

# A machine-learning guided approach to explore the cis-regulatory code involved in neuronal differentiation

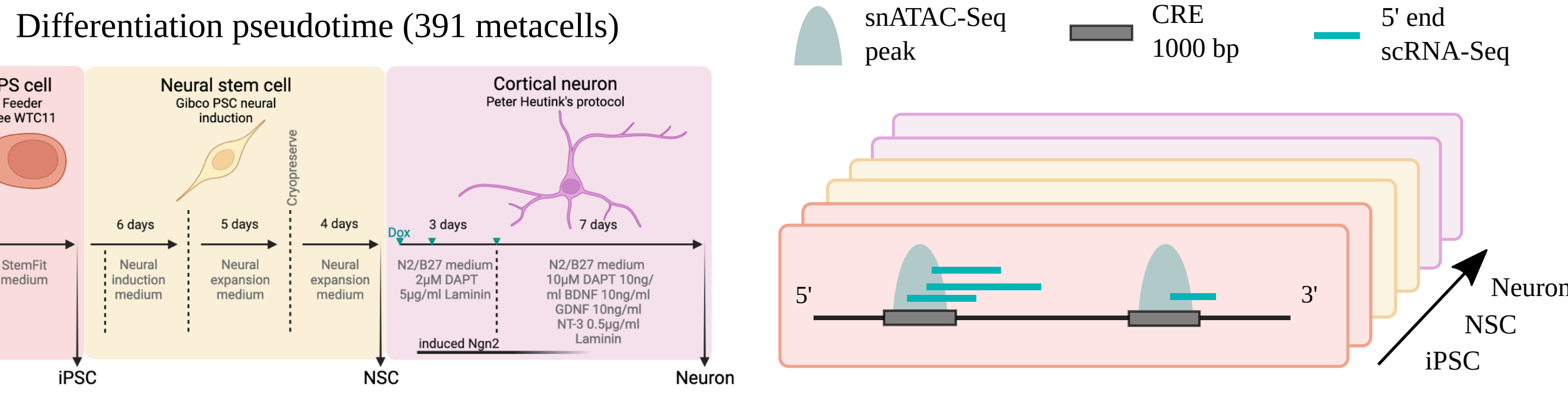
Océane Cassan<sup>1\*</sup>, Christophe Vroland<sup>1,2</sup>, Julien Raynal<sup>1,2</sup>, FANTOM consortium, Masaki Kato<sup>3</sup>, Hazuki Takahashi<sup>3</sup>, Takeya Kasukawa<sup>3</sup>, Piero Carninci<sup>3</sup>, Chi Wai Yip<sup>3\*</sup>, Laurent Bréhélin<sup>1\*</sup> & Charles-Henri Lecellier<sup>1,2\*</sup>

## Context

Gene expression is controlled by proximal and distal **cis-regulatory elements (CREs)**, containing DNA motifs bound by various transcription factors (TFs). Other sequence features, such as specific k-mers or low complexity regions, have also been implicated<sup>1-3</sup>.

However, in a dynamic biological process such as cell differentiation, we lack an understanding of how the transcriptional activity of CREs progressively change and what sequence features underlie these transitions, which may reflect common and/or coordinated regulatory processes.

