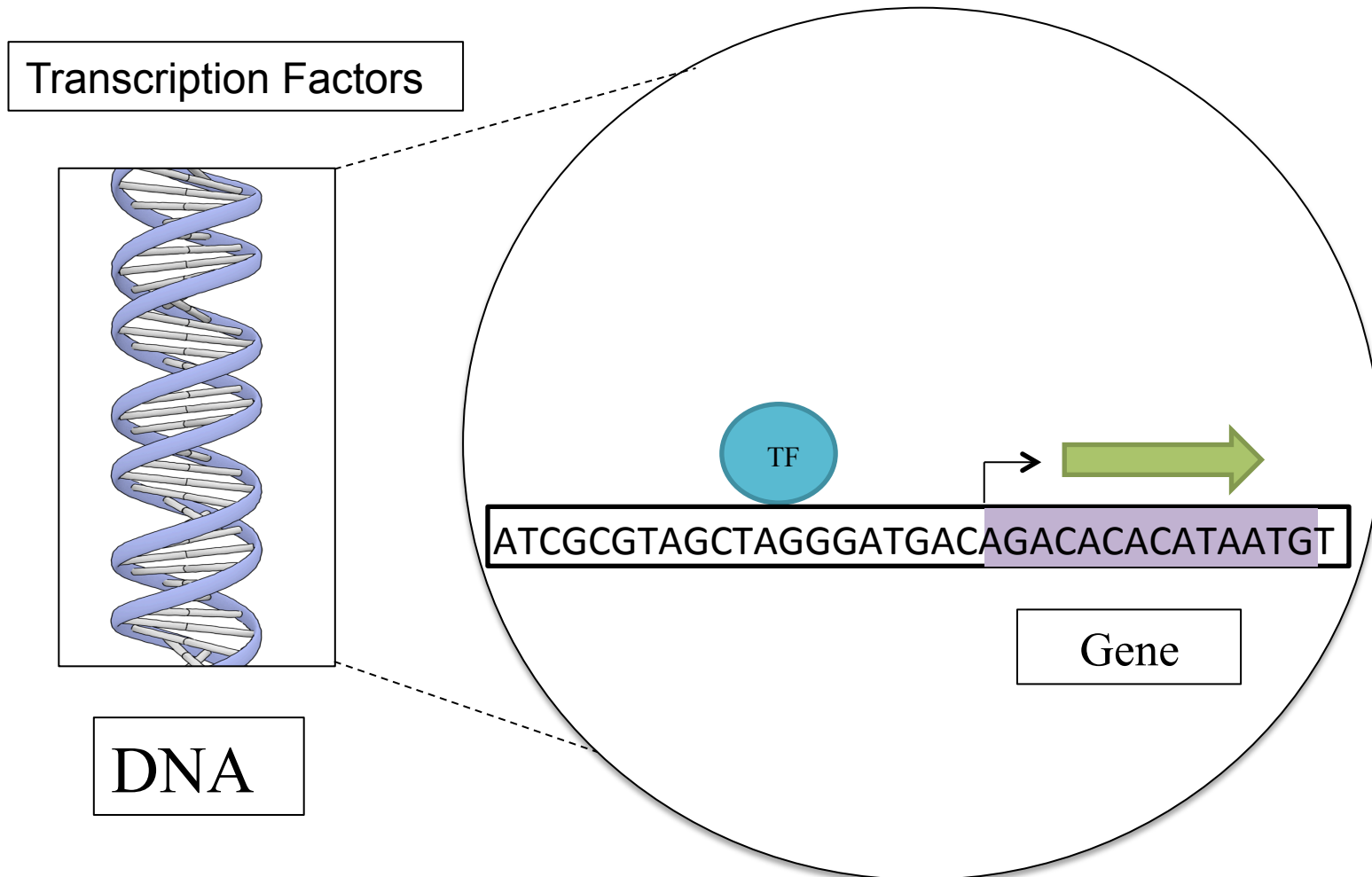
Epigenetic Factors

- Encyclopedia of DNA Elements (ENCODE)
- Roadmap Epigenetics Project (REMC)



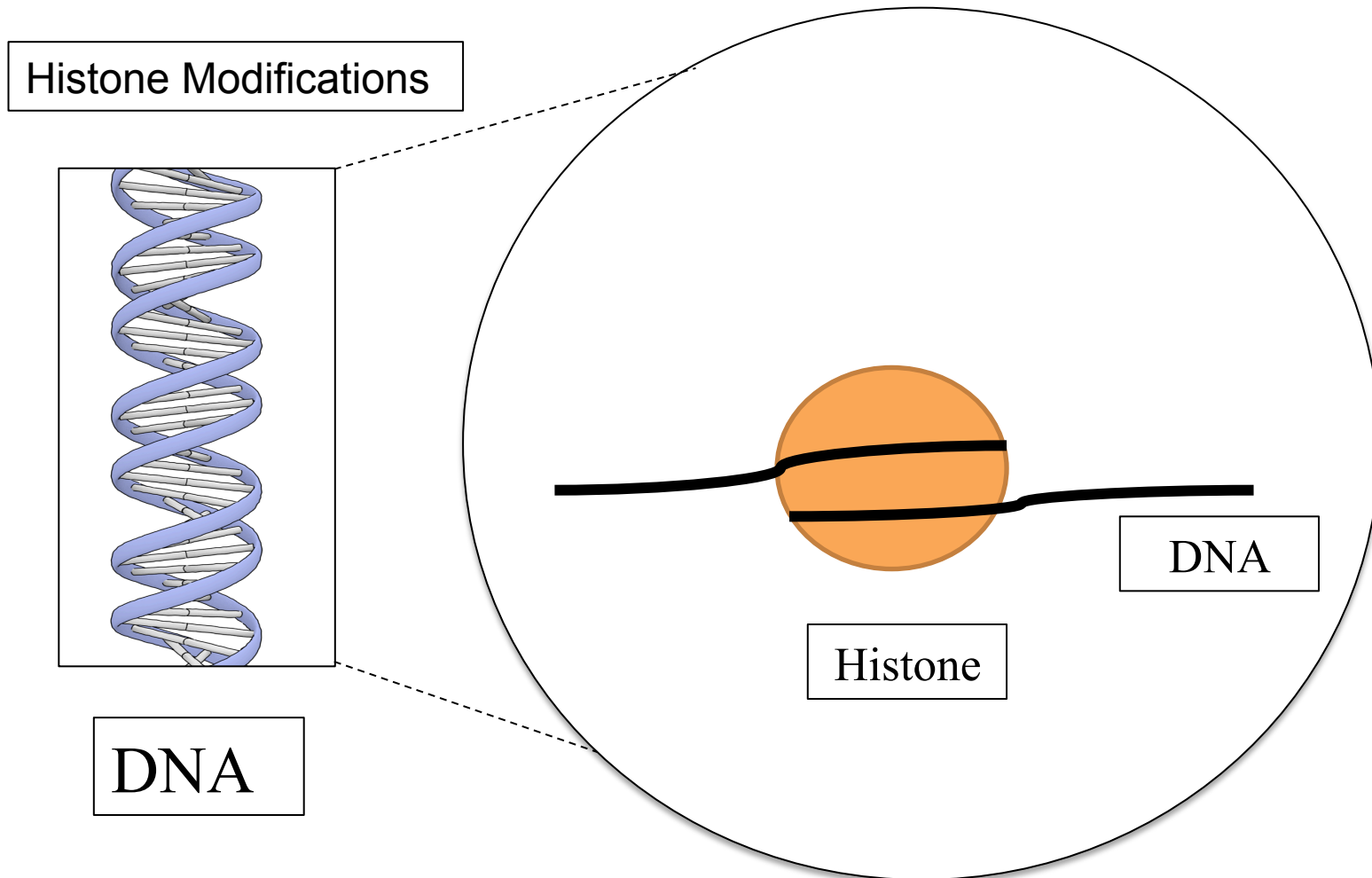
Epigenetic Factors

- Encyclopedia of DNA Elements (ENCODE)
- Roadmap Epigenetics Project (REMC)



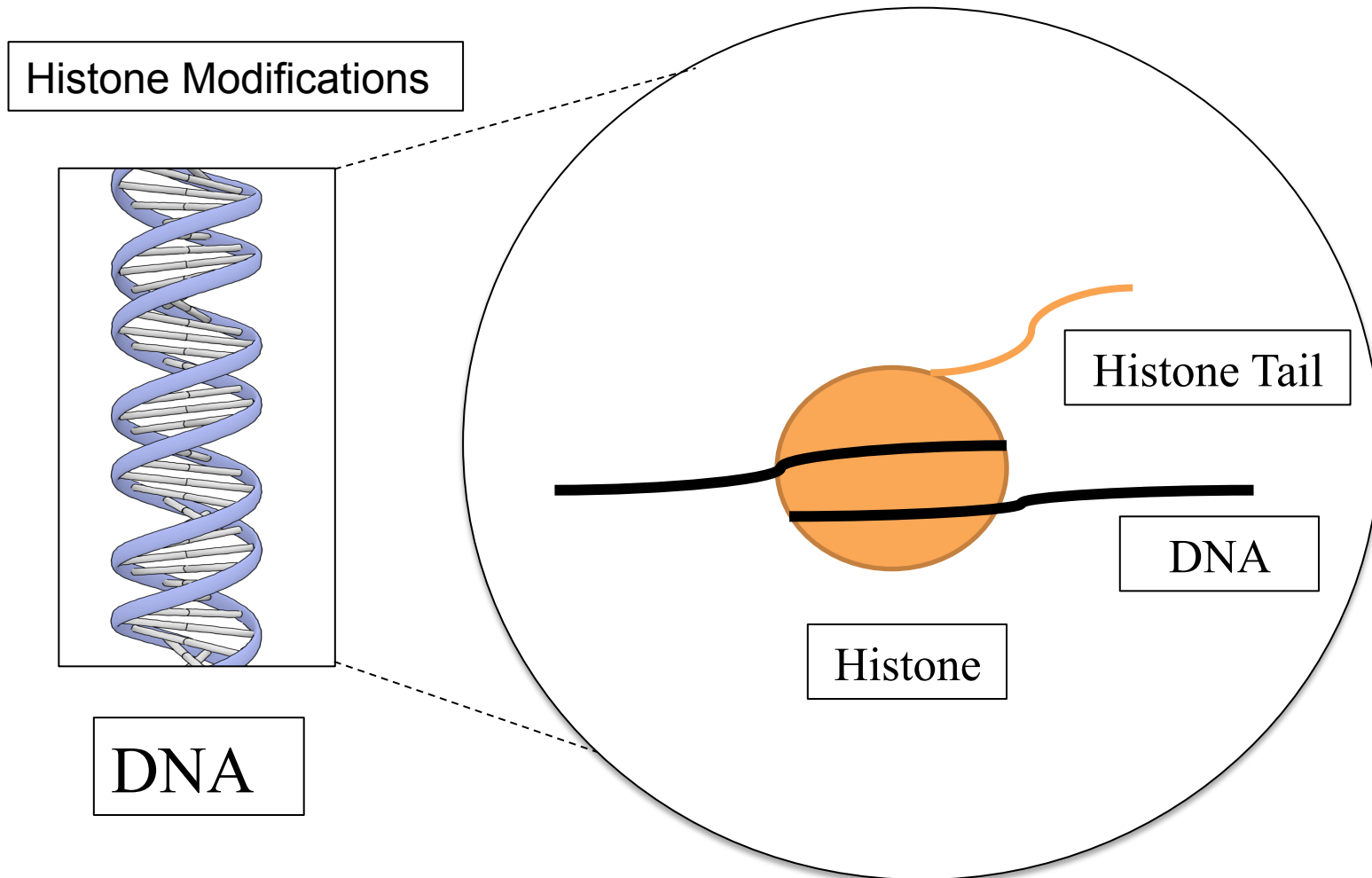
Epigenetic Factors

- Encyclopedia of DNA Elements (ENCODE)
- Roadmap Epigenetics Project (REMC)



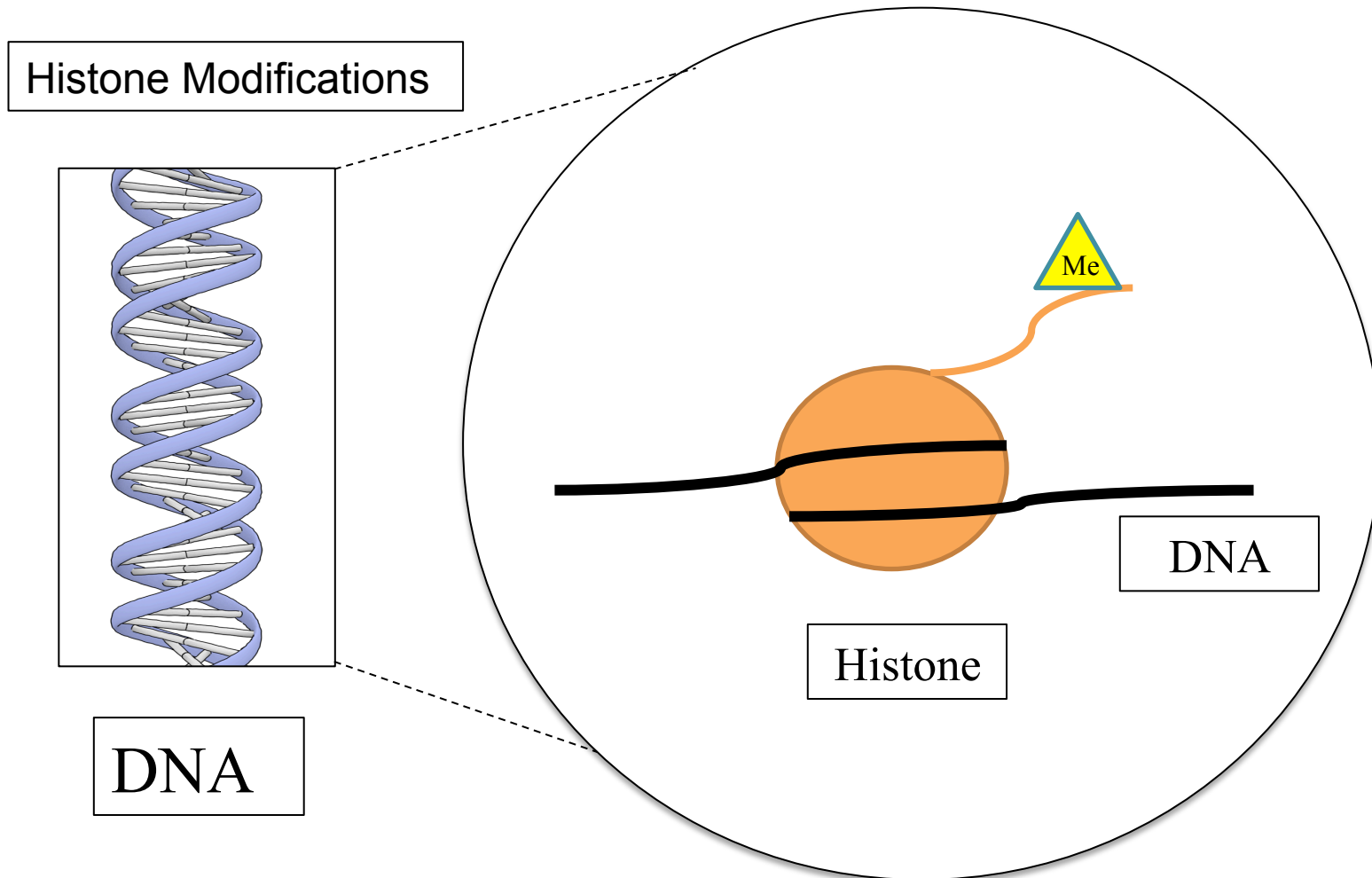
Epigenetic Factors

- Encyclopedia of DNA Elements (ENCODE)
- Roadmap Epigenetics Project (REMC)



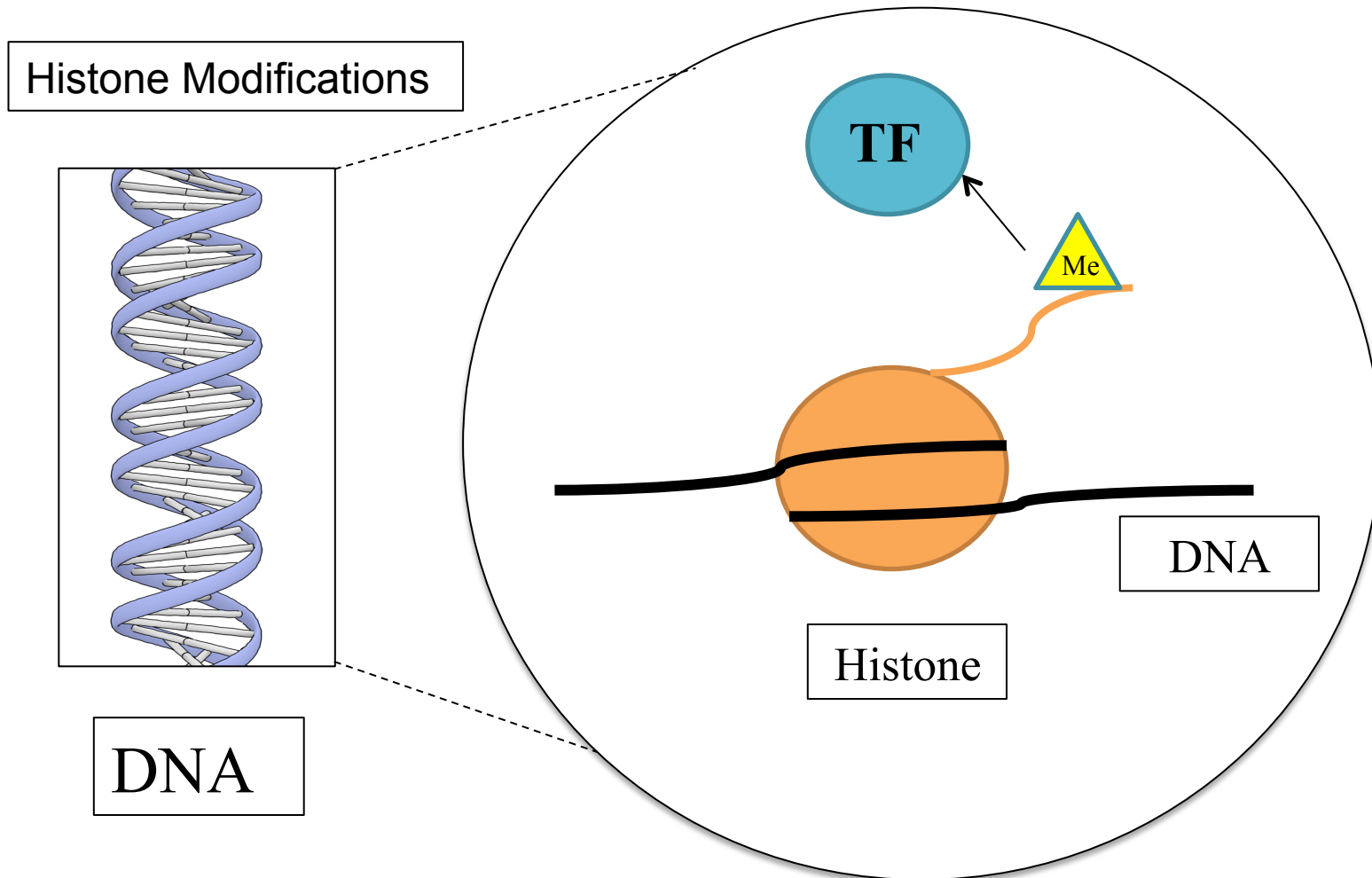
Epigenetic Factors

- Encyclopedia of DNA Elements (ENCODE)
- Roadmap Epigenetics Project (REMC)

