# Application of Machine Learning Methods to Genome-wide Maps of Histone Methylations
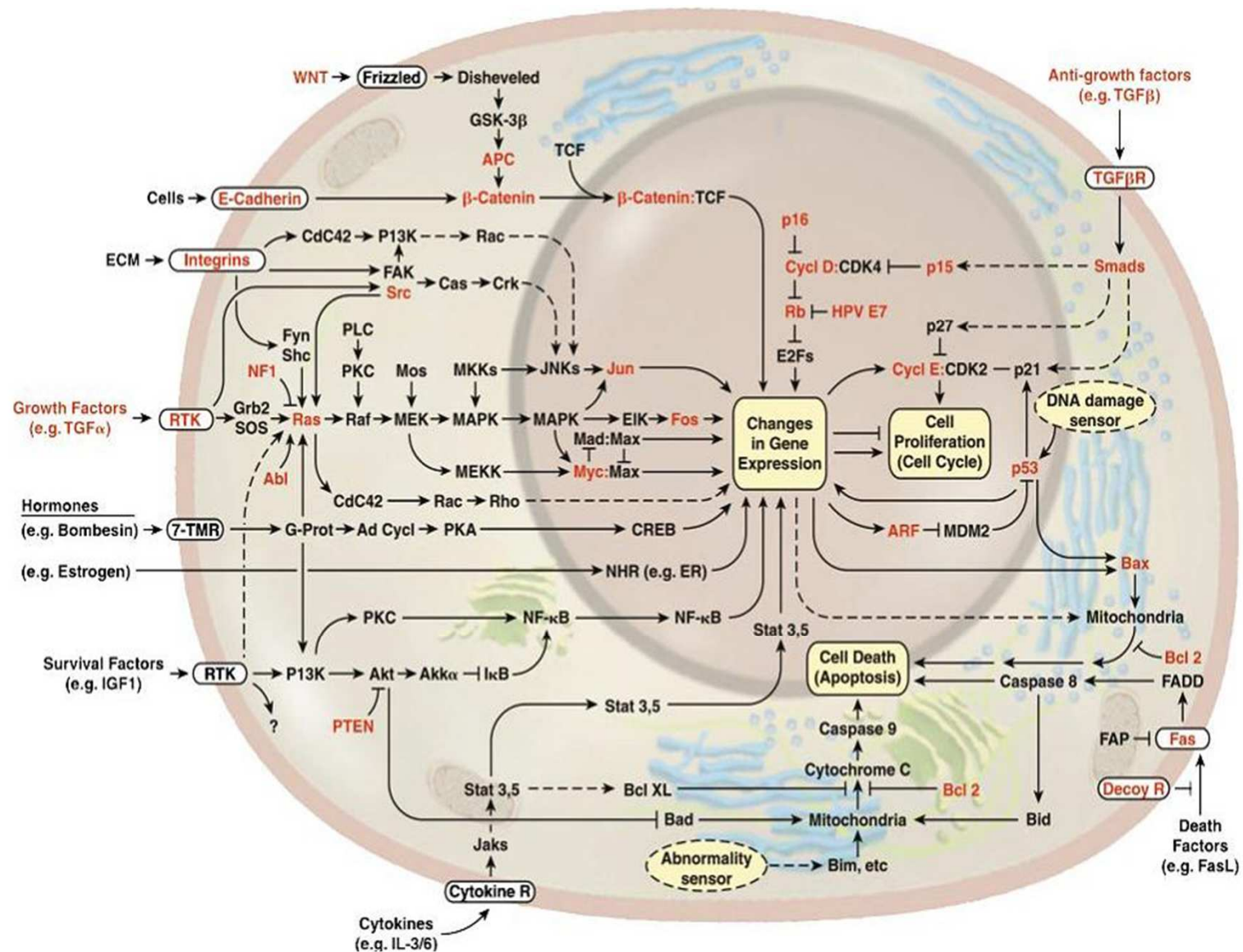
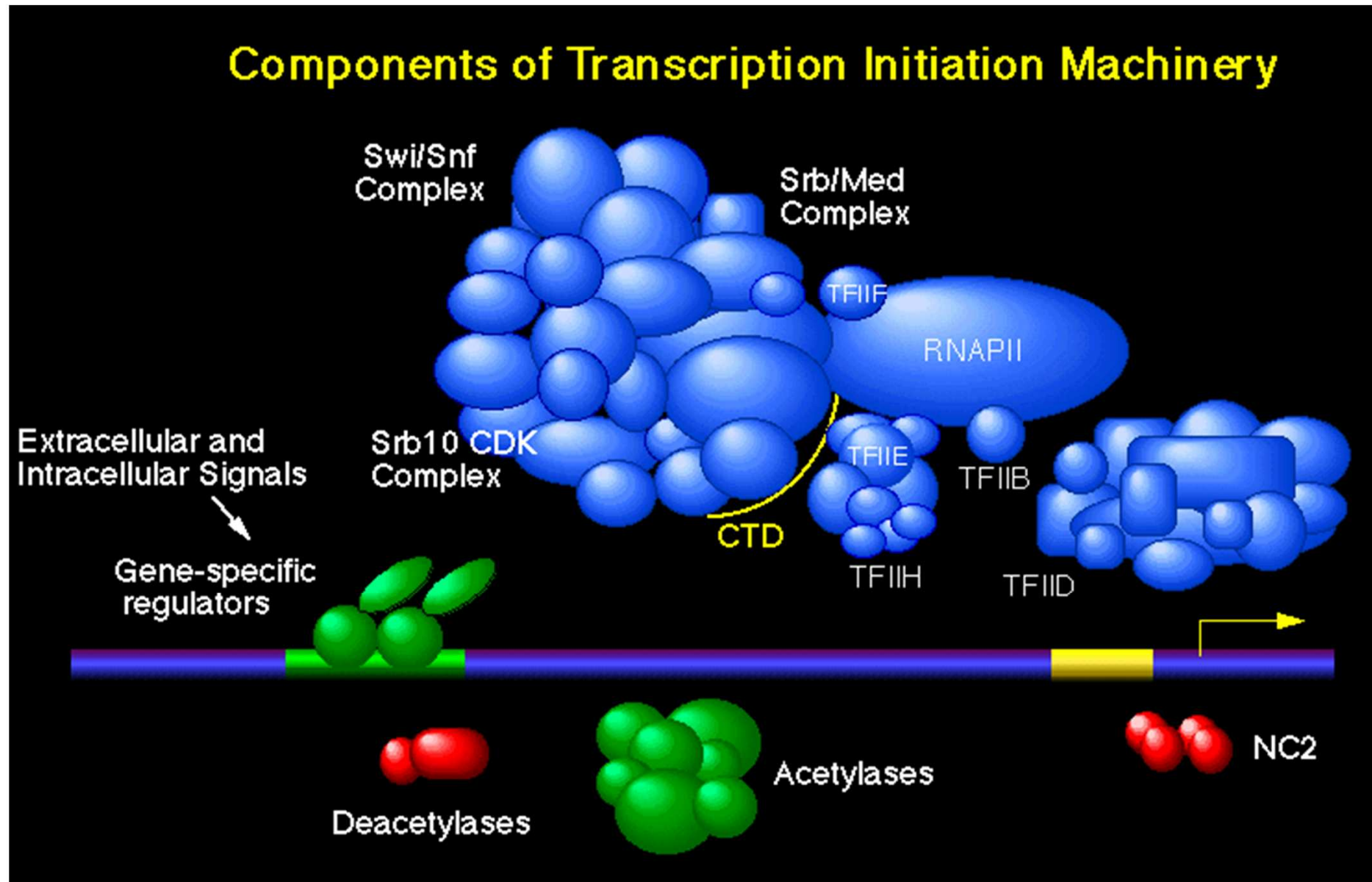Stefan Bekiranov

University of Virginia

# Outline

- Molecular Biology Introduction/Overview
- High Throughput Genomics Technologies
  - Gene Expression Microarrays
  - High Throughput Sequencing: ChIP-Seq
- Machine Learning Methods & Our Results
  - Multilinear Regression
  - Multivariate Adaptive Regression Splines (MARS)
  - Prediction: H4R3me2s is globally repressive
- Experimental Studies on H4R3me2s
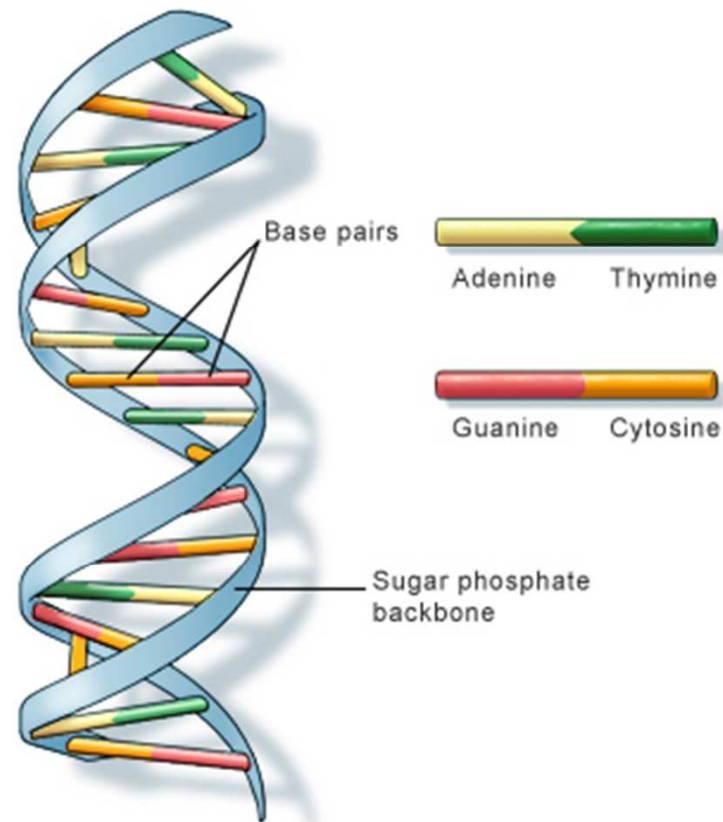
# Molecular Biology Introduction/Overview