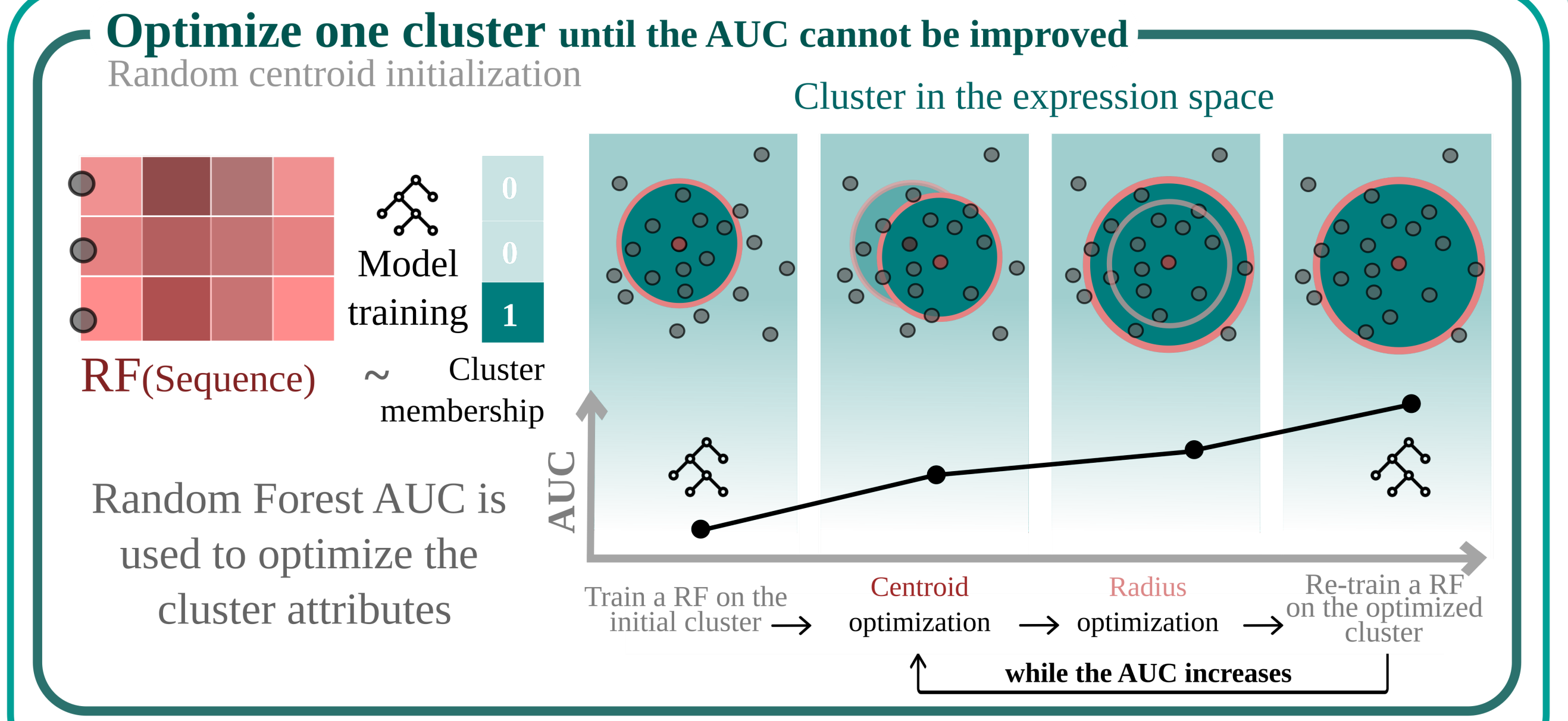
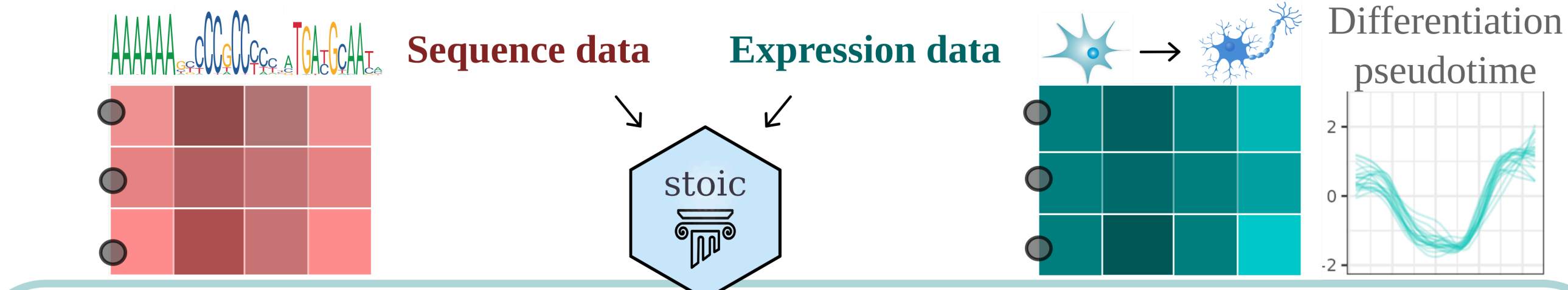
## Single cell dataset



