# A weighted two-stage sequence alignment framework to identify DNA motifs from ChIP-exo data

**Cankun Wang**
Biomedical Informatics Specialist
Department of Biomedical Informatics
The Ohio State University

July 18, 2023

THE OHIO STATE UNIVERSITY
WEXNER MEDICAL CENTER

BMBL
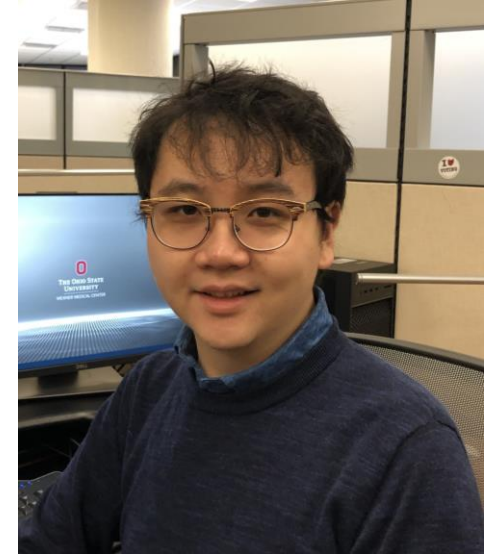Bioinformatics and Mathematical Biosciences Lab

**Research interest:**

**Inference of gene regulatory mechanisms across various organisms**
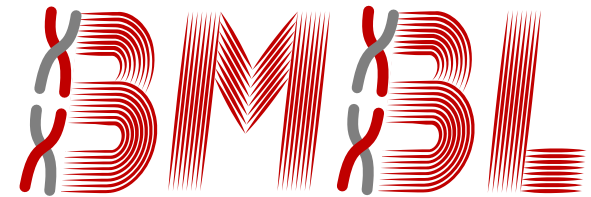
**Developing tools and benchmarking pipelines for next-generation sequencing data**

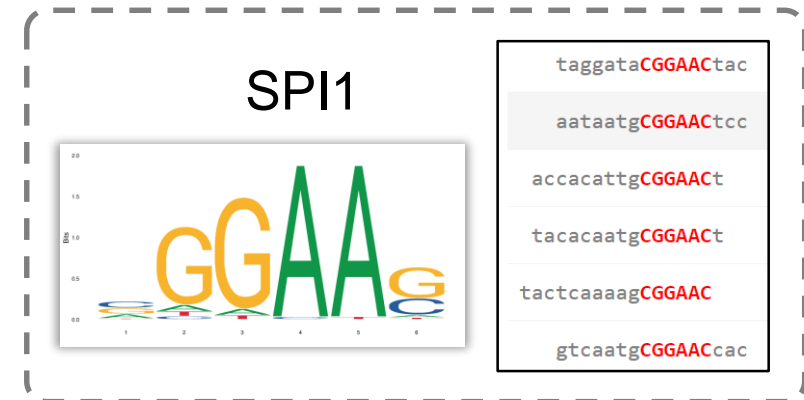**Developing cloud-native biomedical applications for webservers and databases**
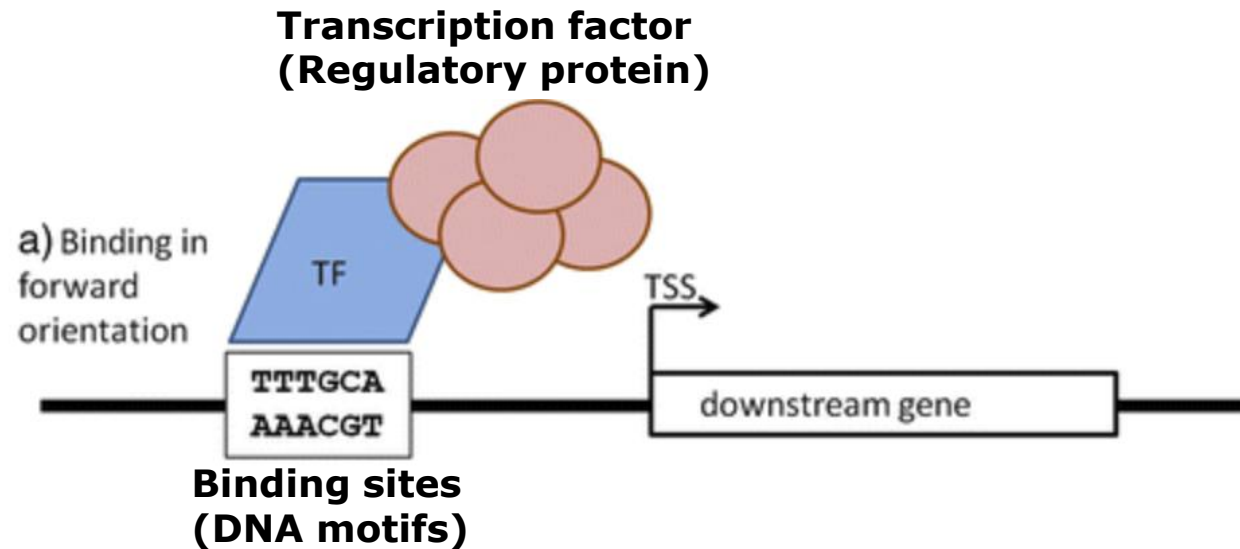
Cankun Wang
Biomedical Informatics Specialist