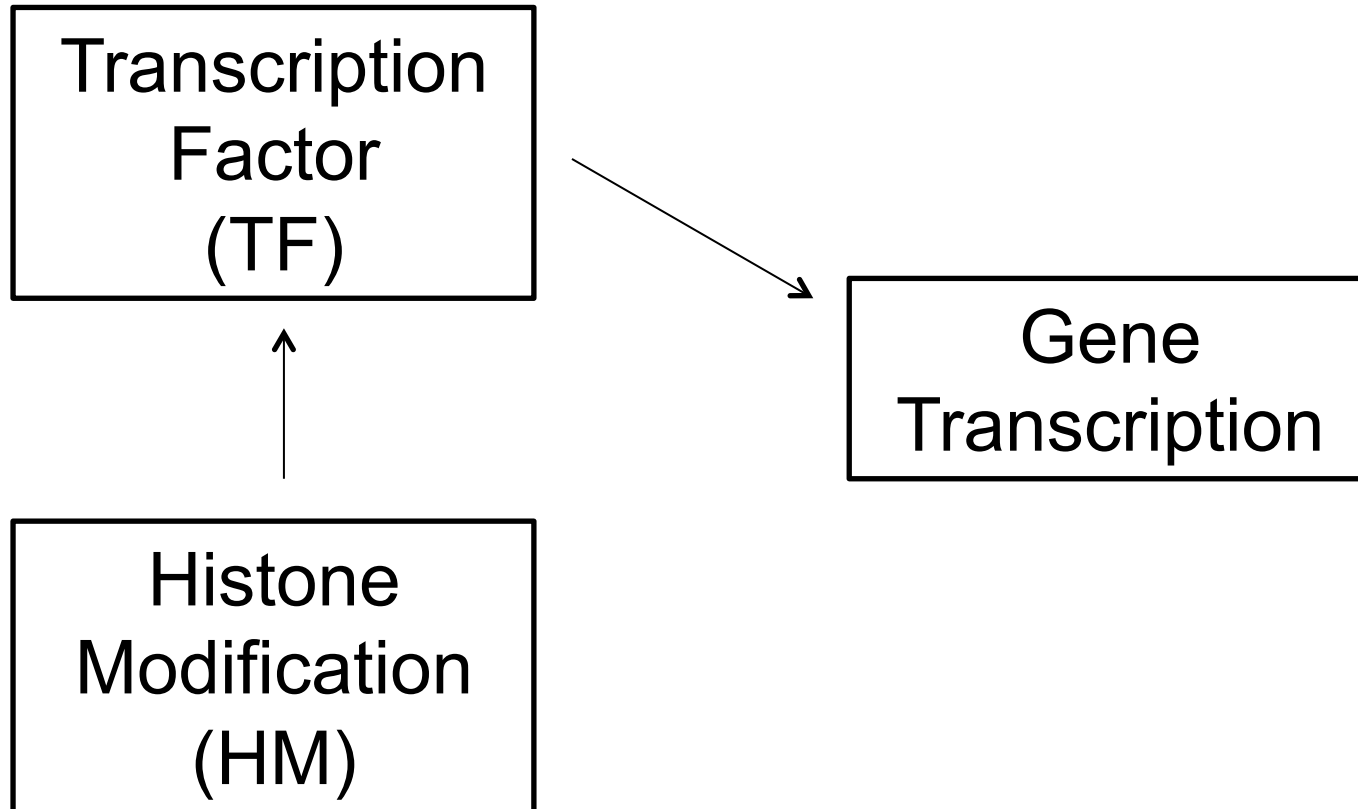


Epigenetic Factors

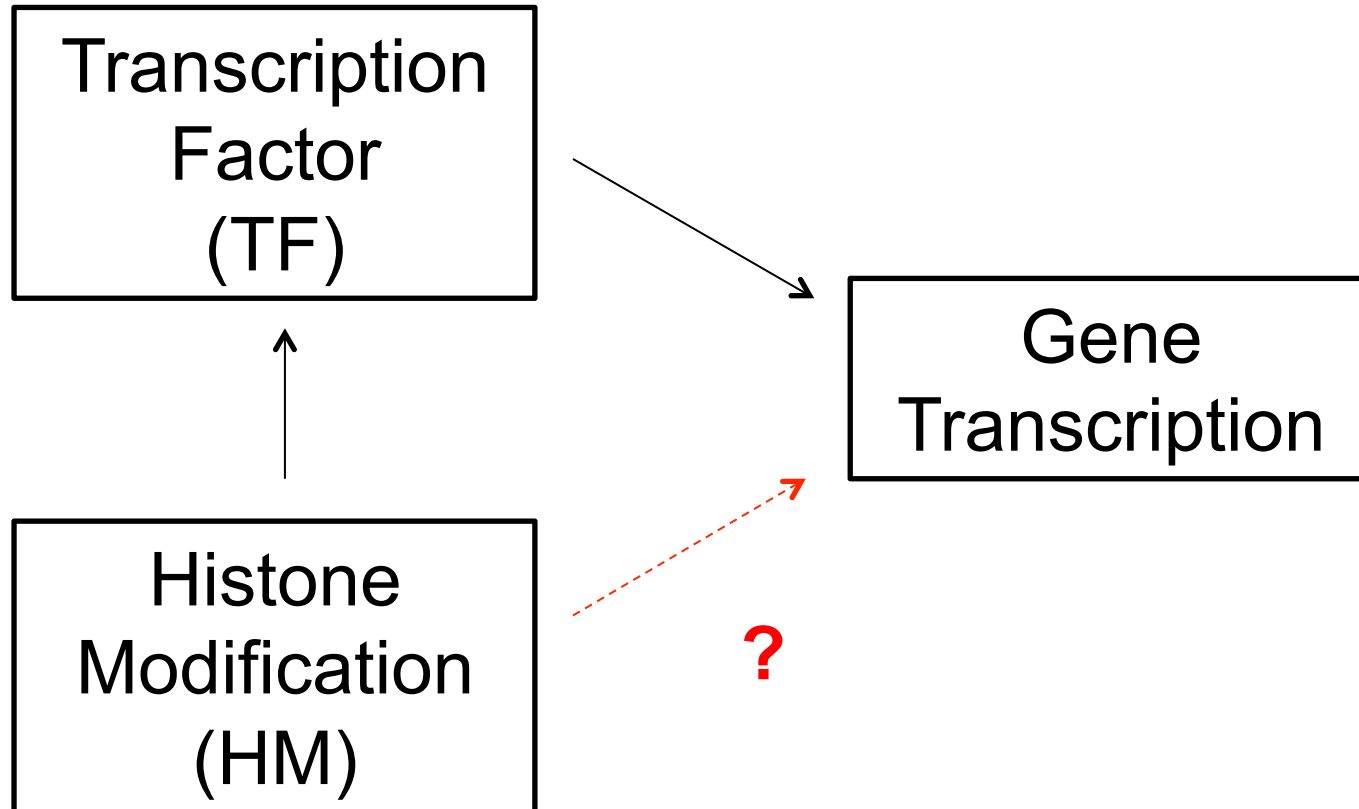
- Encyclopedia of DNA Elements (ENCODE)
- Roadmap Epigenetics Project (REMC)



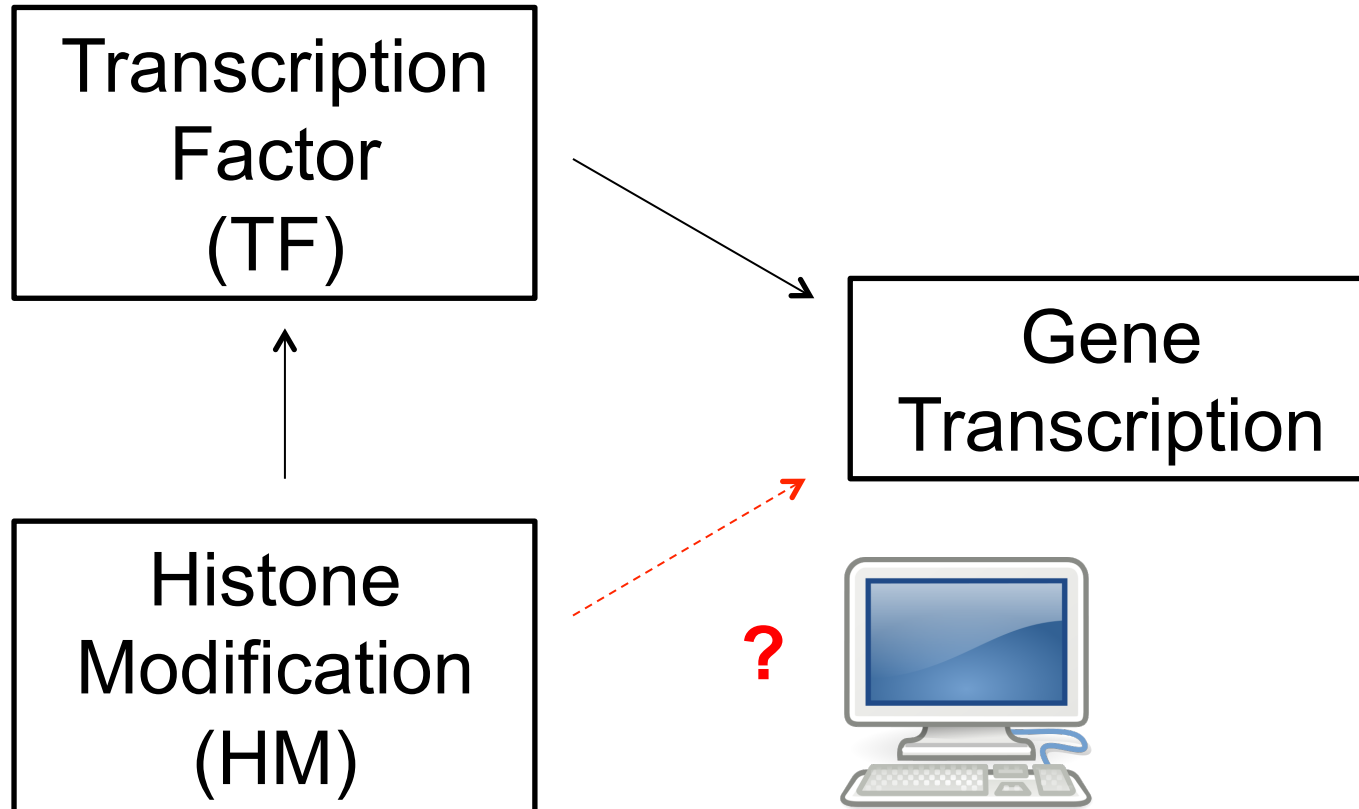
Histone Modification and Gene Transcription



Histone Modification and Gene Transcription



Histone Modification and Gene Transcription

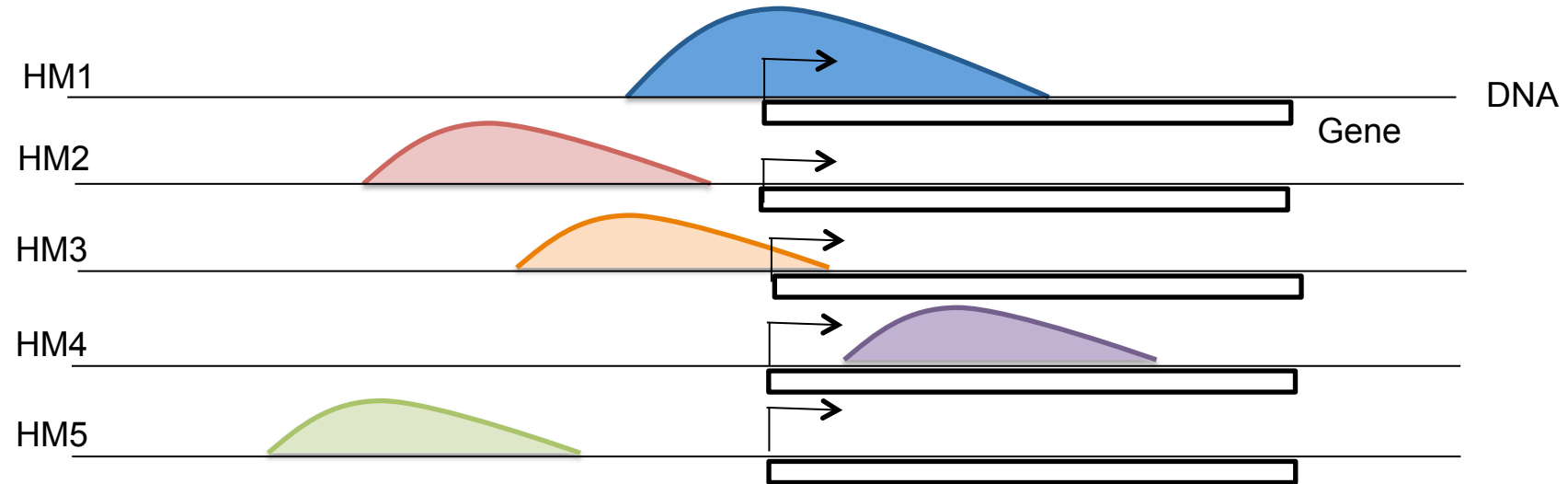


Task Formulation

Prediction Task

Input :