# Signal Transduction
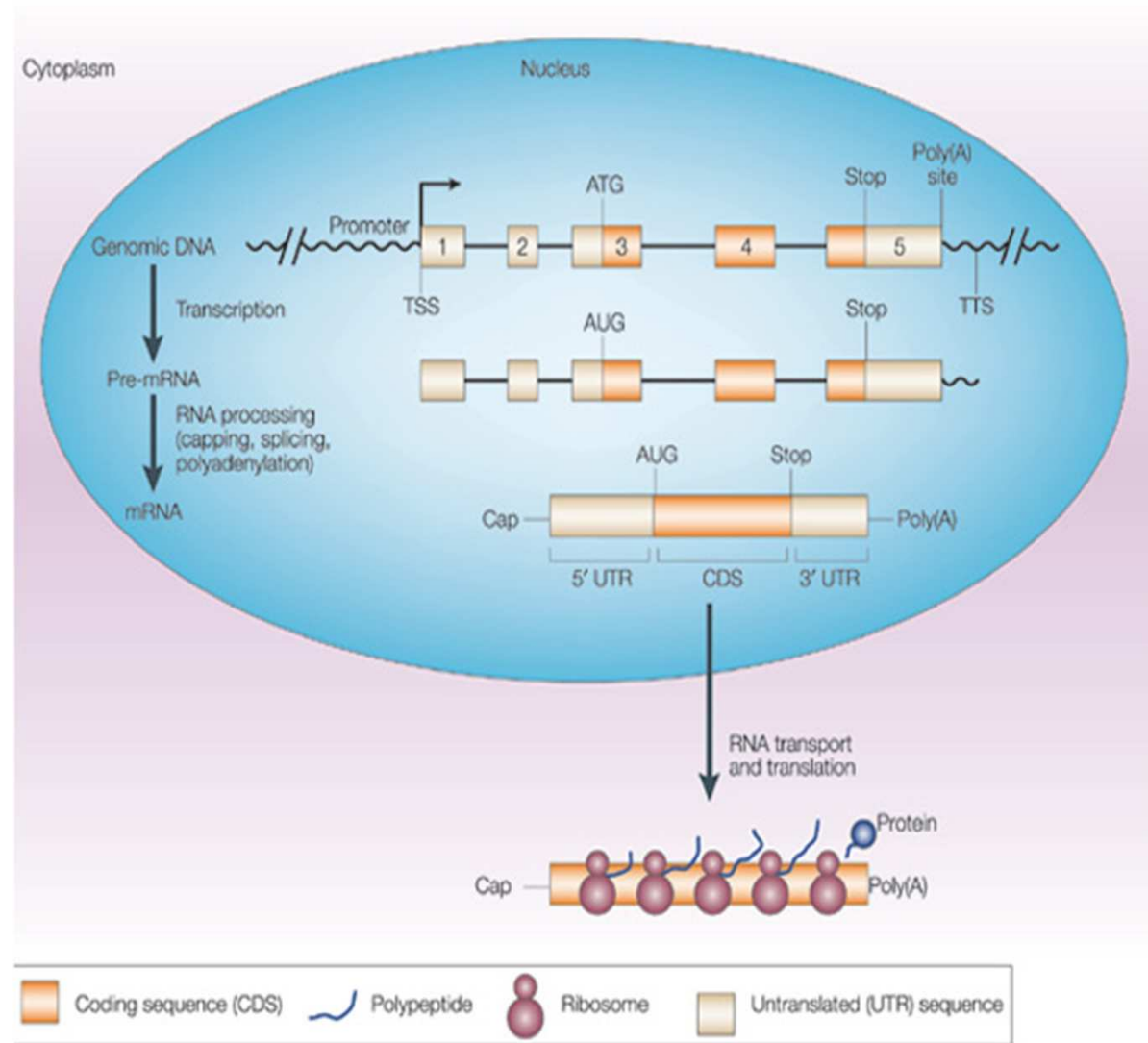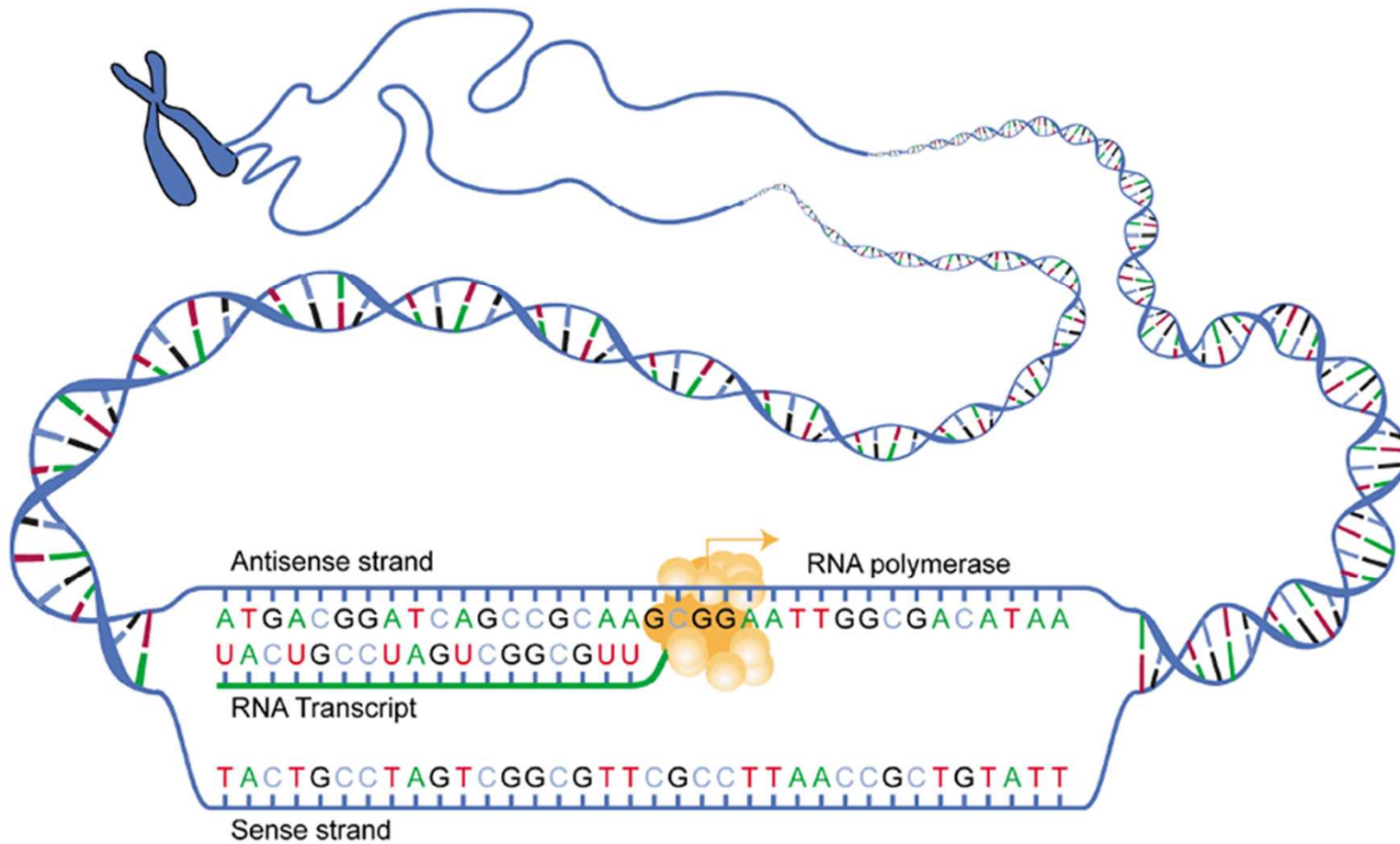
# Transcription Initiation: Turn "on" a gene



Components of Transcription Initiation Machinery

# DNA



Base pairs

Adenine    Thymine

Guanine    Cytosine

Sugar phosphate backbone

U.S. National Library of Medicine

# Gene -> mRNA -> Protein



Nature Reviews | Genetics

# mRNA -> Protein

# Chromatin



Short region of DNA double helix — 2 nm

"Beads on a string" form of chromatin — 11 nm

30-nm chromatin fibre of packed nucleosomes — 30 nm

Section of chromosome in an extended form — 300 nm

Condensed section of chromosome — 700 nm
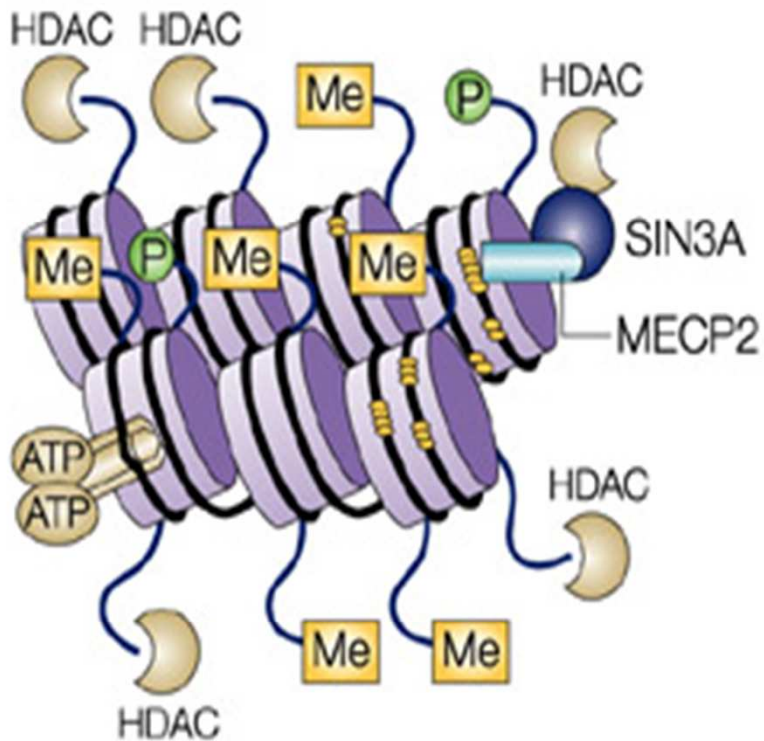
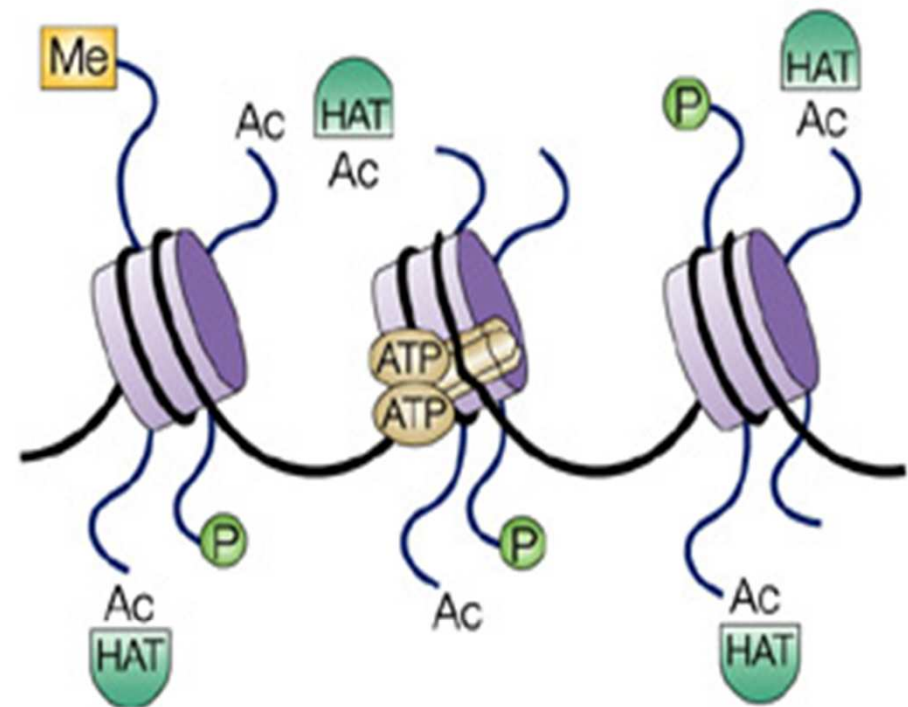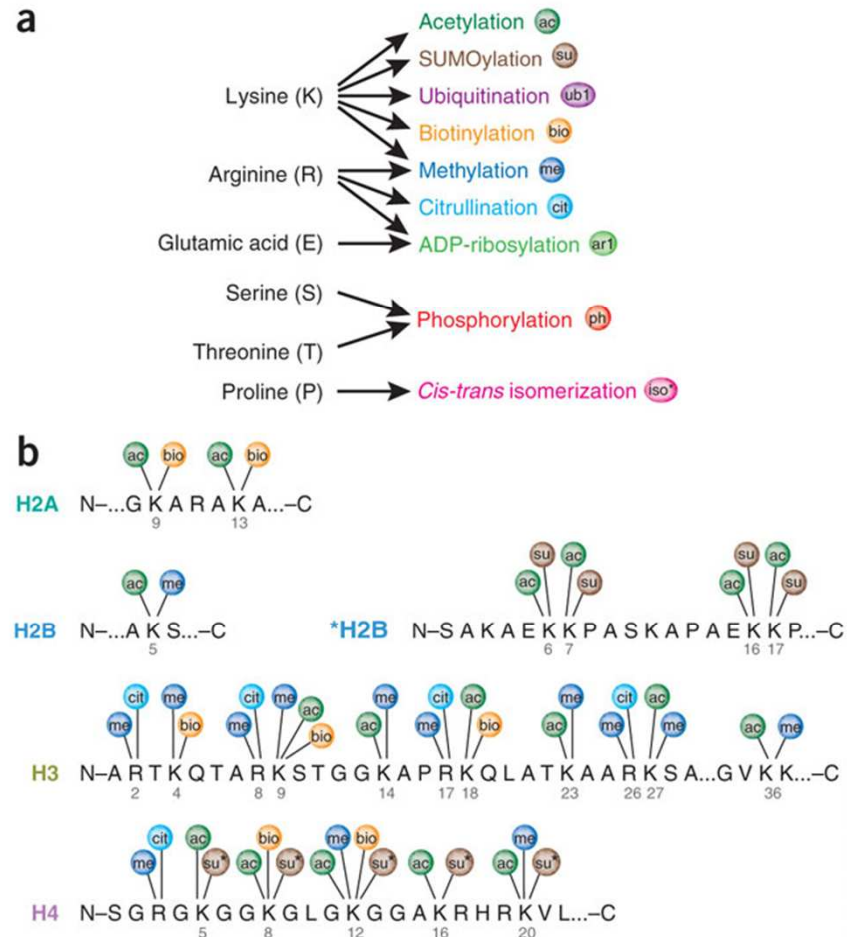Entire mitotic chromosome — Centromere — 1,400 nm

