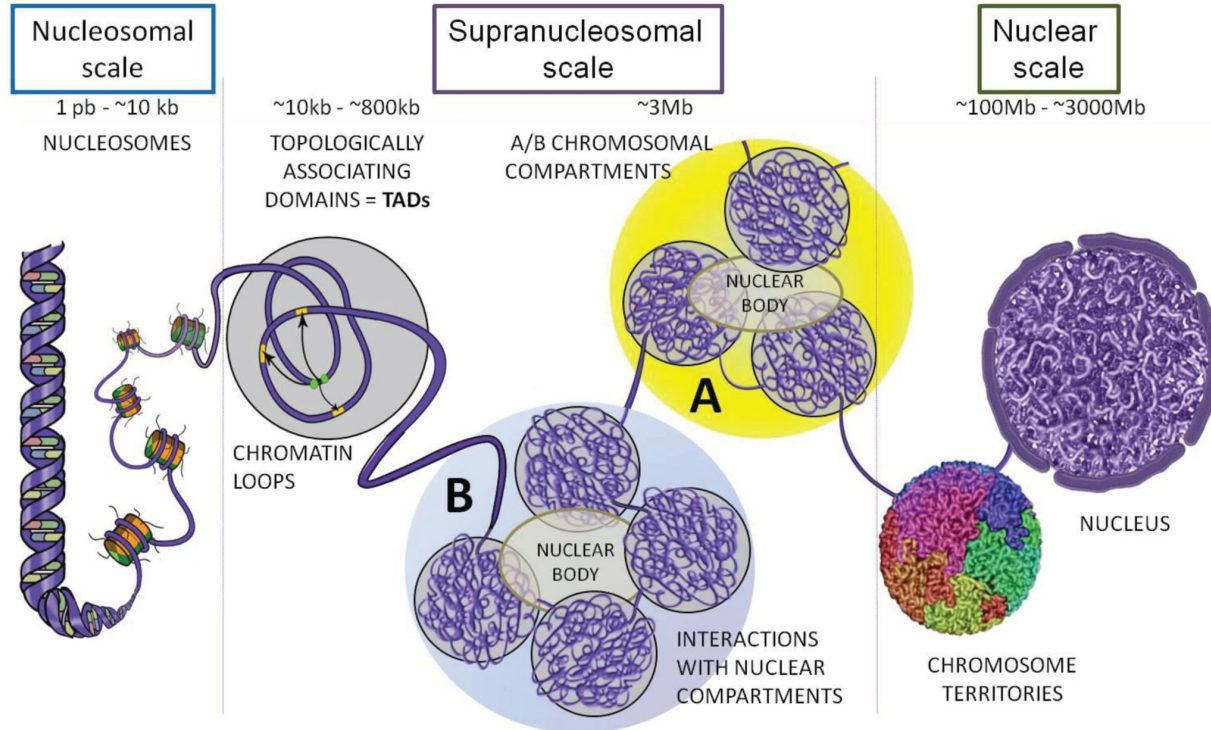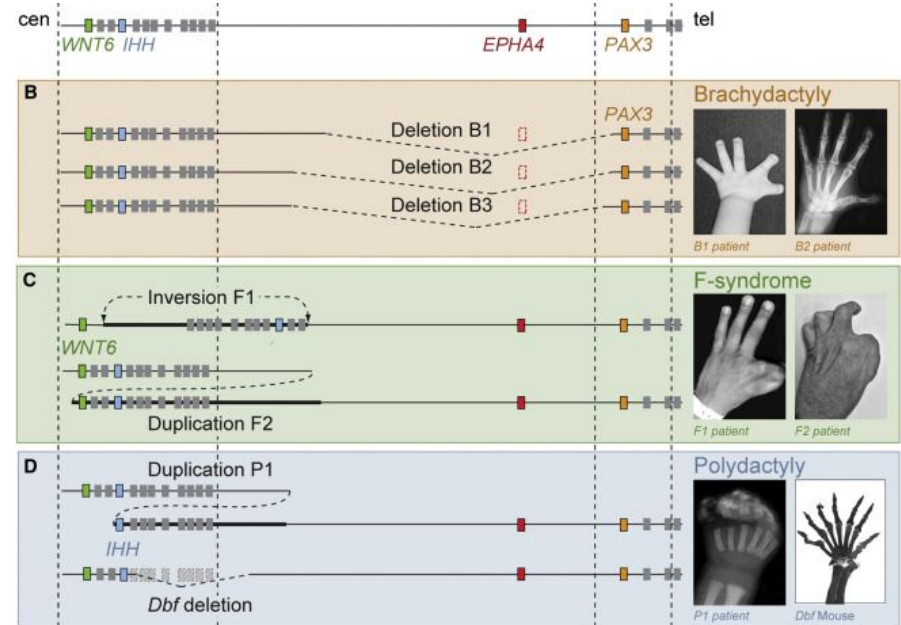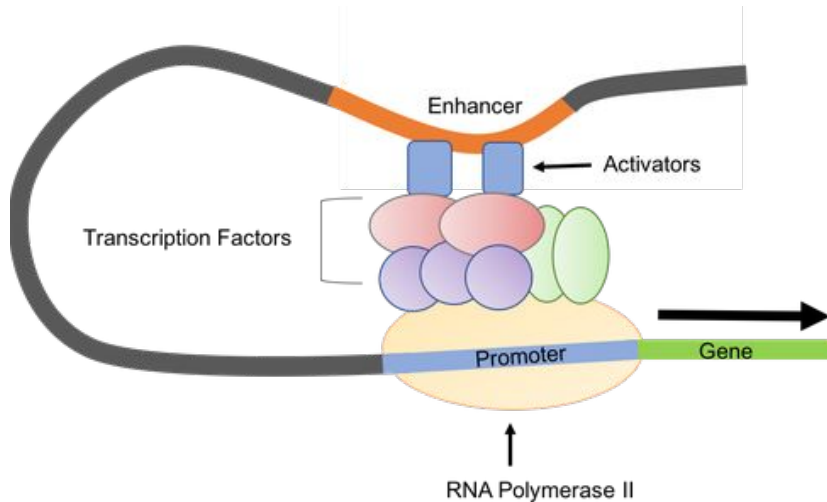# Analysis of Hi-C data using numerical linear algebra approaches

