

Deep Learning Reading Group

MPII, August 17, 2017

DeepCpG: accurate prediction of single-cell DNA methylation states using deep learning

Christof Angermueller, Heather J. Lee, Wolf Reik, Oliver Stegle

Genome Biology, 18(1), 67. doi:10.1186/s13059-017-1189-z

Presented by Fabian Müller



Outline

- Biological background
- DeepCpG model
 - Architecture
 - Validation
 - Interpretation
- Summary



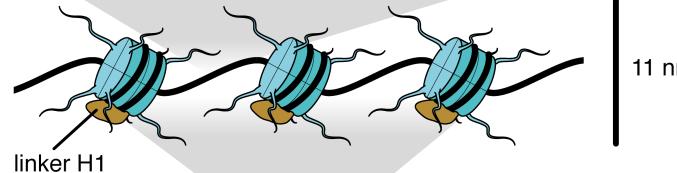
Chromatin Organization

DNA double helix



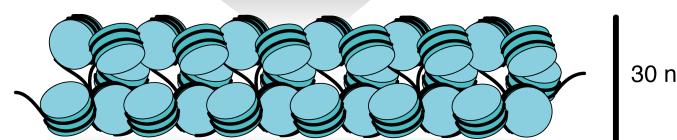
2 nm

Beads-on-a-string



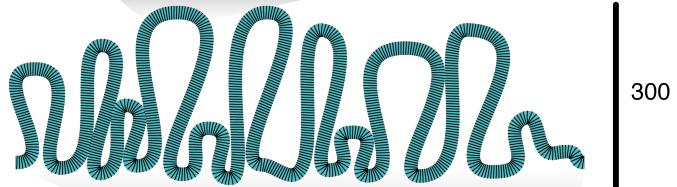
11 nm

30-nm fibre



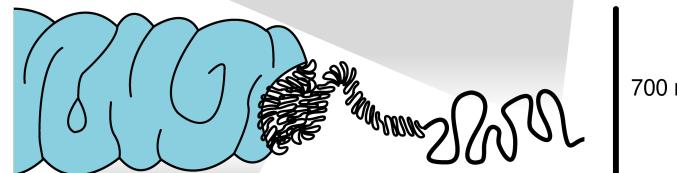
30 nm

Extended section of a chromosome



300 nm

Condensed section of a chromosome



700 nm

Metaphase chromosome

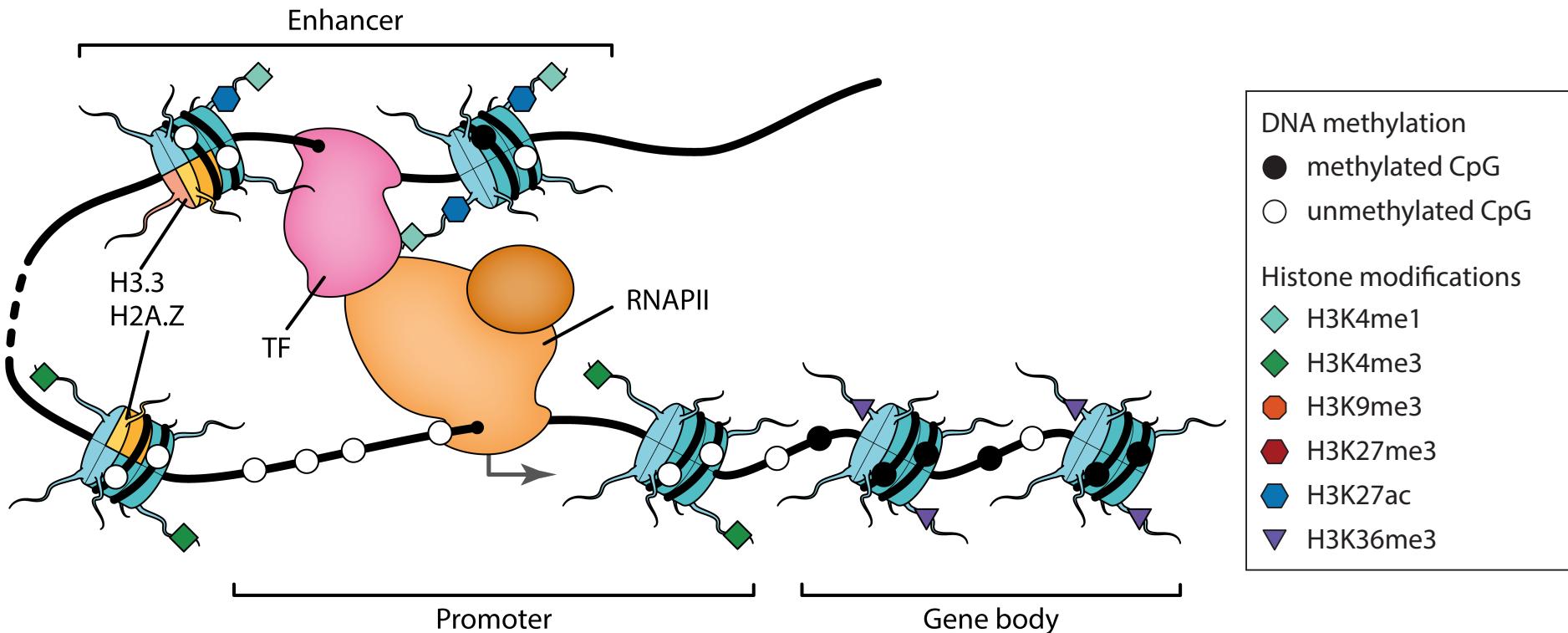


1400 nm

- Hierarchical organization
- Protein scaffold
 - Nucleosome: DNA wrapped around histone proteins
- Accessibility regulates gene expression
- Tightly regulated by a plethora of molecular factors

