## Deep learning has potential for harmonising multi-omics data to discover weak regulatory features

## **Tyrone Chen<sup>1</sup>, Jason D. Rigby<sup>2</sup>, Matthew D. McGee<sup>1</sup>, Sonika Tyagi<sup>1</sup>**

1 School of Biological Sciences, Monash University, Australia; 2 ASPREE National Co-ordinating and Clinical Trial Centre, Monash University, Australia











Reads containing the CTCF motif



## The CTCF motif is an example of a signal present within data



Use sequence data directly to minimise information loss
Work on any combination of omics data types

Expected problems during processing:

- Sequence length distribution differ in datasets
  Signals can be sparse and discontiguous
- Multiple signals may be at the same location
- Category imbalance skews model training

Strategies to resolve expected problems:

- Sequence vectorisation & pyramid pooling can resolve uneven sequence lengths
- Segmentation shows multiple signals in data
- Interleaving and shuffling data redistributes data more evenly

We detect CTCF with an accuracy of 86% in:

• ATAC-Seq

- Chip-Seq
- RIP-Seq data
- 3D genome data integration is in progress

## School of Biological Sciences

Genetics, Genomics and Health





We acknowledge high performance computing support provided by the Monash eResearch Centre.