## Stoic: a machine-learning guided method to identify the sequence features associated with specific CRE expression profiles

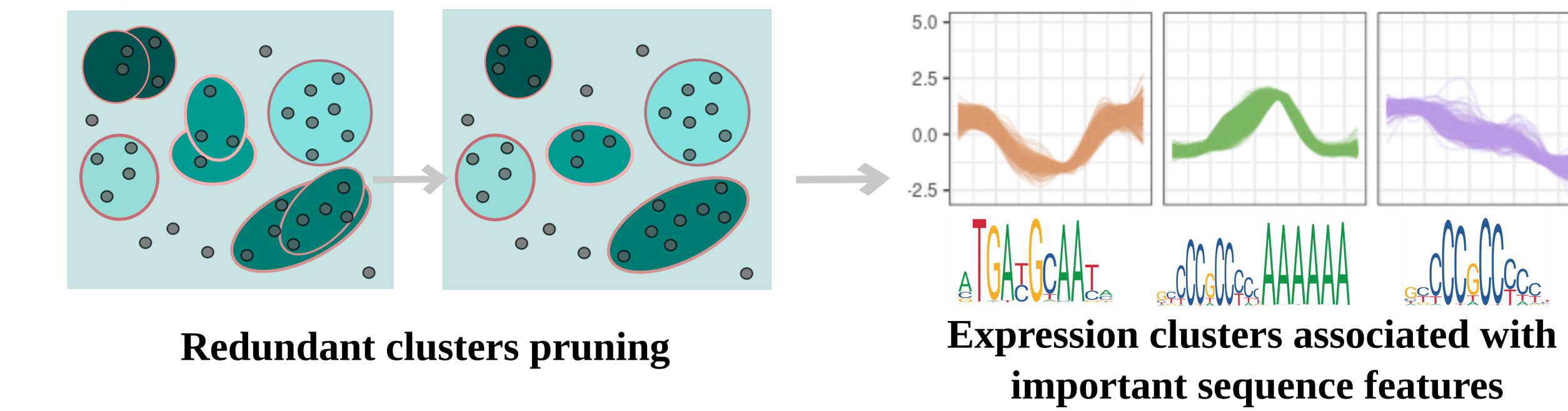
Stoic is an **original approach that aims to identify the sequences features responsible for specific CRE expression profiles** observed in the data. For this, Stoic explores the expression space and delineates the CRE clusters iteratively in order to optimize the performance of a supervised classifier predicting CRE cluster membership using only DNA sequence features.

~400 sequence features for each CRE (k-mers & TF motifs)      ~11000 transcriptionally variable CREs