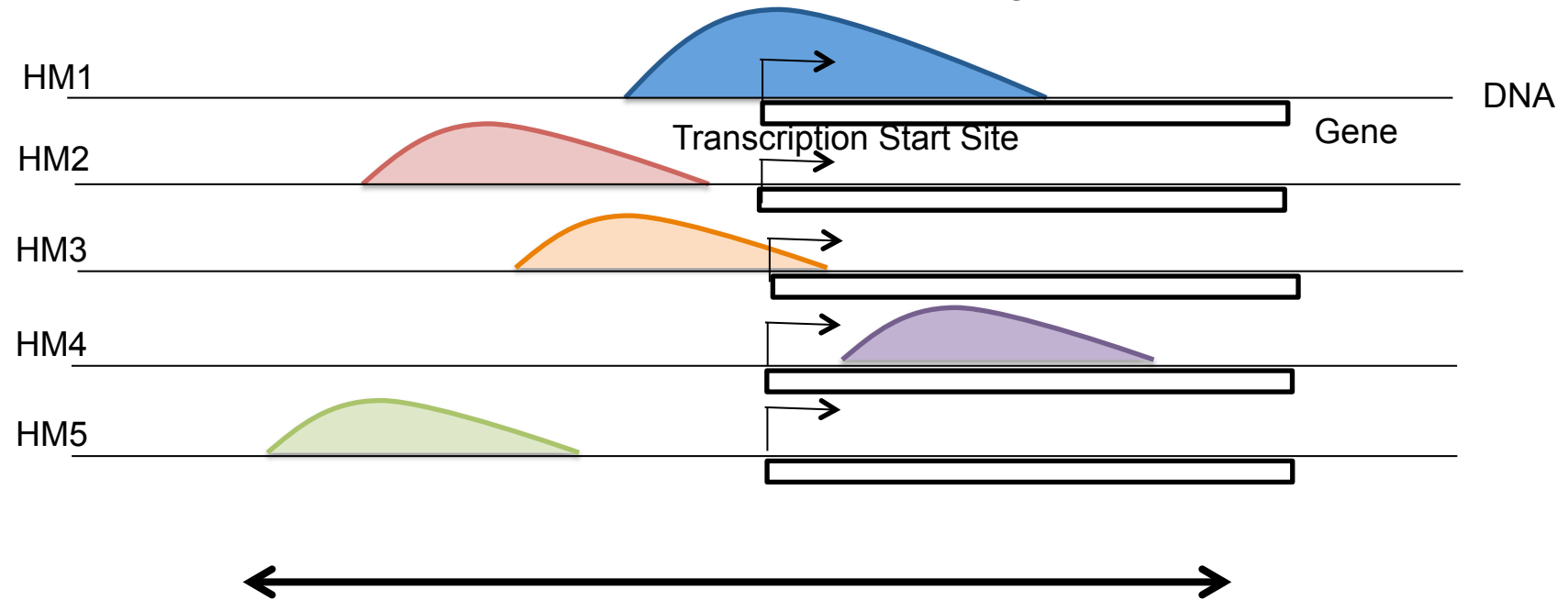
Histone Modification Signals



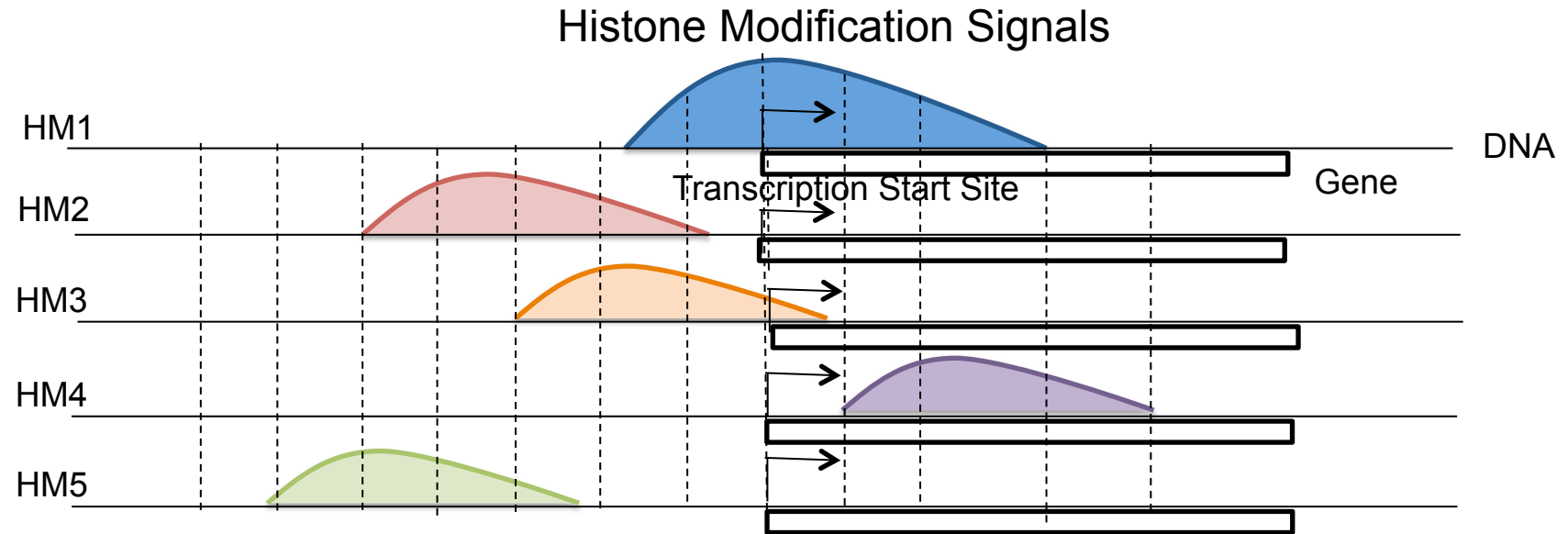
Output : Gene ON/OFF

Input

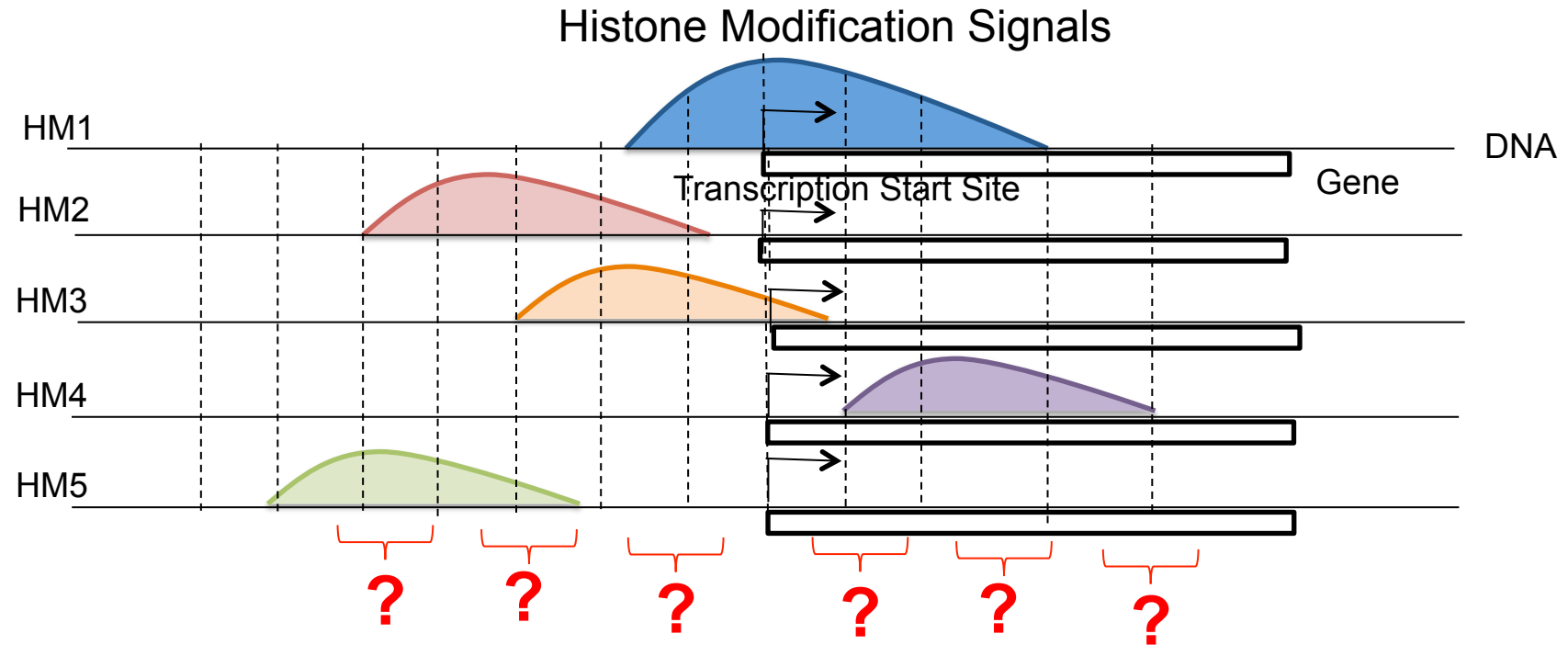
Histone Modification Signals



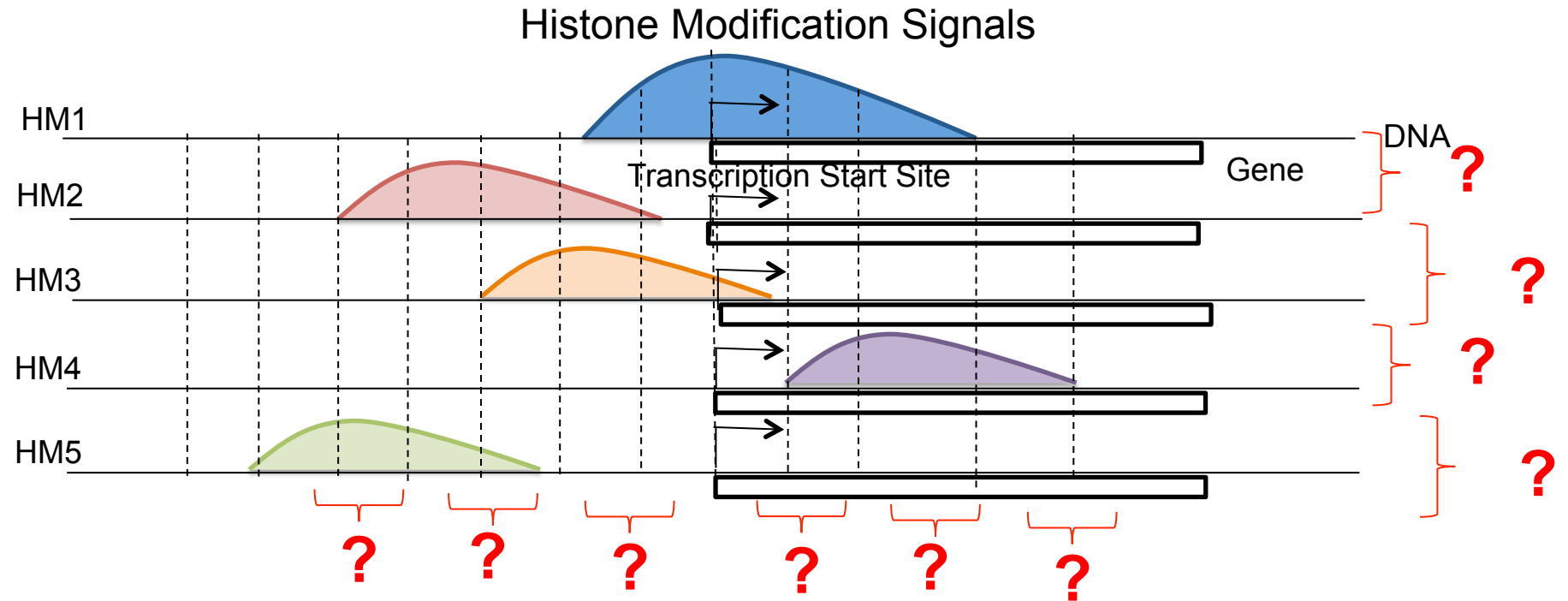
Input



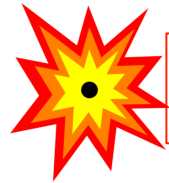
Challenge



Challenge



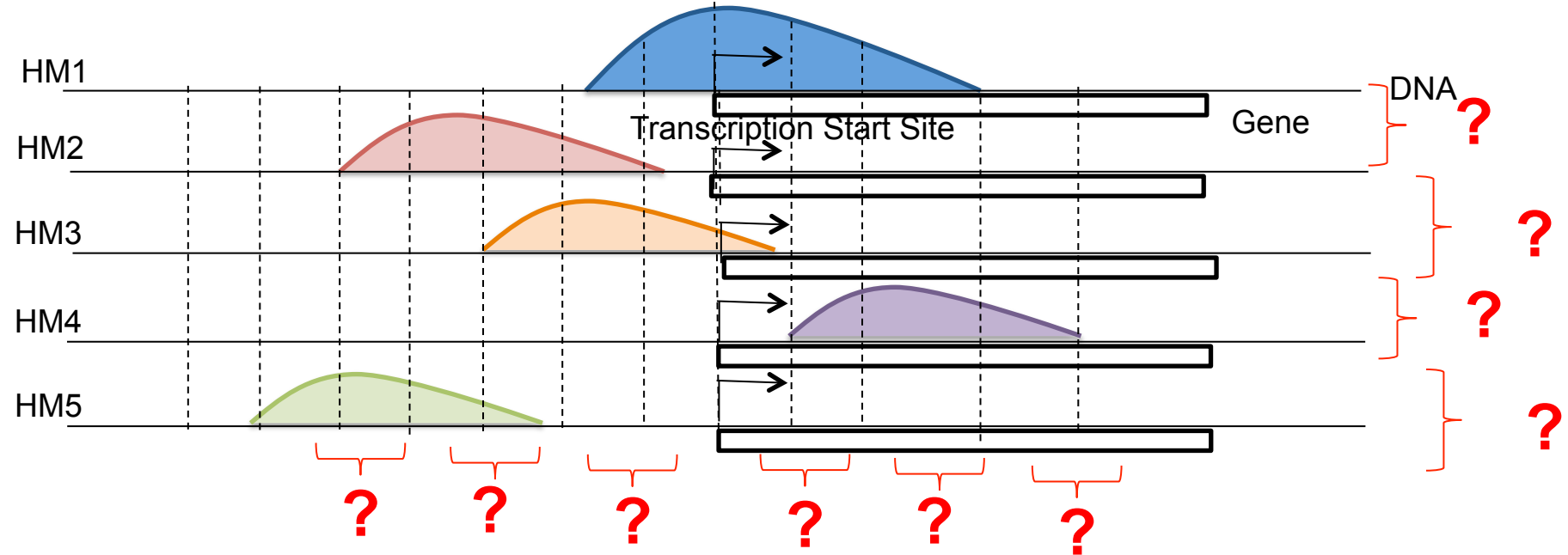
Challenge



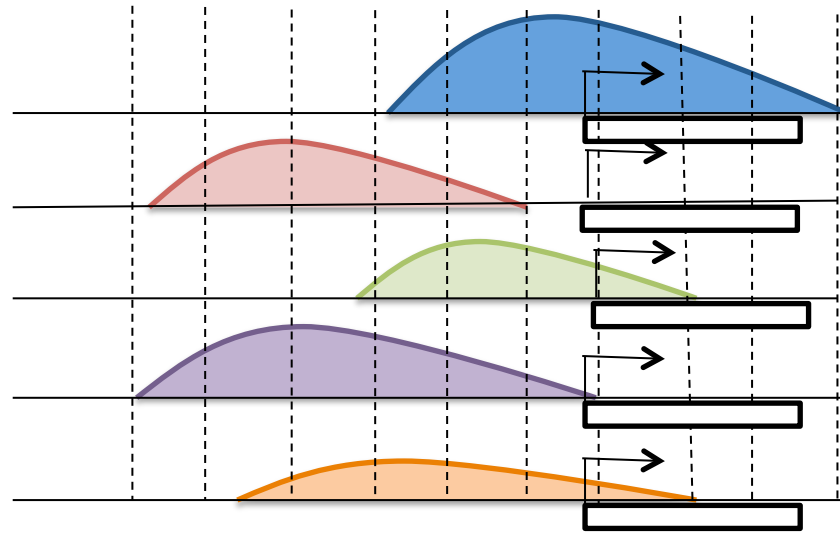
Search Space = 100^5



Histone Modification Signals



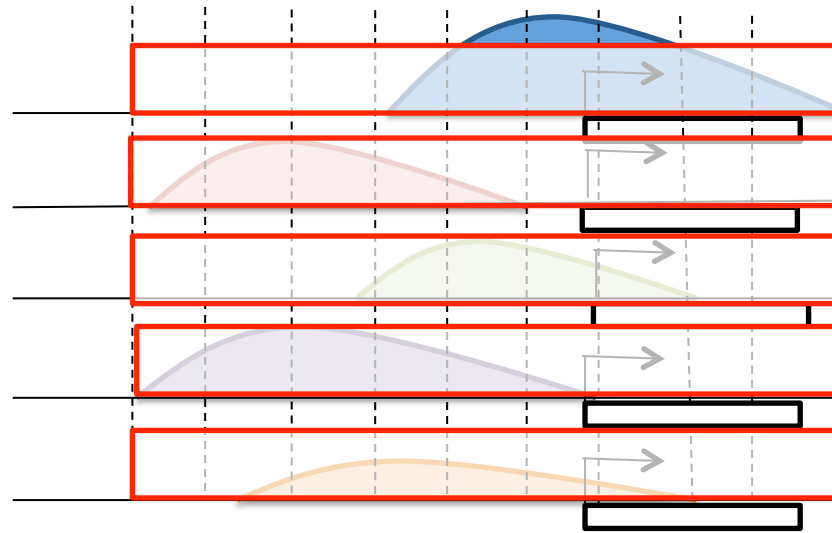
Related Work



**Linear Regression,
SVM,
Random Forest**

Gene ON/OFF

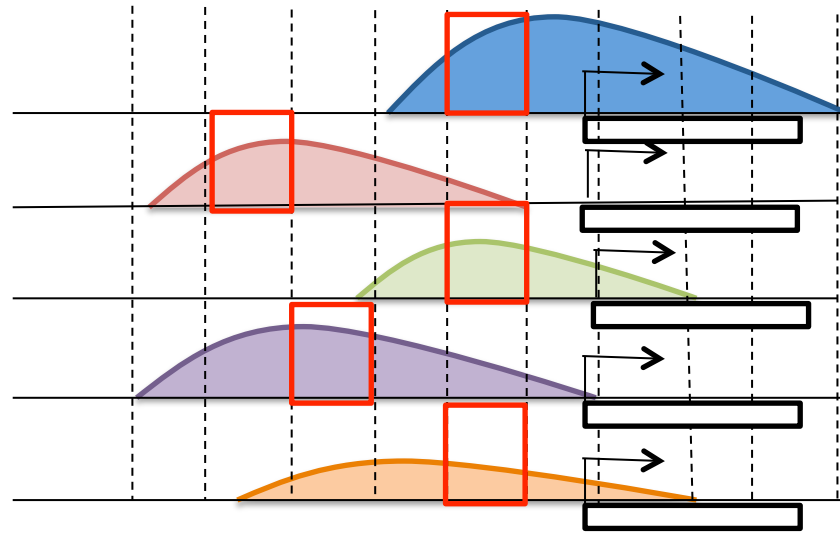
Related Work