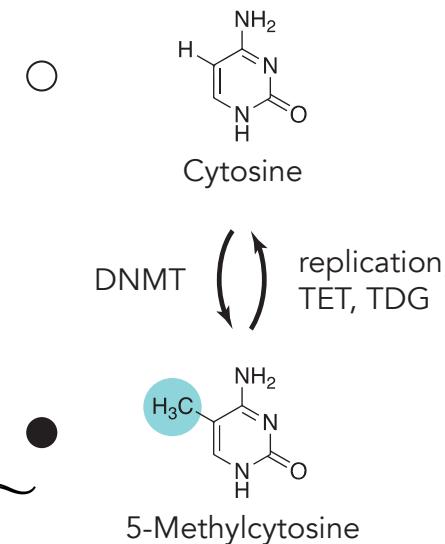
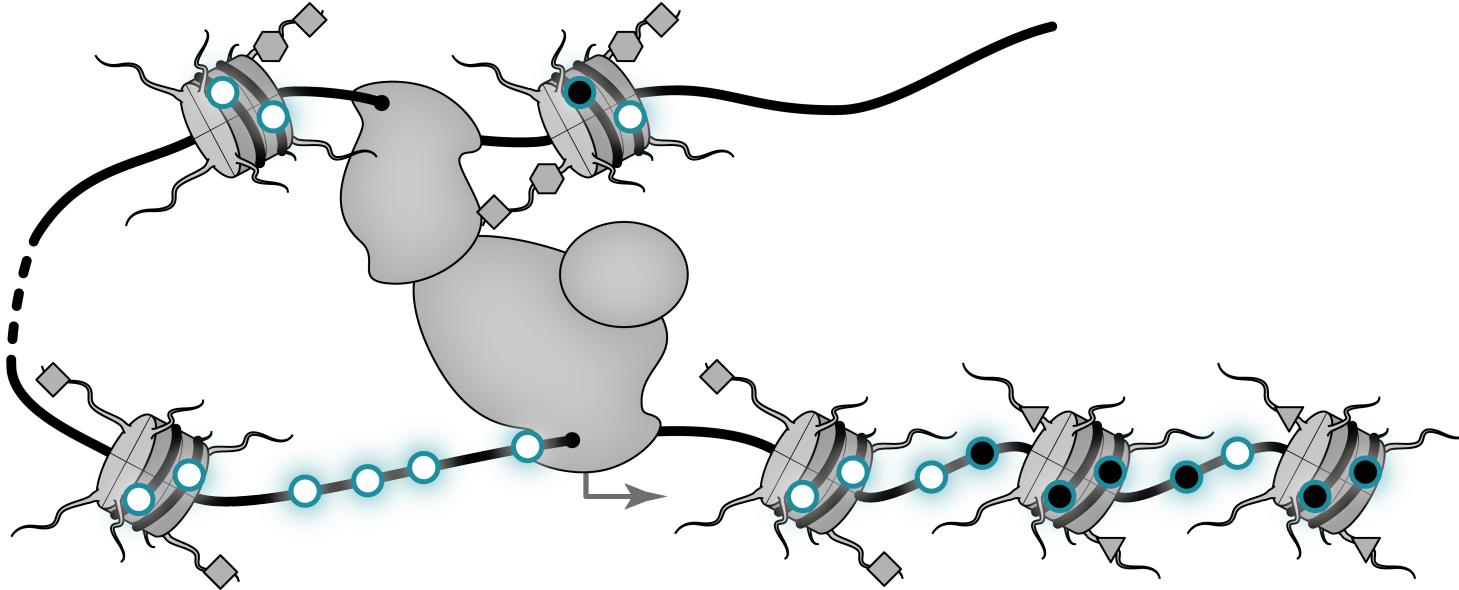


Epigenetic Regulation



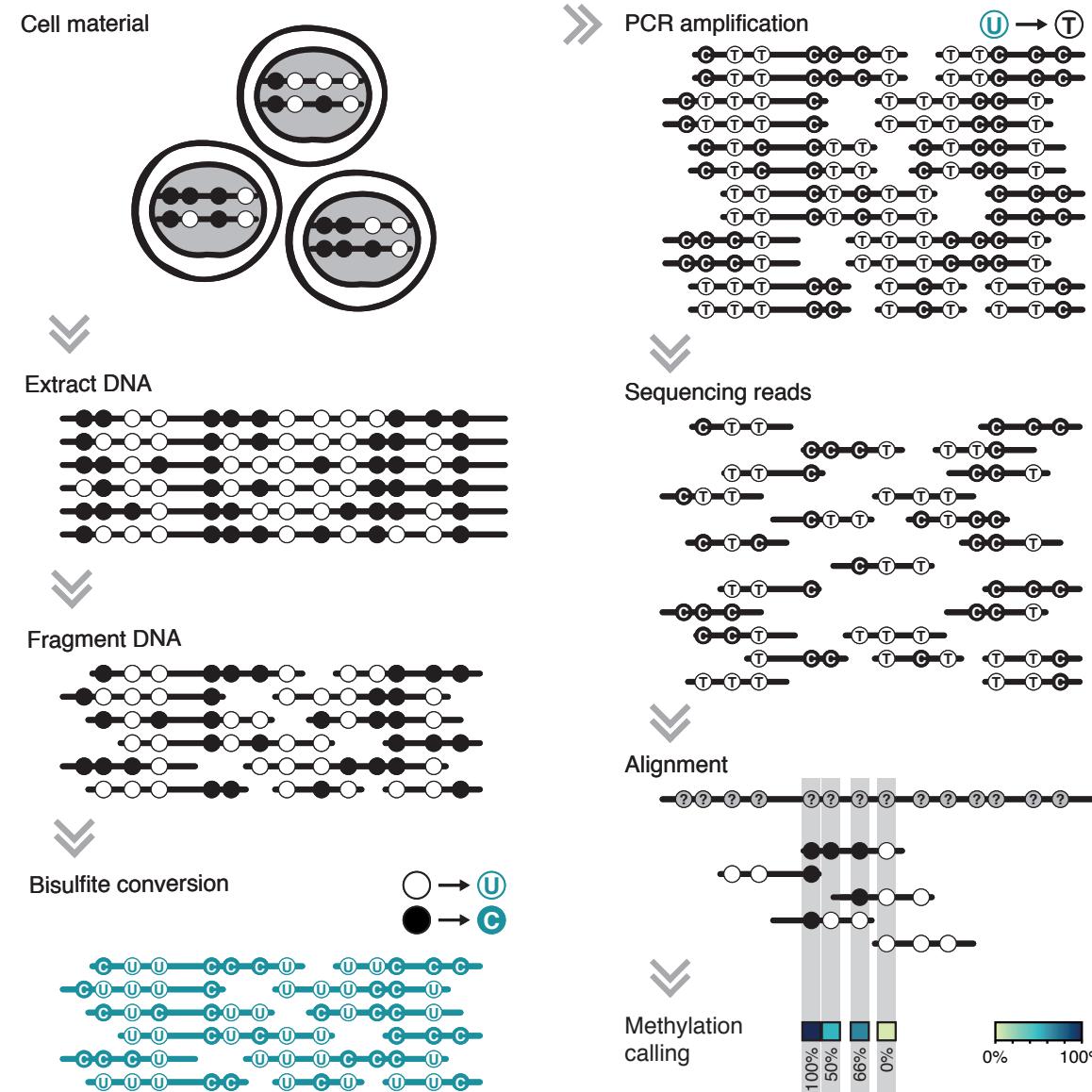
- Accessibility regulates gene expression
- Regulation beyond promoter regions
 - Enhancers, gene bodies, ...
- 80.4% of the genome is attributed with regulatory roles [1]

DNA Methylation



- Occurs predominantly in CpG dinucleotides in mammalian genomes
- Associated with regulation of gene expression
- Deregulated in disease
- Genome-wide quantification:
 - Bisulfite sequencing
 - Methylation arrays
 - Enrichment-based sequencing
 - ...