Kirill Ulianov
Marina Morozova
Maksim Grigoryan
Shamil Magomedov

# A few words about DNA…
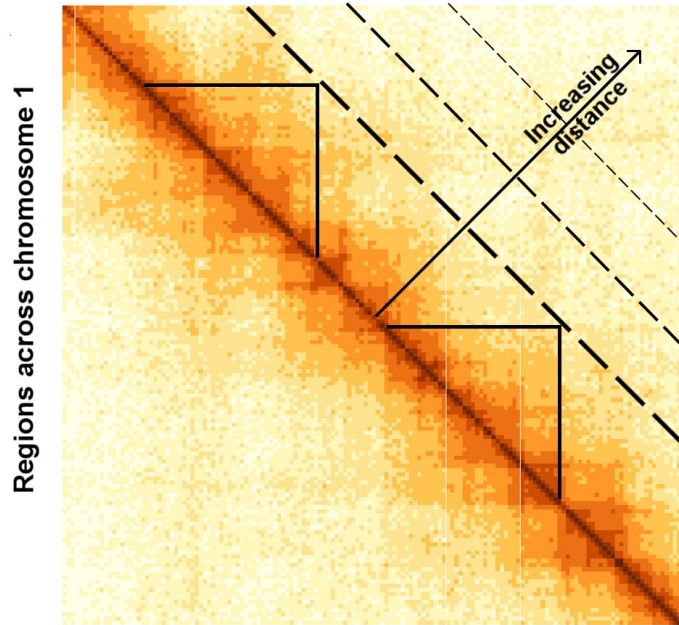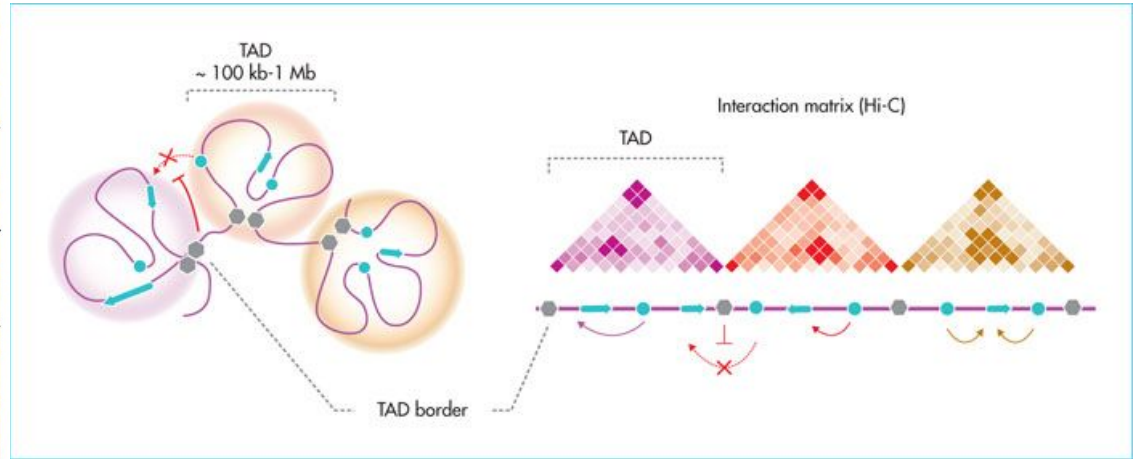
# Why does it matter?



3

# What is a Hi-C matrix?

# Example of a related task

One is common purpose in Hi-C analysis is detecting of the large conglomerates of active and repressed genome parts - compartments.

The conventional approach (Lieberman-Aiden, 2009) utilises only the first eigenvector of preprocess Hi-C matrix as the main predictor of DNA compartments distribution

Some methods use all eigenvectors to restore compartment, e.g. Gaussian Network Model (GNM).

Our **goal** was to use both approaches for Hi-C data and compare the results. Additionally, we were validated by gene activity patterns inferred from independent experiments

# Data

For the project we selected two Hi-C matrix from a public database derived for human and yeast cell cultures. From each map we chose only one chromosome for the analysis.

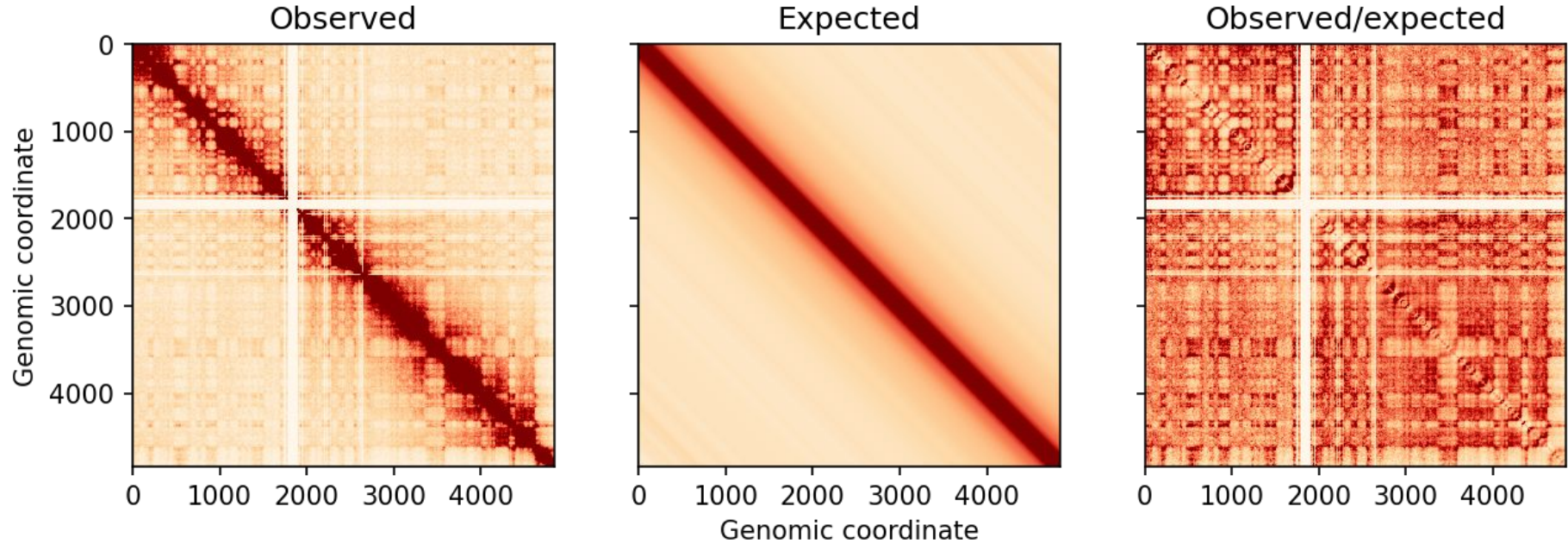The human 2nd chromosome was splitted on 4844 equally sized bins.