# Histone Modifications



**a** Closed chromatin: transcriptional repression

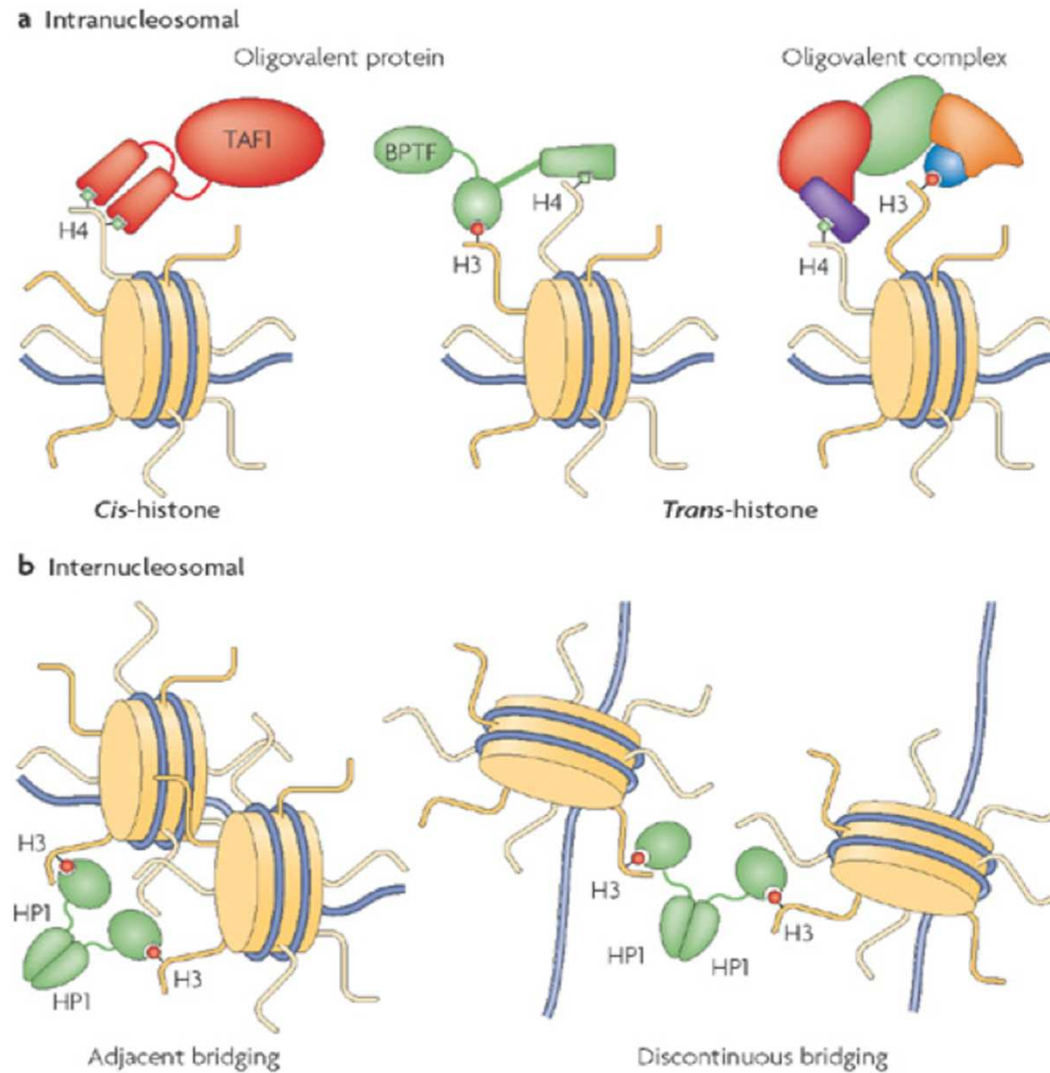**b** Open chromatin: transcriptional activation

# Histone Modifications

# Cross-regulation of Histone Modifications

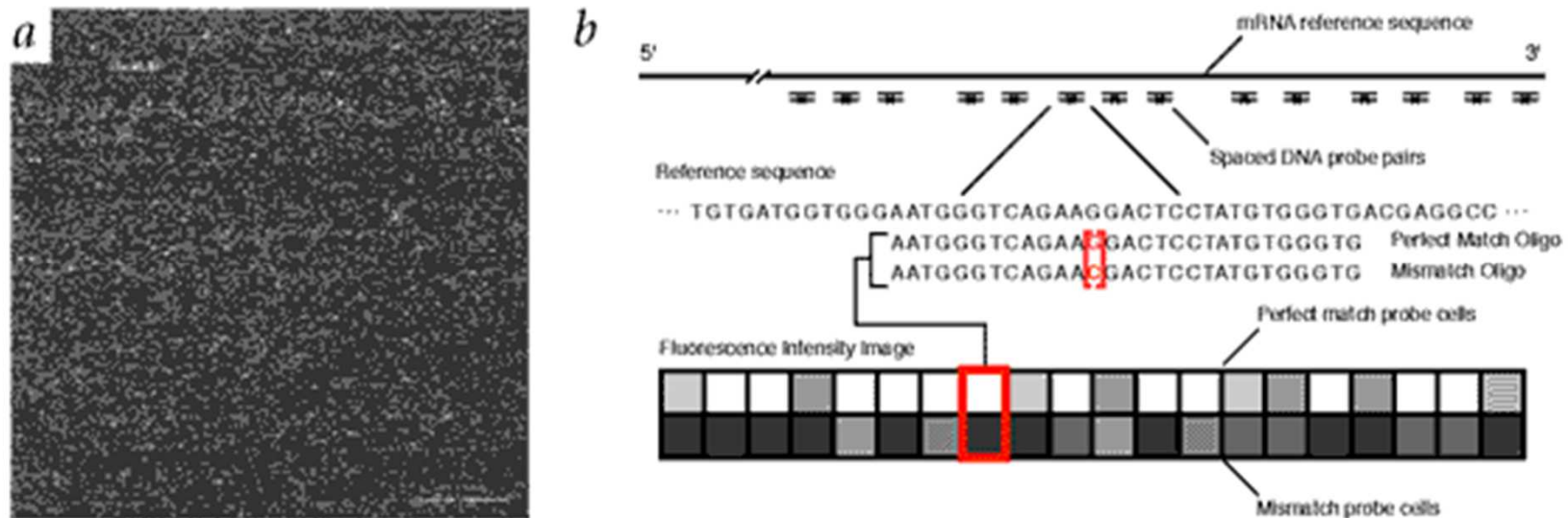# Multivalent Chromatin Engagement



Nature Reviews | Molecular Cell Biology

# High Throughput Genomics Technologies

# Affymetrix Probe Set: Gene Expression Arrays (Human U133 Array ~2002)

# Disruptive Technology: High Throughput Sequencing

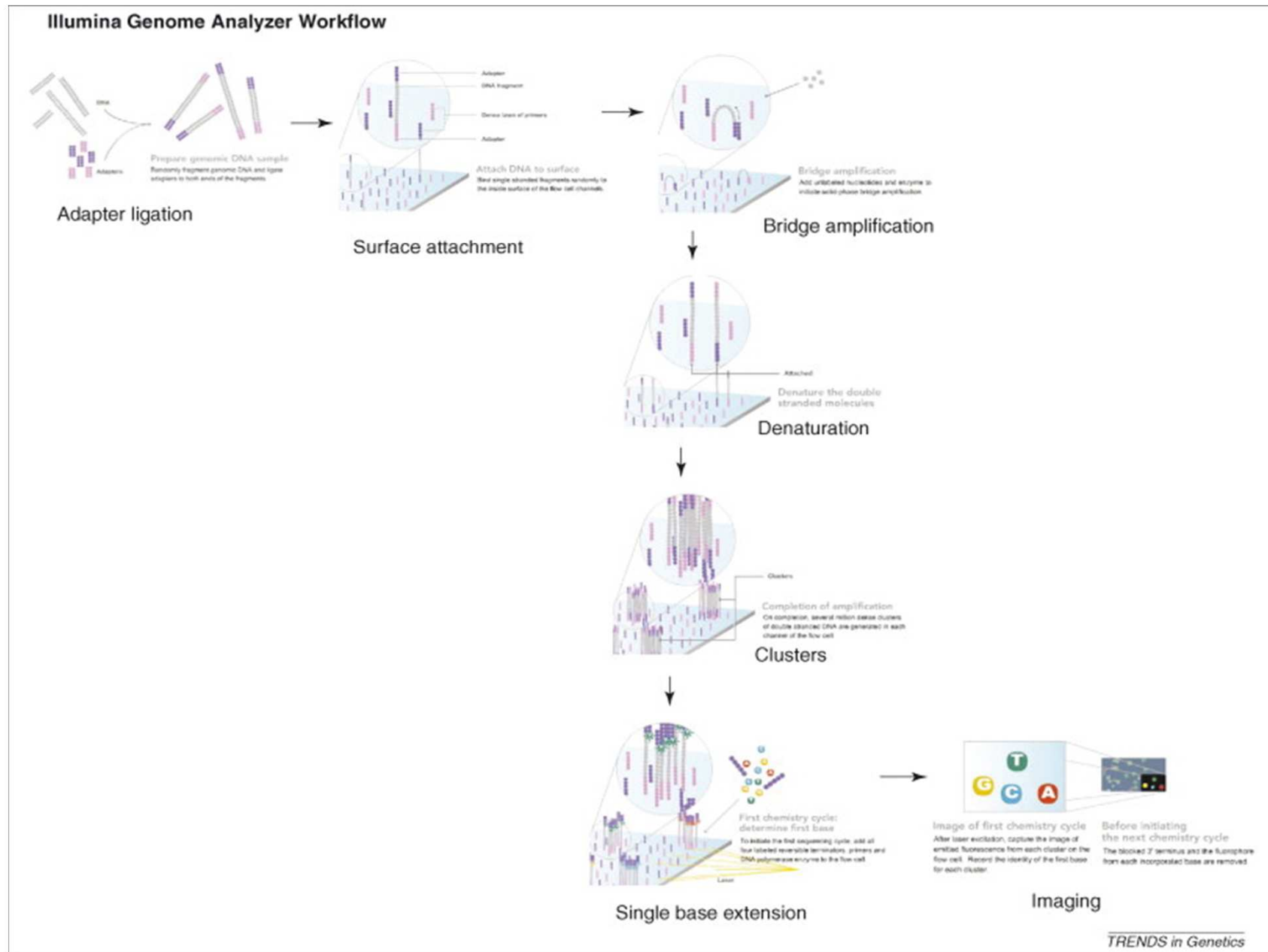Prepare gDNA Library
3 hours hands-on, 6 hours total time