**BMBL**
Bioinformatics and Mathematical Biosciences Lab
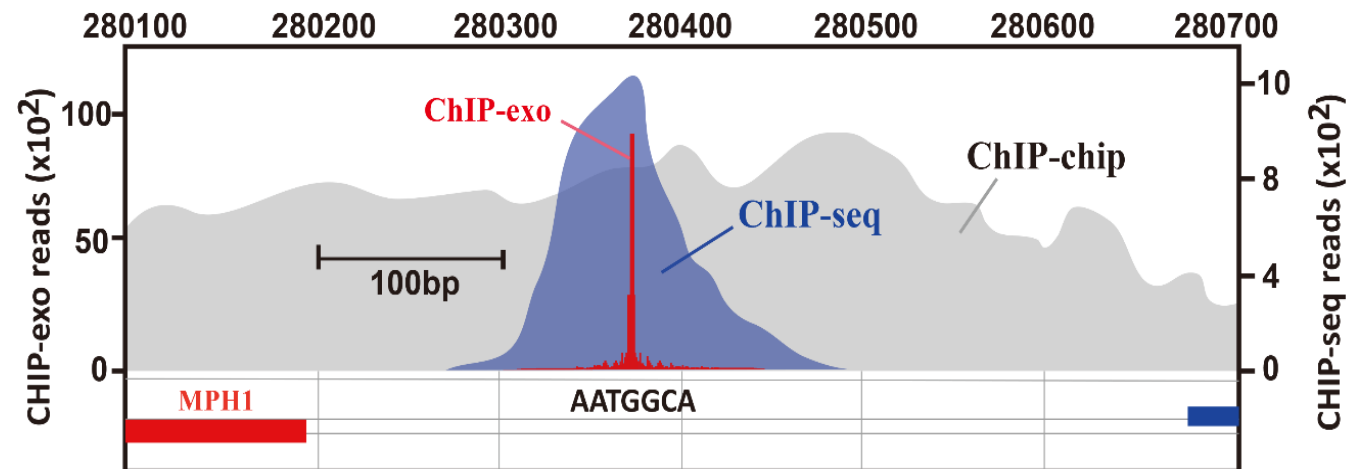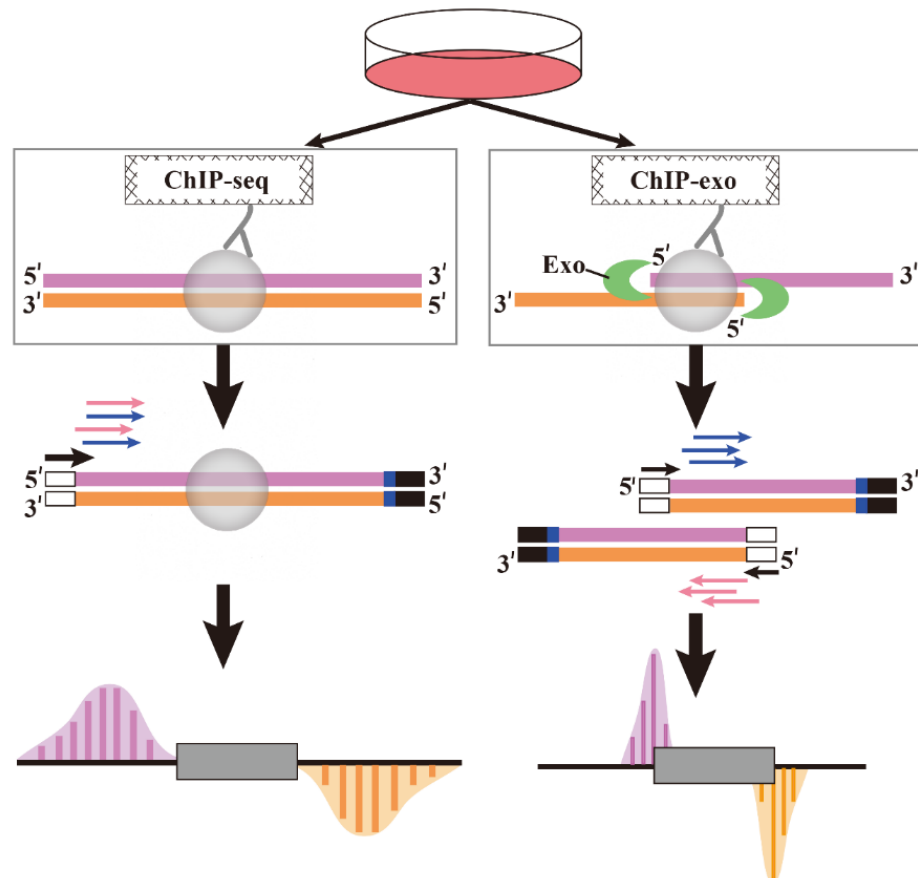
Advisor: Dr. Qin Ma