The yeast 5th chromosome was splitted on 722 equally sized bins.

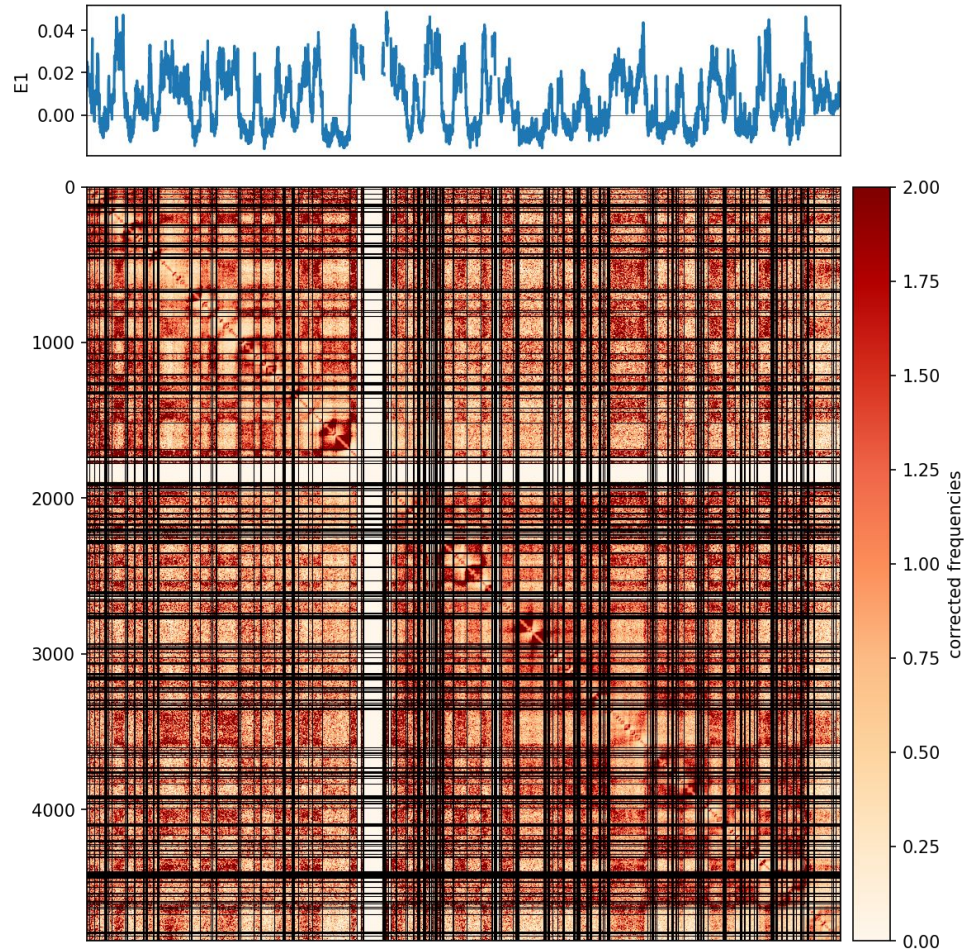# Eigendecomposition
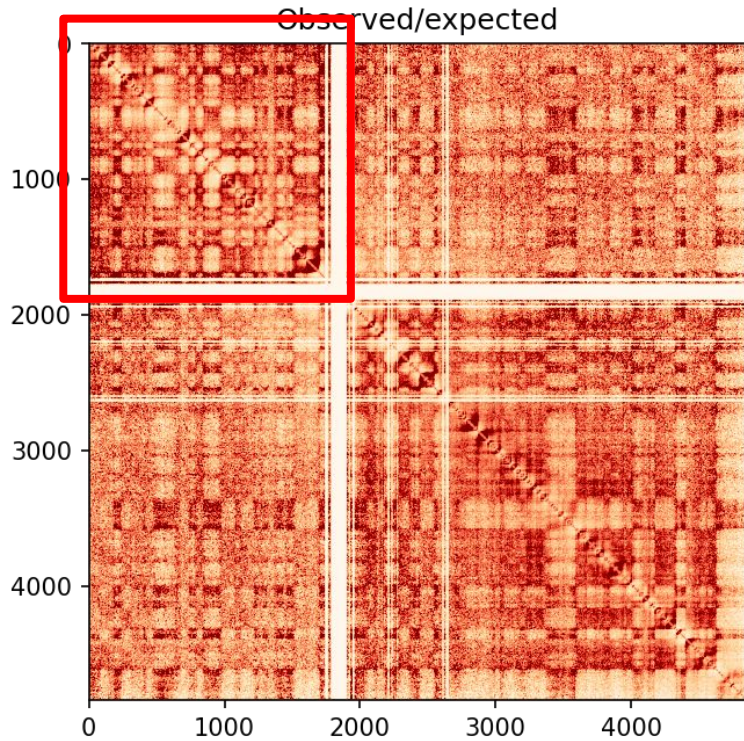
First eigenvector and properties of Hi-C matrix