Generate Clusters
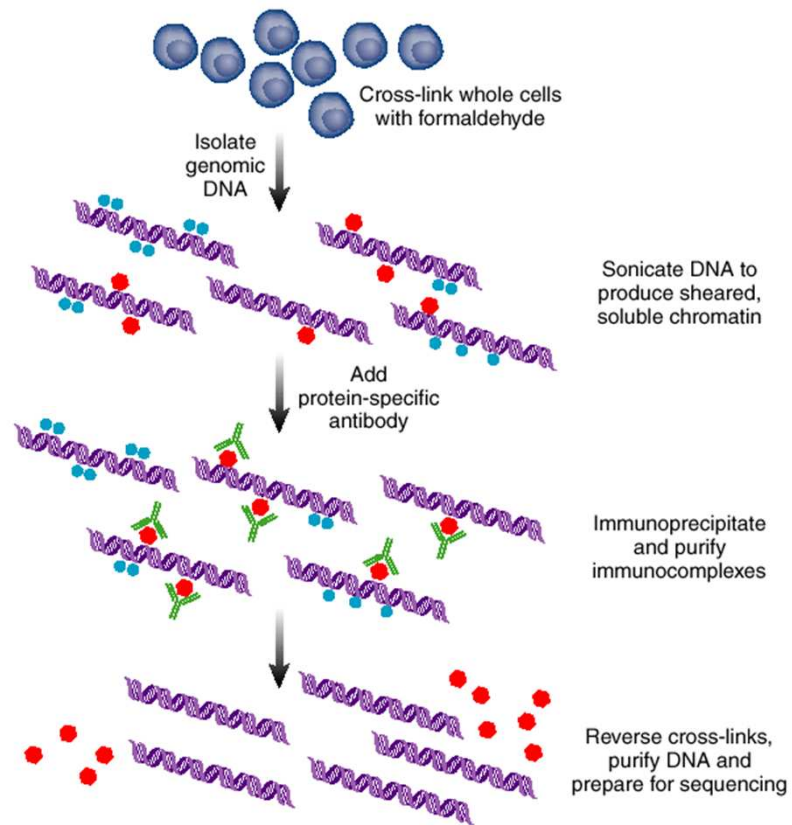<1 hour hands-on time, 5 hours total time
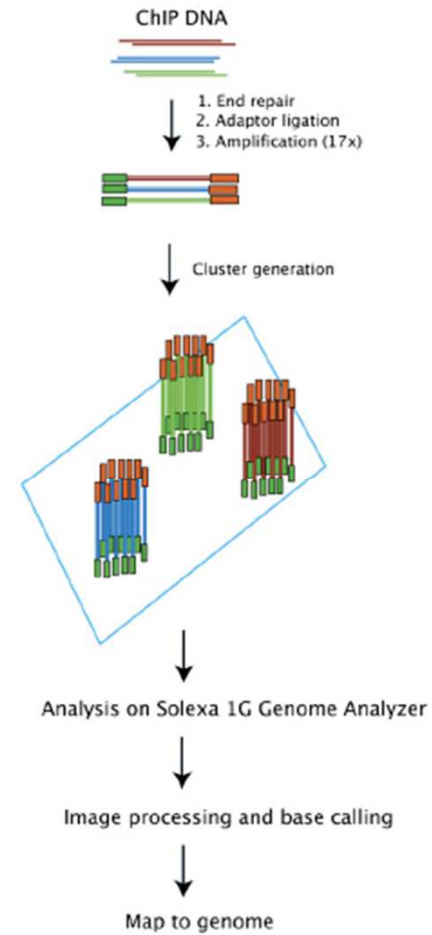
Sequence Clusters
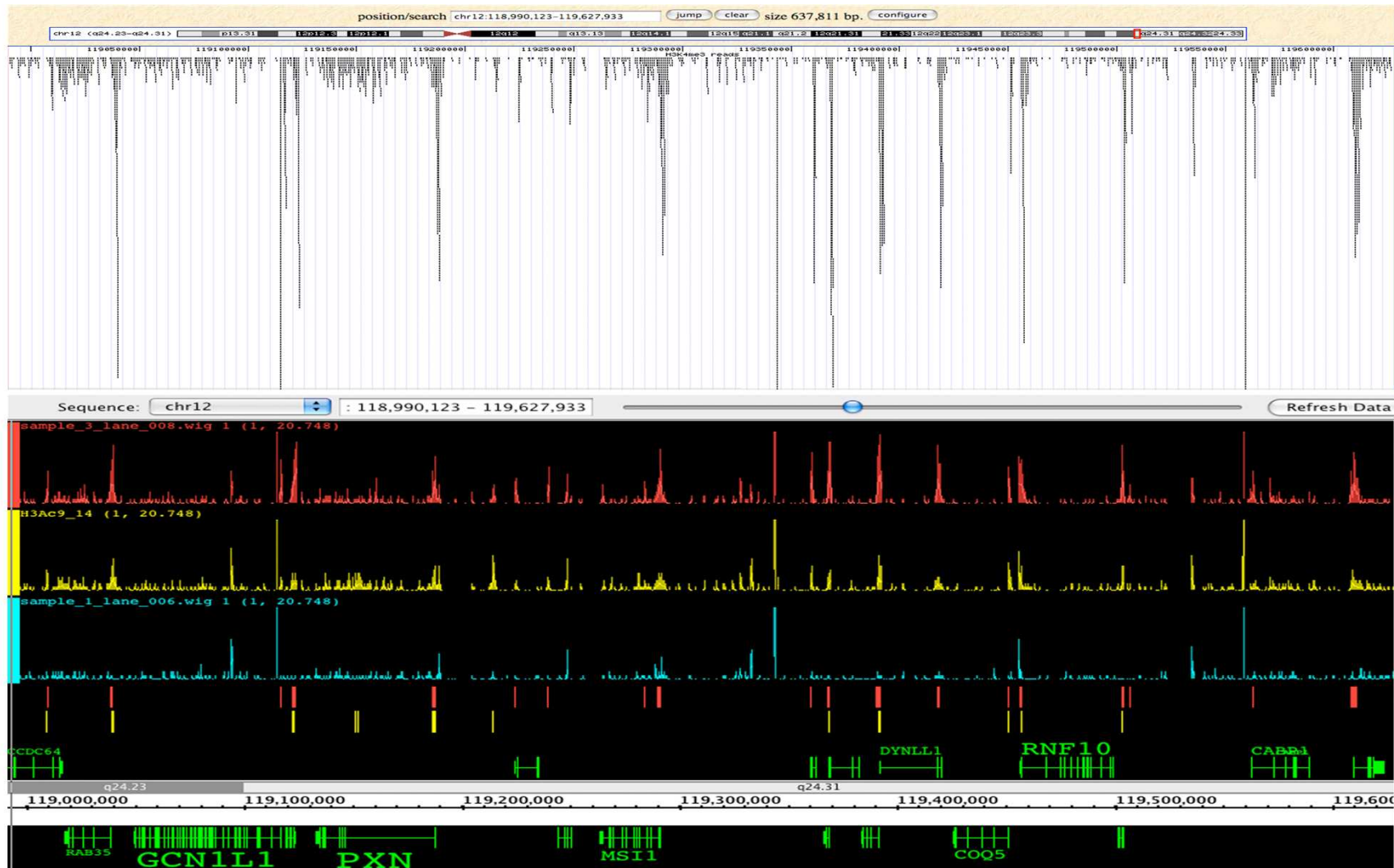2.5 days single read (36 bases)

# Illumina (Solexa) Workflow (~2006)

# ChIP-Seq

# Map Reads; Calculate Enrichment Profile

# ChIP-Seq Data Sets

- In CD4$^+$ T cells, Keji Zhao Lab (NIH) Generated Maps of:
  - 20 Histone Lysine and Arginine Methylations
  - 18 Histone Acetylations
  - Pol II, CTCF and H2A.Z
  - Nucleosome Positions
  - Effector Proteins

- Epithelial to Mesenchymal Transition (EMT)
  - ~15-20 histone modifications/variants
  - ~6 Master-Switch Transcription Factors including Twist, Snail, Slug, SIP1, SMAD, NF-κB

# Read Enrichments Stratified By Gene Expression Levels

# Rogues Gallery

# Machine Learning Methods and Results

# Questions

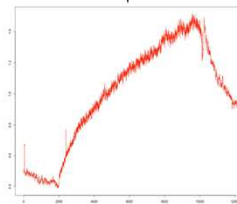- Why does the cell require ~100 histone modifications to maintain two states (open or closed chromatin)?
  - Histone Code Hypothesis (David Allis)
- Are two-body, three-body,… interaction terms strong, dominant given evidence of multivalent readers and writers of histone modifications?
- Can predictive modeling recapitulate known and uncover unknown dependencies among histone modifications?

# (Machine) Learning the Histone Code



ChIP-Seq Data

(1) Calculate average enrichment profile (i.e., template) for each histone mark from ChIP-Seq data

$$X_j^k = \frac{\sum_i c_{i,j}^k t_{i,j}}{\sum_i t_{i,j}^2}$$

(2) Calculate an amplitude value for each gene for each mark using the templates to spatially weight the enrichment

Input = $X_j^k = \dfrac{\sum_i c_{i,j}^k t_{i,j}}{\sum_i t_{i,j}^2}$
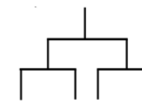
Output = log2 Gene Expression

(3) Use mark amplitudes as inputs and log2 gene expression as outputs for regression models

Nonlinear MARS Model
- Greedy algorithm
- Fit using GCV score

Multilinear Regression Model
- Stepwise regression
- Fit using 10-fold cross validation

Regression Tree
- Greedy Algorithm
- Fit using 10-fold cross validation

# Estimating Amplitudes

- Normalize Templates: $\dfrac{1}{N}\displaystyle\sum_{i=1}^{N} t_{i,j} = 1$

- Least Squares Sum: $Q_j^k = \displaystyle\sum_i \left(c_{i,j}^k - X_j^k t_{i,j}\right)^2$

- Least Squares Solution: $X_j^k = \dfrac{\displaystyle\sum_i c_{i,j}^k t_{i,j}}{\displaystyle\sum_i t_{i,j}^2}$

- Solution for Uniform Template: $X_j^k = \dfrac{1}{N}\displaystyle\sum_i c_{i,j}^k$

# Stepwise Linear Regression

Model: $Y^k = \beta_0 + \sum_{j} \beta_j X_j^k + \sum_{j<l} \beta_{j,l} X_j^k X_l^k + \sum_{j<l<m} \beta_{j,l,m} X_j^k X_l^k X_m^k + \varepsilon^k$

1. Fit the initial (seed) model
2. Use F-test to calculate p-values of potential additional term from pool of candidate terms.
   - If min(p) < 0.05, add term.
   - Repeat until no terms with p < 0.05.
   - If min(p) > 0.05, go to step 3.
3. If any terms in model have p > 0.05, remove term with max p-value and go to step 2; otherwise, end.