# Background of prediction of transcription factor binding sites



**Transcription factor (TF) ---** a protein that controls the rate of transcription of genetic information from DNA to messenger RNA, by binding to a specific DNA sequence, i.e., TF binding sites (TFBSs).

Bailey, *et al., Nucleic Acids Res.* (2015)
Liu, *et al., Briefs. Bioinformatics* (2018)

# Why ChIP-exo Over ChIP-seq for Motif Finding

- ChIP-seq provides genome-wide map of protein-DNA interactions, yet <u>resolution could be insufficient</u> for accurate binding site identification
- Compared with ChIP-seq, Chromatin Immunoprecipitation combined with lambda exonuclease digestion followed by high-throughput sequencing (ChIP-exo) has relatively <u>low noise</u> and achieves <u>near-base pair resolution</u>.

Peter J. Park, *Nature Reviews Genetics.* (2009)
Rossi, *et al., Nature Communications* (2018)

# Why new motif discovery tool is needed for ChIP-exo data?

1.  **High Sensitivity of ChIP-exo**: Traditional tools may not handle the increased sensitivity to experimental conditions in ChIP-exo data adequately.

2.  **Peak Position Variance:** Traditional motif finding tools often assume the binding site is at the peak center, which does not hold true for ChIP-exo data, necessitating a more flexible algorithm.

3.  **Sharpness of Peaks**: ChIP-exo data generates sharper peaks than ChIP-seq, making it difficult for traditional tools to distinguish between adjacent binding events.