1) Observed/Expected values
2) Eigendecomposition
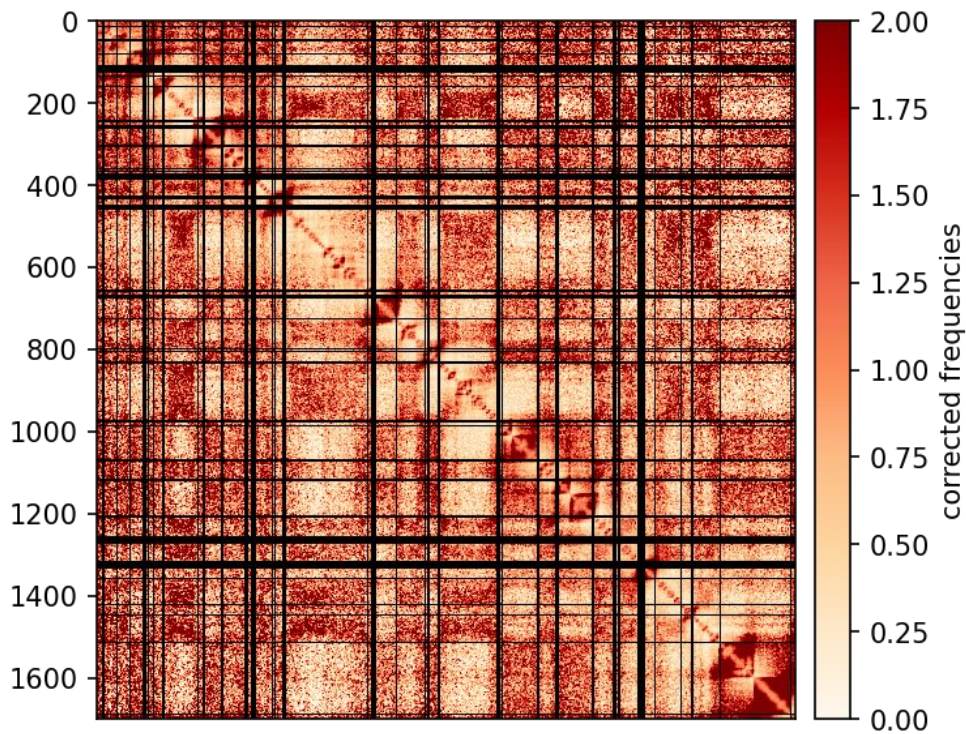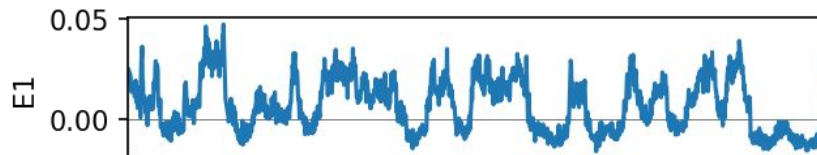
# Human 2nd chromosome
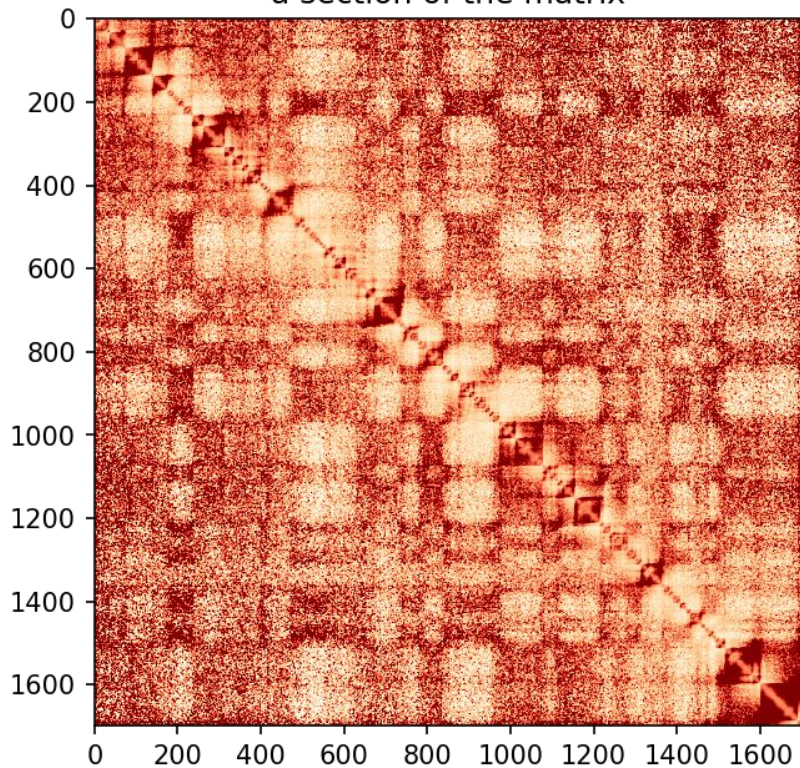


Observed     Expected     Observed/expected

Averaged diagonals of Observed
(Toeplitz matrix)

8

# First eigenvector (E1)



Observed/expected

Imakaev, M., Fudenberg, G., McCord, R. P., Naumova, N., Goloborodko, A., Lajoie, B. R., … Mirny, L. A. (2012). Iterative correction of Hi-C data reveals hallmarks of chromosome organization. Nature Methods, 9(10), 999–1003. doi:10.1038/nmeth.2148 (https://doi.org/10.1038/nmeth.2148)
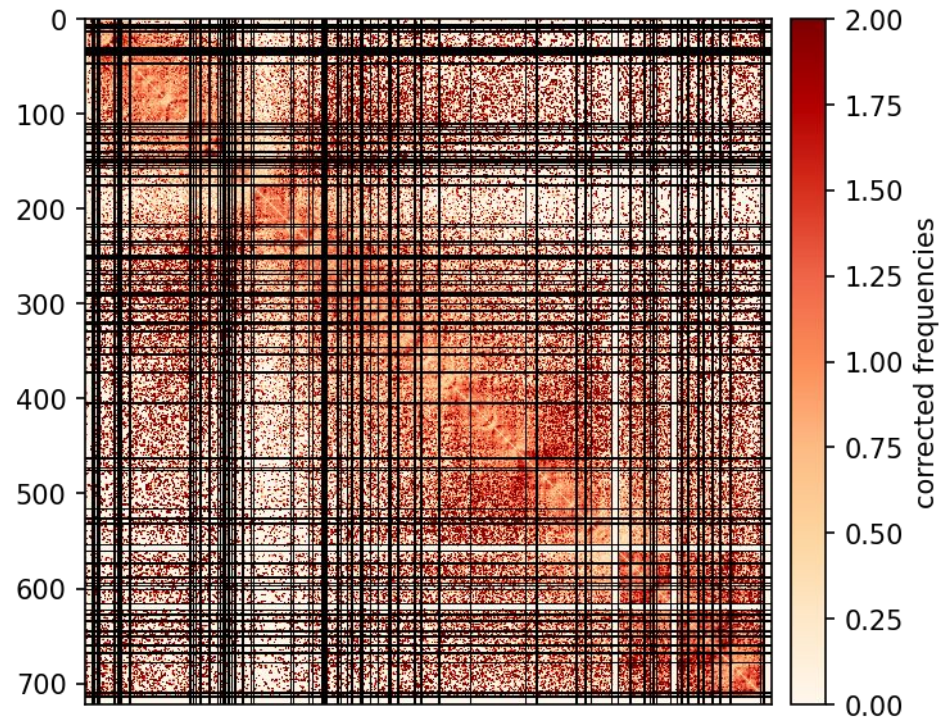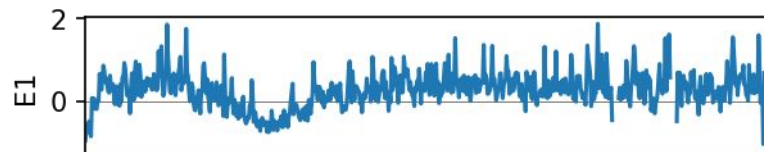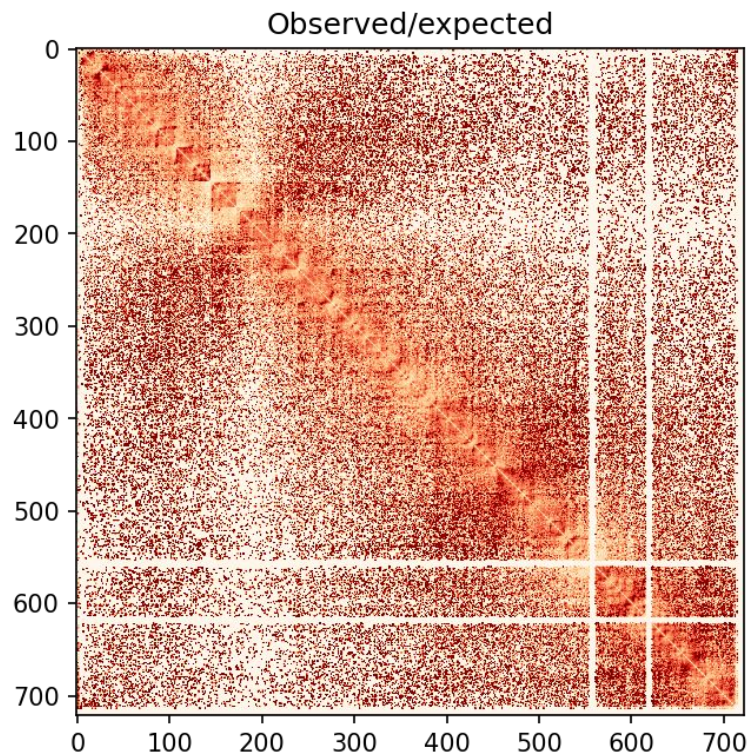
# Closer look