**Linear Regression,
SVM,
Random Forest**

Gene ON/OFF

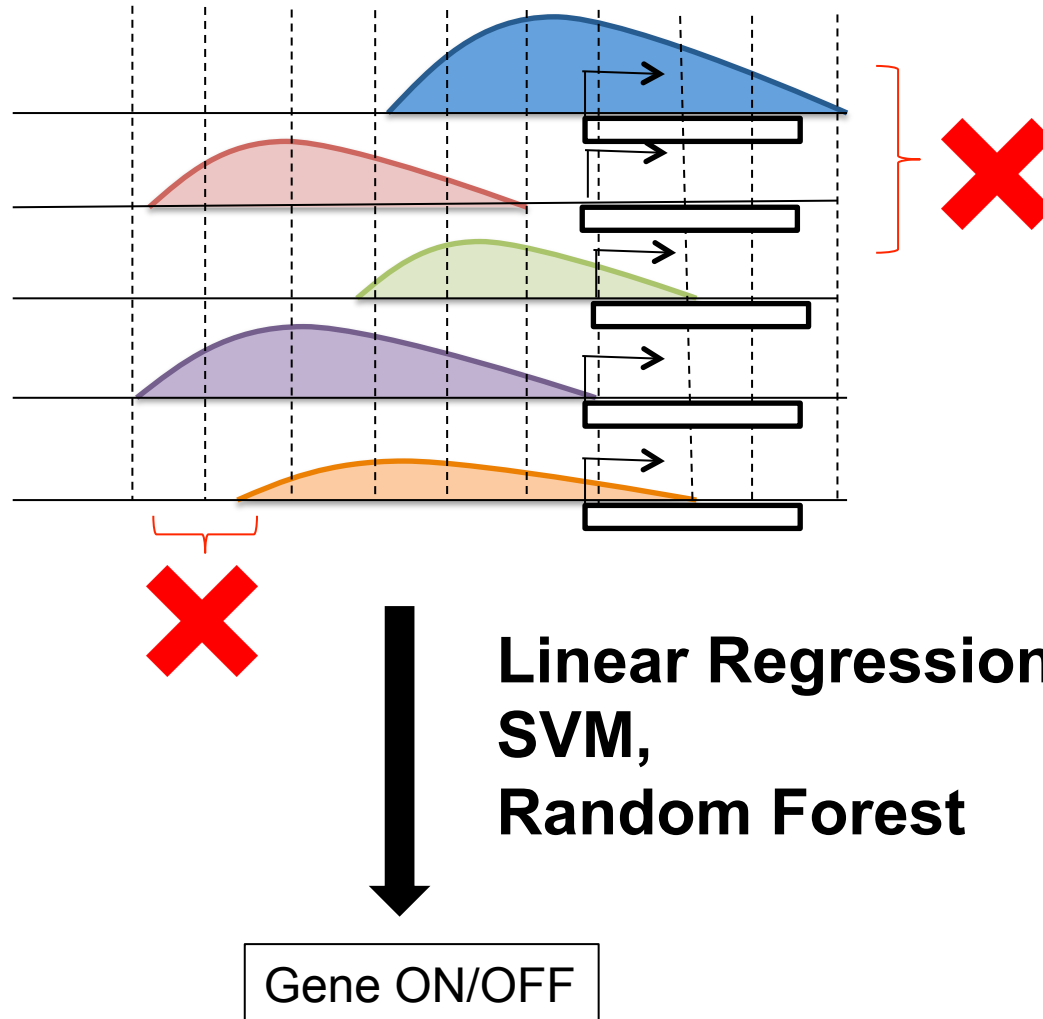
Related Work



**Linear Regression,
SVM,
Random Forest**

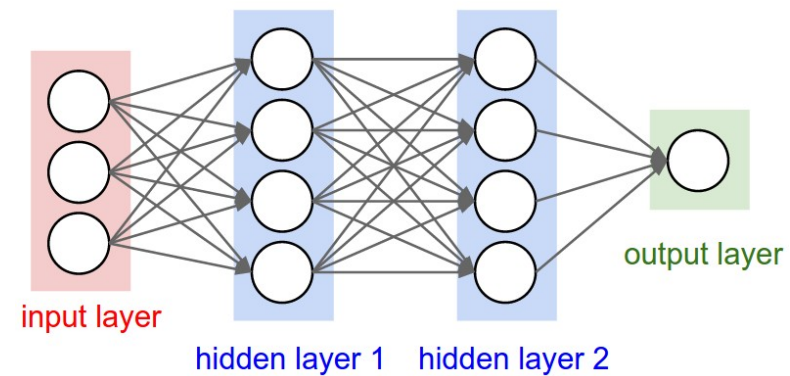
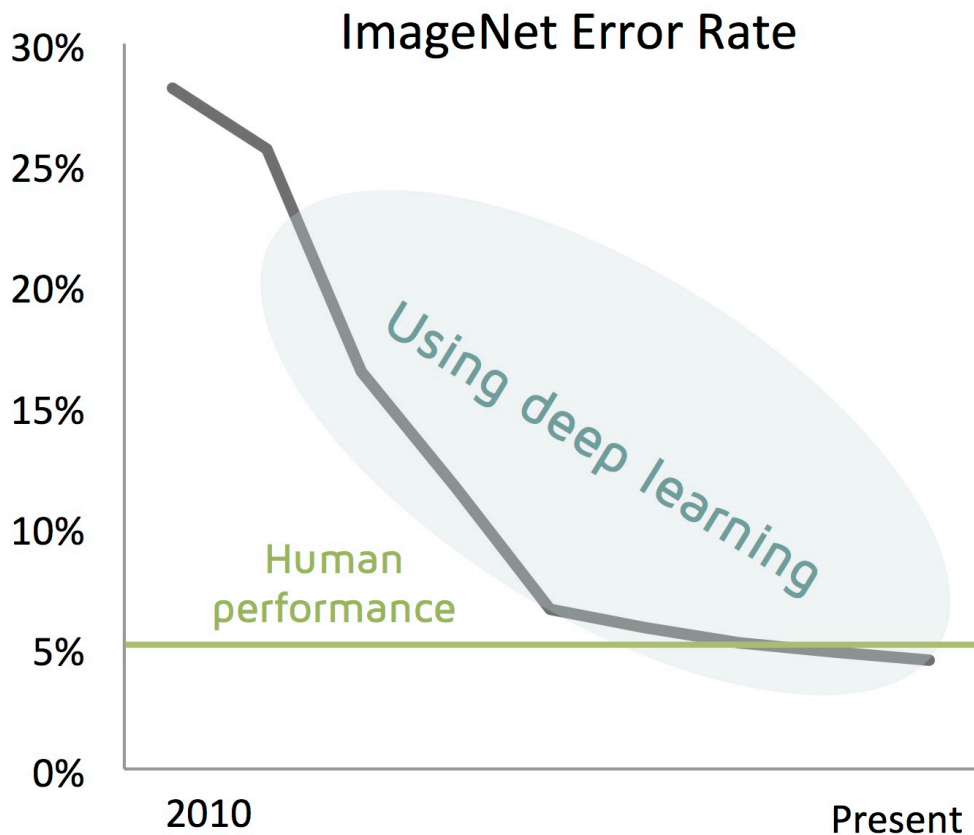
Gene ON/OFF

Drawback



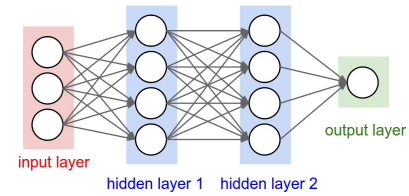
Solution

Convolutional Neural Network (CNN)



Solution

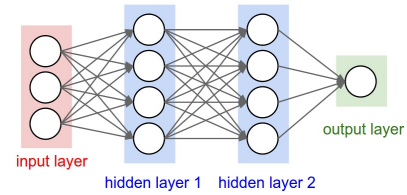
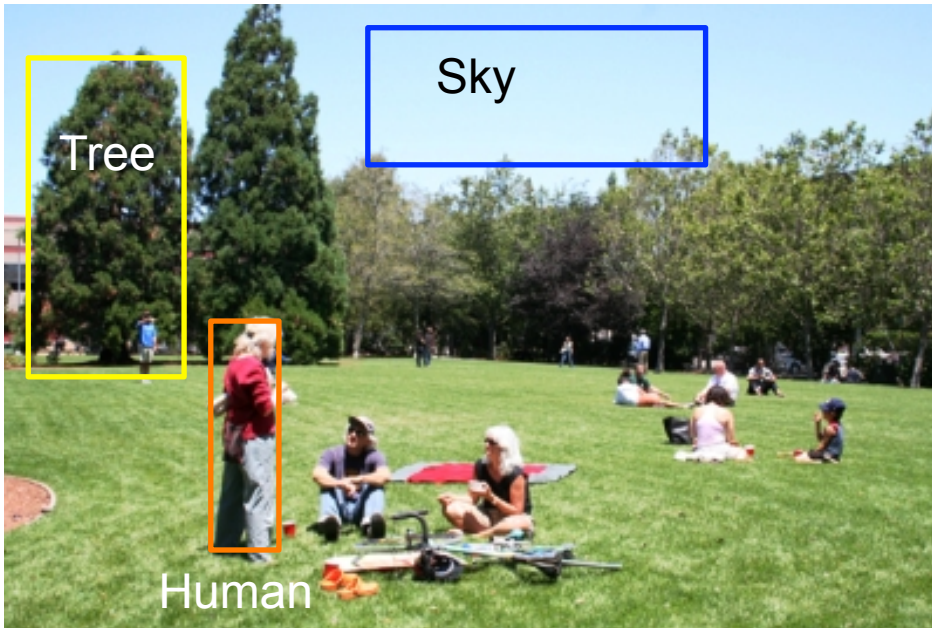
Convolutional Neural Network (CNN)



Park

Solution

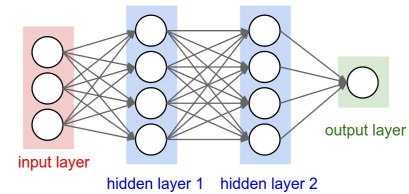
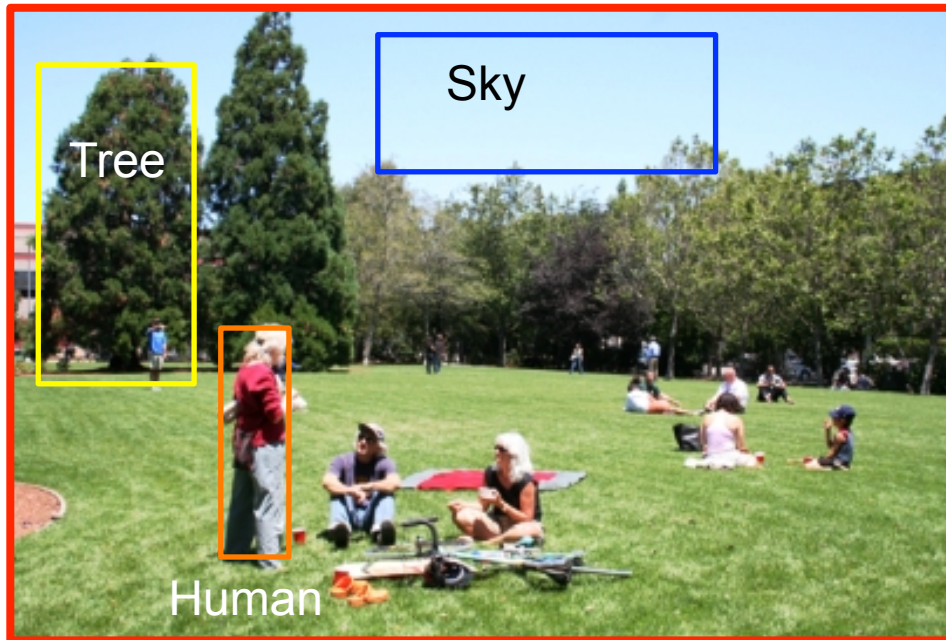
Convolutional Neural Network (CNN)



Park

Solution

Convolutional Neural Network (CNN)