4.  **Need for Advanced Techniques**: New tool is required to account for these intricacies, enhancing signal-to-noise ratio and providing a more precise motif discovery.
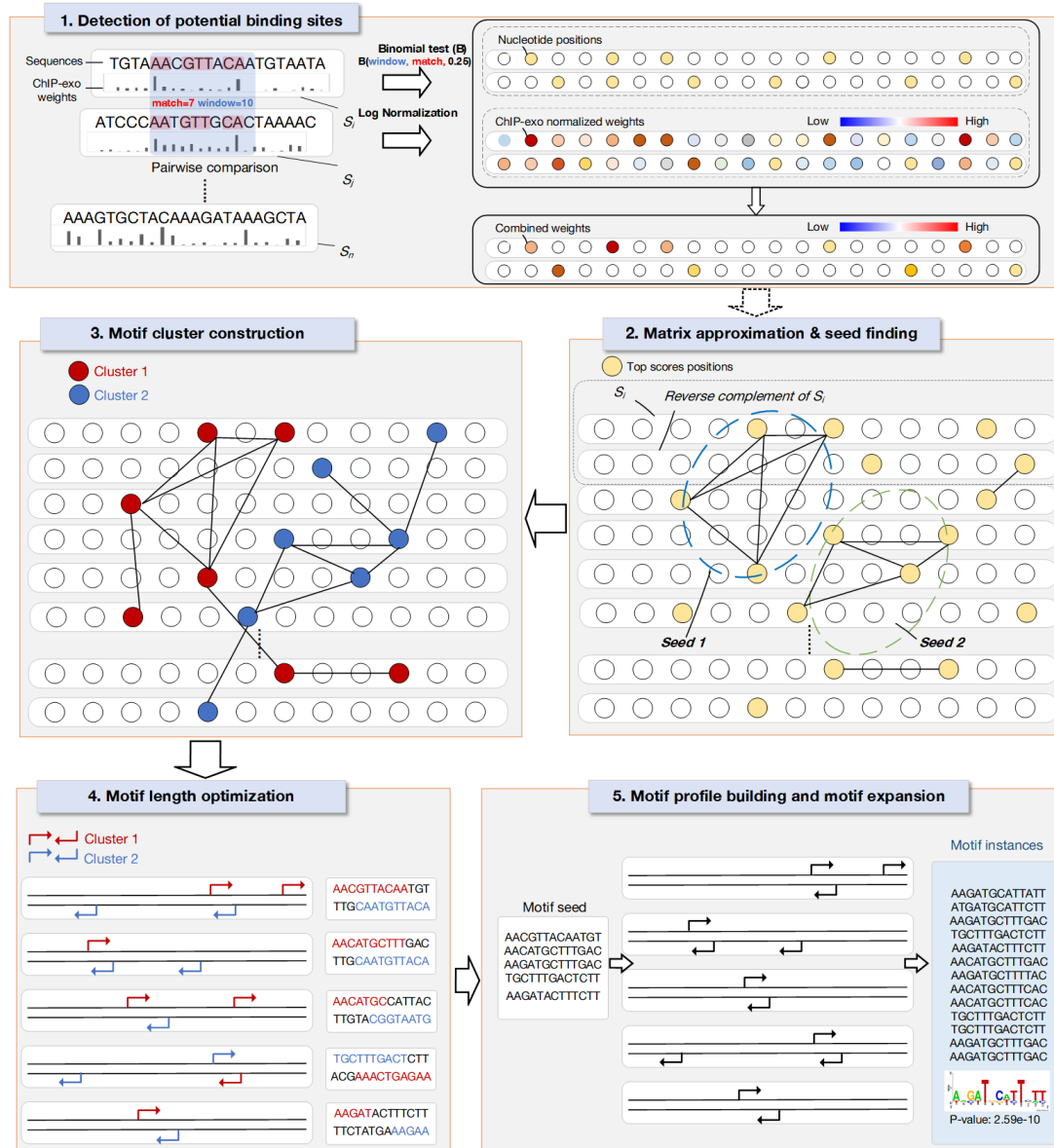
# TESA is a novel motif finding tool designed for ChIP-exo data

TESA (**A weighted two-stage sequence alignment framework to identify DNA motifs from ChIP-exo data** ) was designed specifically to address the unique challenges presented by ChIP-exo data.