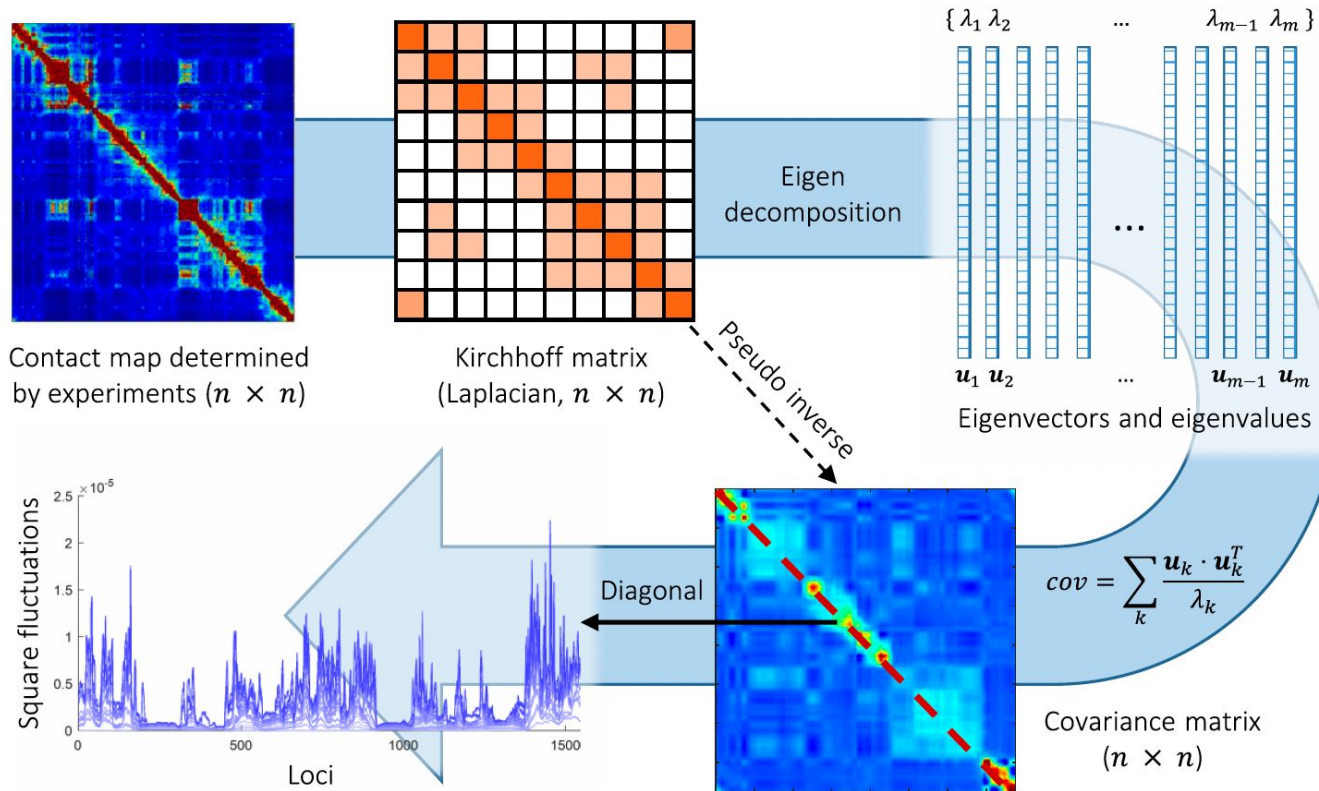Observed/expected;
a section of the matrix

