# DeepDiff: DEEP-learning for predicting DIFFerential gene expression from histone modifications

*A. Sekhon, R. Singh, and Y. Qi, Bioinformatics 34 (17)*

**UnConference 2018**

**[@_alesssia](#) - [@AIClubGenderMinority](#) - [@H2Oai](#)**
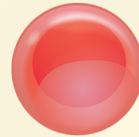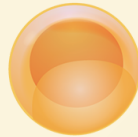
# What? Why?

**Totipotent embryonic stem cell**

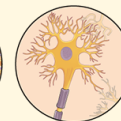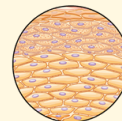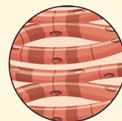**Pluripotent embryonic stem cells**
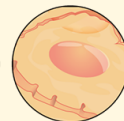
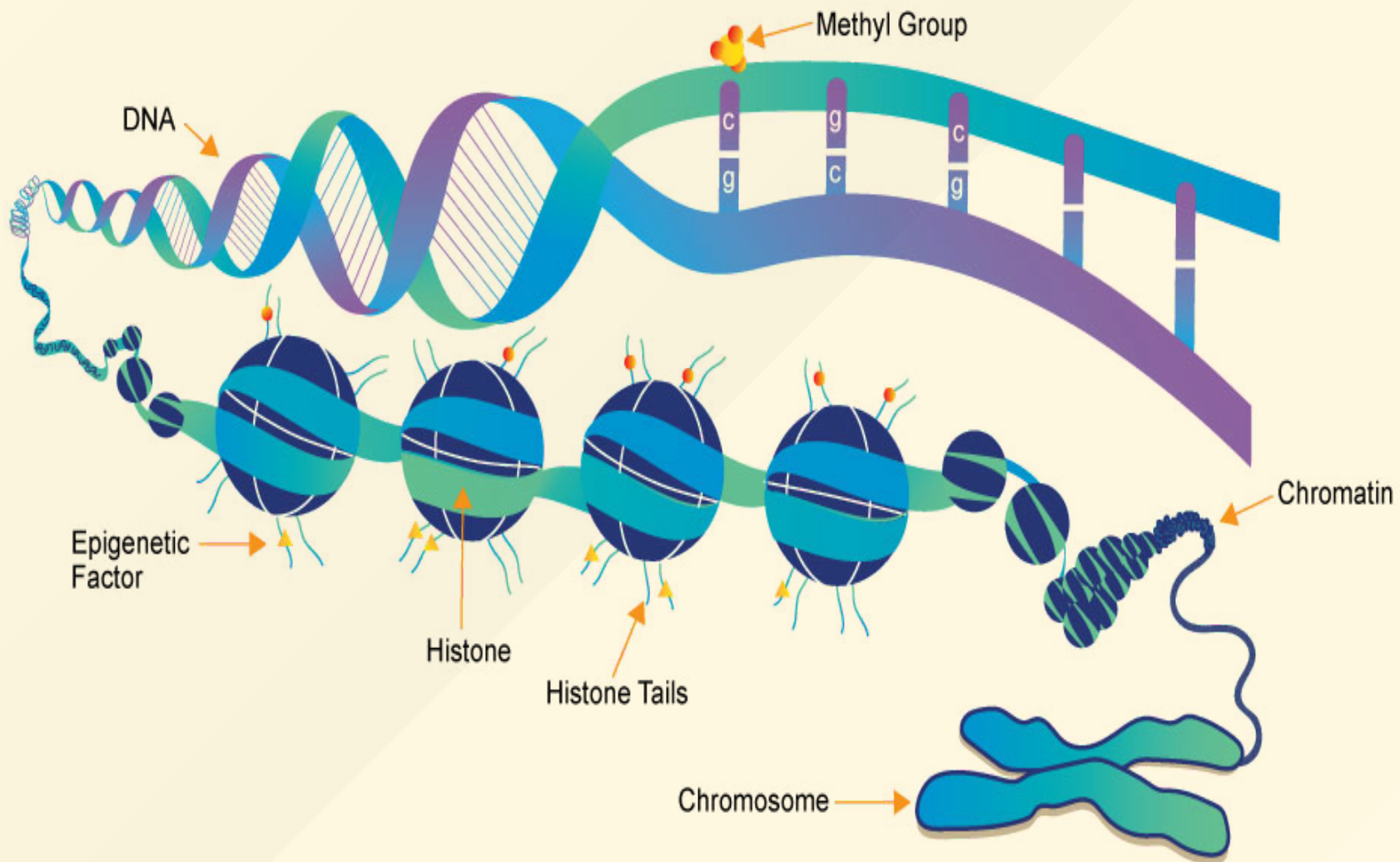Endoderm line      Mesoderm line      Ectoderm line