Unique features:

- Base-level Signal Extraction: TESA leverages ChIP-exo's high-resolution, base-level data, enhancing signal-to-noise ratio and providing precise motif discovery

- Dynamic Motif Length Optimization: Through a bookend model, TESA automatically adjusts motif lengths based on TFBS clustering

- Precision in Distinguishing Adjacent Binding Events: Using a binomial test, TESA determines whether potential TF binding site clusters should be combined or treated separately
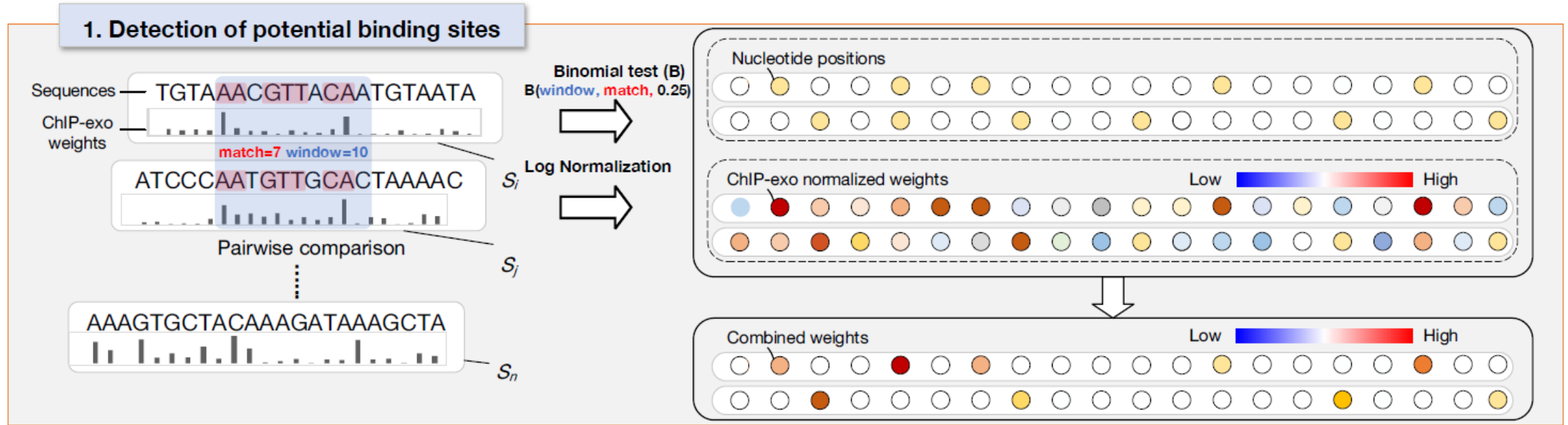
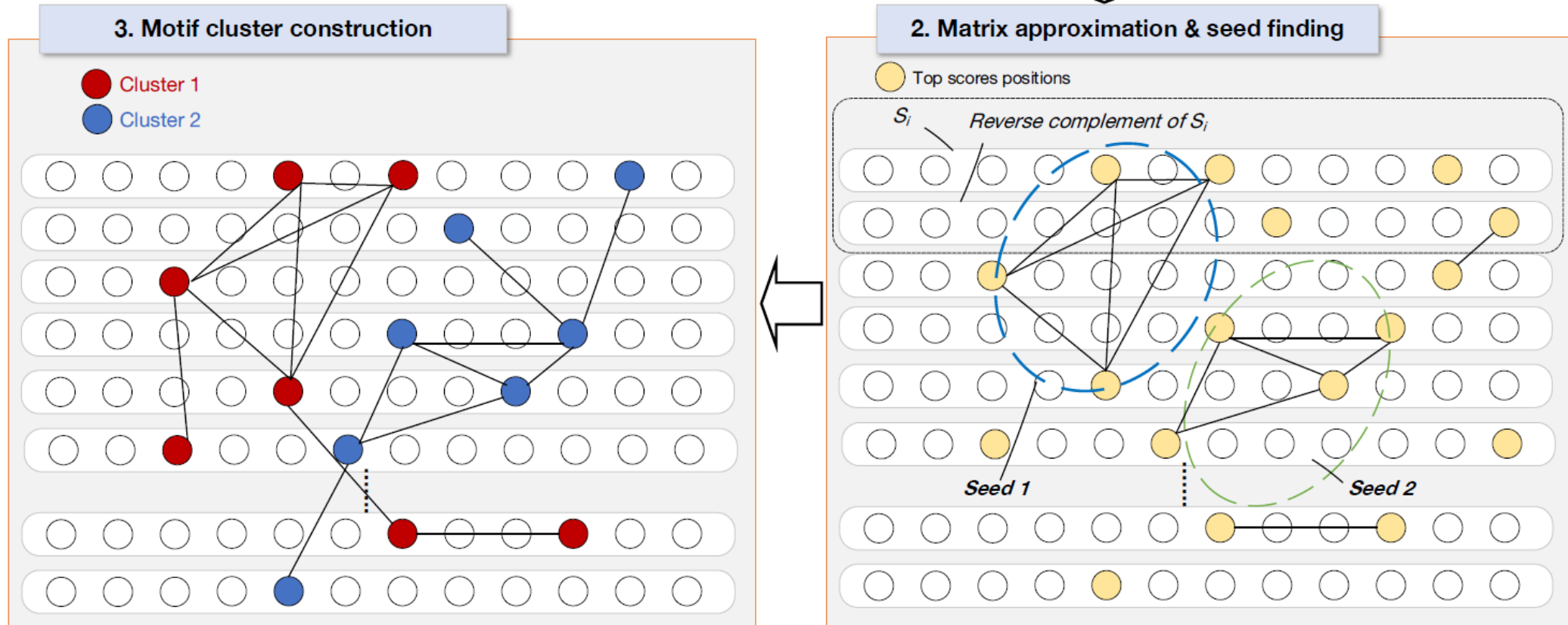# Overview of TESA (weighTEd two-Stage Alignment tool)