# Yeast 5th chromosome

# Gaussian network model

Eigendecomposition
with extra steps
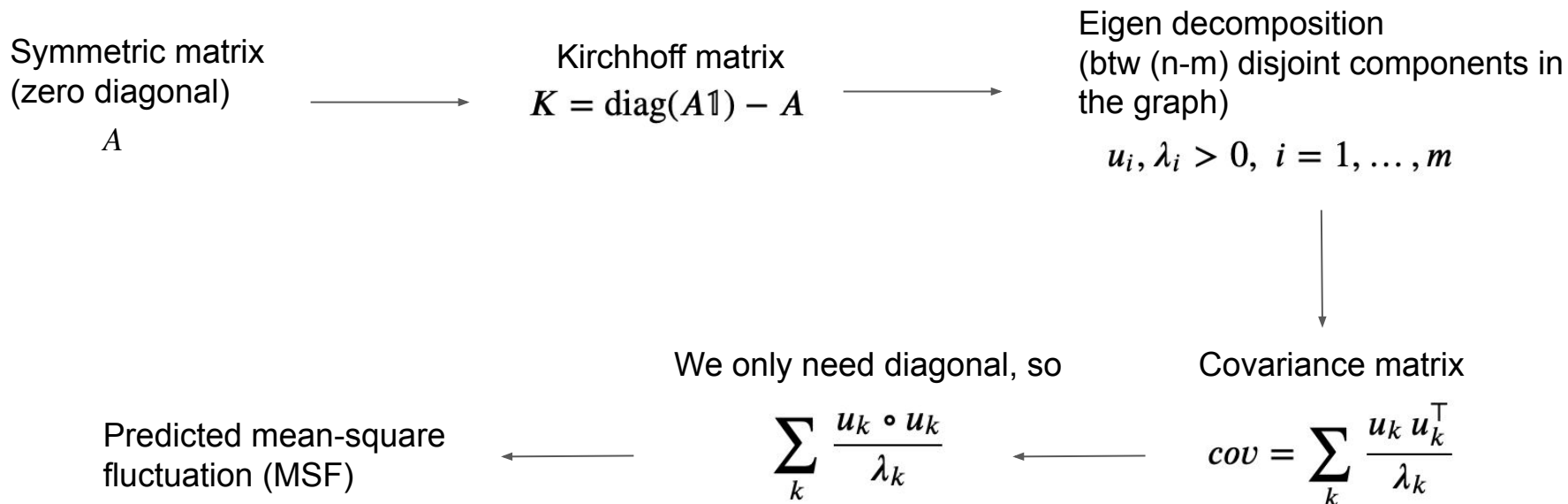
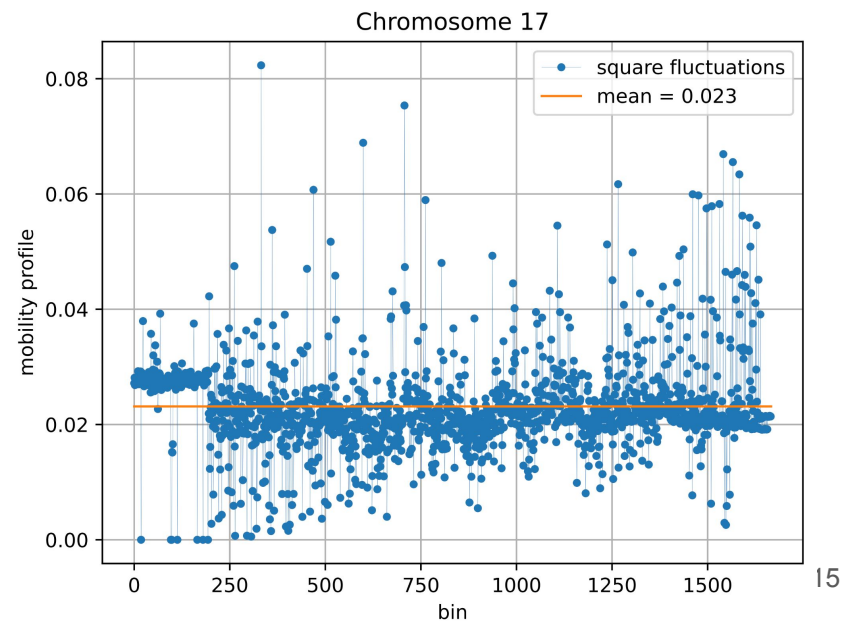1) Kirchhoff matrix
2) Covariance matrix
3) Mean-square fluctuation