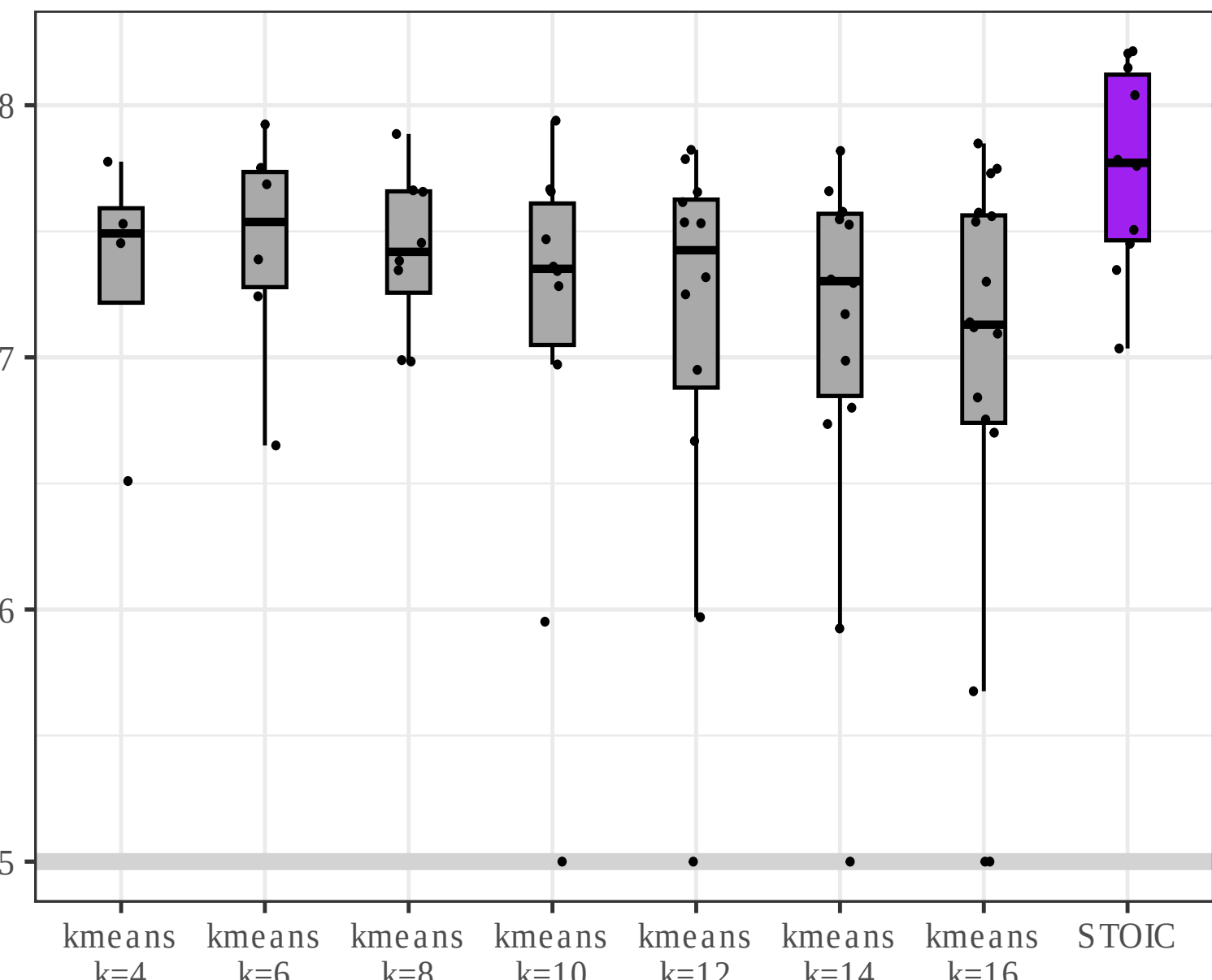
**Add clusters until the expression space is explored**  
New cluster centroids drawn as far away as possible from existing ones

**Run several times with different initializations**



## Performances

Stoic was compared to a non-guided approach that first runs a clustering of the CRE expression profiles with a standard k-means algorithm, and then trains k RFs to predict cluster membership using sequence features.