**Multipotent stem cells**

Lung     Pancreas     Heart muscle     Red blood cell     Skin     Neuron

Methyl Group

DNA

c
g

g
c

c
g

Epigenetic
Factor

Histone

Histone Tails

Chromatin

Chromosome

3' Antisense strand

RNA polymerase

5'

ATGACGGATCAGCCGCAAGCGGAATTGGCGACATAA
UACUGCCUAGUCGGCGUU

RNA Transcript

TACTGCCTAGTCGGCGTTCGCCTTAACCGCTGTATT

5'

Sense strand

3'

# Differential gene expression: cell differentiation and <u>diseases</u>

# Challenges

# 1. Genome-wide histone modification (HM) signals are spatially structured and may have long-range dependency
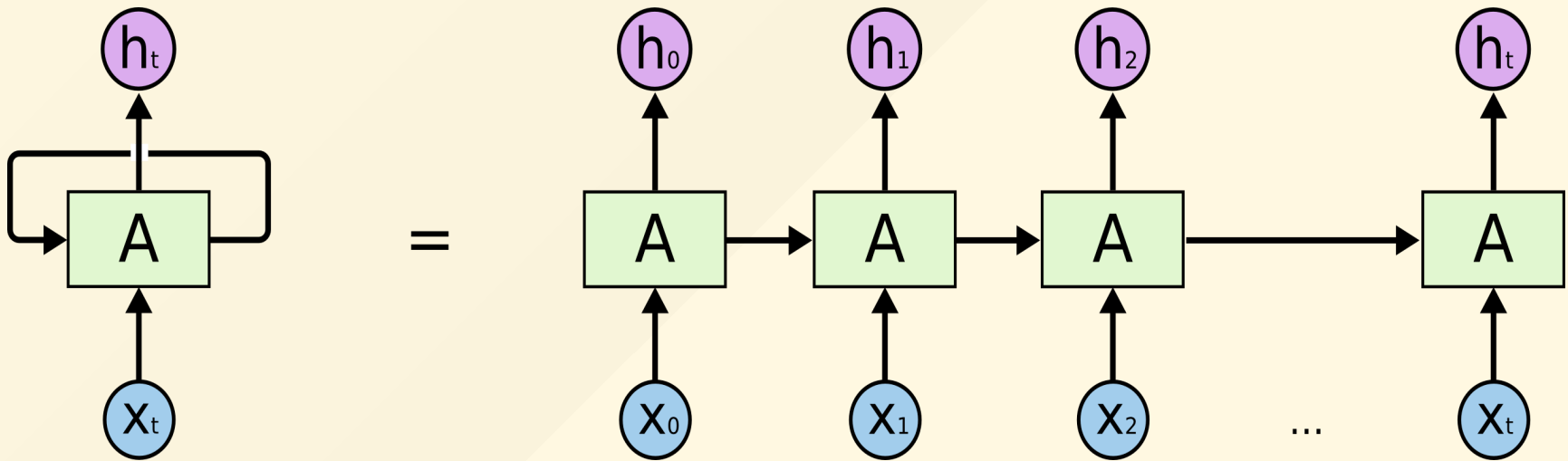
**2. We want to understand what the relevant HM factors are, and how they work together**