Bisulfite Sequencing

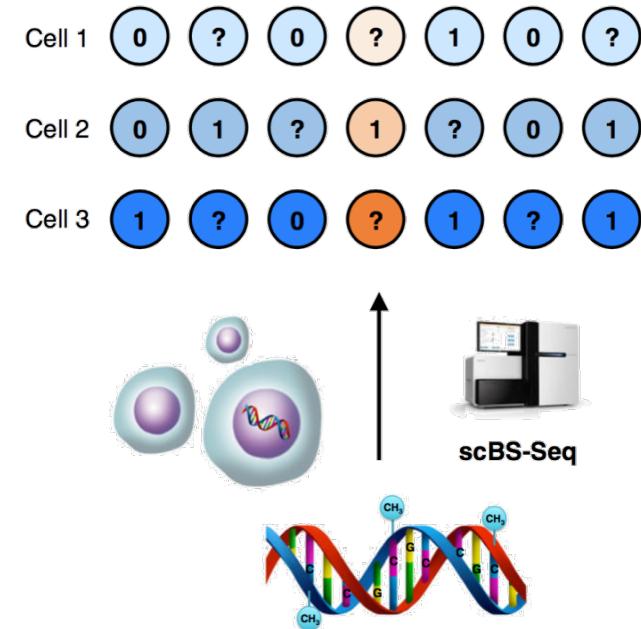


Prediction Task

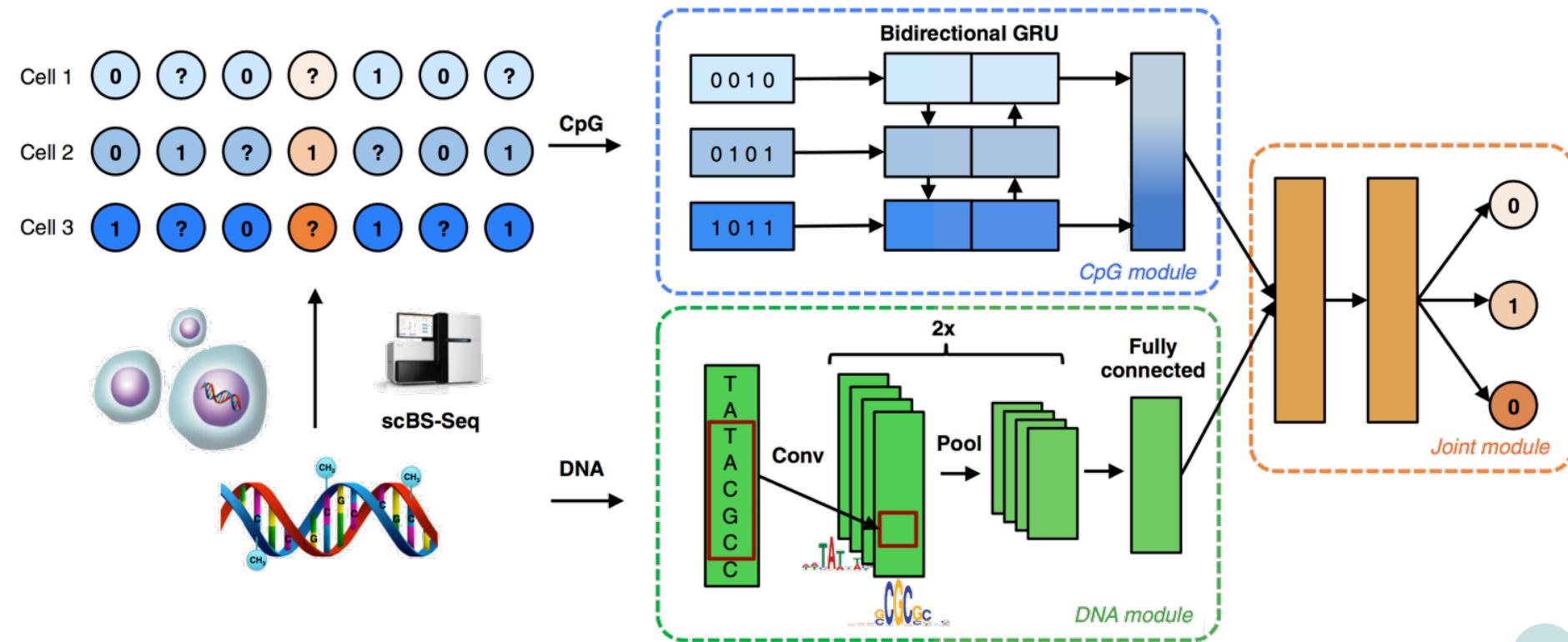
- Problem: single-cell bisulfite sequencing data is sparse

⇒ **Predict missing CpG methylation states in single cells**

- Output: methylation state (binary) for each CpG and cell
- Features:
 - Local DNA sequence
 - Observed CpG methylation



DeepCpG Model Architecture



- Combining CpG methylation and sequence features using different modules

DNA module

- Convolutional network using position weight matrices
- Input: 1001bp DNA sequence windows, centered at each target CpG (n): :

$$s_n \in \{0, 1\}^{1001 \times 4}$$

A = [1, 0, 0, 0]
 T = [0, 1, 0, 0]
 G = [0, 0, 1, 0]
 C = [0, 0, 0, 1]

- 1st convolution:

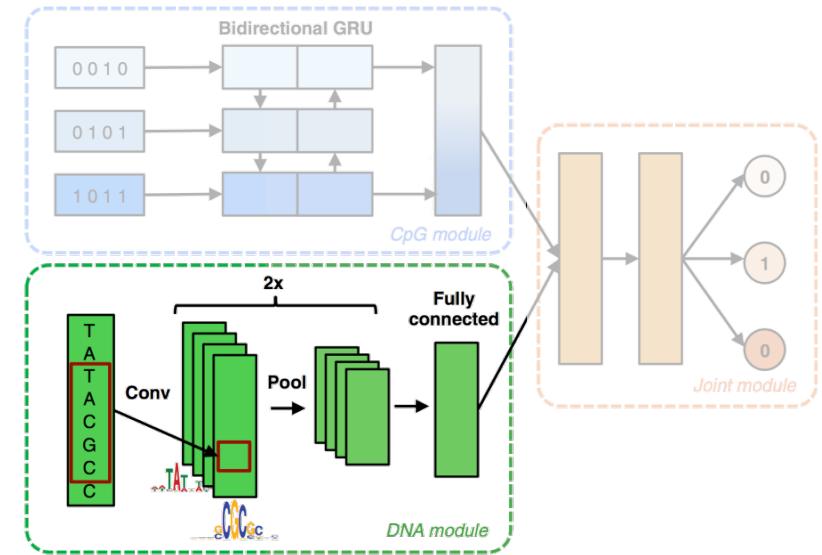
filter length $\overbrace{\quad\quad\quad}$
 $a_{nfi} = \text{ReLU} \left(\sum_{l=1}^L \sum_{d=1}^D w_{fld} s_{n,i+l,d} \right)$
 \uparrow multiple filters f

with $\text{ReLU}(x) = \max(0, x)$

- Max-pooling layer:

$$p_{nfi} = \max_{|k| < P/2} (a_{nf,i+k})$$

\uparrow
 step size for non-overlapping pooling



- Multiple convolution-pooling layers capture sequence-motif interactions
 - Number selected on validation set
- One final fully connected layer (e.g. 128 nodes, ReLU activation)

CpG Module

- Bidirectional gated recurrent network (GRU)
 - Input: methylation states and relative distances for the 25 left and right neighbor CpGs of each target CpG in each of T cells:
- $$x_1, \dots, x_T \in (\{0,1\}^{50}, [0,1]^{50}) \subset [0,1]^{100}$$
- Embedding layer combines input and captures interactions for each cell:
- $$\bar{x}_t = \text{ReLU}(W_{\bar{x}} \cdot x_t + b_{\bar{x}}).$$
- $\uparrow \in \mathbb{R}^{256}$
- GRU layers (forward+backward) capture cell interactions :

- Reset gate:

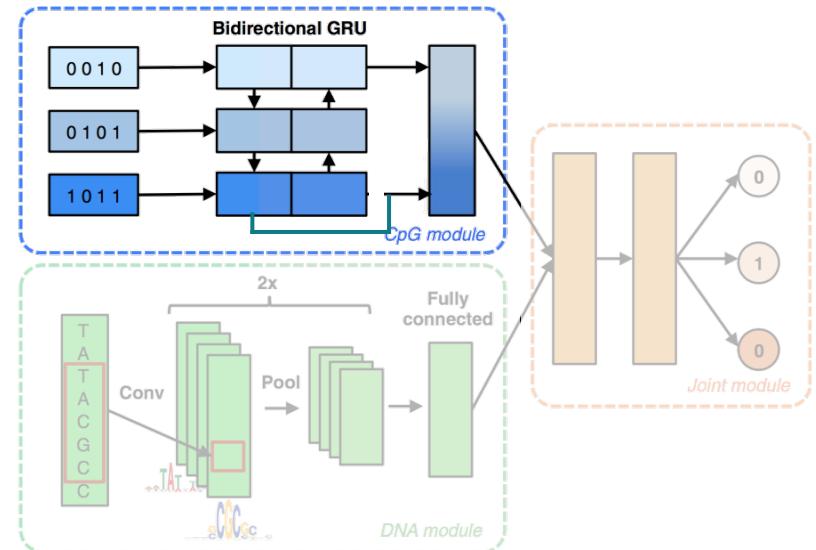
$$r_t = \text{sigmoid}(W_{r\bar{x}} \cdot \bar{x}_t + W_{rh} \cdot h_{t-1} + b_r)$$

- Update gate:

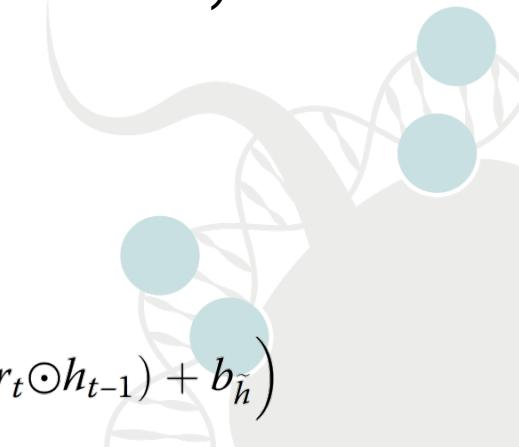
$$u_t = \text{sigmoid}(W_{u\bar{x}} \cdot \bar{x}_t + W_{uh} \cdot h_{t-1} + b_u)$$

- Hidden state:

$$h_t = (1-u_t) \odot h_{t-1} + u_t \odot \tilde{h}_t.$$



- Output: concatenation of forward and backward GRUs' last hidden state (512-dimensional)

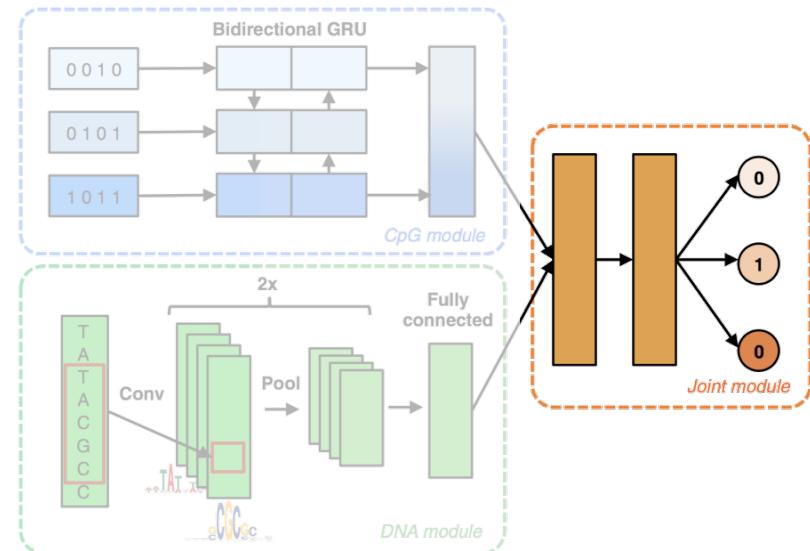


Joint Module

- Combine outputs of DNA and CpG modules
- Two fully-connected layers (512 nodes, ReLU activation)
- Output: methylation state for each CpG n and cell t :

$$\hat{y}_{nt}(x) \in [0, 1]$$

$$\hat{y}_{nt}(x) = \text{sigmoid}(x) = \left(\frac{1}{1 + e^{-x}} \right)$$



Model Training

- Minimize loss:

$$L(w) = \text{NLL}_w(\hat{y}, y) + \lambda_2 \|w\|_2$$
$$\text{NLL}_w(\hat{y}, y) = -\sum_{n=1}^N \sum_{t=1}^T o_{nt} [y_{nt} \log(\hat{y}_{nt}) + (1-y_{nt}) \log(1-\hat{y}_{nt})]$$

weight decay ↓
 ↑
 meth. State
 observed?

- Dropout with different rates for DNA, CpG and Joint module
- Optimization using mini-batch stochastic gradient descent
- Learning rate decay
- Early stopping: terminate if validation loss did not improve

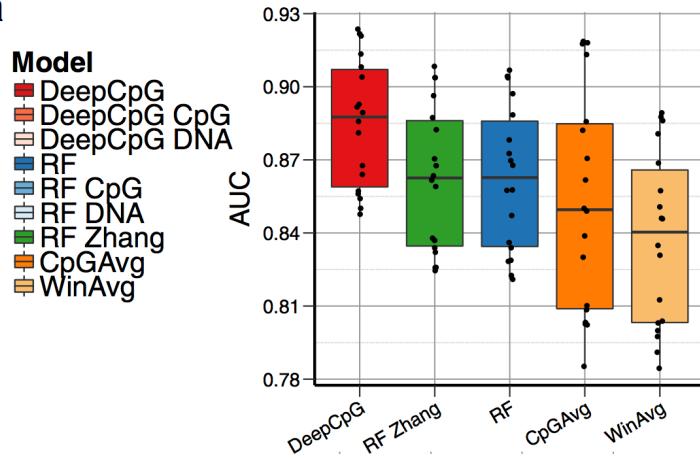
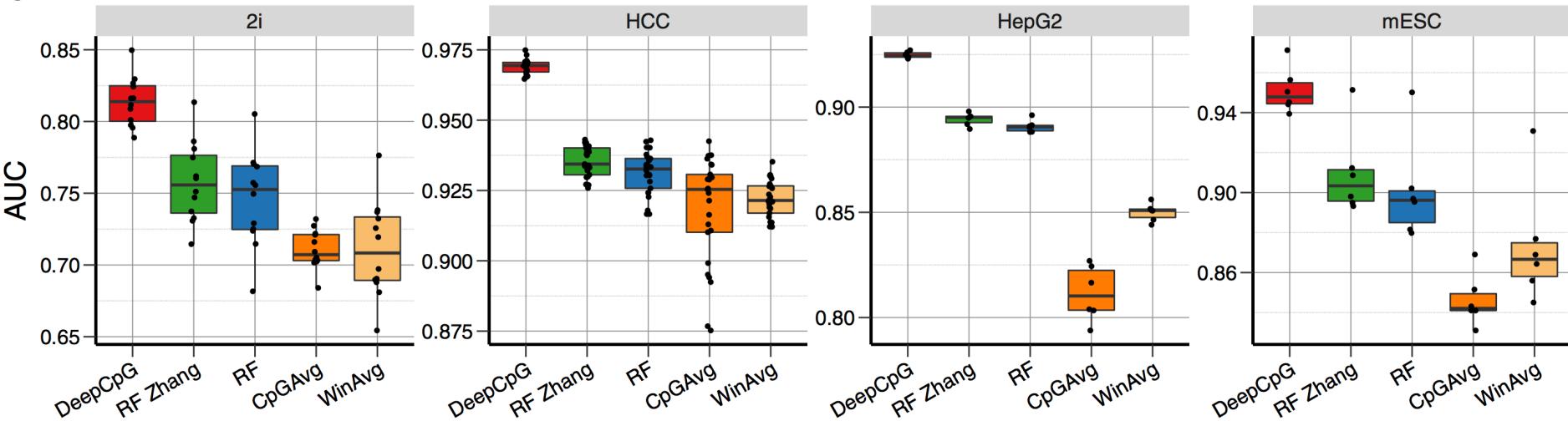
- Hyper-parameter optimization: random sampling on validation set
- Training time:
 - DNA module: 24h
 - CpG module: 12h
 - Joint module: 4h
- Implemented in Python using Theano and KERAS [1]