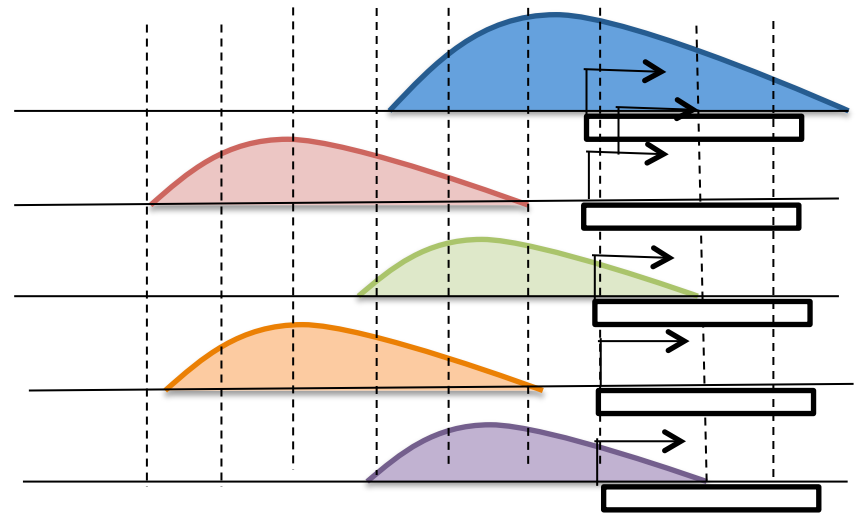
Park

Solution

Analogy to our task

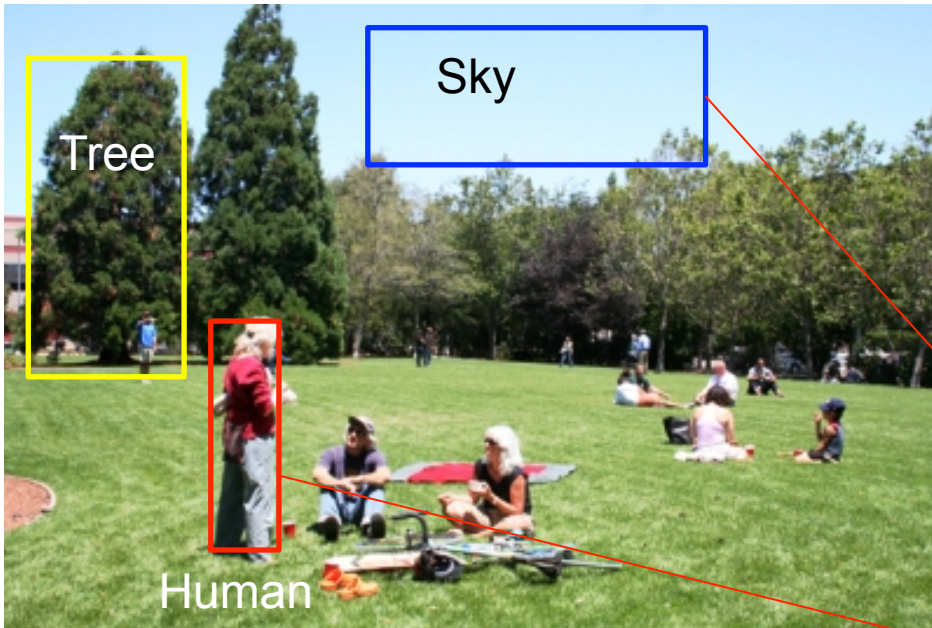


Histone Modification Signals

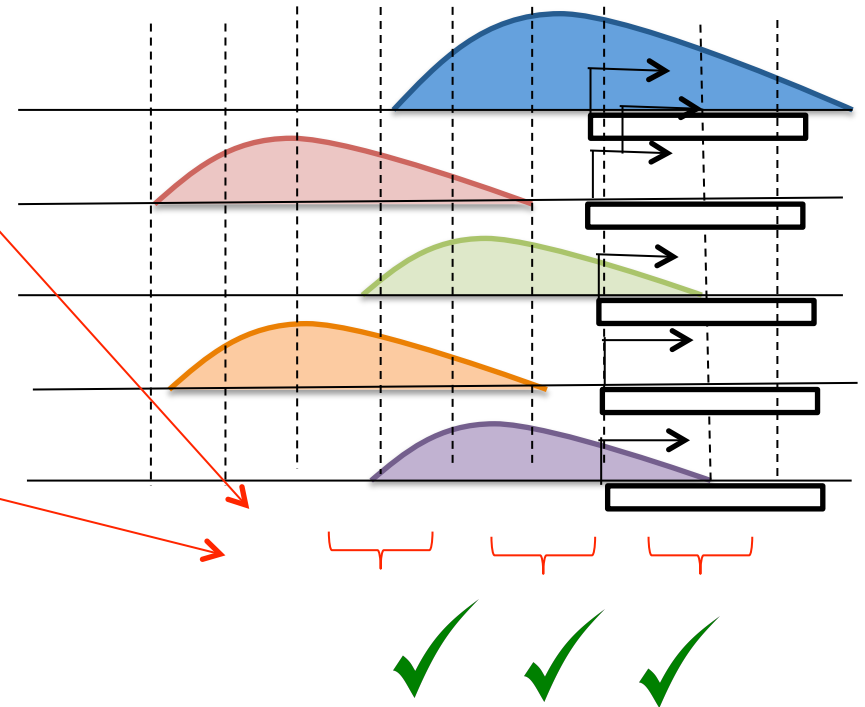


Solution

Analogy to our task

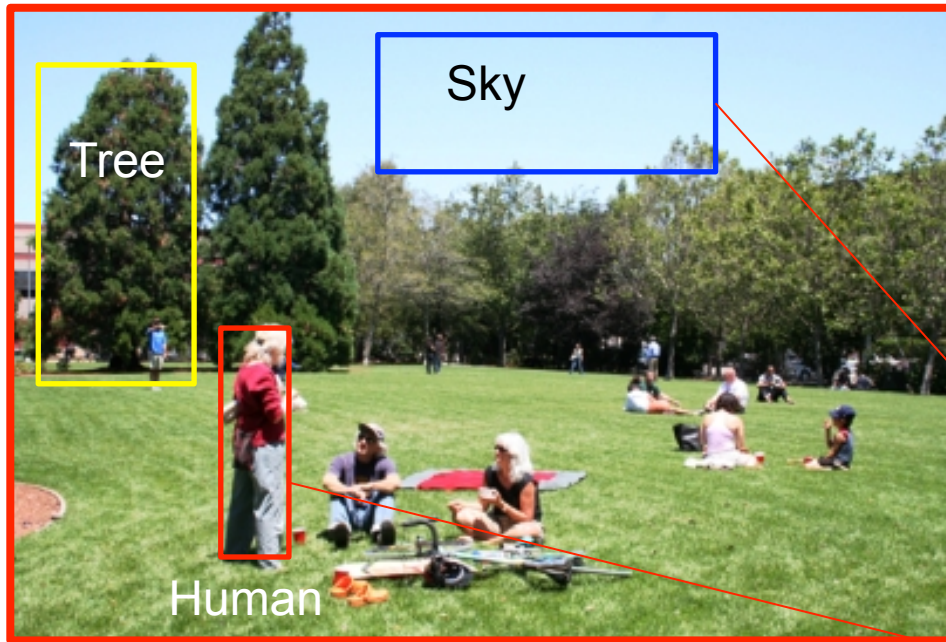


Histone Modification Signals

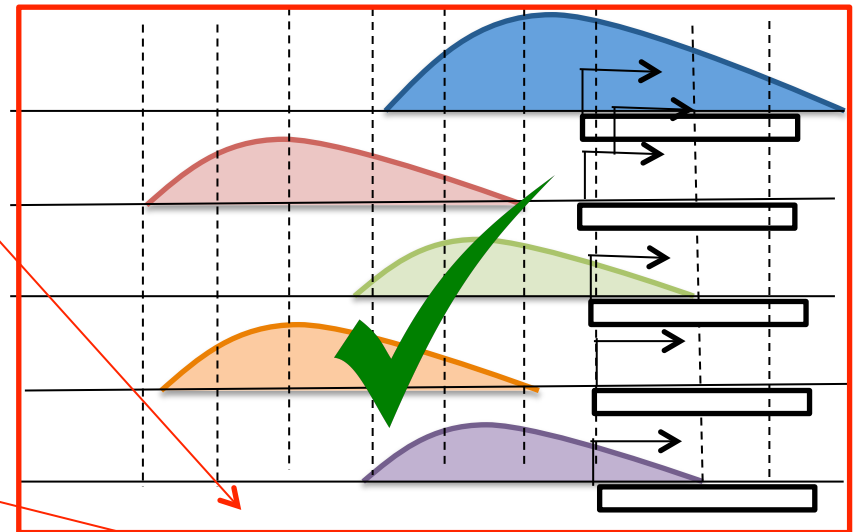


Solution

Analogy to our task

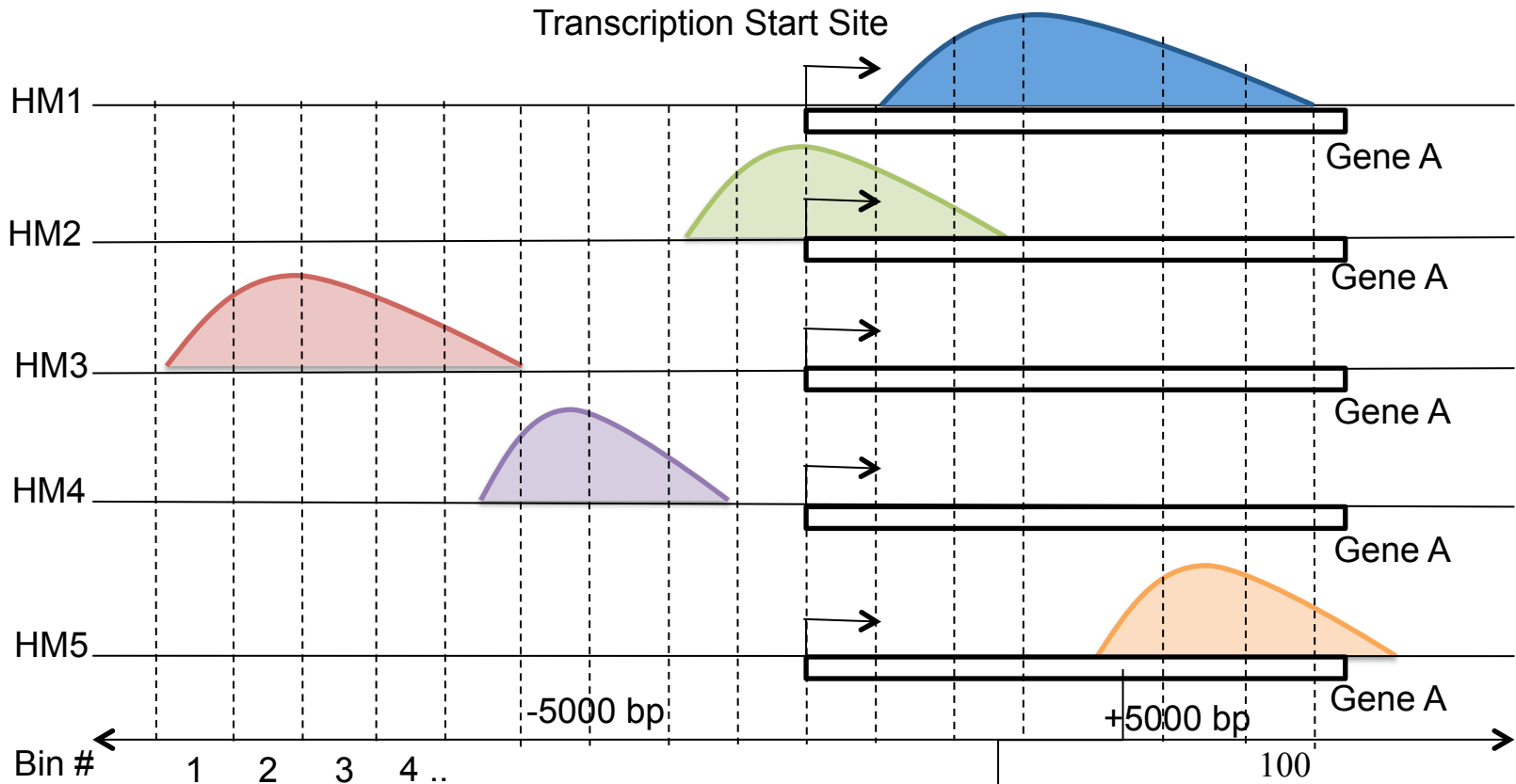


Histone Modification Signals



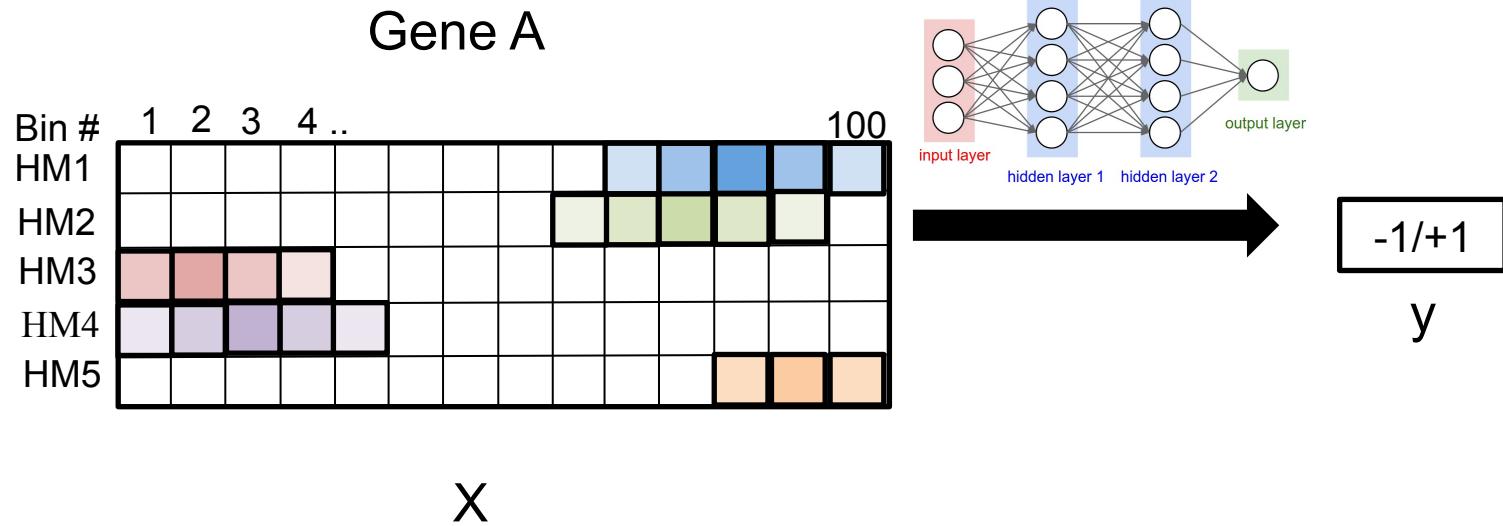
DeepChrome

Data

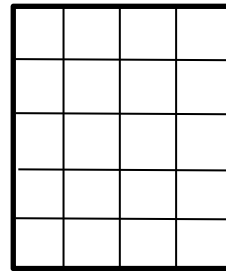
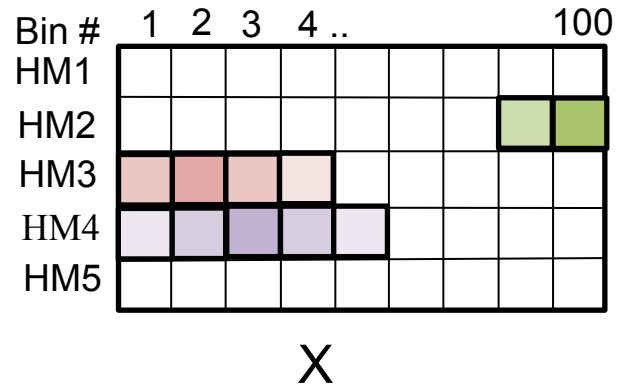


Bin #	1	2	3	4	..	100
HM1						Blue
HM2						Green
HM3	Red	Red	Red	Red		
HM4	Purple	Purple	Purple	Purple		
HM5						Orange