**3. We want to understand how HMs affect gene regulation, therefore we require a model with a degree of interpretability**
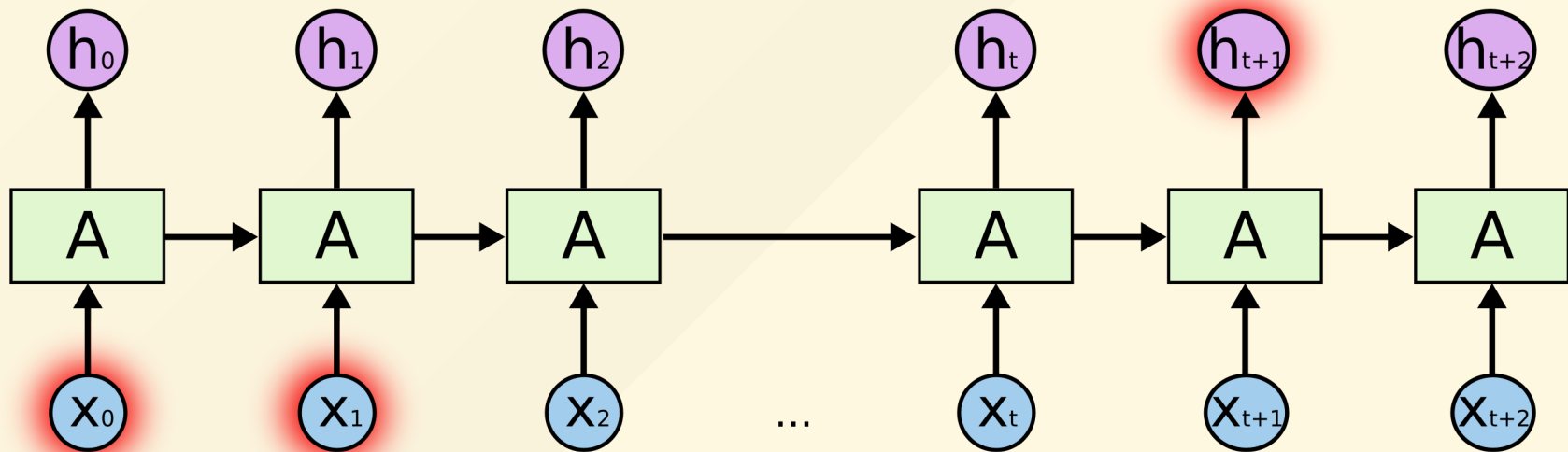
# 4. We are dealing with multiple cells
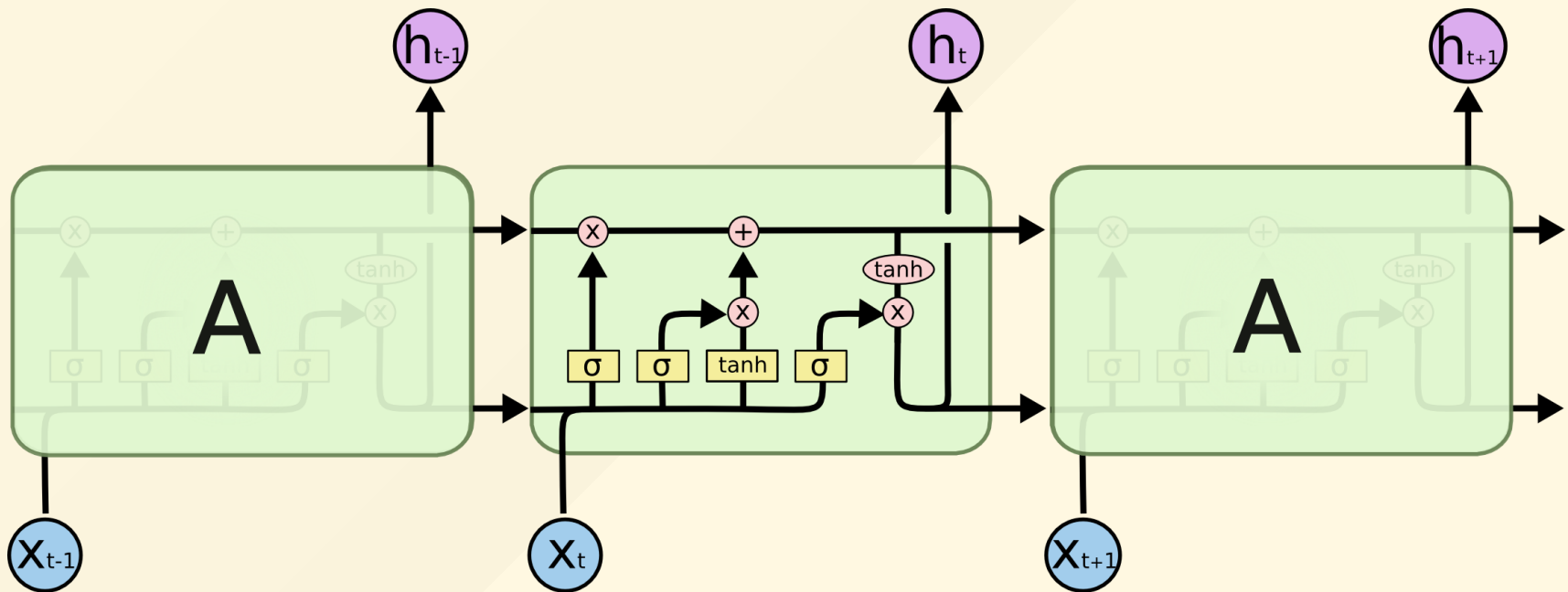
# How?