Performance Evaluation

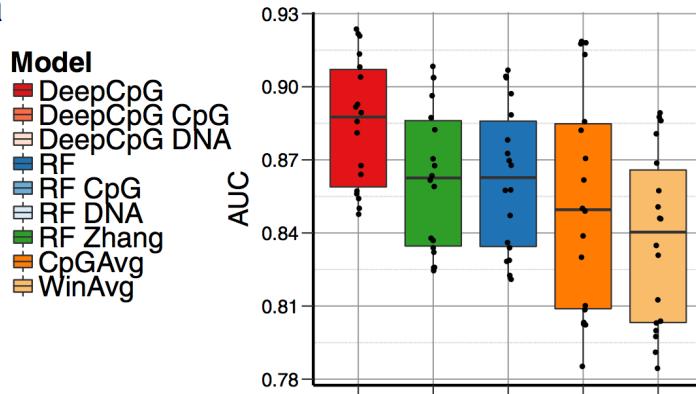
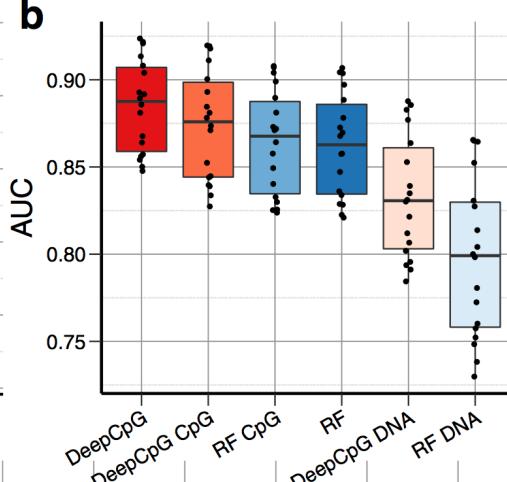
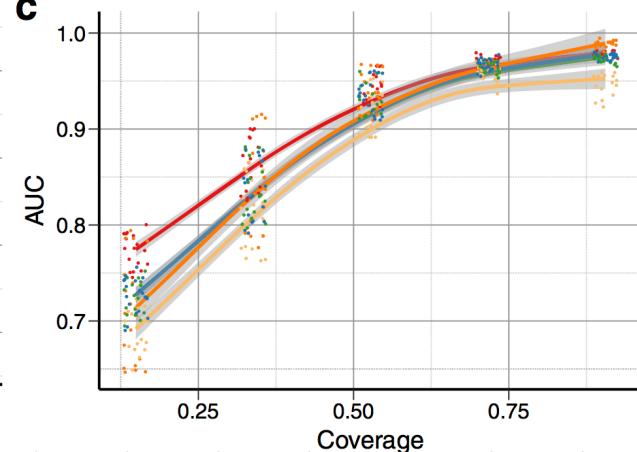
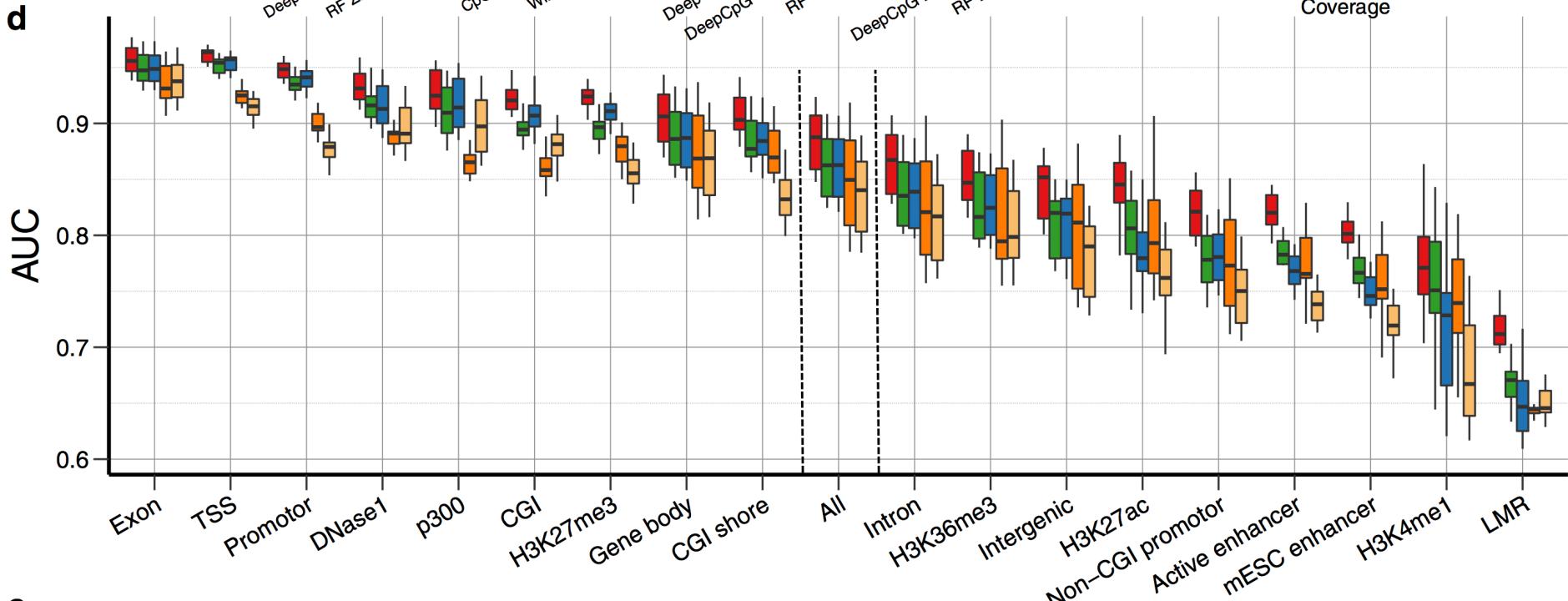
- Comparison with
 - Genomic window averaging (*WinAvg*)
 - Averaging across cells (*CpGAvg*)
 - Random forest classifier using sequence and neighboring CpG states (*RF*)
 - Random forest classifier using comprehensive annotations (*RF Zhang*)
- Holdout validation
 - Different chromosomes in training, test and validation set
- Datasets:
 - Mouse Embryonic Stem Cells (mESCs; 18 serum, 12 2i medium; scWGBS)
 - 6 mESCs (scRRBS)
 - 25 human hepatocellular carcinoma (HCC; scRRBS)
 - 6 HepG2 (scRRBS)