Overview

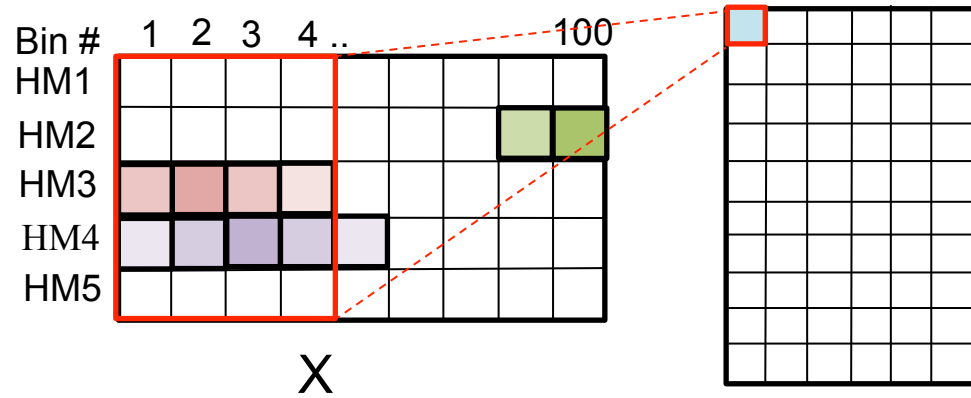


CNN Model



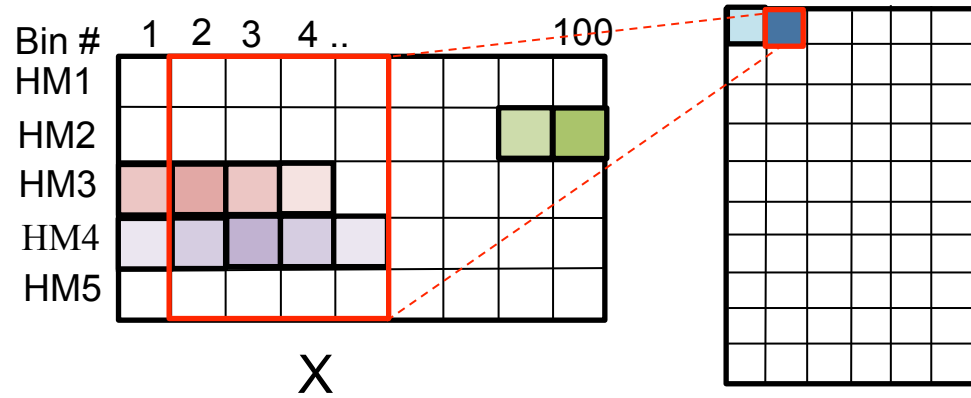
1. Convolution

CNN Model



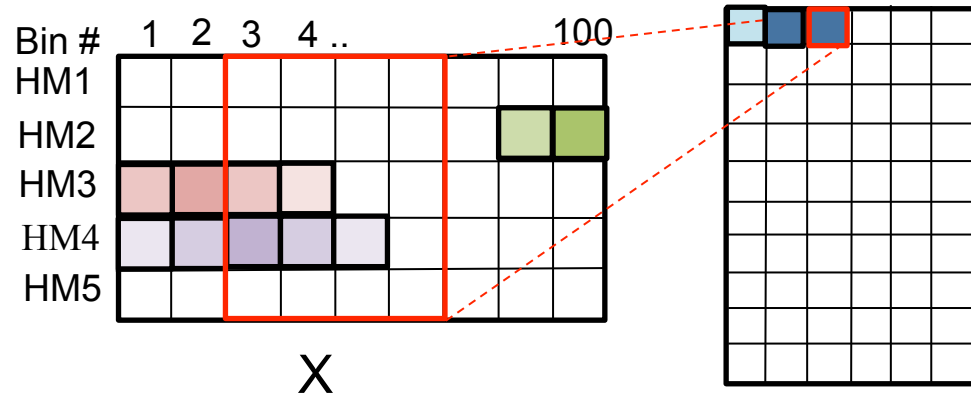
1. Convolution

CNN Model



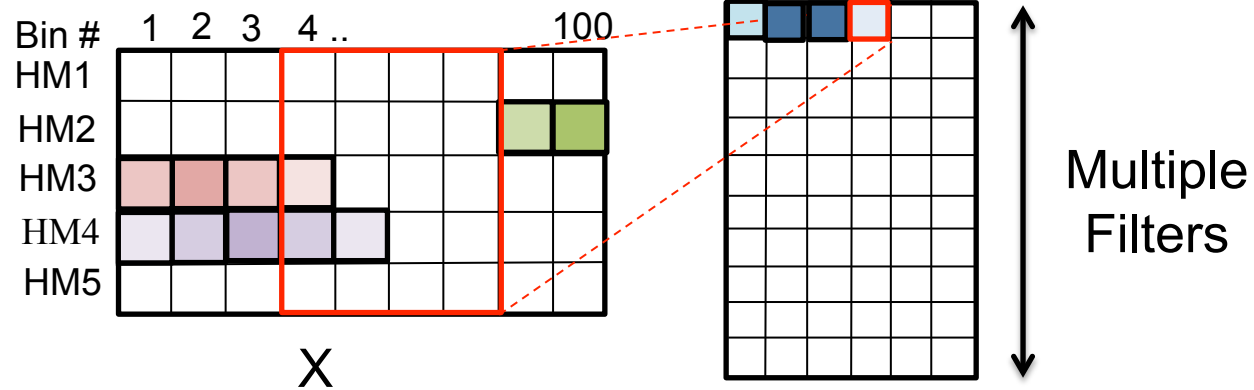
1. Convolution

CNN Model



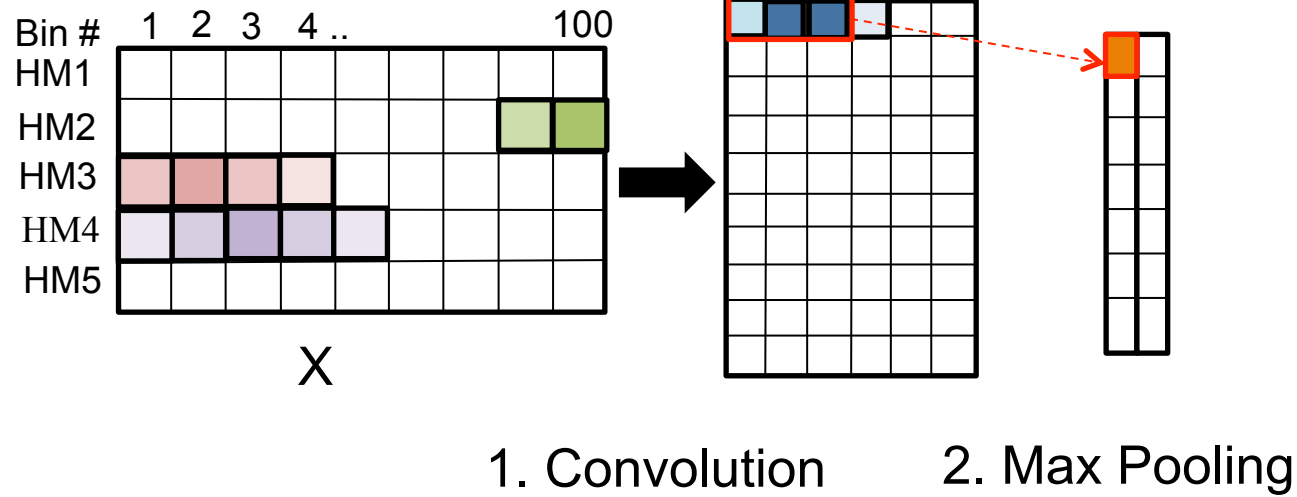
1. Convolution

CNN Model

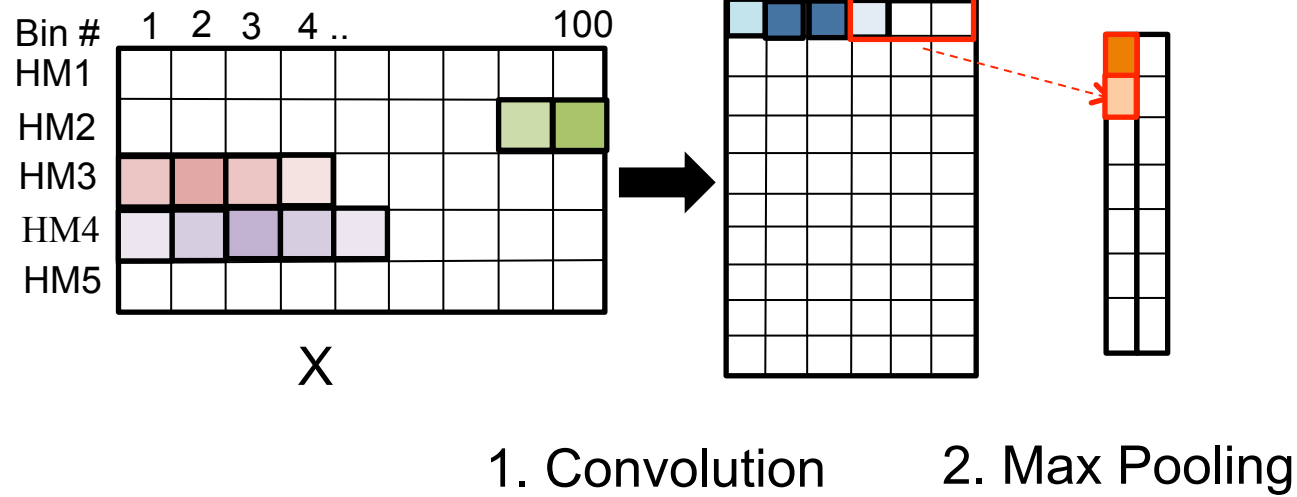


1. Convolution

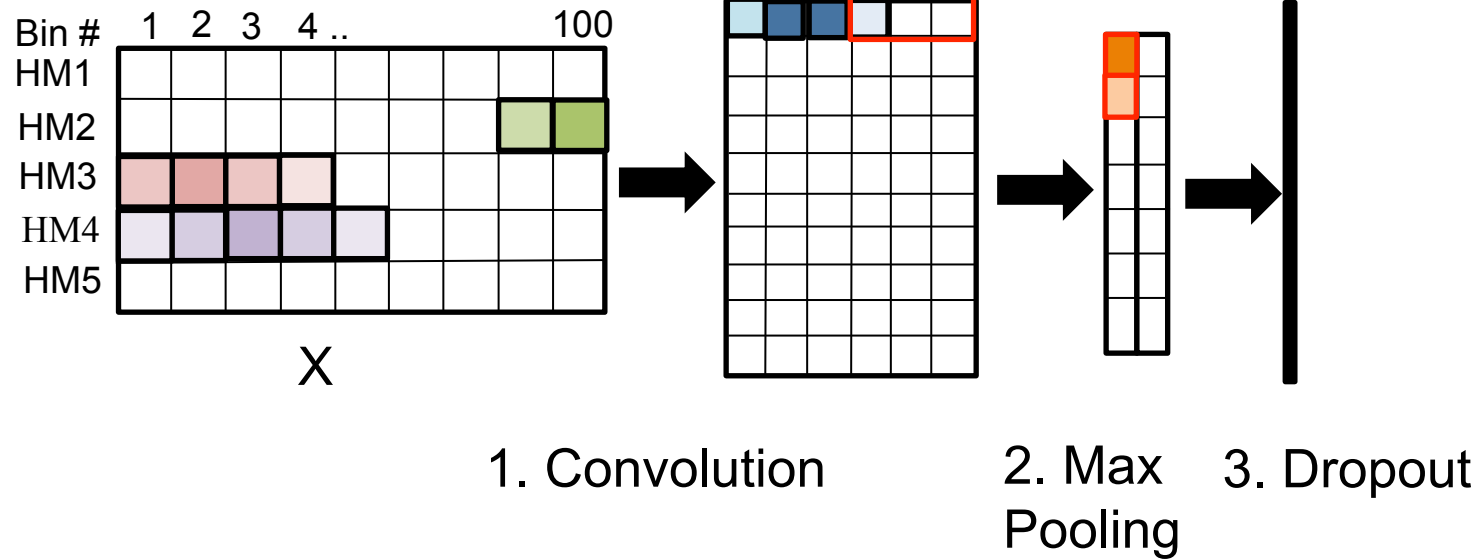
CNN Model



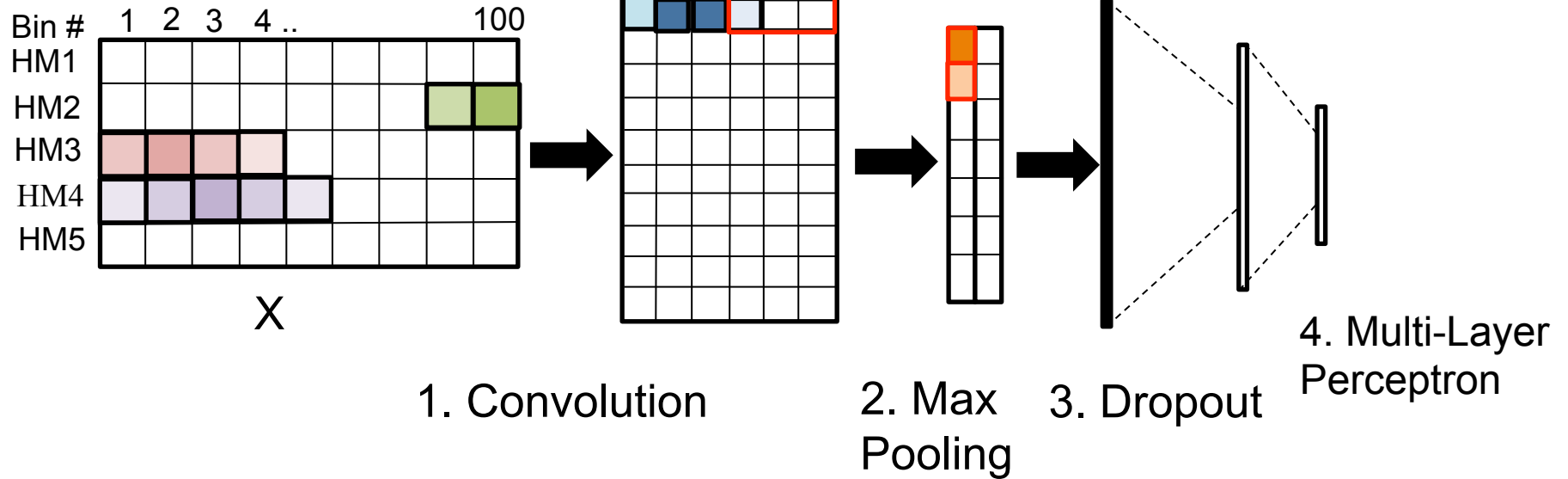
CNN Model



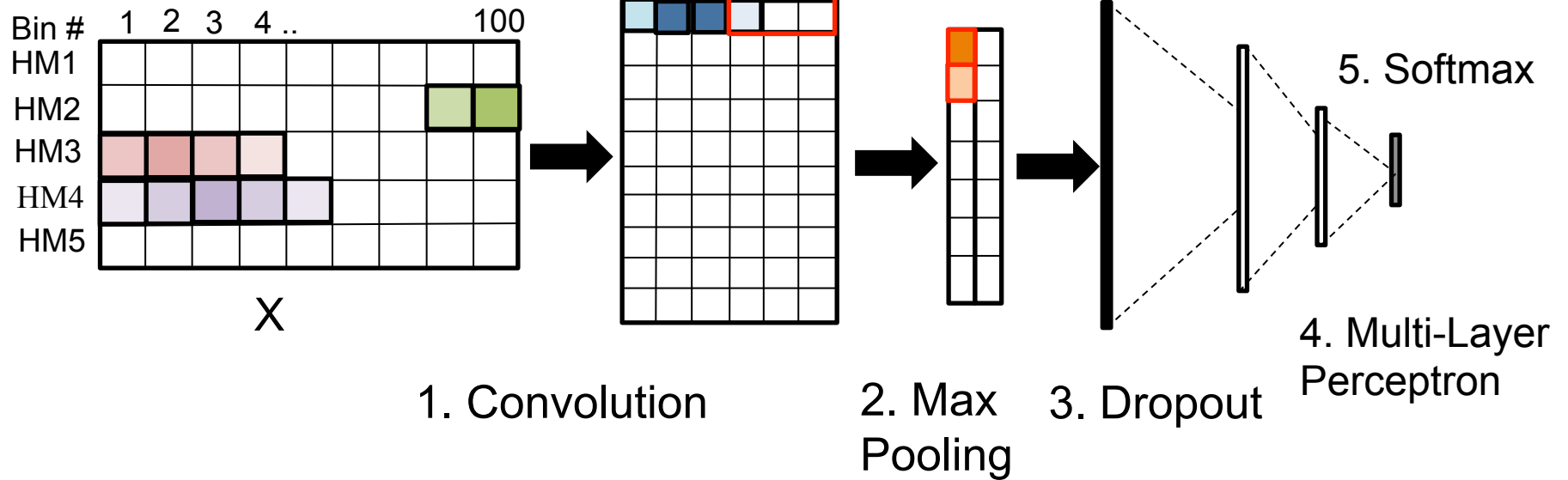
CNN Model



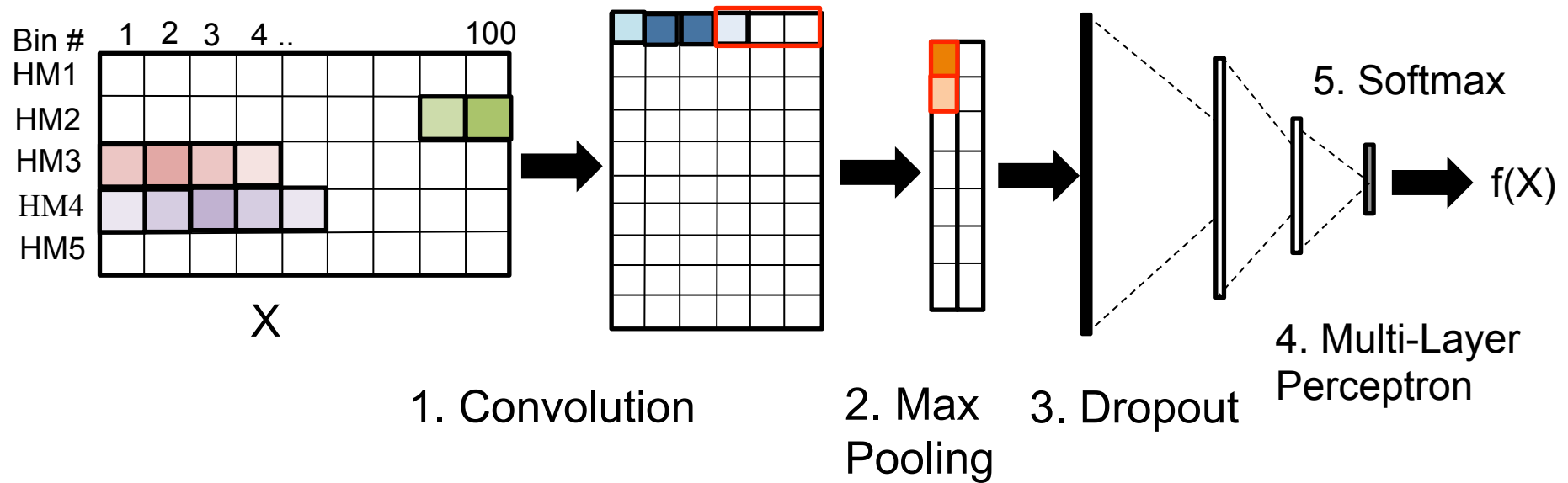
CNN Model



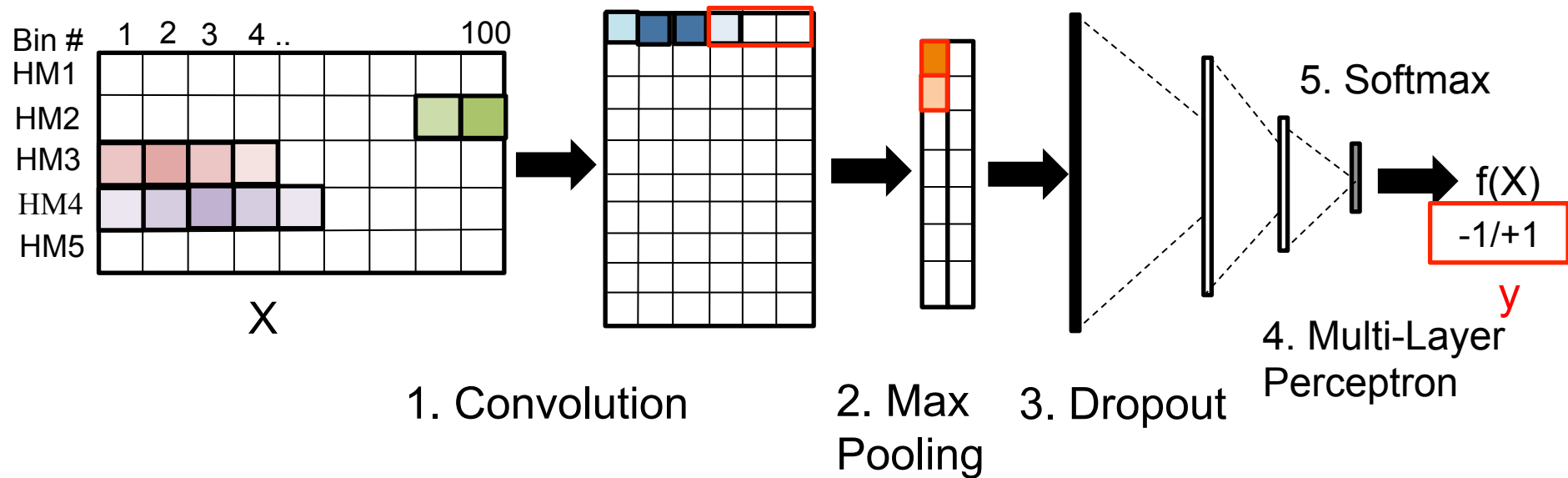
CNN Model



CNN Model

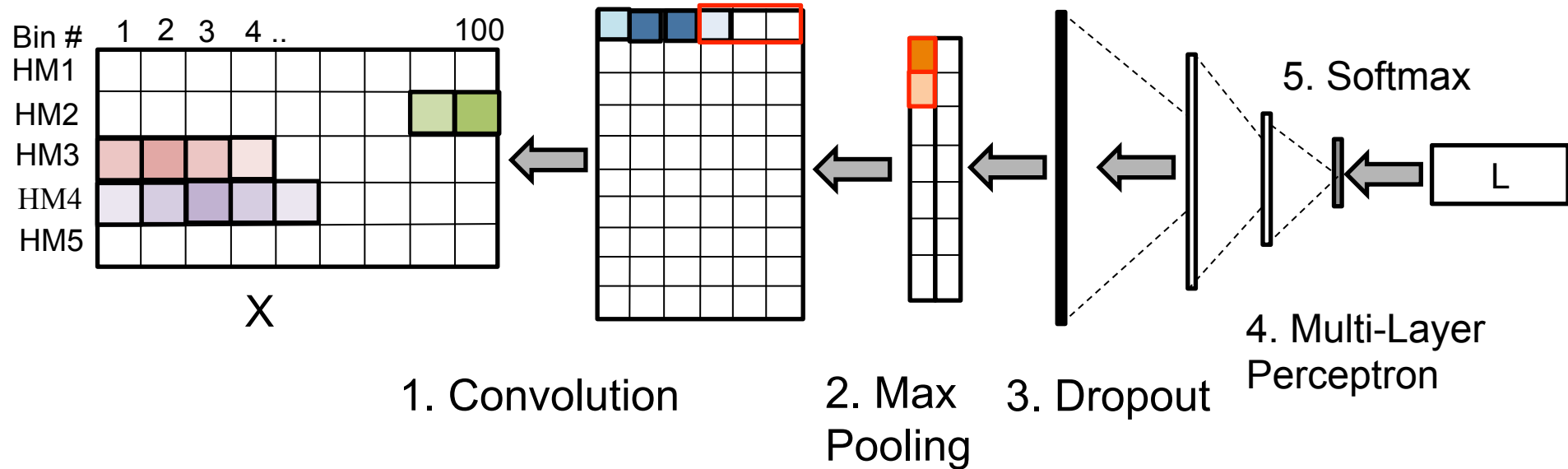


CNN Model



$$L = \sum_{n=1}^{N_{\text{samp}}} \text{loss}(f(X^{(n)}), y^{(n)})$$

CNN Model



Back-propagation:

$$\Theta \leftarrow \Theta - \eta \frac{\partial L}{\partial \Theta}$$

Experimental Setup

- **Cell-types:** 56
- **Input (HM):** ChIP-Seq Maps (REMC)
- **Output (Gene Expression):** Discretized RNA-Seq (REMC)

Experimental Setup

- **Cell-types:** 56
- **Input (HM):** ChIP-Seq Maps (REMC)
- **Output (Gene Expression):** Discretized RNA-Seq (REMC)

Histone Mark	Functional Category
H3K27me3	Repressor
H3K36me3	Promoter
H3K4me1	Distal Promoter
H3K4me3	Promoter
H3K9me3	Repressor

Experimental Setup

- **Cell-types:** 56
- **Input (HM):** ChIP-Seq Maps (REMC)
- **Output (Gene Expression):** Discretized RNA-Seq (REMC)

Histone Mark	Functional Category
H3K27me3	Repressor
H3K36me3	Promoter
H3K4me1	Distal Promoter
H3K4me3	Promoter
H3K9me3	Repressor