# Recurrent Neural Networks (RNNs)



Understanding LSTM Networks

# Recurrent Neural Networks (RNNs)

# Long Short Term Memory networks (LSTMs)

# Attention-based deep-learning methods



A woman is throwing a <u>frisbee</u> in a park.
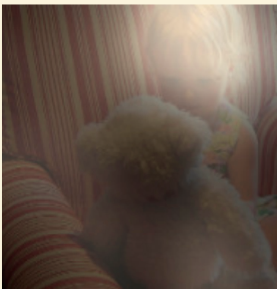
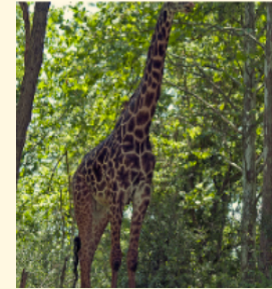A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

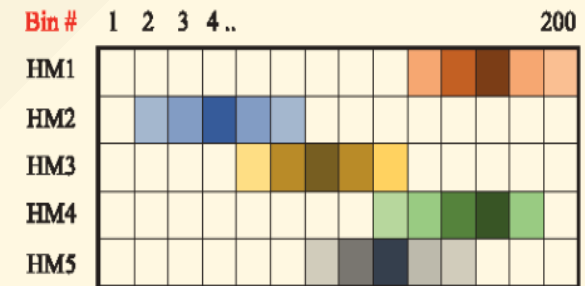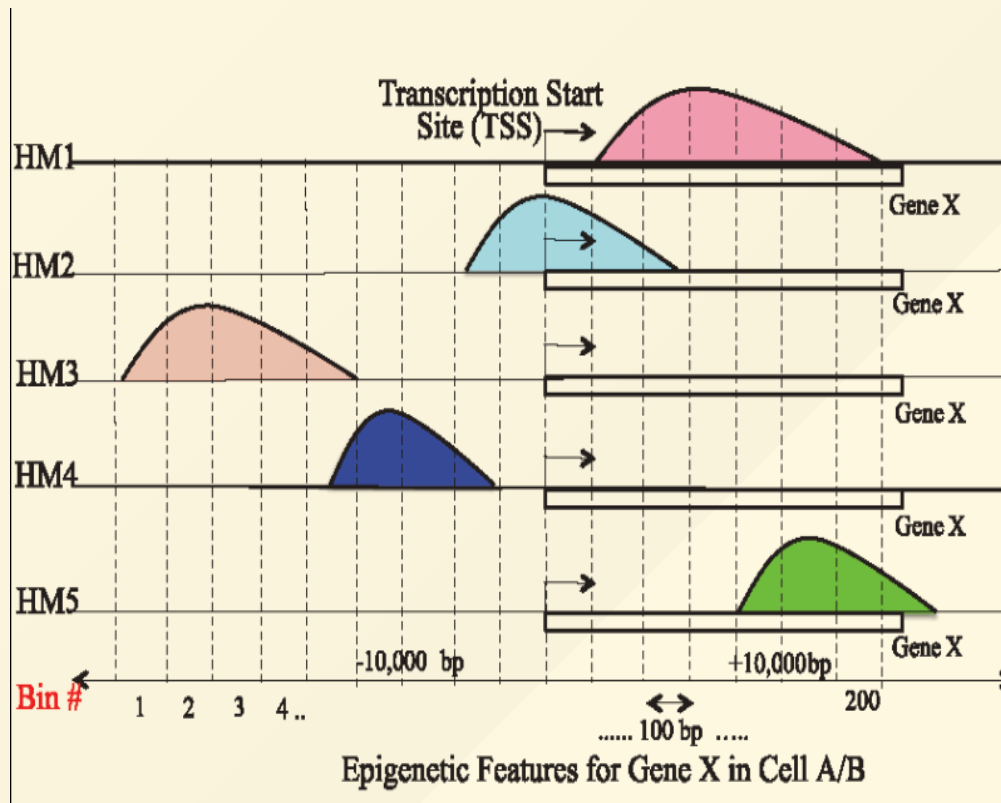A giraffe standing in a forest with <u>trees</u> in the background.

