

# Clustering reveals ubiquitous heterogeneity and asymmetry of genomic signals at functional elements.

Anshul Kundaje<sup>1</sup>, Sofia Kyriazopoulou-Panagiotopoulou<sup>1</sup>, Max W. Libbrecht<sup>1</sup>, Cheryl L. Smith<sup>2</sup>, Debasish Raha<sup>3</sup>, Elliott E. Winters<sup>4</sup>, Steven M. Johnson<sup>4</sup>, Michael Snyder<sup>5</sup>, Serafim Batzoglou<sup>2</sup>, Arend Sidow<sup>5</sup>

<sup>1</sup> Stanford CS, <sup>2</sup> Stanford Pathology, <sup>3</sup> Yale Molecular, Cellular, and Developmental Biology,

<sup>4</sup> Brigham Young University Department of Microbiology and Molecular Biology,

<sup>5</sup> Stanford Genetics

## ABSTRACT

The advent of high-throughput DNA sequencing has enabled the measurement of many types of chemical phenomena in the human genome [2]. These methods have enabled the location of both punctate phenomena, such as transcription factor binding sites (TFBSs), and broad phenomena, such as the location of chemical modification of the histone proteins that wrap DNA. Experiments measuring broad phenomena are typically processed into real-valued signal tracks defined across every base pair.

A popular and highly effective method for visualizing and quantifying the relationship between a punctate phenomenon and a genomic signal is the so-called aggregation plot (AP). In a typical AP, the signal around several predefined anchor sites (such as TFBSs) is averaged for each position within a window around the sites. If the signal behaves similarly across the anchor sites, the AP will reveal these common signal patterns.

For an AP to display signals that are asymmetric about the anchor, the alignment of features has to be robust, and some other data is utilized to provide the correct orientation. Features that can be aligned, but for which there exists no external information regarding their orientation (e.g. TFBSs), can produce APs with strong but obligatorily symmetric signals. We developed a hierarchical agglomerative clustering strategy for this data which may reverse the orientation of sites in order to merge clusters that are mirror images of one another. This analysis revealed that asymmetries of chromatin signals are a pervasive feature at TFBSs, not only within genes but, surprisingly, equally strongly at gene-distal sites [1].

## BODY

*Aggregation plots cannot represent asymmetrical signals around sets of un-oriented sites. Clustering is necessary for such analysis.*

## REFERENCES

- [1] A. Kundaje, S. Kyriazopoulou-Panagiotopoulou, M. Libbrecht, C. L. Smith, D. Raha, E. E. Winters, S. M. Johnson, M. Snyder, S. Batzoglou, and A. Sidow. Ubiquitous heterogeneity and asymmetry of the chromatin environment at regulatory elements. *Genome Research*, 22(9):1735–1747, 2012.
- [2] The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature*, 489:54–74, 2012.

*Volume 2 of Tiny Transactions on Computer Science*

This content is released under the Creative Commons Attribution-NonCommercial ShareAlike License. Permission to make digital or hard copies of all or part of this work is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. CC BY-NC-SA 3.0: <http://creativecommons.org/licenses/by-nc-sa/3.0/>.