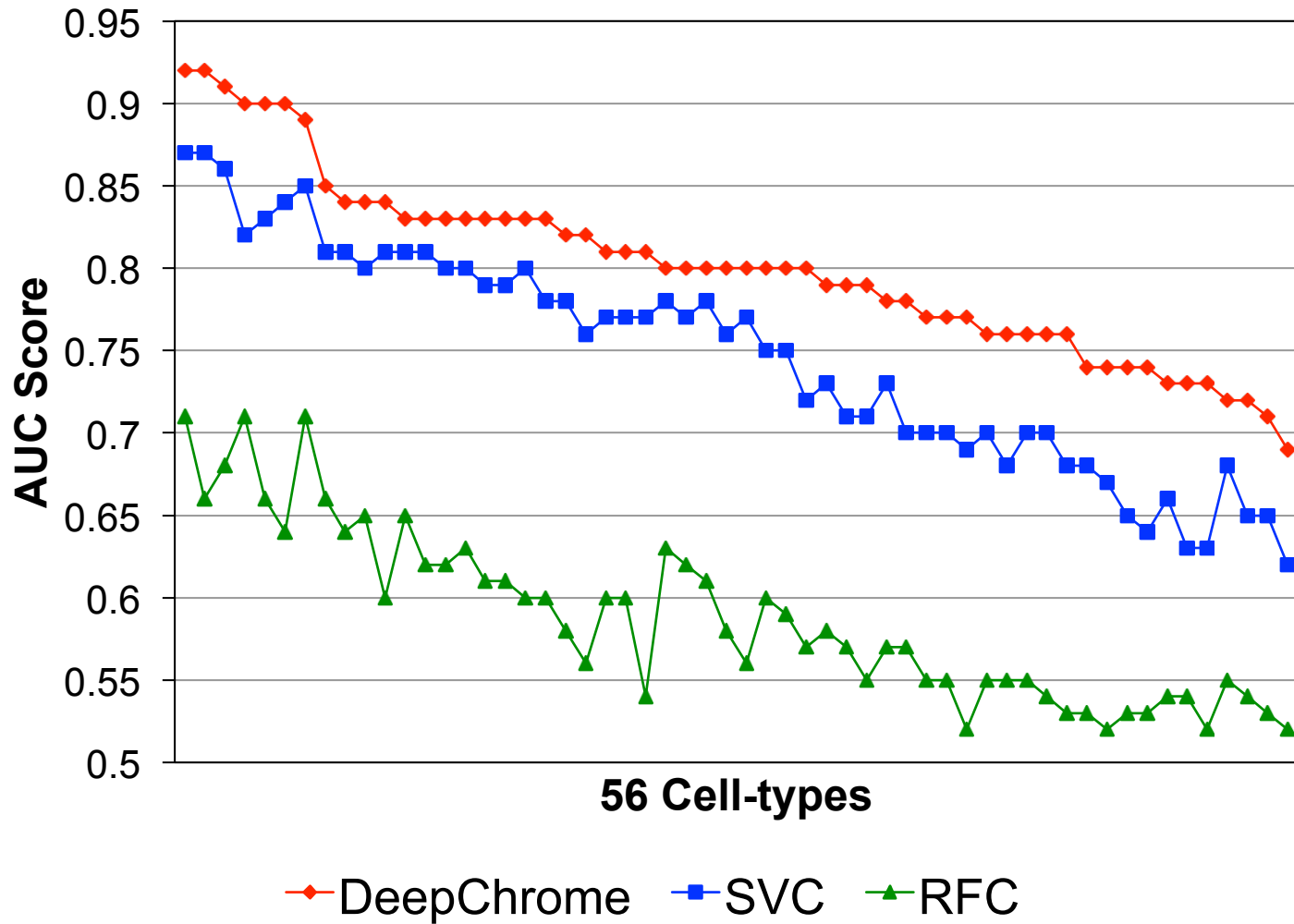
- **Baselines:** Support Vector Classifier (SVC) and Random Forest Classifier (RFC)

Training Set
6601 Genes

Validation Set
6601 Genes

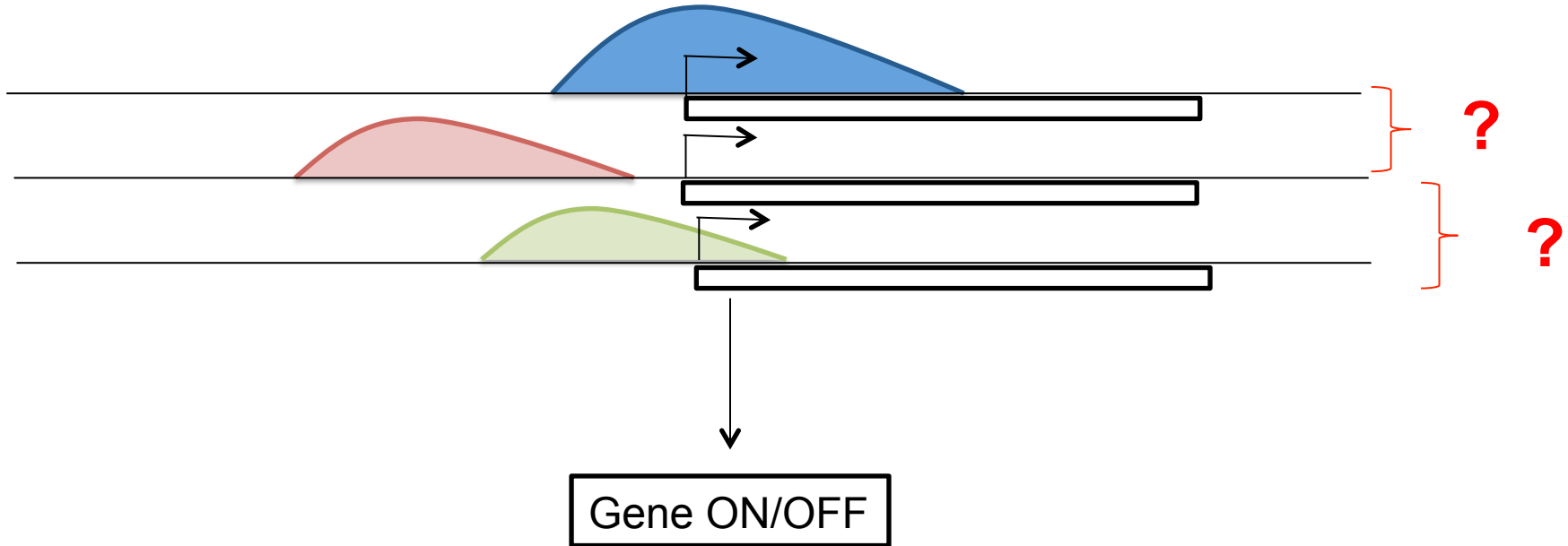
Test Set
6600 Genes

Results: Accuracy



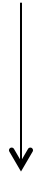
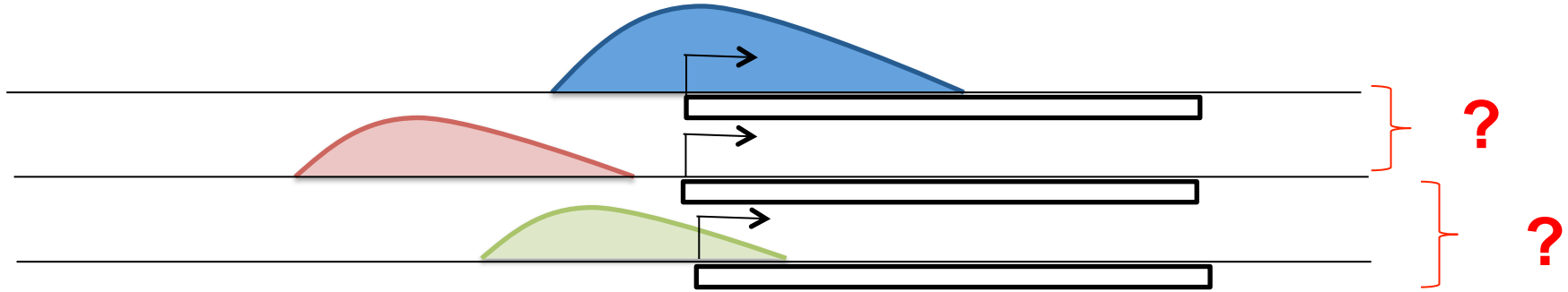
Visualization

Histone Modification Signals



Visualization

Histone Modification Signals



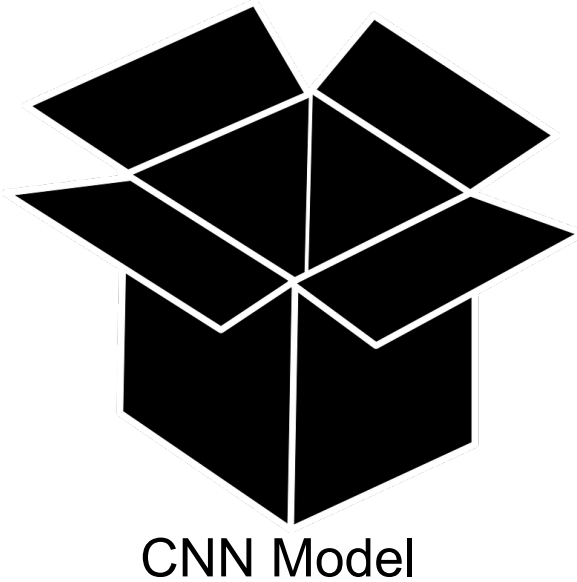
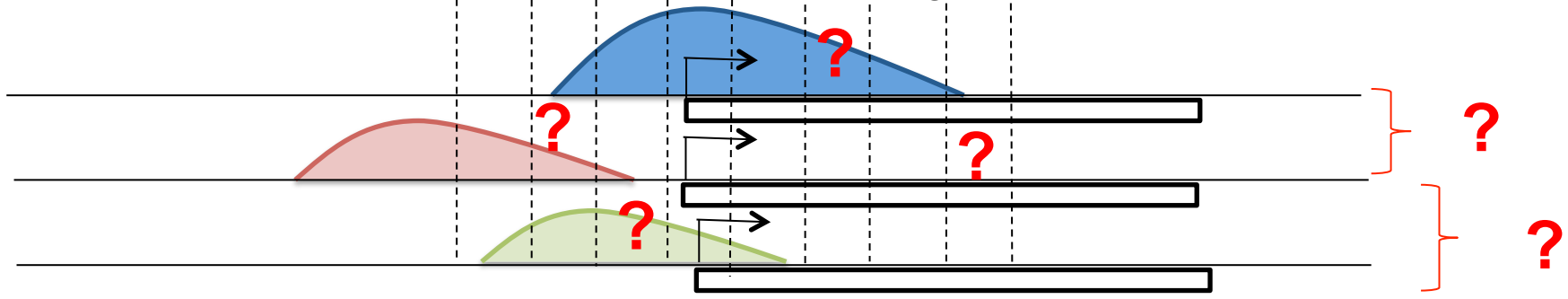
Gene ON/OFF



CNN Model

Visualization

Histone Modification Signals



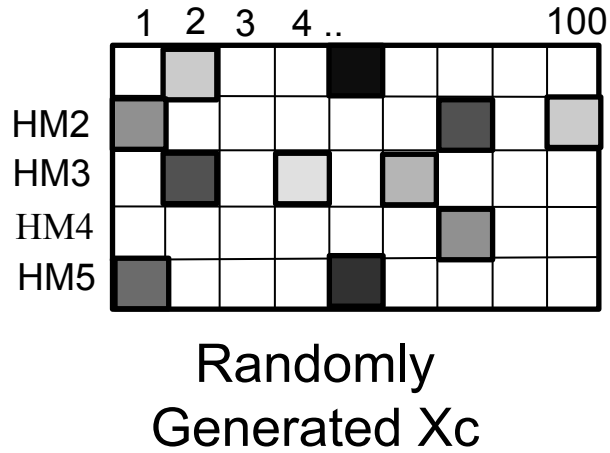
CNN Model

Gene ON/OFF

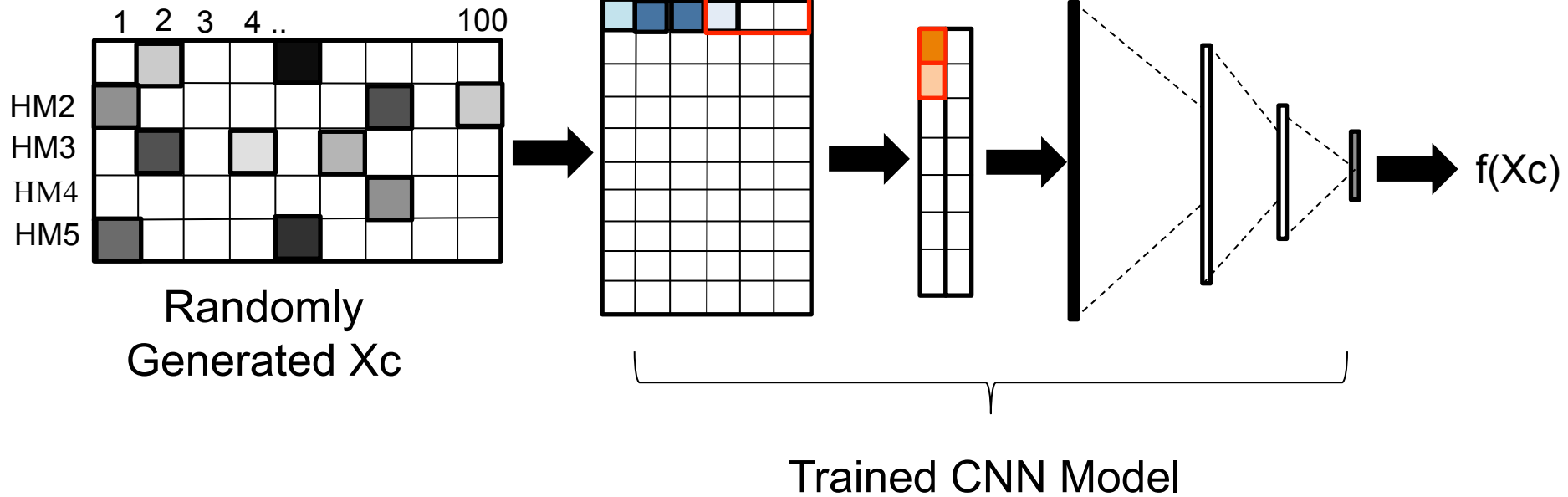
Bin #	1	2	3	4	..										
HM1									?					?	
HM2		?		?											
HM3				?					?					?	
HM4										?					
HM5				?					?					?	