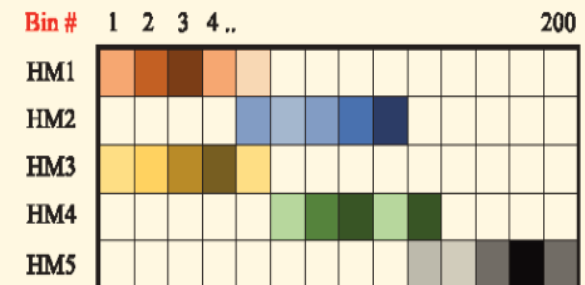[How to Automatically Generate Textual Descriptions for Photographs with Deep Learning](#)

# DeepDiff

# Input generation
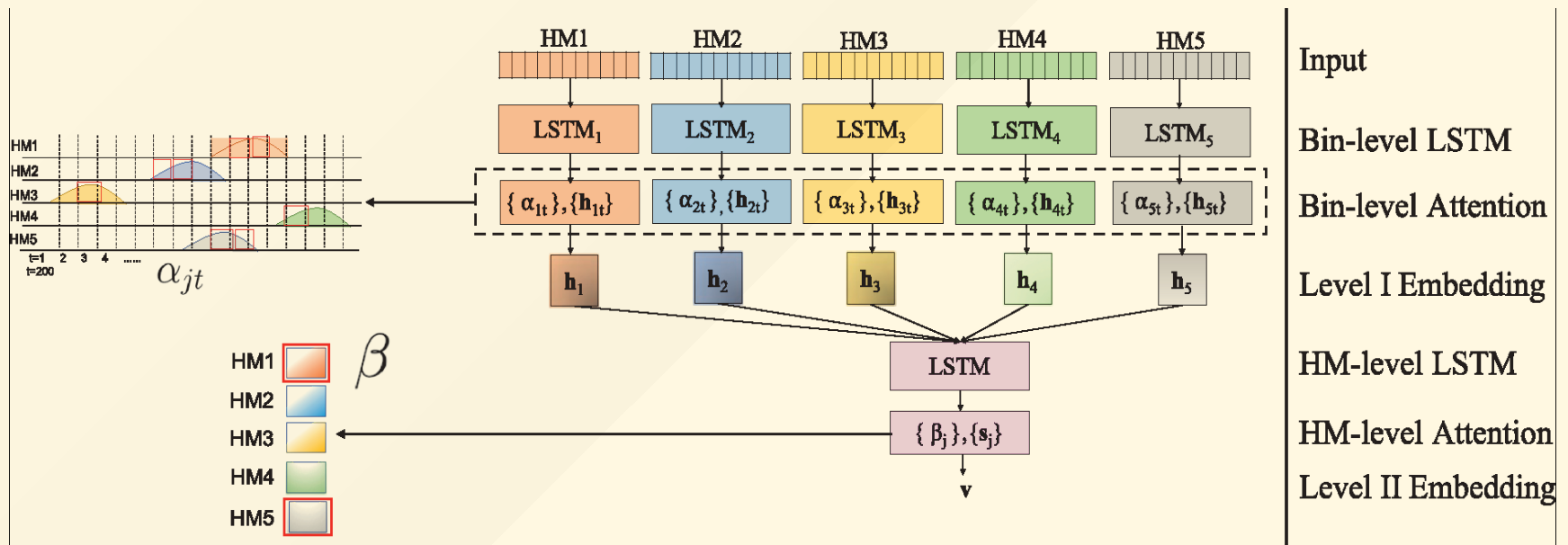


Epigenetic Features for Gene X in Cell A/B

Raw Input Matrix for Gene X in Cell-type A: $X^A$

Raw Input Matrix for Gene X in Cell-type B: $X^B$

# DeepDiff network architecture
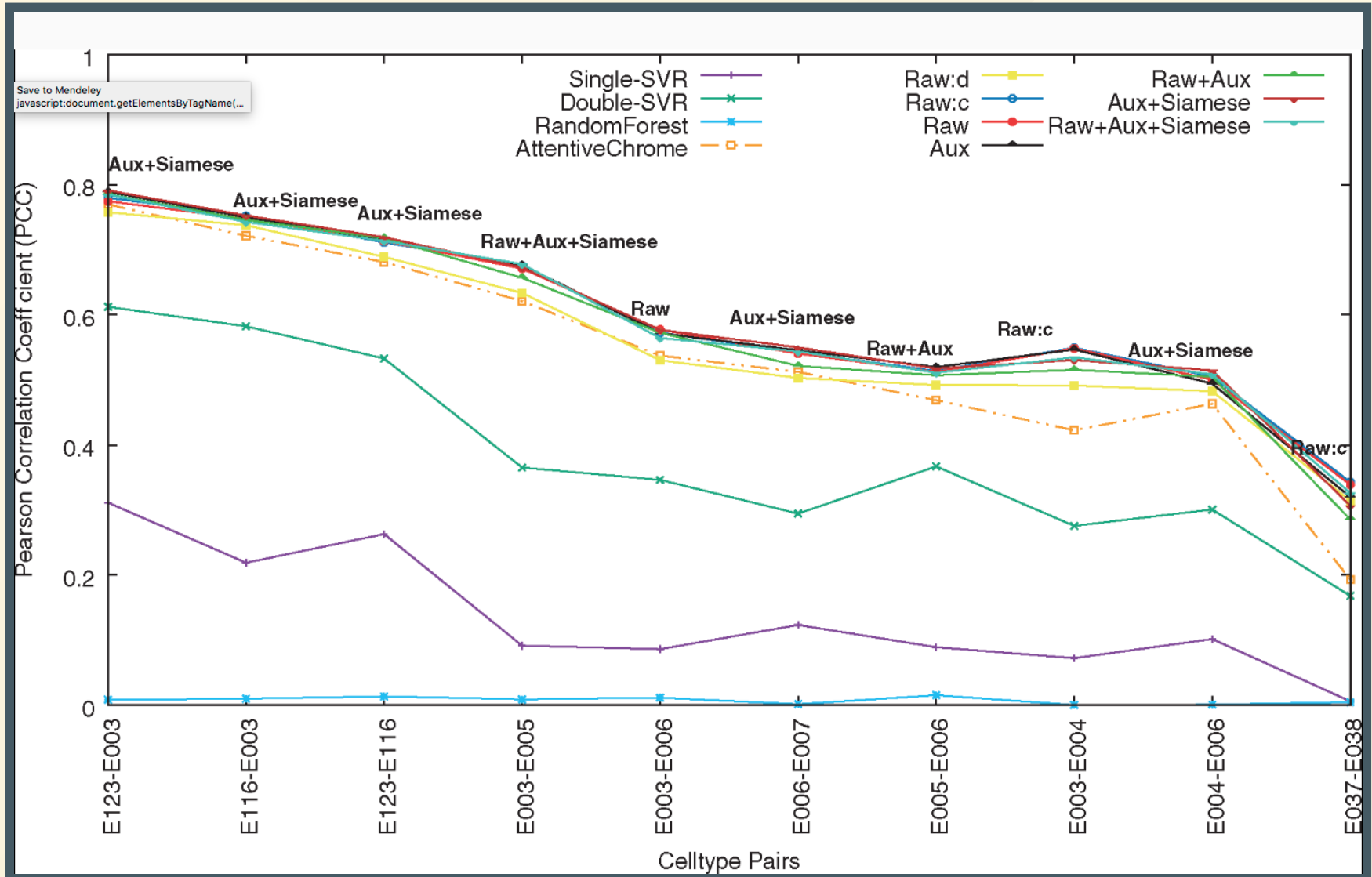