**Steps of TESA:**
1. Detection of potential binding sites
2. Matrix approximation and seed finding
3. Motif cluster construction
4. Motif length optimization
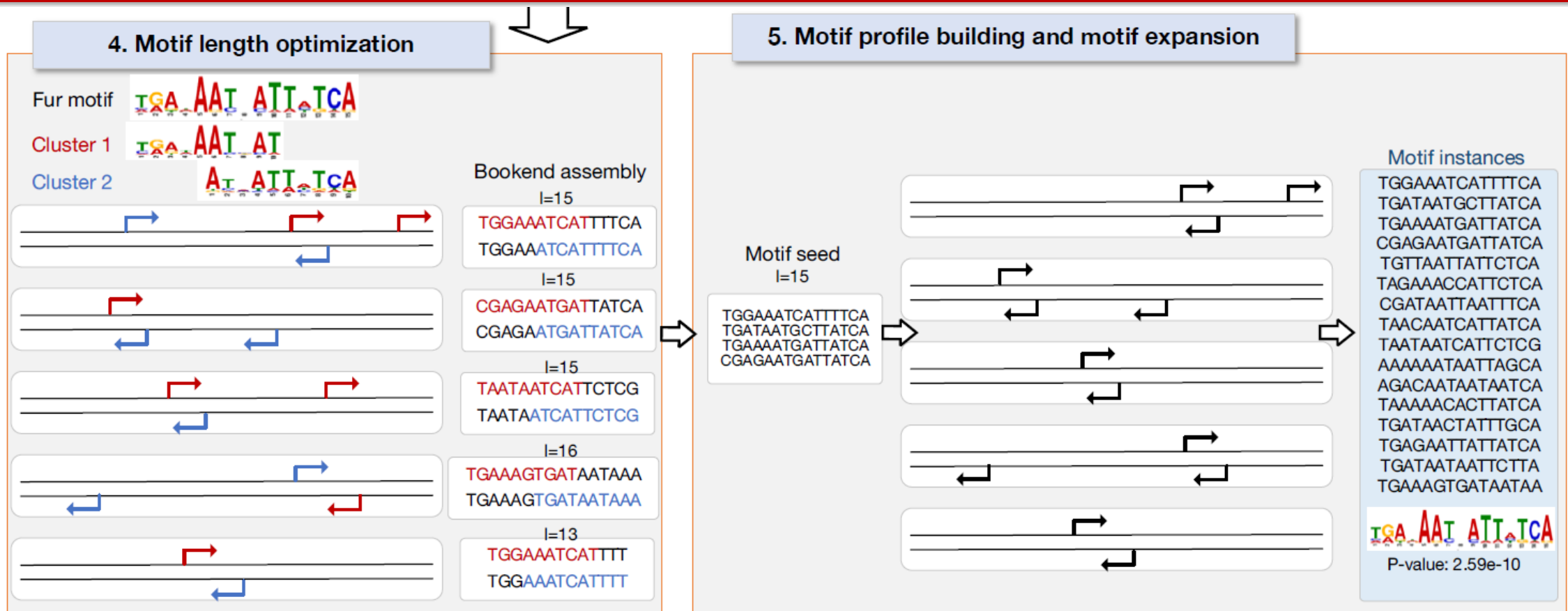5. Motif profile building and motif expansion