Performance Evaluation

a**e**

Performance Evaluation

a**b****c****d**

Model Interpretability

- Motif analysis of DNA module
- For each filter: select sequence windows with

$$a_{nfi} > 0.5 \max_{ni}(a_{nfi})$$

$$a_{nfi} = \text{ReLU} \left(\sum_{l=1}^L \sum_{d=1}^D w_{fld} s_{n,i+l,d} \right)$$

filter length ↘
↑ multiple filters f ↗
sequence indicator matrix

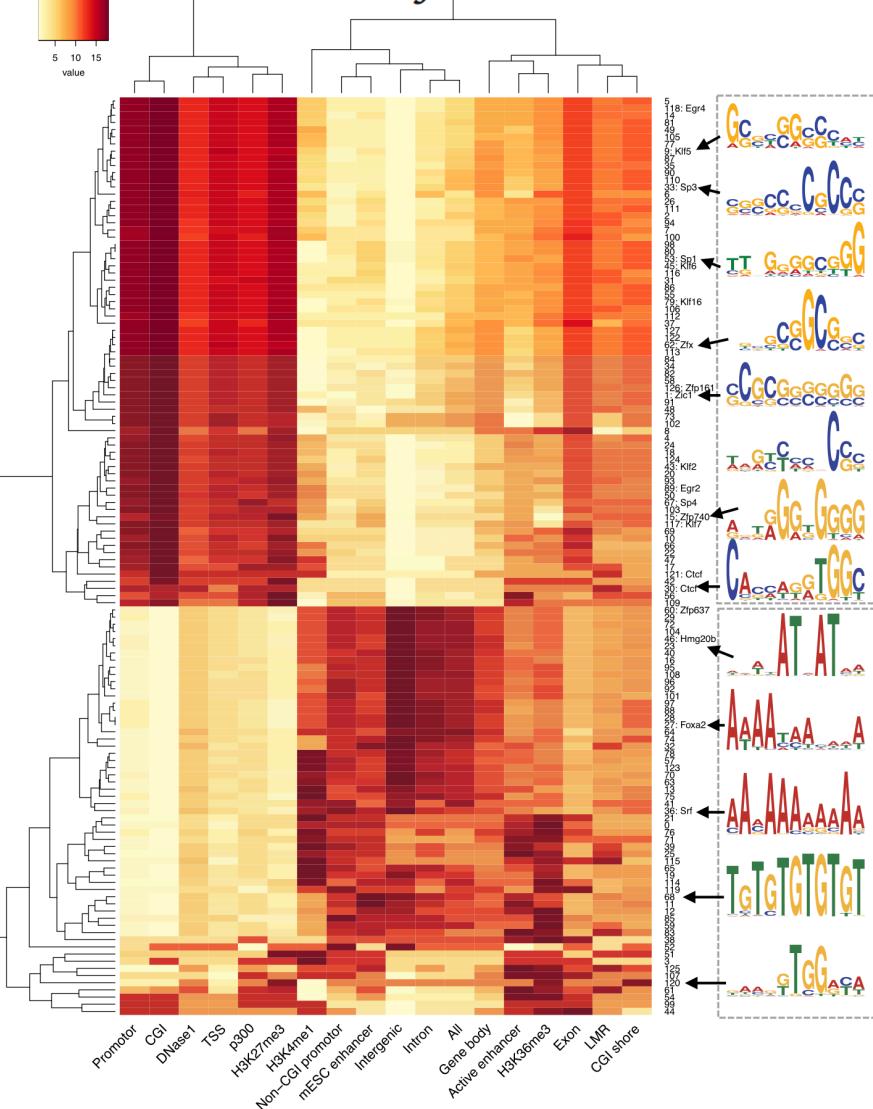
- Filter activity: average weighted mean sequence activity \bar{a}_{nf}

- Allows for the quantification of
 - Association with DNA methylation state:
 $r_{ft} = \text{cor}_n(\bar{a}_{nfs}, \hat{y}_{nt})$
 - Sequence motifs