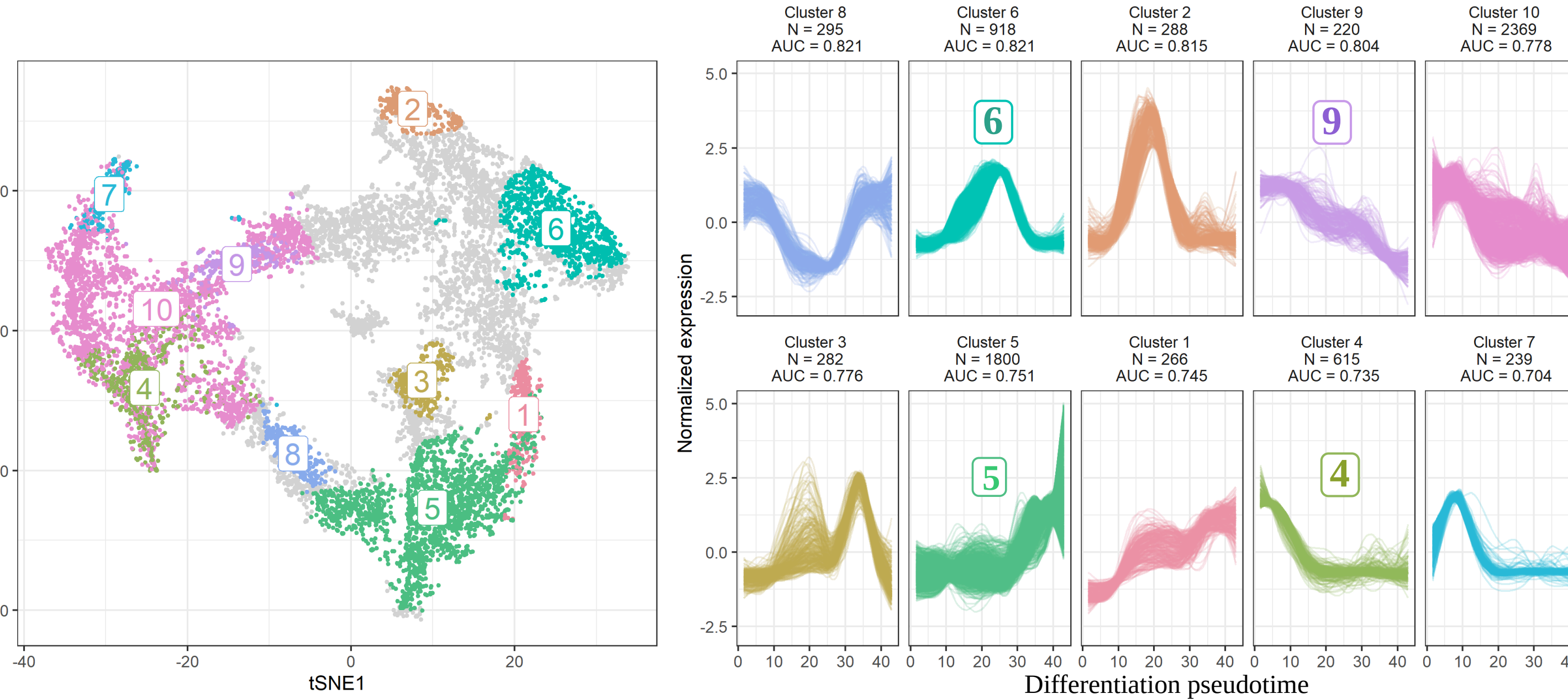


Studies of clustering stability showed that Stoic is comparable to the k-means approach for the same number of clusters and starting points.

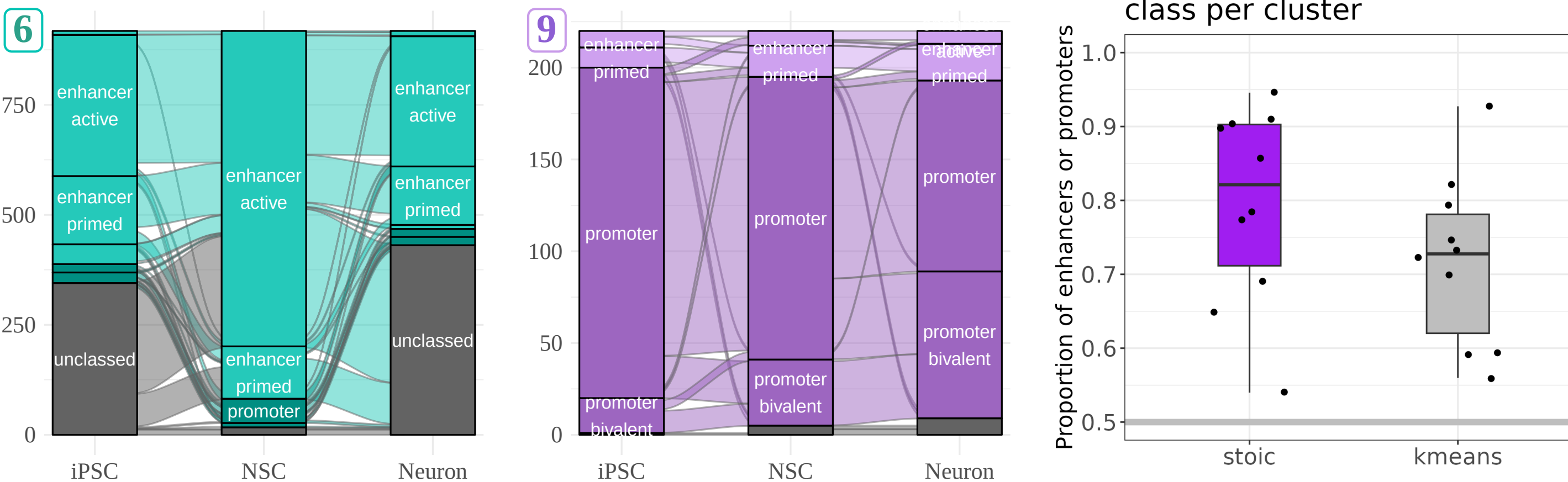
However, Stoic provides higher AUC values for discriminating cluster membership using sequence features, showing that **it recovers stronger sequence to expression associations** than the standard k-means.