# Multitasking with auxiliary tasks

1. Cell-type specific prediction

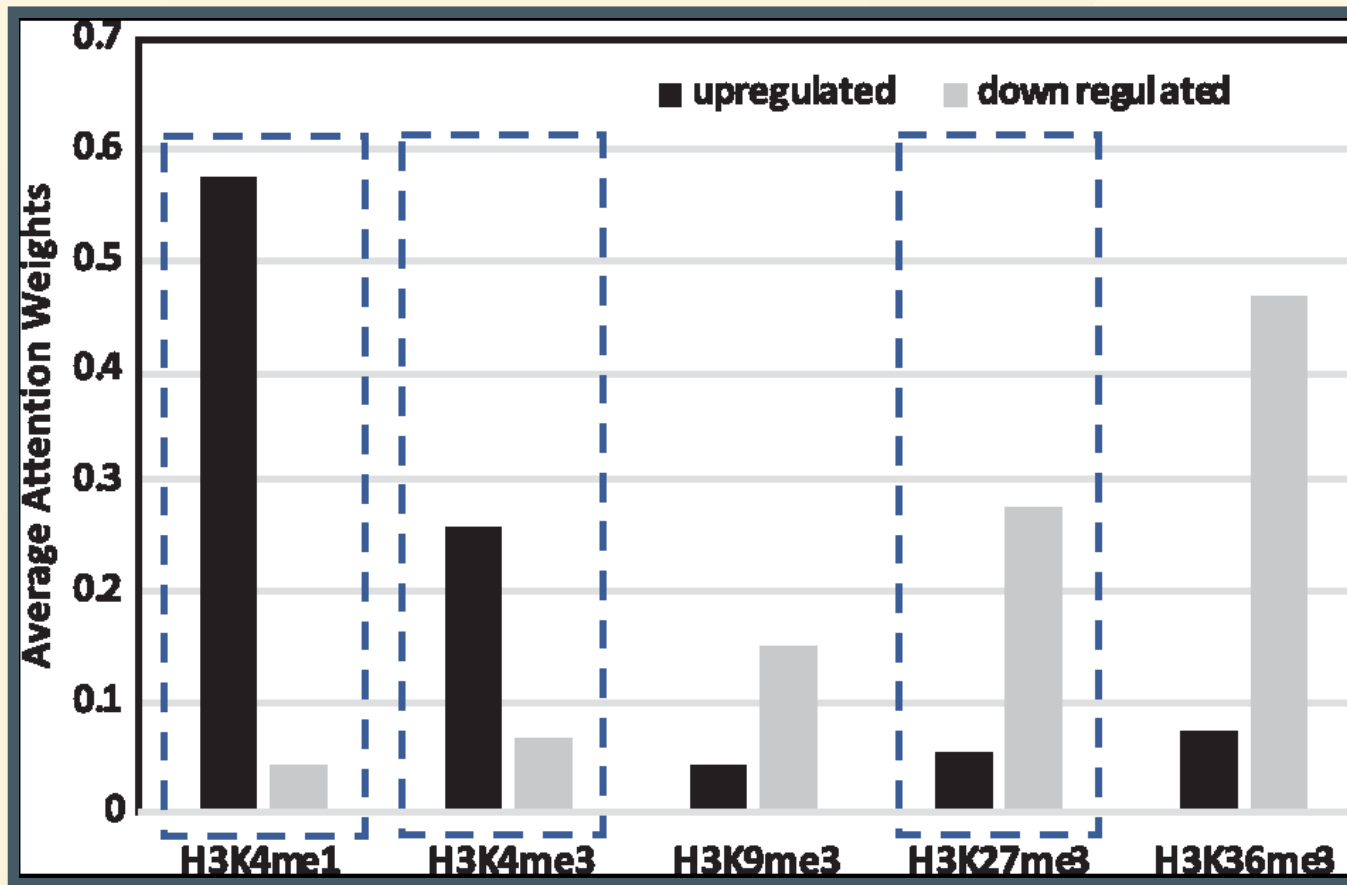2. Siamese architecture formulation

# Results

# Dataset

- 10 cell lines from the RoadMap project (18,460 genes)

  - 10,000 genes used for training

  - 2,360 genes used for validation (and tuning)

  - 1,100 genes used for testing

- Performances as PCC between DeepDiff variants and *baselines* (aka, alternative approaches)

# Performances

# Interpretation via Attention

# Discussion

# Points for discussion

1. Any thoughts on DeepDiff architecture?

2. Will the DeepDiff generalise with new cell lines?

3. Is this the best way to evaluate DeepDiff?

4. Is attention helpful?

Thanks for surviving until the end of the last session 😎

You deserve 🍕

H₂O.ai