# Multilinear Model Terms

- Ran stepwisefit 100 times with randomly seeded models.
- Required terms to appear in 35% of the models: 167 core terms.
- 10-fold cross validation with seed model including 167 core + 60 random terms.
- Ran each fold 10 times and selected the model with best test MSE.
- Pruned model using F-test on test data.
- Required a term to appear in ≥ 5 folds: 24 terms total in final model.

| Term | B | Z-score | p-value | Impact |
|---|---|---|---|---|
| H3K79me1 | 6.74 | 18.23 | 0 | 1.33 |
| H3K36me3 | 4.09 | 17.80 | 0 | 0.92 |
| H3K79me3 | 3.08 | 23.91 | 0 | 0.60 |
| H4K20me1 | 0.977 | 21.44 | 0 | 0.45 |
| H3K4me2-H3R2me1 | 18.27 | 7.85 | 1.66e-15 | 0.44 |
| H3K4me2*-H3K9me1* | -1.58 | -3.34 | 3.5e-4 | -0.12 |
| H4R3me2 | -11.12 | -13.23 | 0 | -0.30 |
| H3K27me2*-H3K36me3 | -31.77 | -8.91 | 0 | -0.31 |
| H3K27me2*-H3K79me1 | -56.54 | -9.53 | 0 | -0.45 |
| H3R2me1 | -11.94 | -16.50 | 0 | -0.60 |

# Multivariate Adaptive Regression Splines

$$Y^k = c_0 + \sum_{i=1}^{n} c_i b_i(\vec{X}_i^k)$$

$b_i(\cdot)$ is a basis function that is made up of either one or a product of two or more hinge functions. Hinge functions are splines that take on the form $h(X_j^k) = \max\left(0, X_j^k - X_j^{k^*}\right)$ or

$h(X_j^k) = \max\left(0, X_j^{k^*} - X_j^k\right)$ where $X_j^{k^*}$ is a special constant known as a *knot*.

Model Built in One Forward and One Reverse Pass:

Forward Pass:
(1) Add Intercept = mean of response variable $Y^k$
(2) Add 1-body terms which maximally reduce RSS at every step
(3) Add 2-body terms (constrained to include one of the 1-body terms) which maximally reduce the RSS at every step
(4) Add 3-body terms (constrained to include one of the 2-body terms) which maximally reduce the RSS at every step

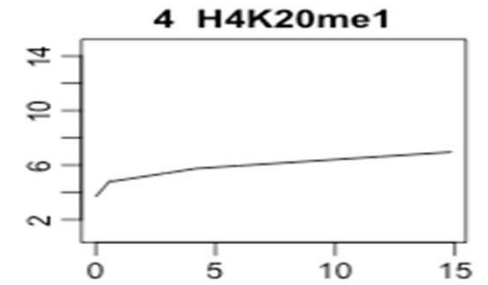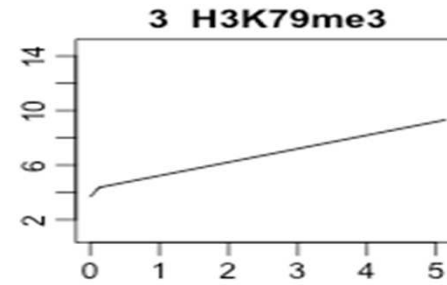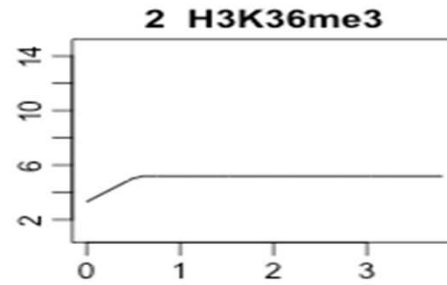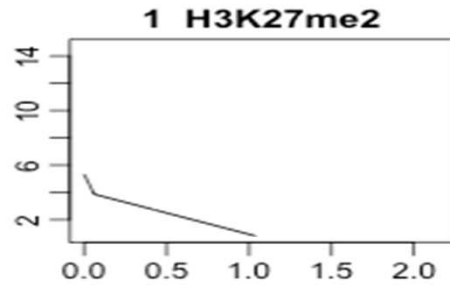Reverse Pass:
(1) Remove terms to optimize a Generalized Cross Validation Score which penalizes model complexity by dividing the RSS by the effectve number of degrees of freedom in the model
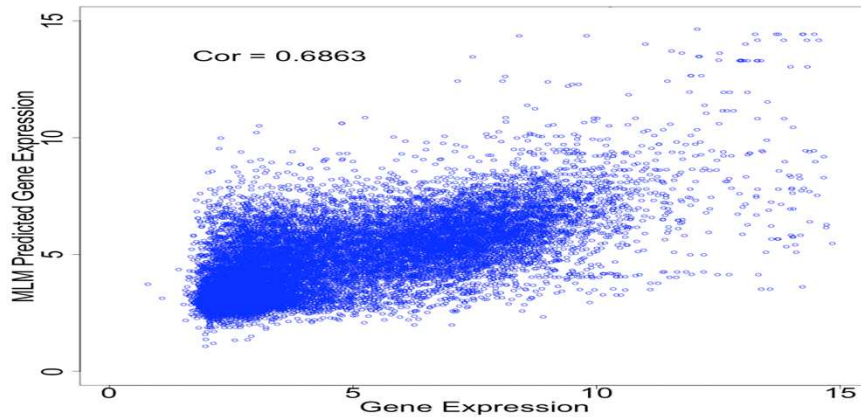
Final Model Contains 24 terms.
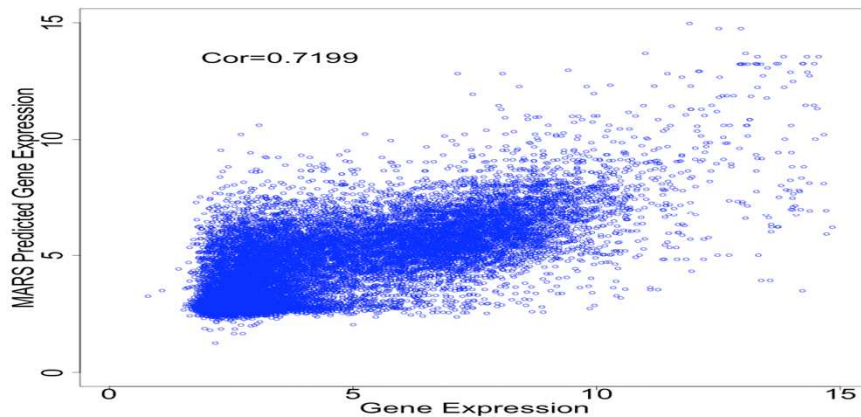
# MARS Model Terms


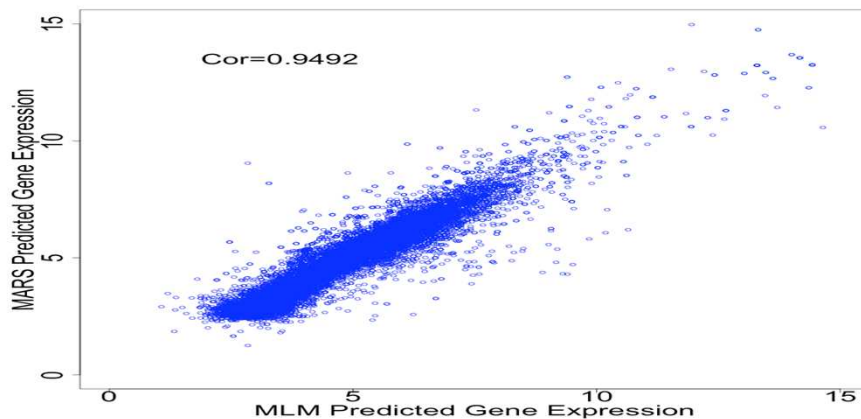
earth(formula=exp_val~.,data=a...

# Scatter Plots



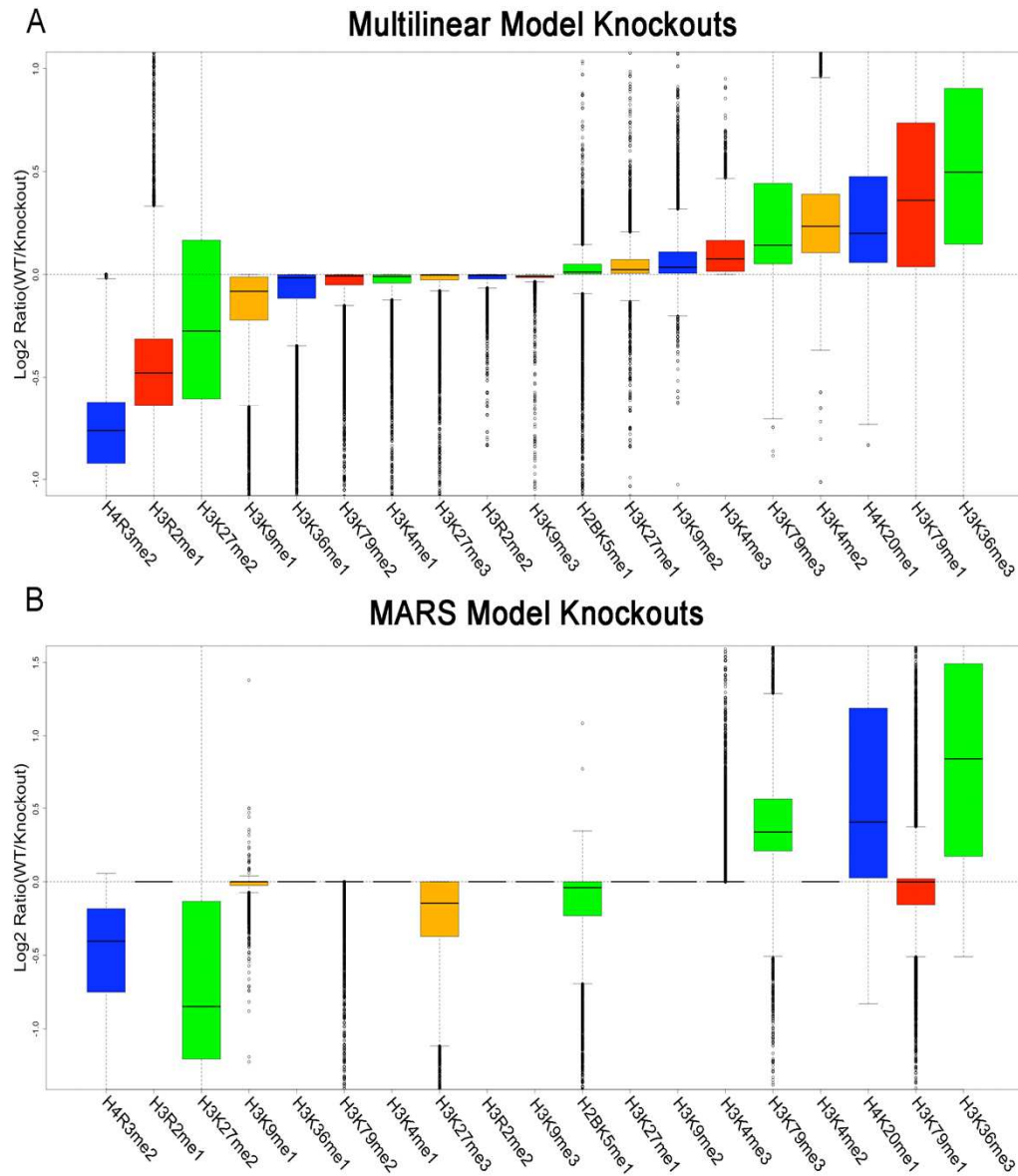Cor = 0.6863

Multilinear Model vs Gene Expression
(24 terms in model)