Compared with ChIP-seq, Chromatin Immunoprecipitation combined with lambda exonuclease digestion followed by high-throughput sequencing (ChIP-exo) has relatively low noise and achieves near-base pair resolution.

# Matrix Approximation and Cluster Finding



- Initial data matrix is refined to highlight potential binding site locations.
- constructed network with points representing possible binding sites.
- Connections are made to identify similar potential binding sites.
- Final clusters are formed by linking these seeds, signifying potential motifs.

# Refinement of motif and length



- Optimized Motif Lengths: TESA prevents premature truncation and enhances accuracy.
- Expanded Motifs: Sequence segments with high similarity scores are added to the motifs.
- Iterative Refinement: Continual reassessment and inclusion ensure highly refined motifs.

# Methods for benchmarking the performance of motif finding

**Methods to compare:**
- **TESA,** BoBro, MEME-ChIP, and Weeder, DiNAMO, Dipartitle, MFMD
- 20 Escherichia coli ChIP-exo datasets downloaded from the proChIPdb database
- Validating the motifs discovered by TESA by comparing them with known motifs from the DPInteract database

BoBro: The segment alignment algorithm as basis of TESA (Li et al. 2011)

MEME-ChIP: This popular tool has been cited 1551 times, widespread use and recognition in the community (Bailey 2011)

Weeder: With 627 citations, Weeder stands as another well-recognized tool (Pavesi et al. 2014)

**Added during revision:**

DiNAMO: An exhaustive and efficient algorithm for motif discovery (Saas et al. 2018)