X 52

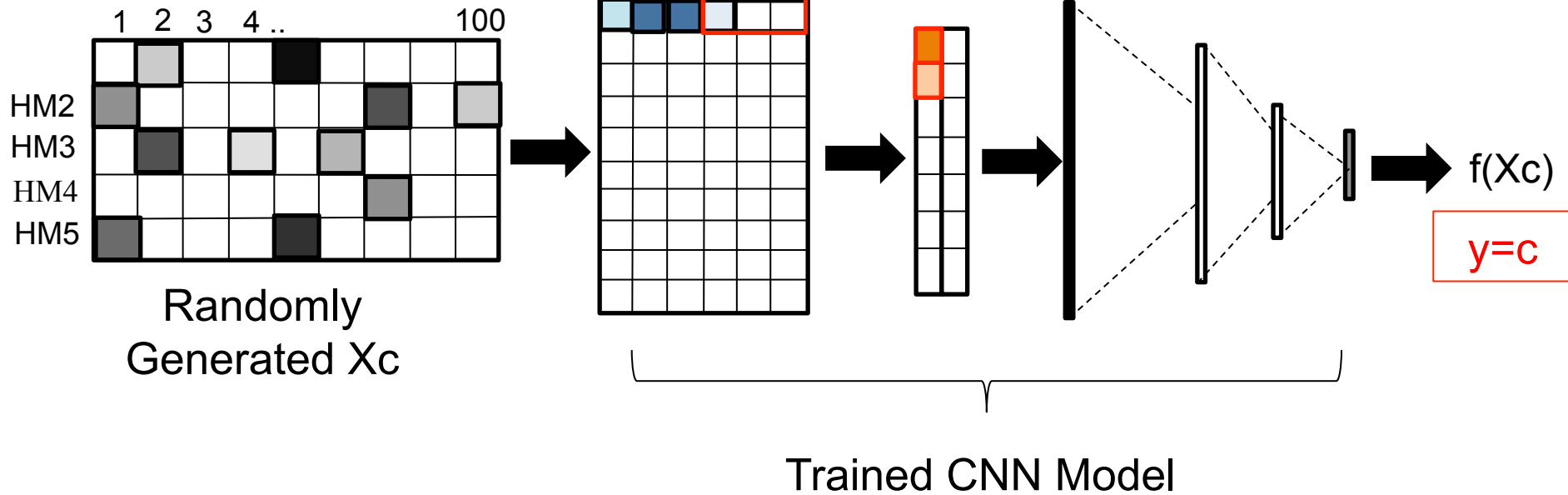
Iterative Most-likely Class Method



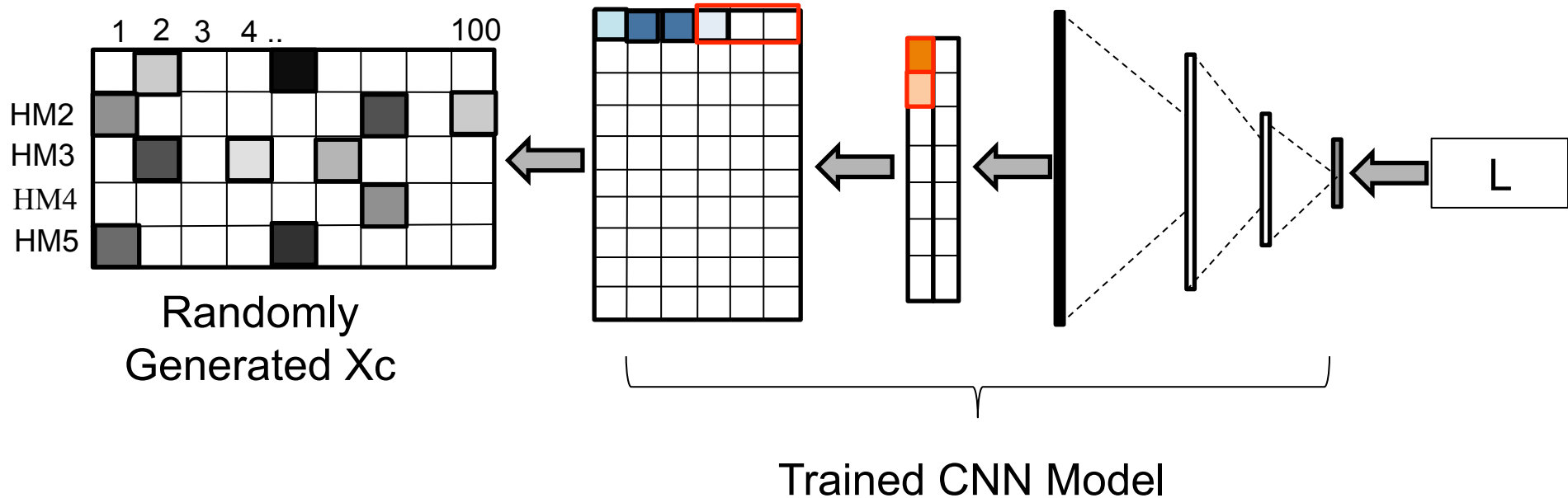
Iterative Most-likely Class Method



Iterative Most-likely Class Method

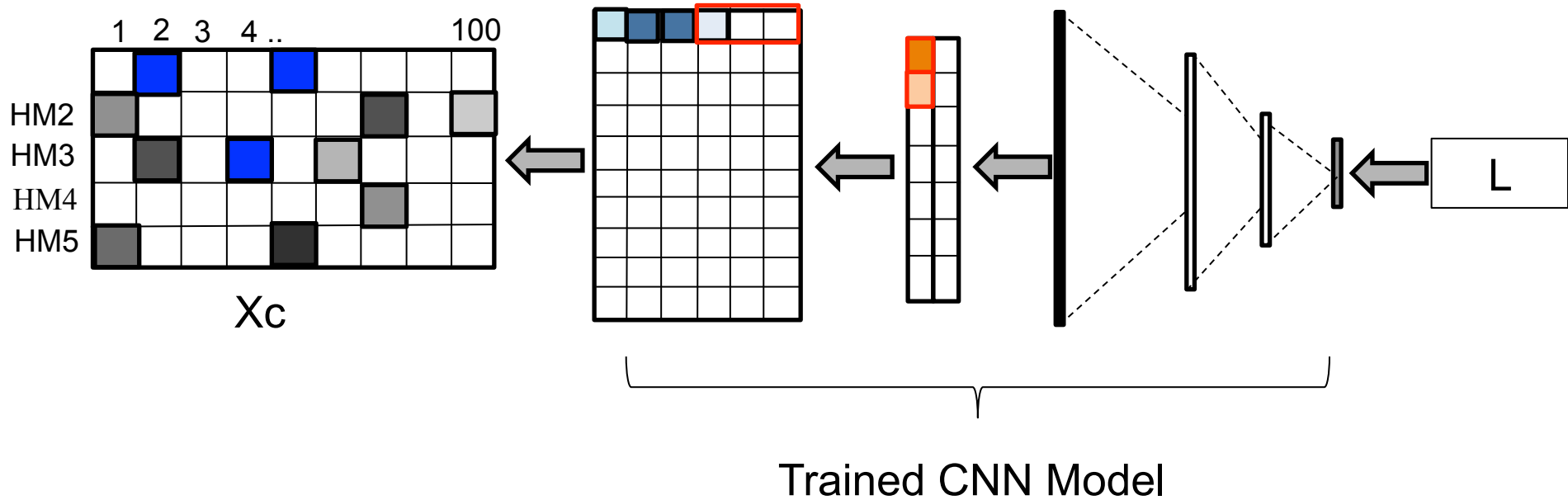


Iterative Most-likely Class Method



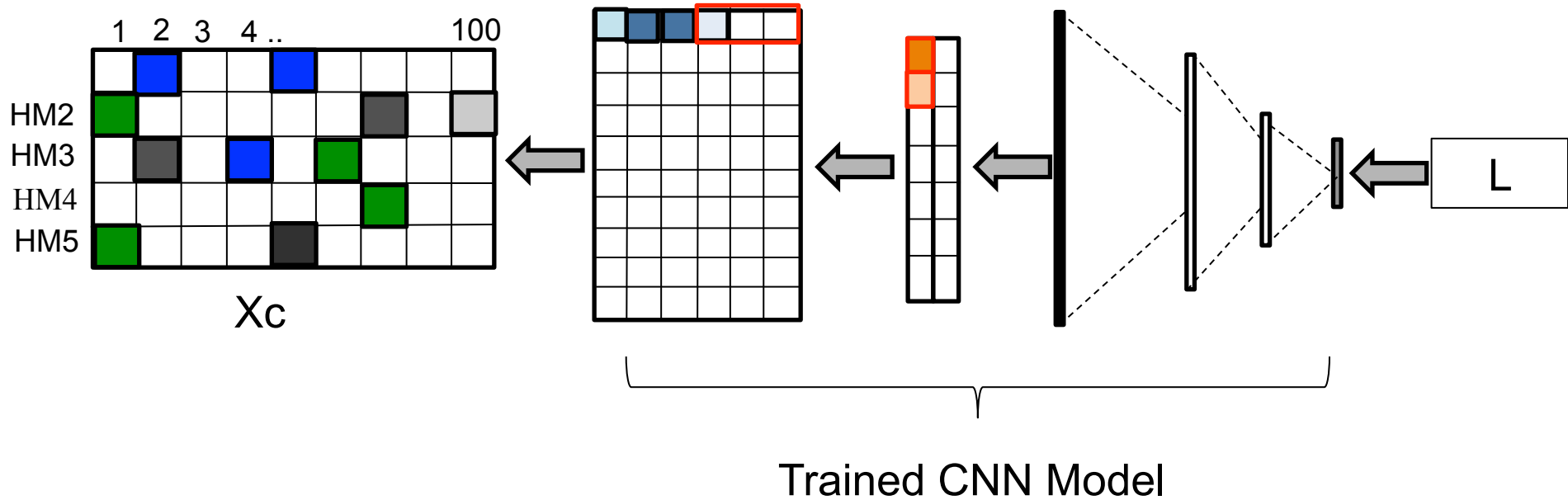
$$\arg \min_{X_c} \{L(f(X_c), y = c)\}$$

Iterative Most-likely Class Method



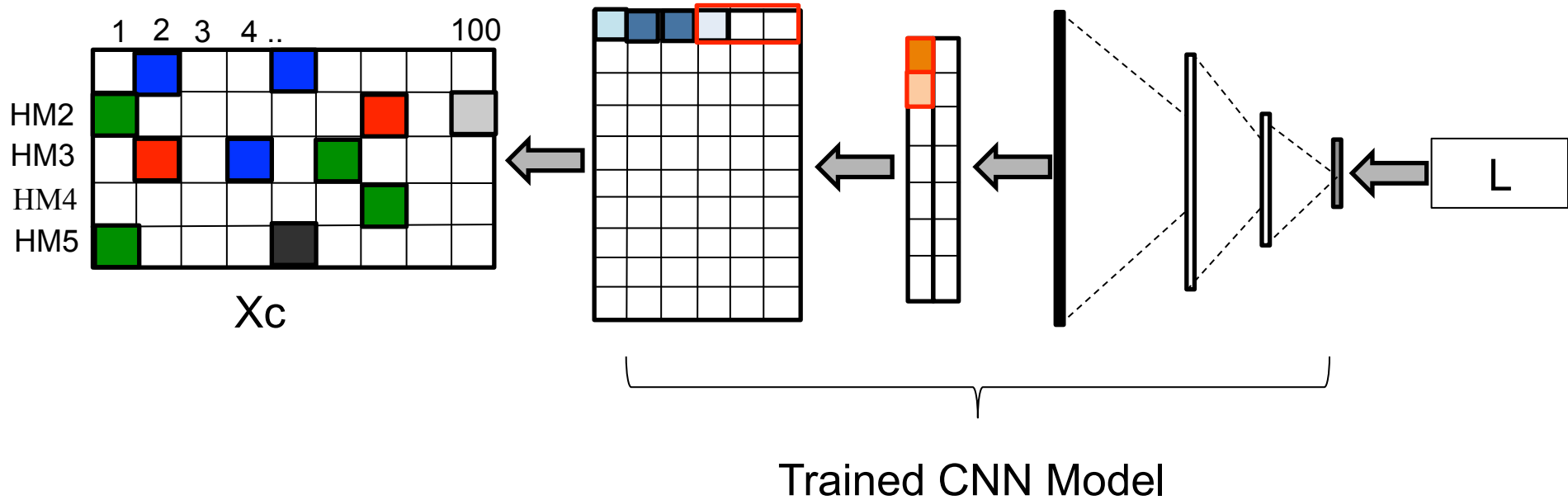
$$\arg \min_{X_c} \{L(f(X_c), y = c)\}$$

Iterative Most-likely Class Method



$$\arg \min_{X_c} \{L(f(X_c), y = c)\}$$

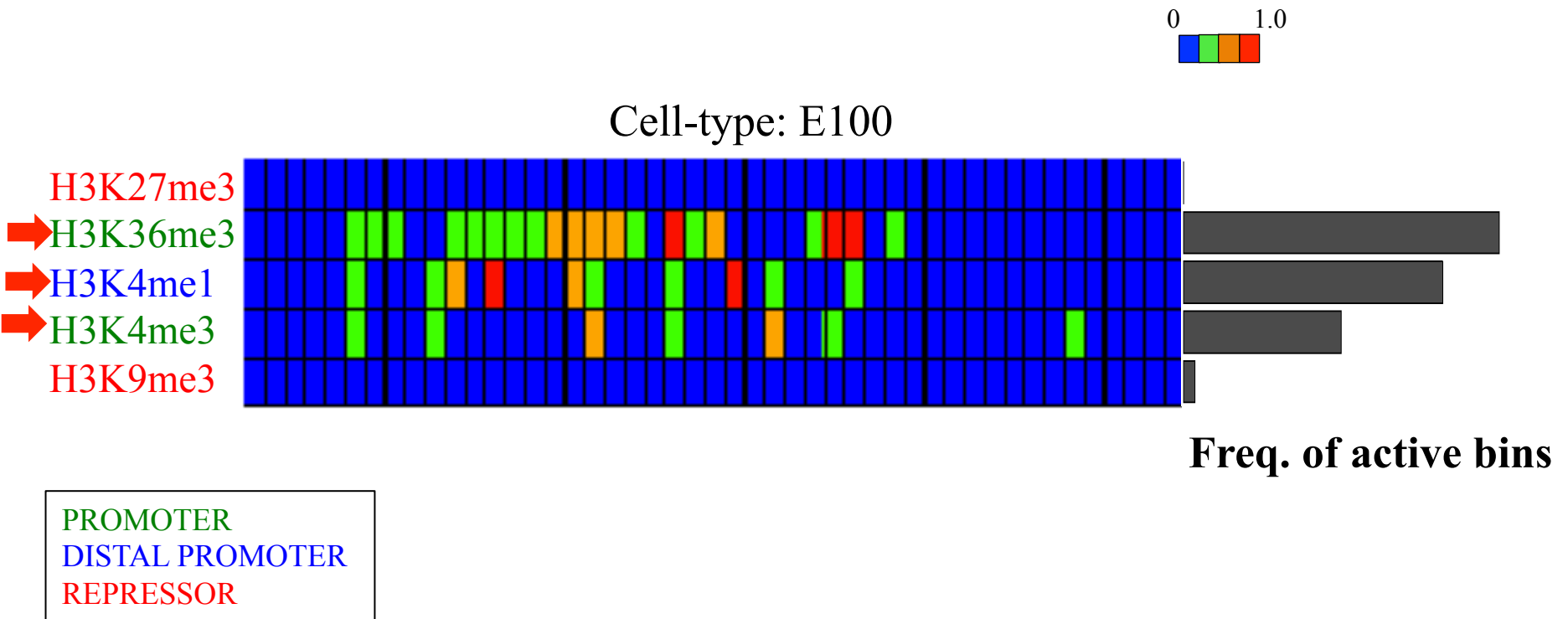
Iterative Most-likely Class Method



$$\arg \min_{X_c} \{L(f(X_c), y = c)\}$$

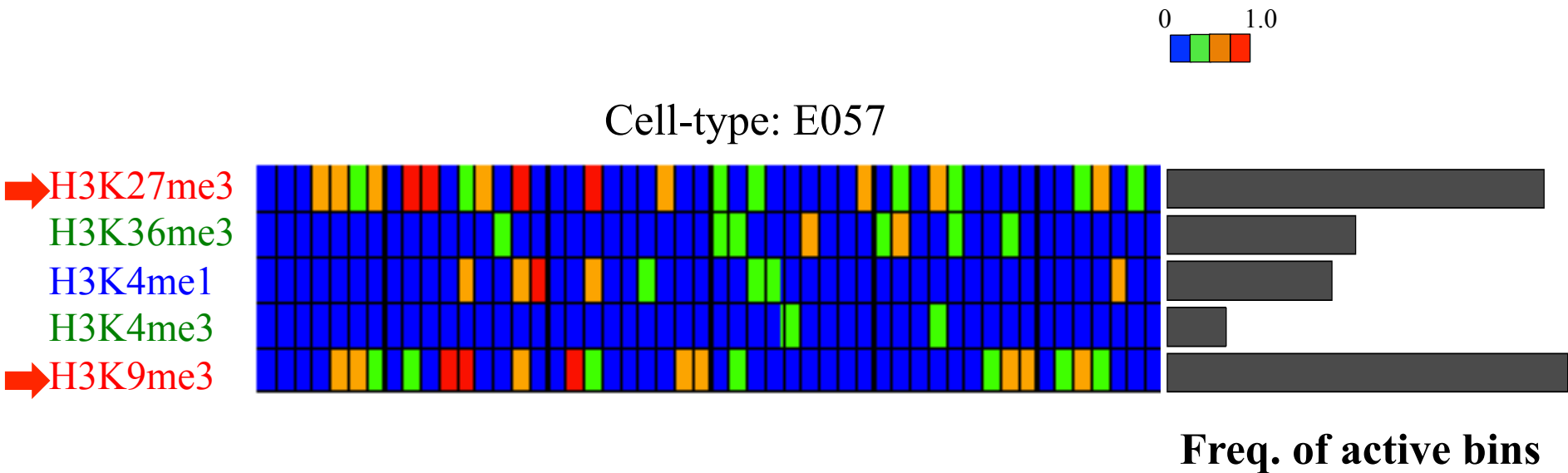
Results: Visualization

Gene : ON ($y=+1$)



Results: Visualization

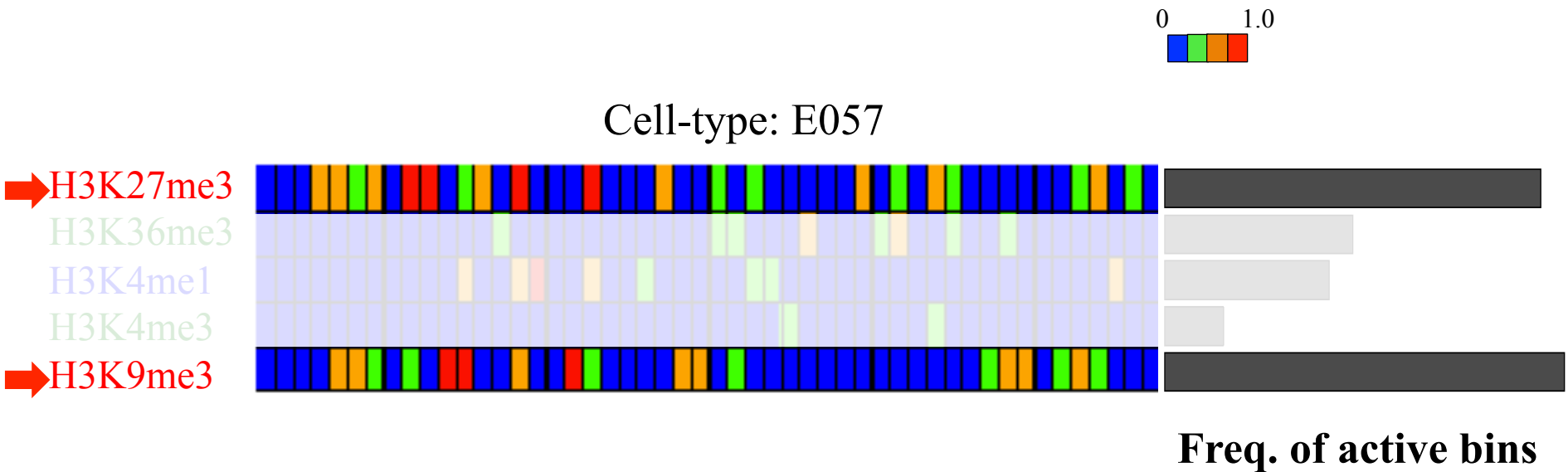
Gene : OFF ($y = -1$)



PROMOTER
DISTAL PROMOTER
REPRESSOR

Results: Visualization

Gene : OFF ($y = -1$)



PROMOTER
DISTAL PROMOTER
REPRESSOR

Conclusion

1. First deep learning implementation for gene expression prediction
2. Unified Framework
 - a. Outperforms state-of-the-art implementations
 - b. Visualization of high-order combinatorial relationships

Available @ www.deepchrome.org

Acknowledgements

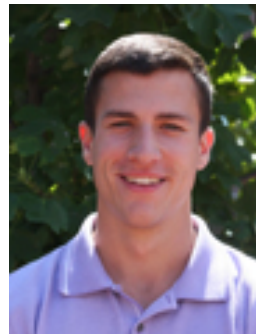


Computer Science
at the UNIVERSITY of VIRGINIA

**Machine Learning and
Bioinformatics Lab**



Dr. Yanjun Qi



Jack Lanchantin

Dr. Gabriel Robins

Marina Sanusi

Beilun Wang

Weilin Xu

Ji Gao

Kamran Kowsari

**Department of
Biochemistry and Molecular
Genetics: Dr. Mazhar Adli**

Travel Fund



Thank you