Cor=0.7199

MARS Model vs Gene Expression
(24 terms in model)

Cor=0.9492

MARS vs Multilinear Model

# *In-silico* Knockout Predictions



A    Multilinear Model Knockouts

B    MARS Model Knockouts
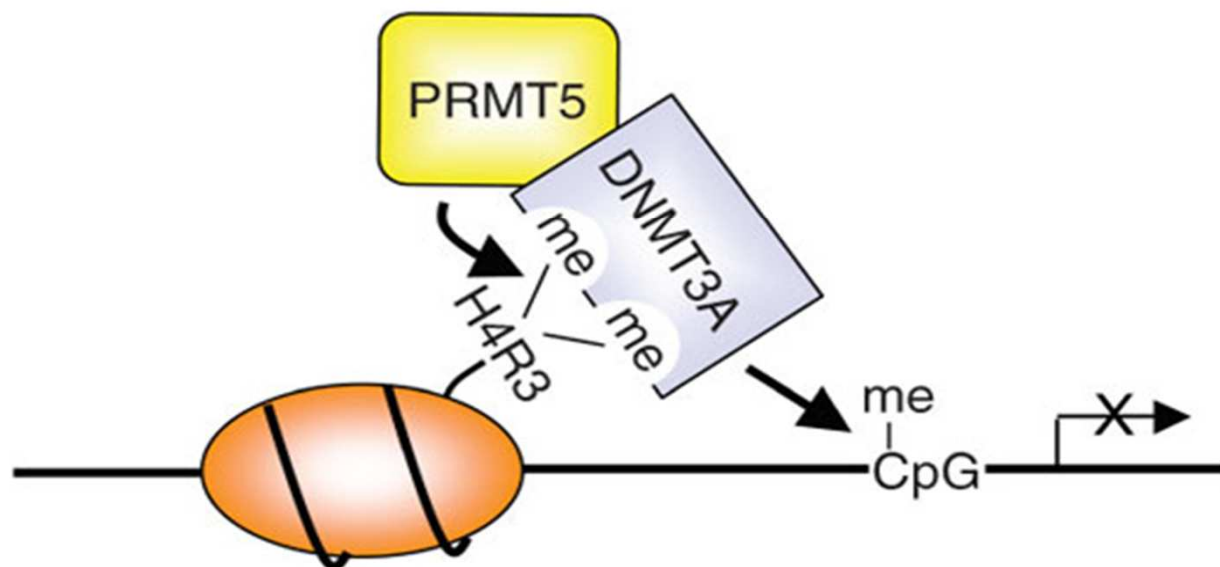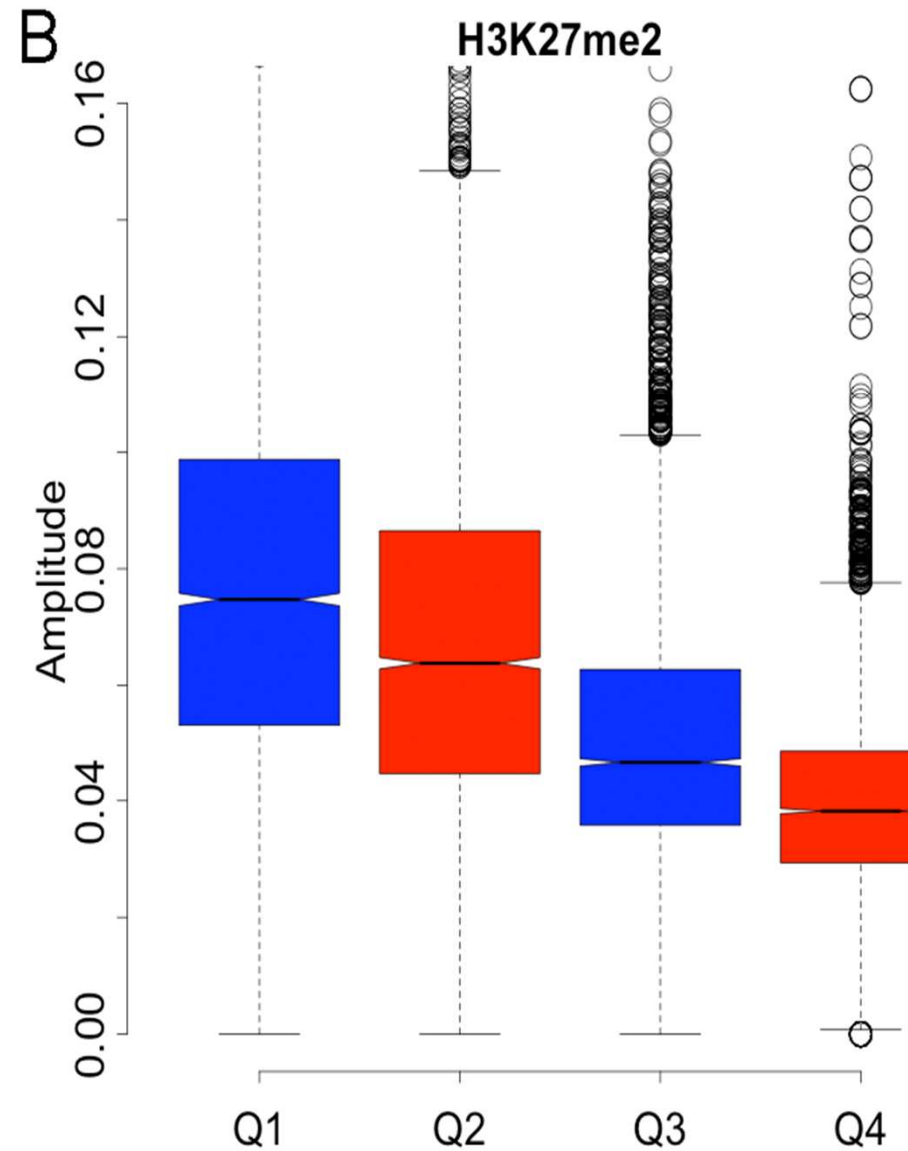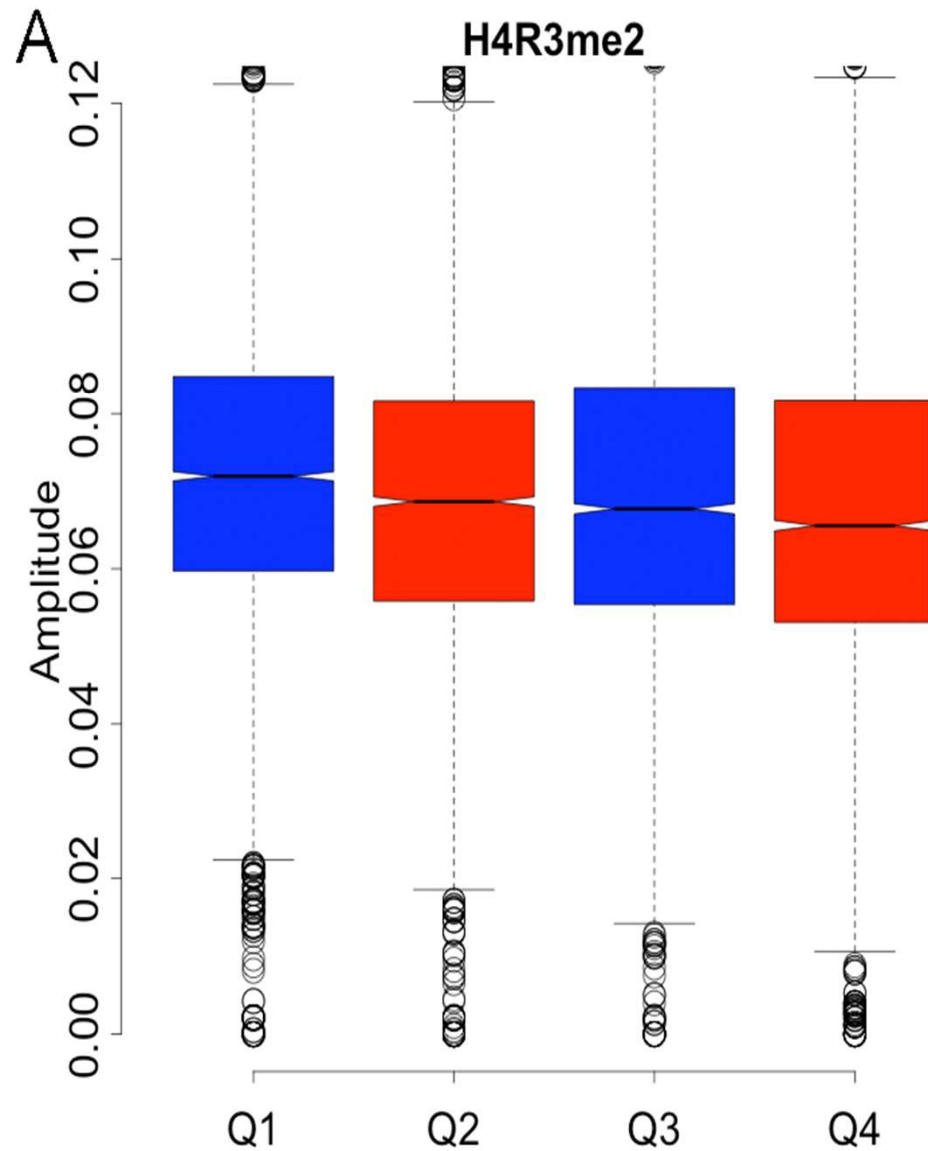
# PRMT5-mediated methylation of histone H4R3 recruits DNMT3A, coupling histone and DNA methylation in gene silencing