# Gaussian network model
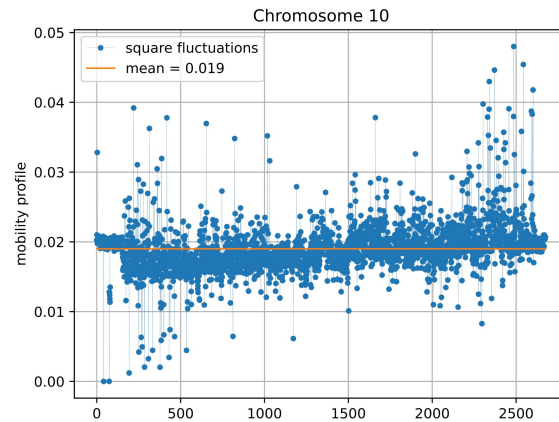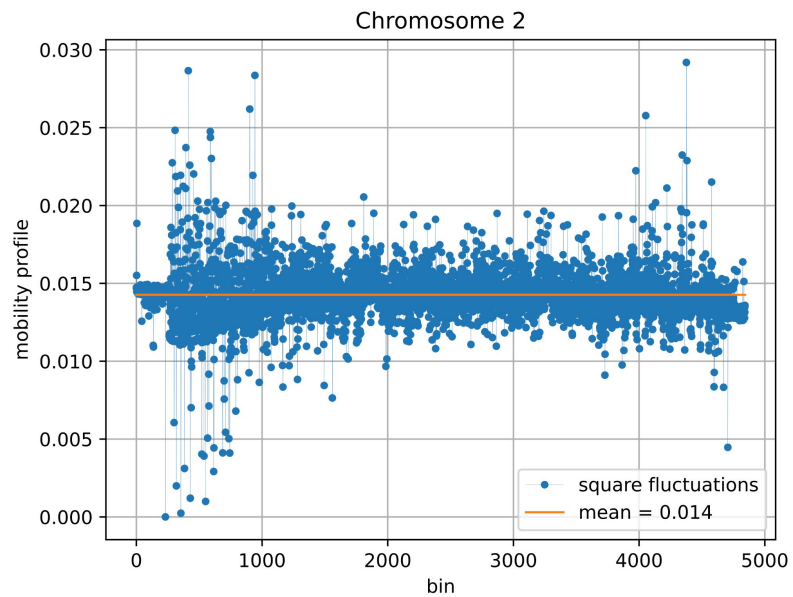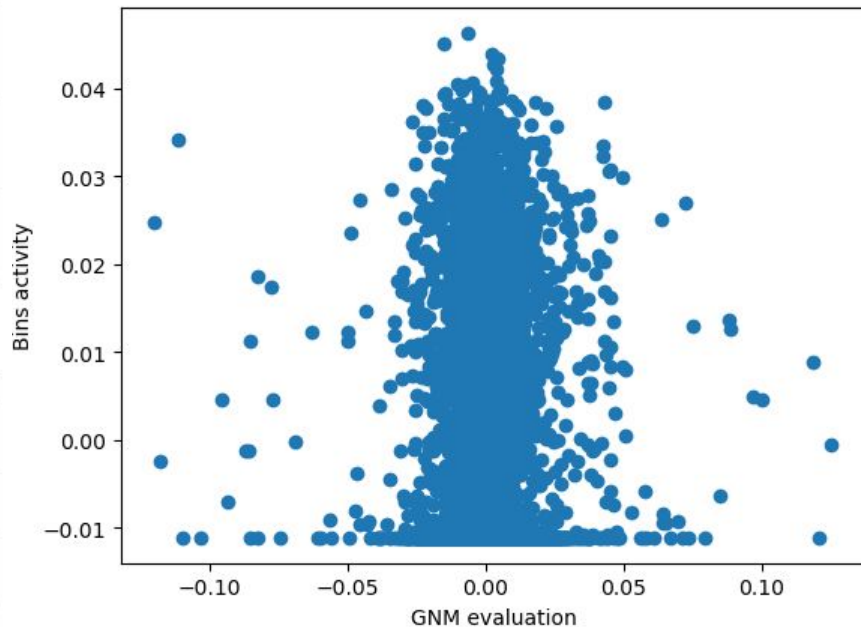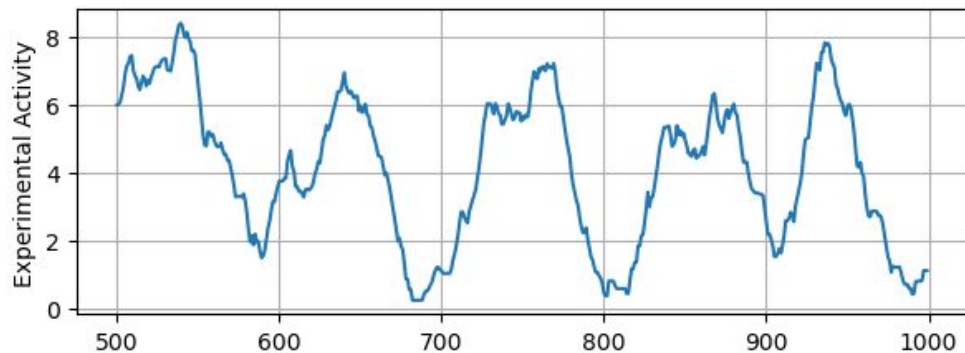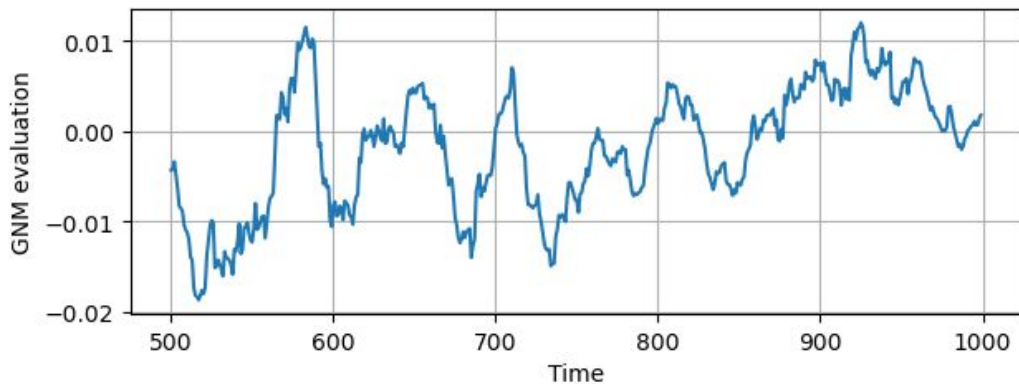


Chromosomal dynamics predicted by an elastic network model explains genome-wide accessibility and long-range couplings. Sauerwald, Natalie et al in Nucleic Acids Research (2017) pp. 3663-3673. doi: 10.1093/nar/gkx172 (https://doi.org/10.1093/nar/gkx172)

# Algorithm

Symmetric matrix
(zero diagonal)

$A$

$\longrightarrow$

Kirchhoff matrix

$K = \text{diag}(A\mathbb{1}) - A$

$\longrightarrow$

Eigen decomposition
(btw (n-m) disjoint components in
the graph)

$u_i, \lambda_i > 0, \ i = 1, \dots, m$

$\downarrow$

Predicted mean-square
fluctuation (MSF)