Further experiments on simulated data also showed good precision and recall values.

## Stoic clusters outline diverse expression profiles and epigenetic signatures

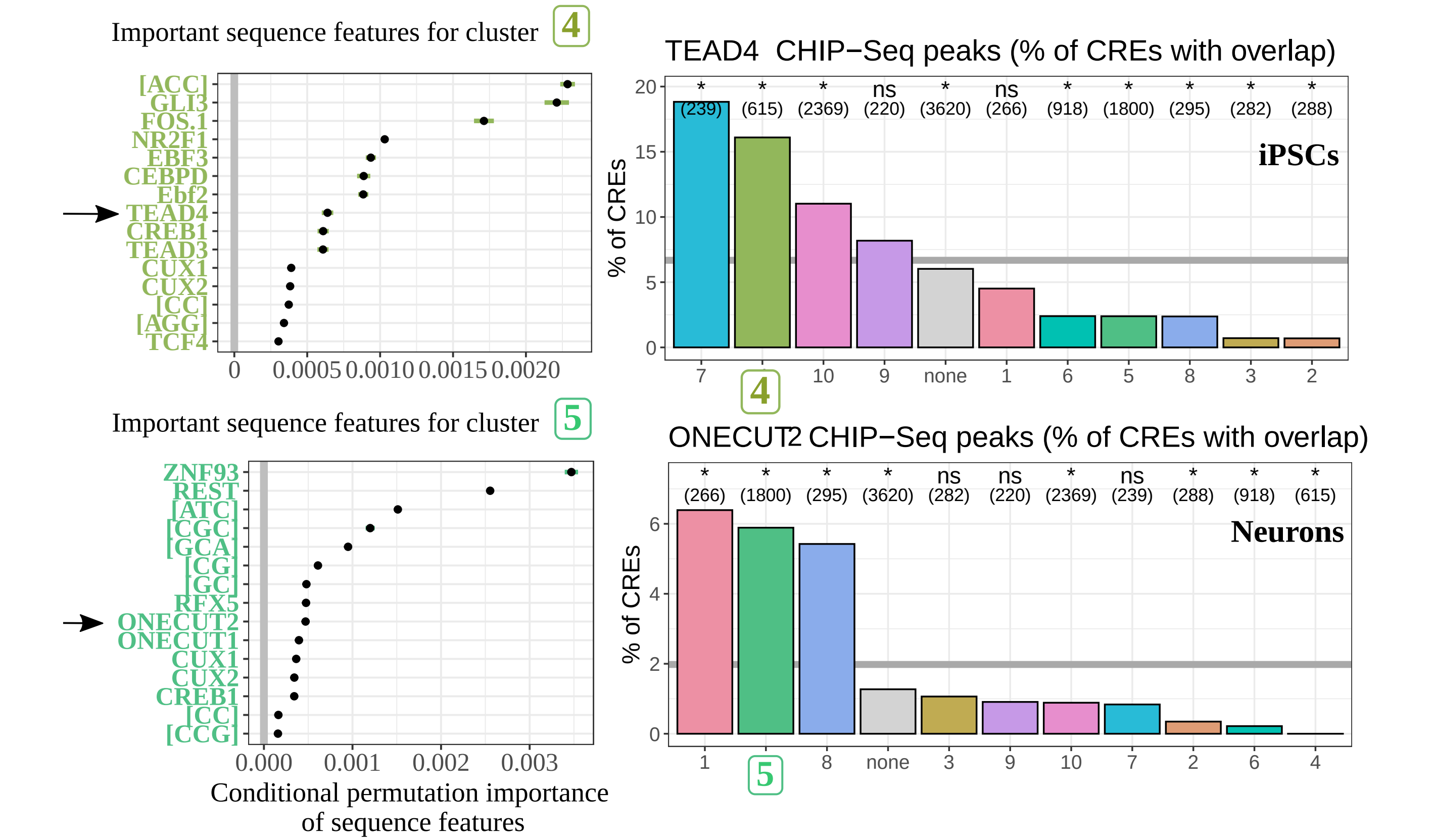


Most inferred clusters are either **enhancer-rich** or **promoter-rich**



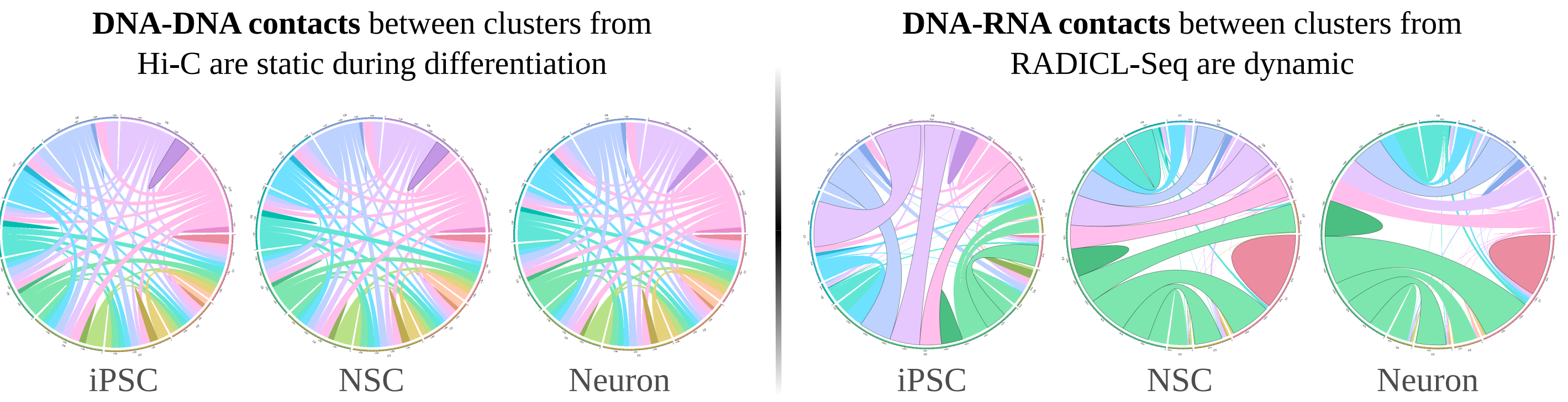
Chromatin states of CREs are predicted from CUT&Tag epigenetic data by chromHMM

## Important sequence features in Stoic clusters are supported by CHIP-Seq experiments



TEAD4 and ONECUT2 are predicted as important sequence features of clusters 4 and 5, and are **enriched in binding events** in the CREs of these clusters.

## Transcriptionally active clusters come into deep contact in a cell type-specific manner



Further interpretations of Stoic clusters based on **eQTLs, repeat elements, or clinically relevant gene sets** provide an updated perspective on the transcriptional regulations at play during neuronal differentiation.

## Stoic R package

Stoic's methodology is available as an R package. The machine-learning guided approach developed in Stoic is applicable to any problem where the clustering of some measurements can be guided by a second matched dataset.

```
library(remotes) # remotes should be installed if it is not
install_github("oceane.cssn/stoic")
```

## Affiliations

- <sup>1</sup>LIRMM, Univ Montpellier, CNRS, Montpellier, France,
- <sup>2</sup>Institut de Génétique Moléculaire de Montpellier, University of Montpellier, CNRS, France
- <sup>3</sup>RIKEN Center for Integrative Medical Sciences, Yokohama, Kanagawa, Japan

## References

1. C. M. et al. Identification of long regulatory elements in the genome of plasmodium falciparum and other eukaryotes. PLOS Computational Biology, 2021.
2. M. G. et al. Discovery of widespread transcription initiation at microsatellites predictable by sequence-based deep neural network. Nat. Communications, 2021.
3. Horton et al. Short tandem repeats bind transcription factors to tune eukaryotic gene expression. Science, 2023.
4. M. Lajoie et al. Computational discovery of regulatory elements in a continuous expression space. Genome biology, 2012.