Dipartitle:  A tool for detecting motifs by considering base interdependencies (Vehed et al. 2018)
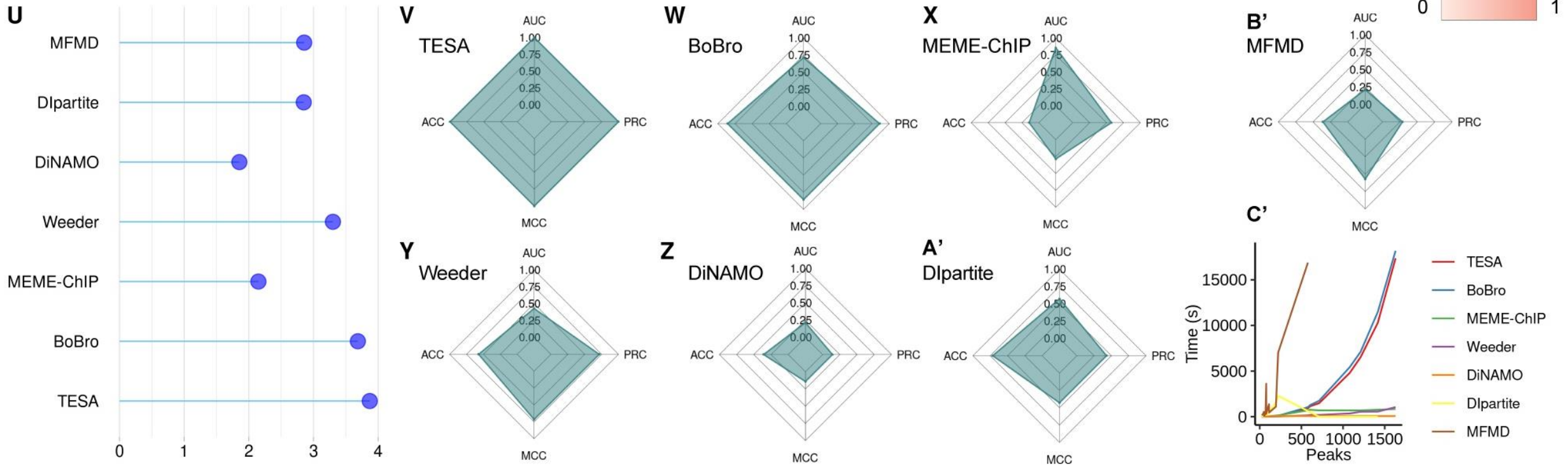
MFMD: Semi-greedy constructive heuristics as a local optimizer for motif finding (Caldonazzo et al. 2019)

# Evaluation metrics of the performance of motif finding

- Utilization of standard performance metrics: PRC, AUC, MCC

- 20 Escherichia coli ChIP-exo datasets downloaded from the proChIPdb database

- Comparison with known motifs from DPInteract, leveraging TOMTOM for similarity computation and significance assessment

# Overall performance of TESA



Lollipop plot: Overall scores of each algorithm summed across 20 datasets

Spider plots: averaged scores of each algorithm

12

**Motif similarity:**

- TESA discovers motifs with significant similarity with target motifs
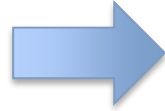- Some TFs are co-factors of the ChIP-ed one

| Dataset | Rank | TF | D | L1 | L2 | Q.value | Logo alignment |
|---------|------|-----|---|----|----|---------|----------------|
| cpxr_EtOH | 2 | cpxR | - | 16 | 15 | 2.61E-05 |  |
| cpxr_EtOH | 6 | gcvA | - | 16 | 20 | 2.39E-02 |  |

...significance of motif alignment (the seventh column), and the eighth column showcases the alignment of motif logos.

# Summary of TESA



**TESA**
A weighted two-stage sequence alignment framework to identify DNA motifs from ChIP-exo data

- Effective Graph Construction from base-resolution ChIP-exo

  - Optimal Motif Length Determination

  - High precision in motif identification

# Limitations of TESA

- Computational resources needs

- Multiple parameter dependency

- Contemplation of alternative negative controls

# Acknowledgement



**Dr. Qin Ma**

**Dr. Bingqiang Liu**

Yang Li

Yizhong Wang

Anjun Ma

**BMBL**
Bioinformatics and Mathematical Biosciences Lab

**https://u.osu.edu/bmbl/**

@Wang-Cankun
@QinMaBMBL

Cankun.Wang@osumc.edu

16

# Thank You

Cankun.Wang@osumc.edu