$\longleftarrow$

We only need diagonal, so

$$\sum_k \frac{u_k \circ u_k}{\lambda_k}$$

$\longleftarrow$

Covariance matrix

$$cov = \sum_k \frac{u_k \, u_k^\top}{\lambda_k}$$

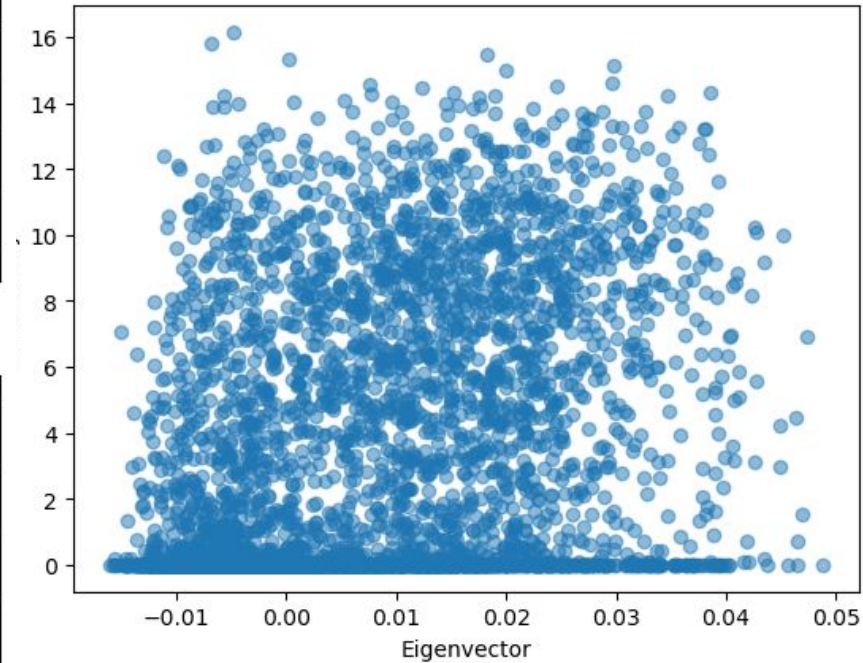Chromosome 2

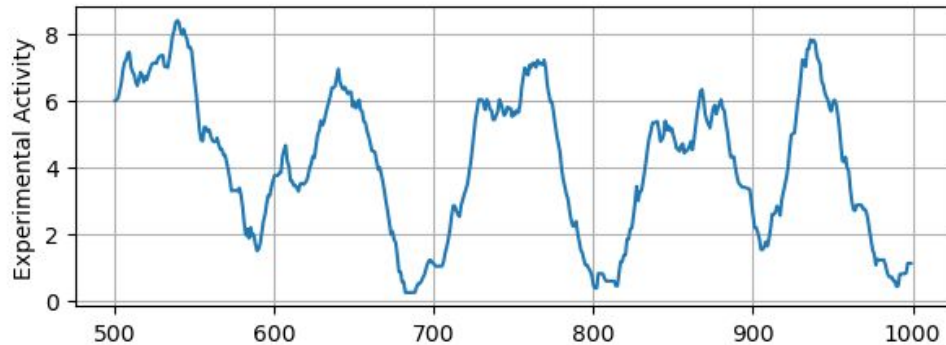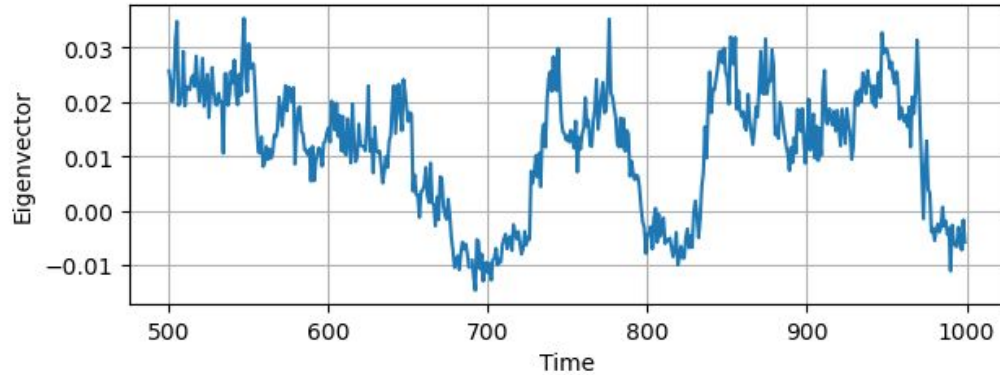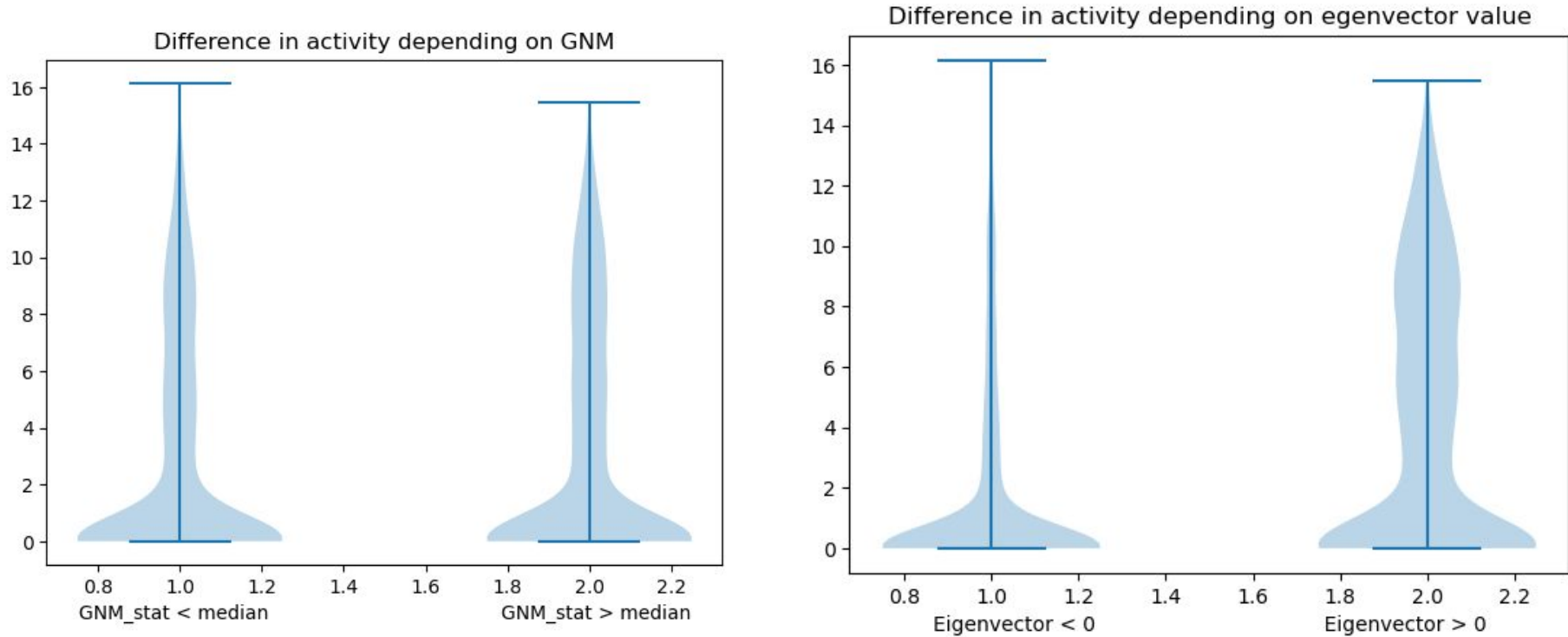Chromosome 10

Chromosome 17

# Gaussian Network Model Evaluation is hardly related with bin activity

# Eigenvector provides valuable information about gene activity

# Eigenvector sign is determined by gene activity



Difference in activity depending on GNM

Difference in activity depending on egenvector value

p< 10^-26

# Conclusions

The First Eigenvector is well-related with gene activity, that is complement with other Hi-C studies

Gaussian Network Model Evaluation shows no capacity to distinguish compartments

Bad performance could be caused by noise from other eigenvectors or some problems within the algorithm