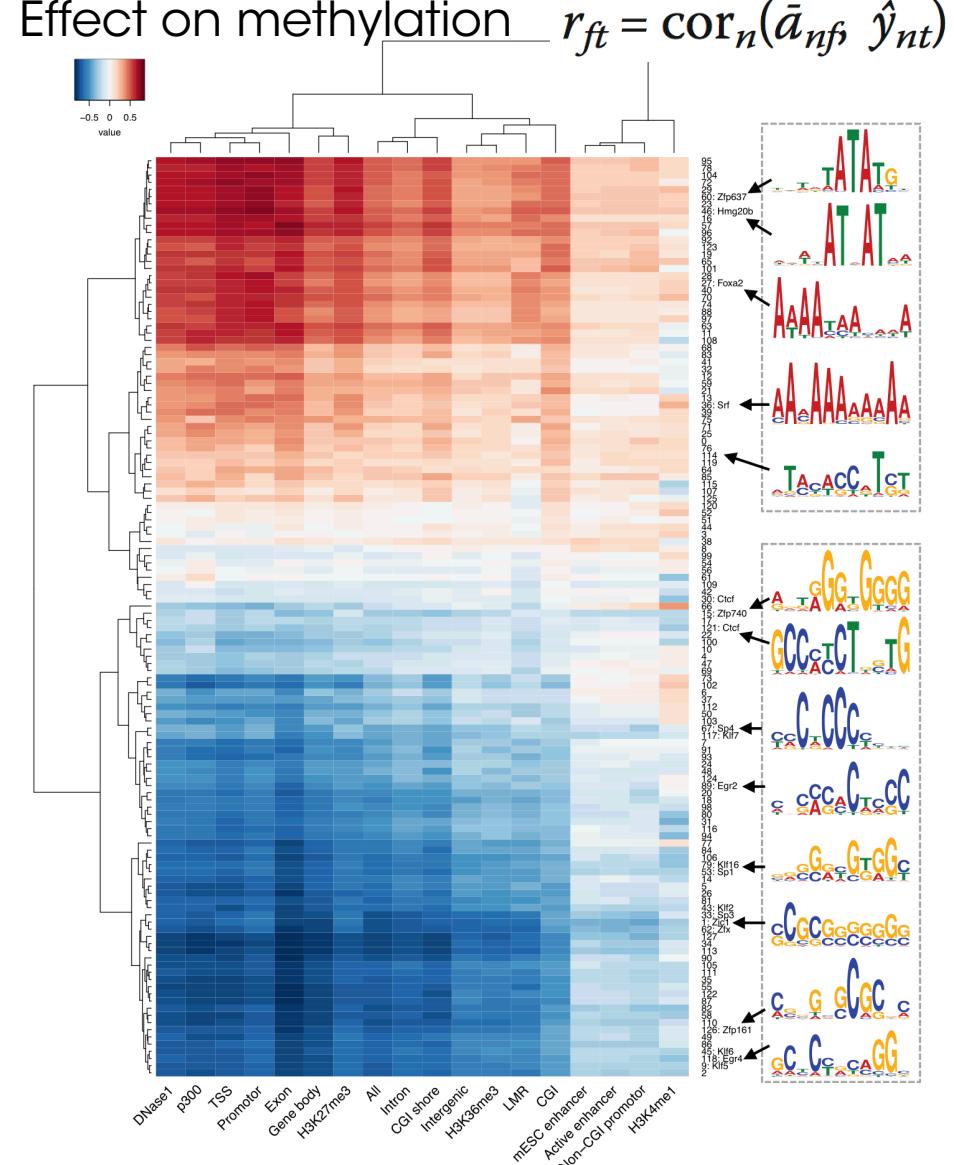


Filter Analysis

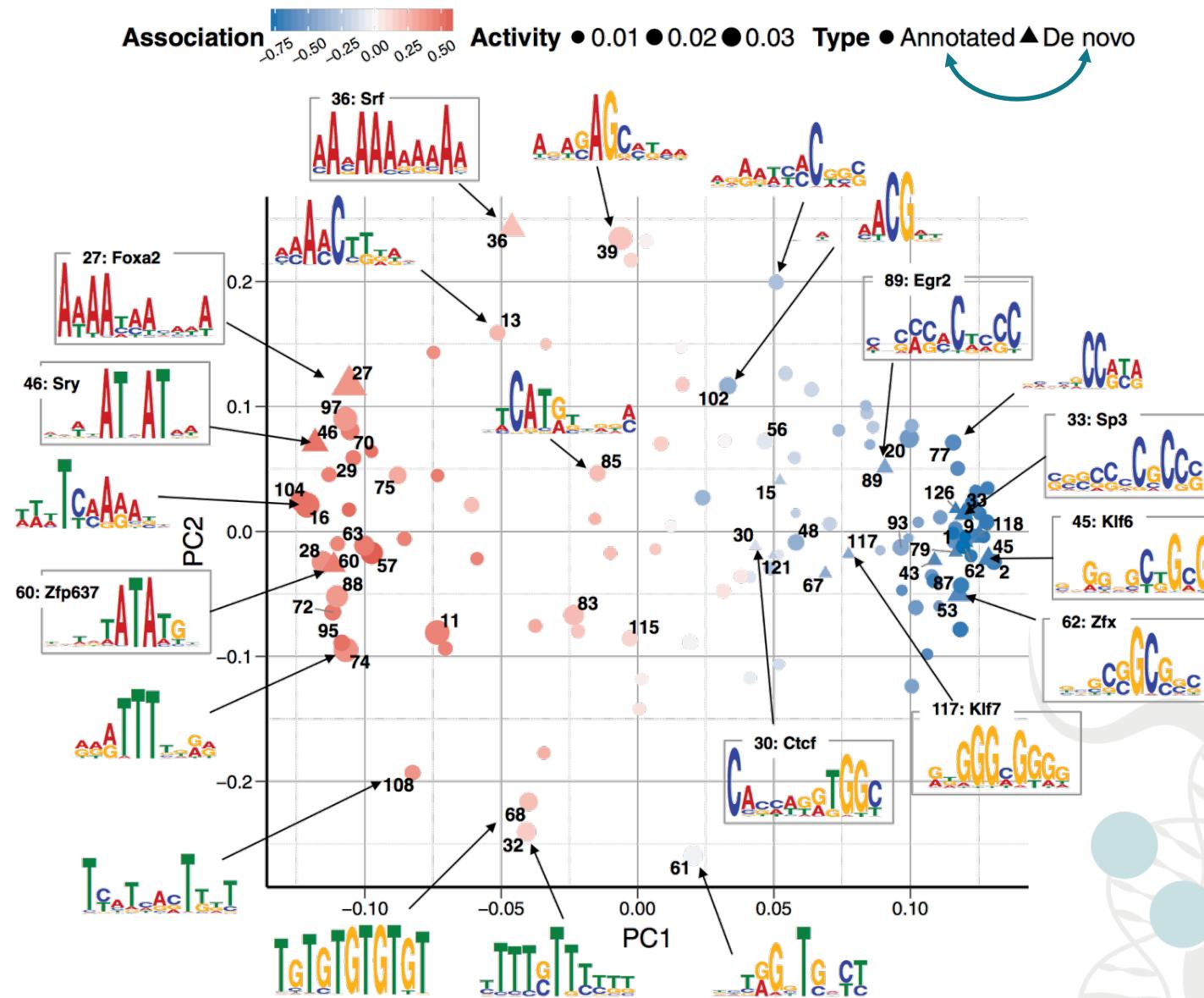
Activity (rank \bar{a}_{nf})



Effect on methylation



Filter Analysis



Effect of Sequence Mutations

- Estimate the effect of SNPs on the predicted methylation:

$$e_{nid}^s = \frac{\Delta \hat{y}_n(s_n)}{\Delta s_{nid}} * (1 - s_{nid})$$

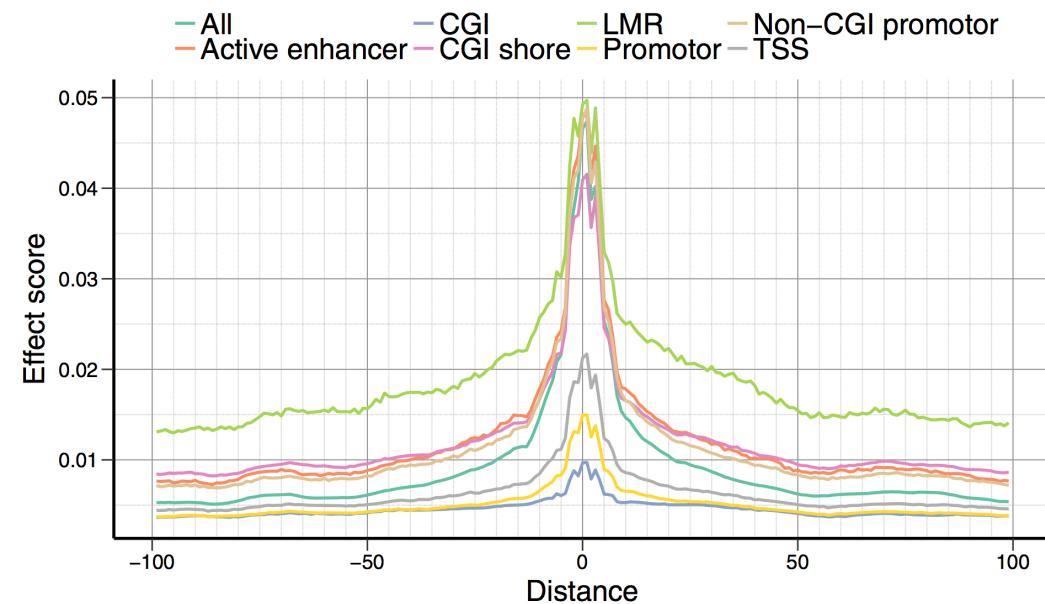
mean predicted methylation across cells

gradient w.r.t. sequence

set effect of WT to 0

$$e_{ni}^s = \max_d |e_{nid}^s|$$

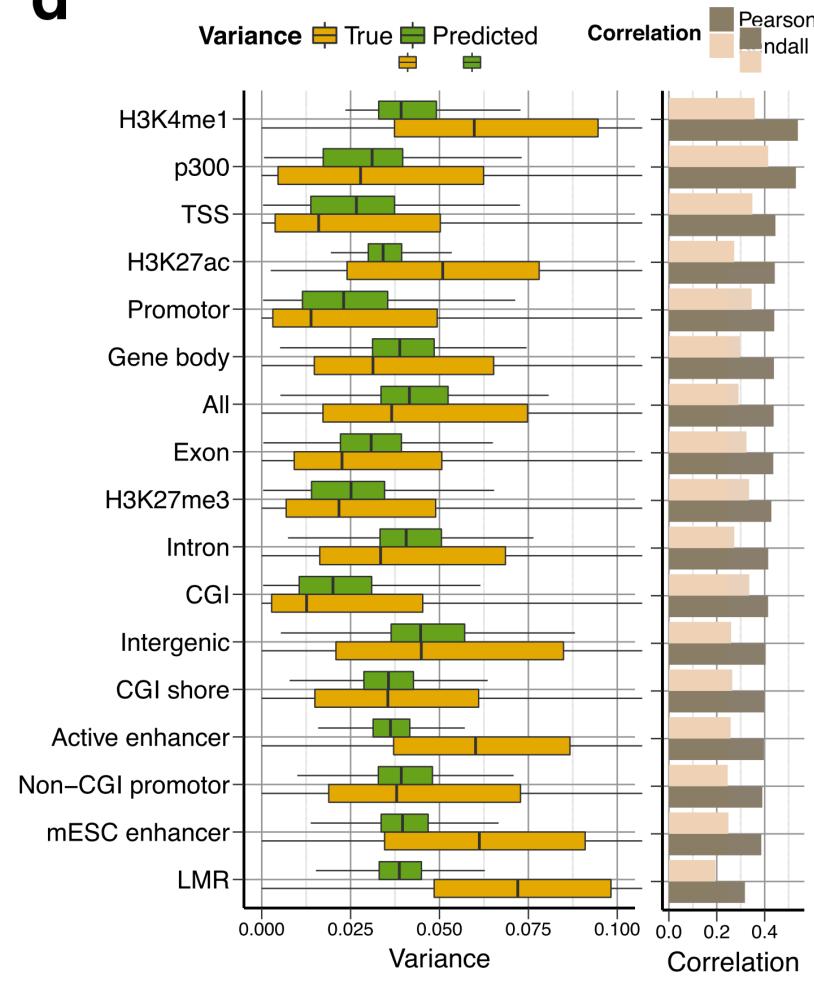
$$e_i^s = \text{mean}_n (e_{ni}^s)$$



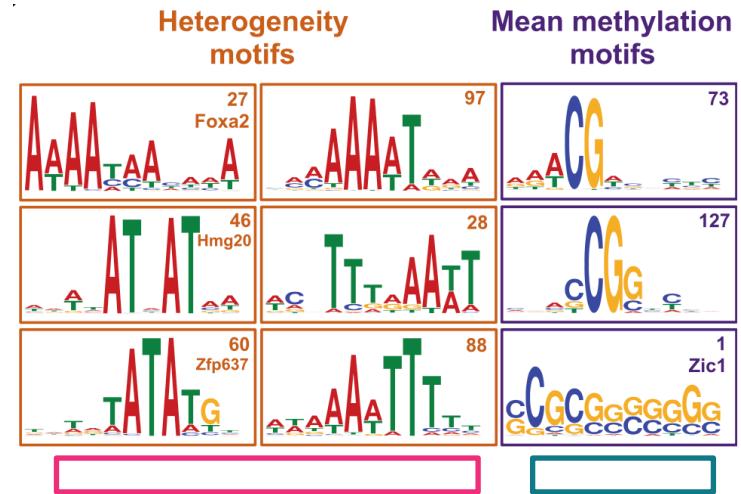
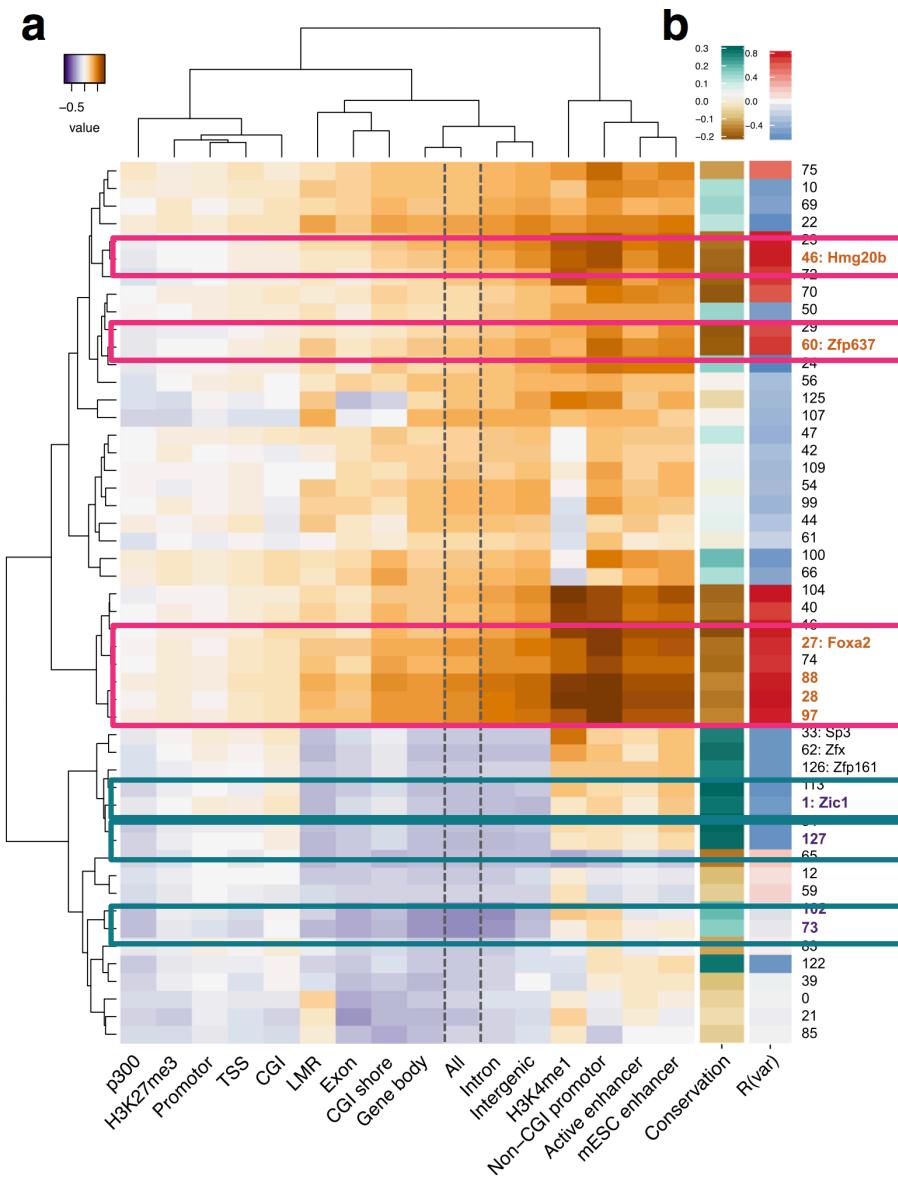
- Known mQTL mutations have significantly larger effect than random

Methylation Variability

- Train another neural net
 - same architecture as the DNA module
 - Output layer:
 - Predict methylation variance and mean methylation across cells, averaged across windows
 - Sigmoid activation
- Compute association (correlation) of motif activity with mean methylation and variance

d

Methylation Variability



- Variance-associated motifs
 - AT-rich
 - Enhancers
 - Associate weakly with changes in gene expression

Summary

- DeepCpG employs deep neural networks to accurately predict missing methylation states based on DNA sequence and observed methylation
- Biologically relevant DNA sequence motifs are learned by the model
- Models can also be used to learn methylation variability



Open Issues

- Can the models be interpreted beyond motifs?
- Sequence-dependent DNA methylation?
- Other epigenetic features?
- Scaling to thousands/millions of single cells?