Quan Zhao[1,2,8], Gerhard Rank[1,8], Yuen T Tan[1], Haitao Li[3], Robert L Moritz[4], Richard J Simpson[4], Loretta Cerruti[1], David J Curtis[1], Dinshaw J Patel[3], C David Allis[5], John M Cunningham[6] & Stephen M Jane[1,7]
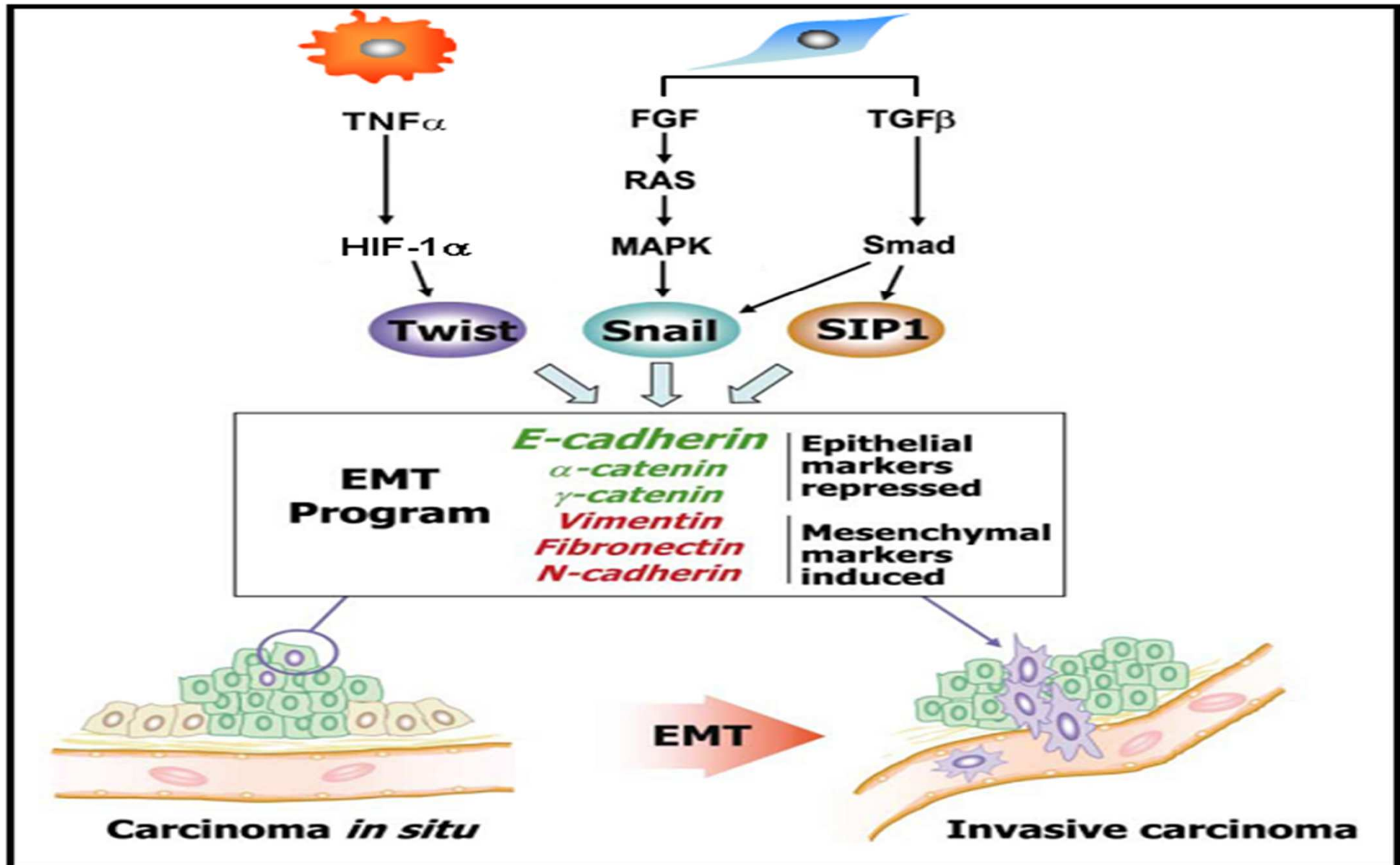
Mammalian gene silencing is established through methylation of histones and DNA, although the order in which these modifications occur remains contentious. Using the human β-globin locus as a model, we demonstrate that symmetric methylation of histone H4 arginine 3 (H4R3me2s) by the protein arginine methyltransferase PRMT5 is required for subsequent DNA methylation. H4R3me2s serves as a direct binding target for the DNA methyltransferase DNMT3A, which interacts through the ADD domain containing the PHD motif. Loss of the H4R3me2s mark through short hairpin RNA–mediated knockdown of PRMT5 leads to reduced DNMT3A binding, loss of DNA methylation and gene activation. In primary erythroid progenitors from adult bone marrow, H4R3me2s marks the inactive methylated globin genes coincident with localization of PRMT5. Our findings define DNMT3A as both a reader and a writer of repressive epigenetic marks, thereby directly linking histone and DNA methylation in gene silencing.
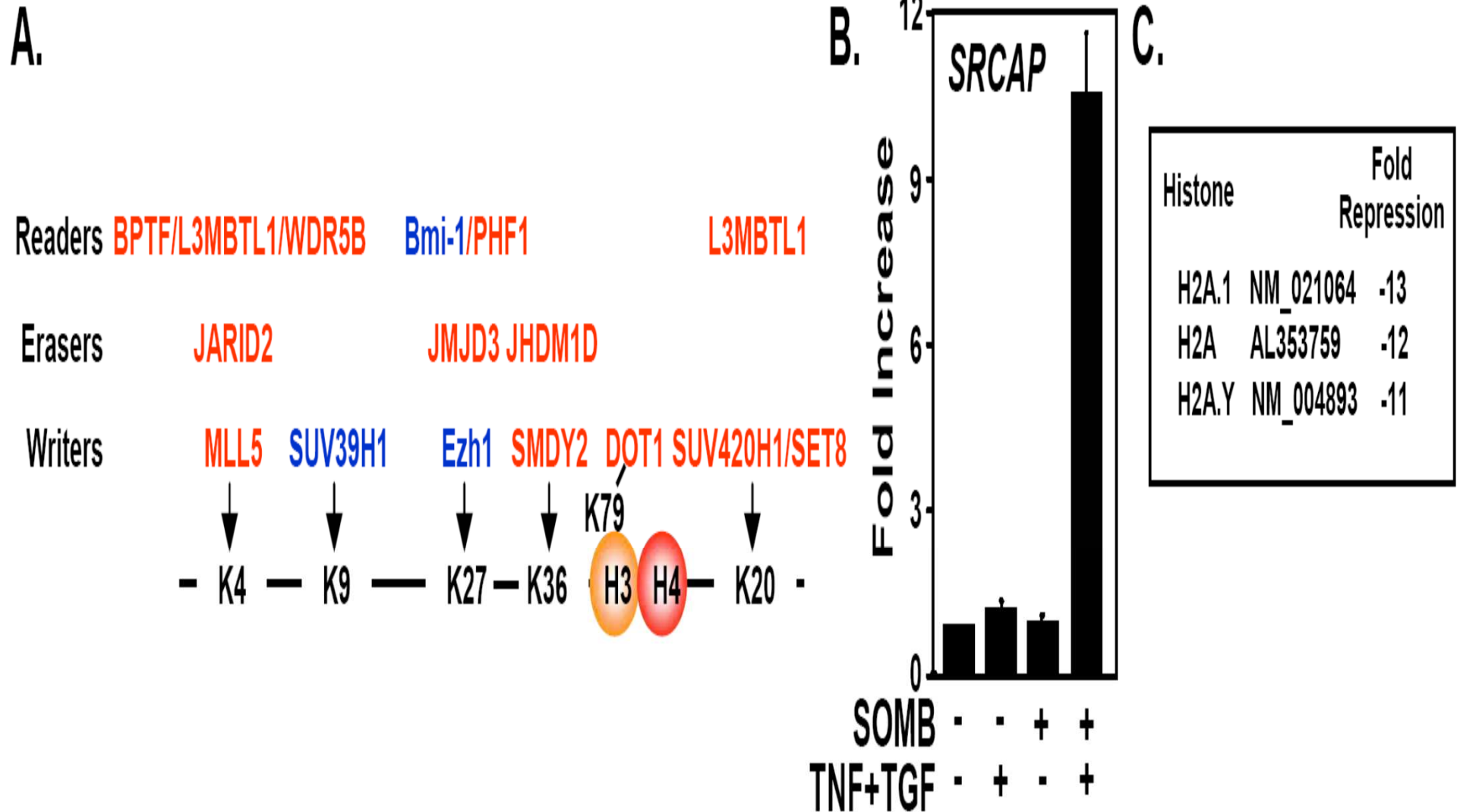
# Amplitude vs Gene Expression  Levels

# Epithelial to Mesenchymal Transition

# ~5,000 Genes Differentially Expressed during EMT including Multiple Readers, Writers and Erasers of Histone Modifications

# Summary

- Multilinear, MARS and Regression Tree Models recapitulate known/univariate trends
  - H3K36me3, H3K79me3, H3K4me3, H4K20me1 activating
  - H3K27me2,3 repressive
- They predict novel synergies:
  - H3K36me3-H3K79me1,3-H4K20me1
- They reveal activating/repressive potential of marks not observed in univariate analysis
  - H3K79me1 activating
  - H4R3me2 repressive
- H4R3me2s (PRMT5 writer) shown to be globally repressive from a number of experimental studies
  - Required for DNMT3A mediated DNA methylation and gene silencing
  - Knockdown of PRMT5 yield more de-repressed than repressed genes

# Acknowledgements

## Bekiranov Lab

Xiaojiang Xu
Stephen Hoang
Kunal Poorey
Veena Valsakumar

## Collaborators

Marty Mayo
Manish Kumar
Natalya N. Baranova
Mayo Lab
Patrick Grant
Mitch Smith
David Auble
Jeff Smith

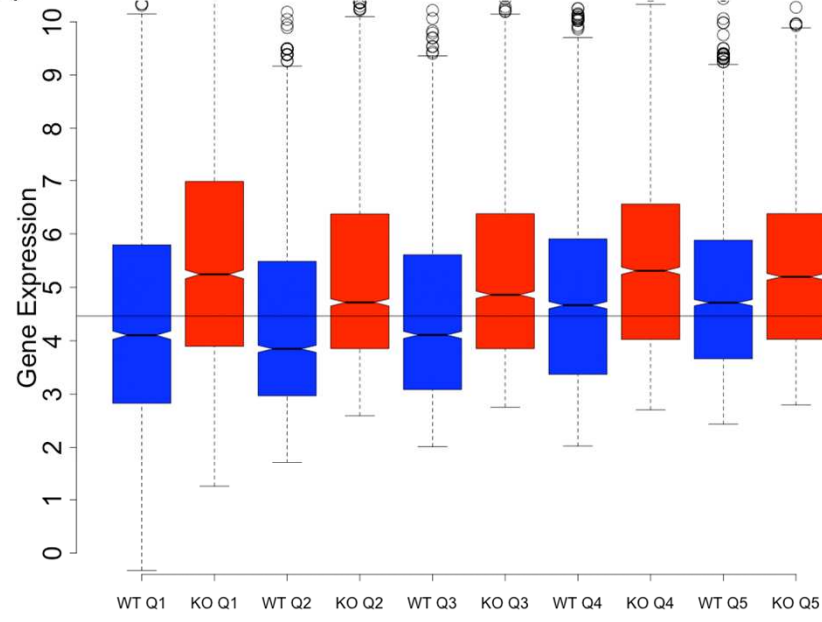## Computational Biologists

Bill Pearson
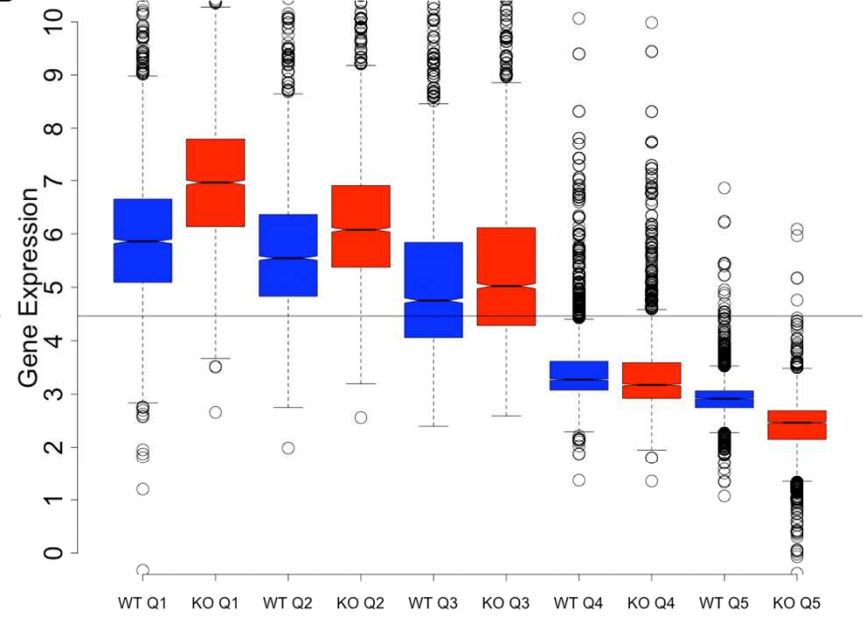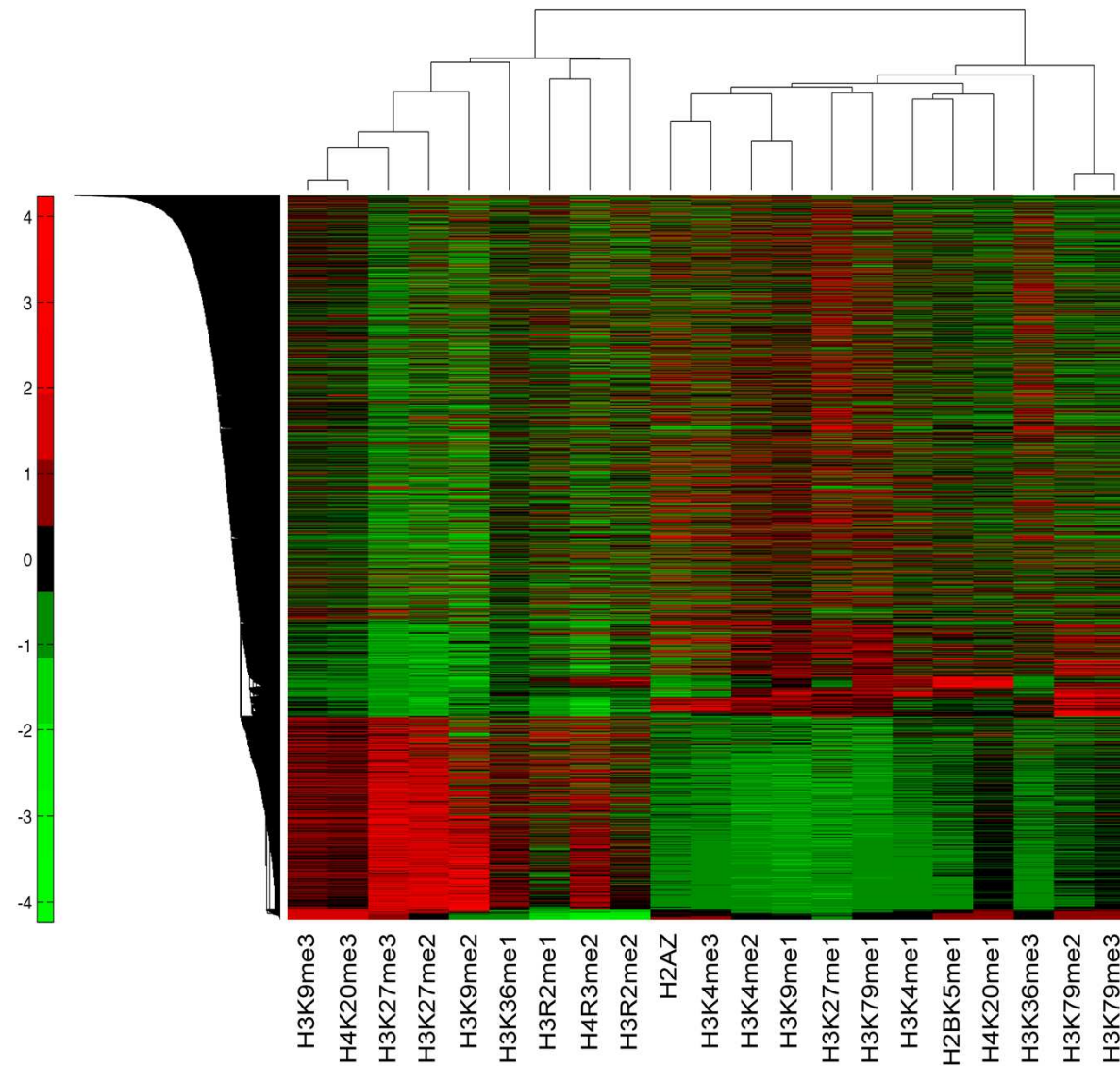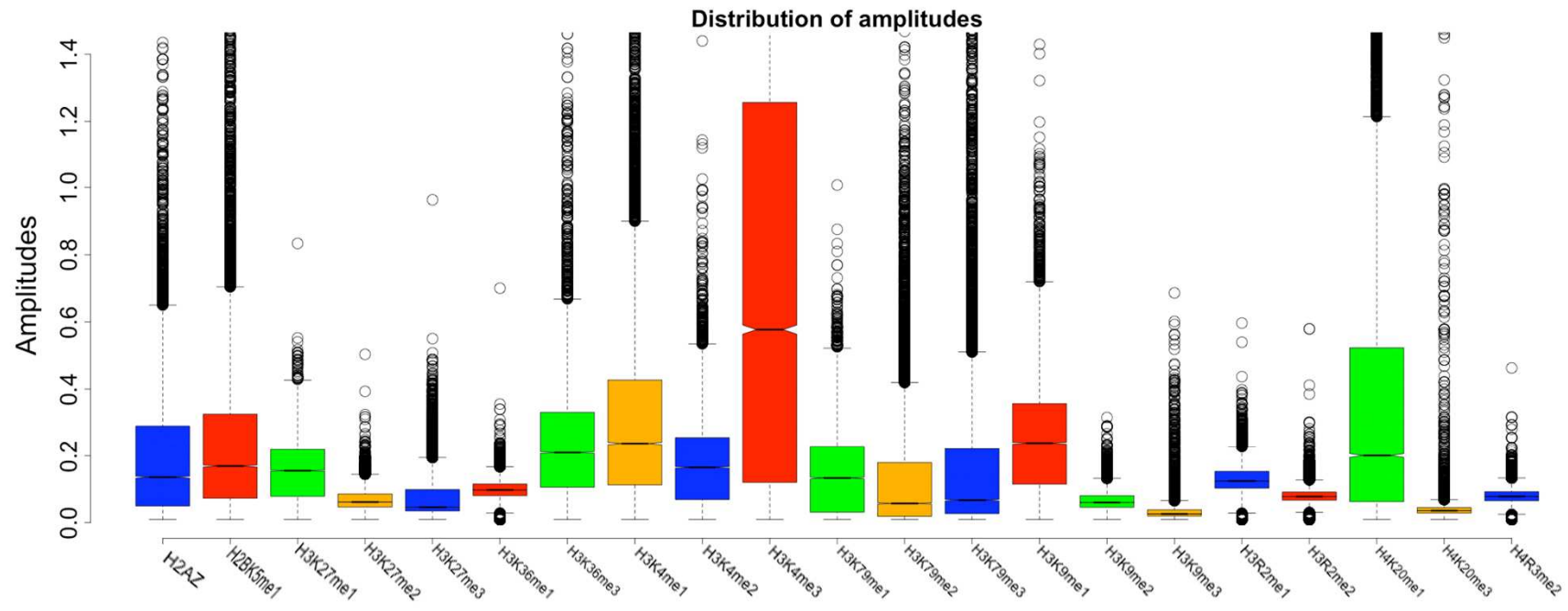Aaron Quinlan
Royden Clark

## Consultation

Sepideh Khorasanizadeh

A  **H4R3me2 KO**

D  **H3K27me2 KO**

Cluster Amplitudes: Genes x Marks

Distribution of amplitudes

# Calculating Enrichment Amplitudes

- Fixed Enrichment Profile Model
  - Average enrichment profile is characteristic of mark/gene boundary
  - Amplitude, not shape, varies across genes
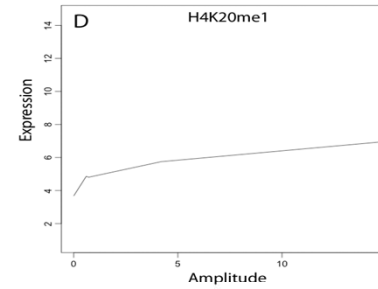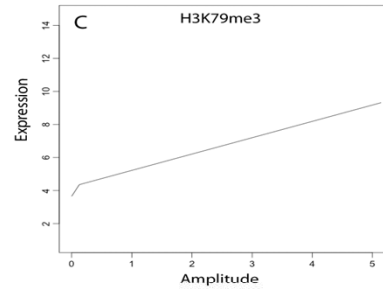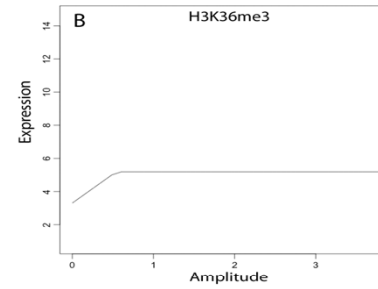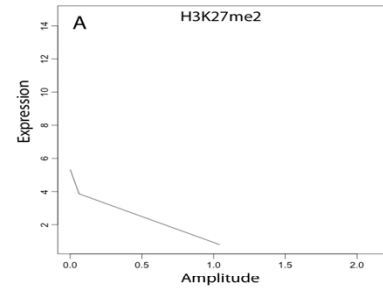  - Calculate amplitude from shape weighted average of read depth
- Gene Dependant Enrichment Profile Model
  - Use Unsupervised Learning Methods (PCA, Clustering Analysis, SOM) to identify enrichment profiles classes
  - Calculate amplitude from shape weighted average of read depth for each profile class
- Annotation Based Estimate
  - Calculate Average read depth in various functional elements: promoters, 5' UTR, exons, introns, 3'UTR
  - Use least squares to weight/summarize vector into one amplitude
- Significantly Enriched Sites Approach
  - Identify significantly enriched sites (at 5% False Discovery Rate)
  - Estimate Amplitude using average, median, or max within site
  - Cases below cutoff: Impute positions of sites from relative position of significant sites within genes & estimate amplitude

# Stepwise Regression

1. Fit the initial model
2. Use F-test to calculate p-values of potential additional term from pool of candidate terms.
   - If min(p) < 0.05, add term.
   - Repeat until no terms with p < 0.05.
   - If min(p) > 0.05, go to step 3.
3. If any terms in model have p > 0.05, remove term with max p-value and go to step 2; otherwise, end.

# MARS Model Terms

# Synthetic Oligonucleotide Arrays (1991)