# Probabilistic generative modeling of multi-mapping reads with mHi-C advances analysis of Hi-C studies

Sündüz Keleş

Department of Biostatistics and Medical Informatics

Department of Statistics

University of Wisconsin, Madison

Google atsnp search

# http://atsnp.biostat.wisc.edu/

Google atsnp search

# High throughput chromatin conformation capture (Hi-C) for studying long-range interactions



Up to 10Mb

Enhancer TF TF Pol II Promoter Gene

ENCODE project generated catalogs of enhancers.

Pombo & Dillon, *Nature Reviews Molecular Cellular Biology*, 2015
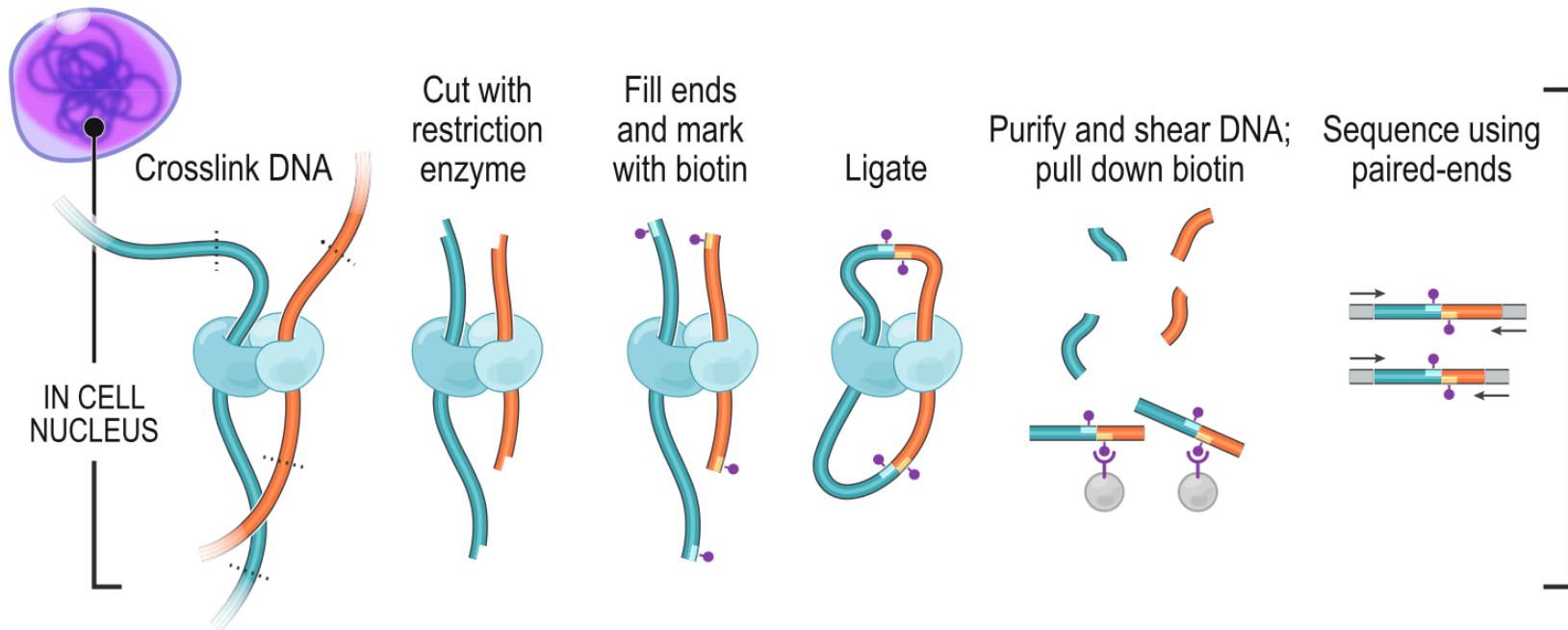
# Hi-C for studying long-range interactions

Up to 10Mb

Looping of DNA



ENCODE project generated catalogs of enhancers.

# Hi-C experimental protocol
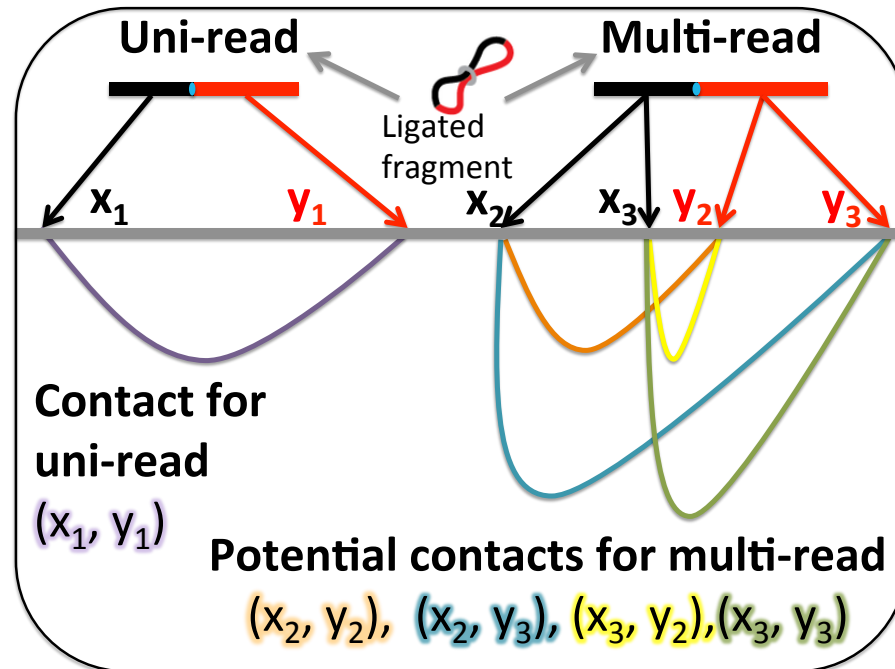


Rao *et al*.,  *Cell*, 2014
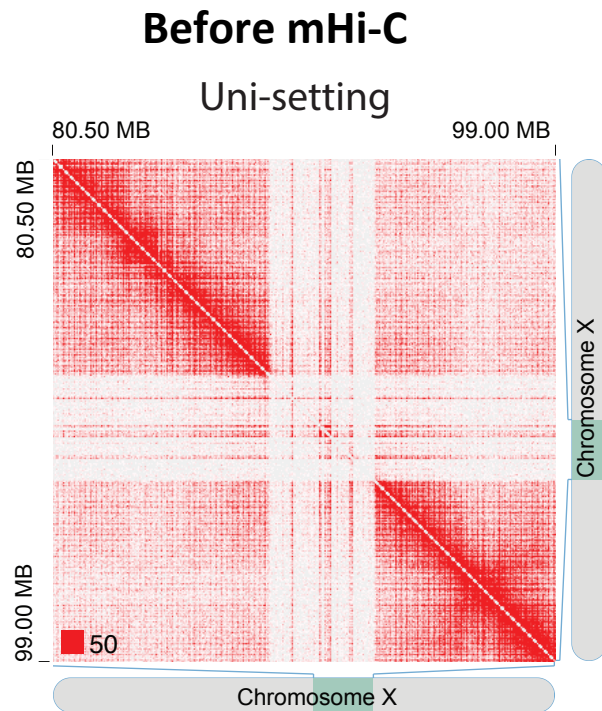
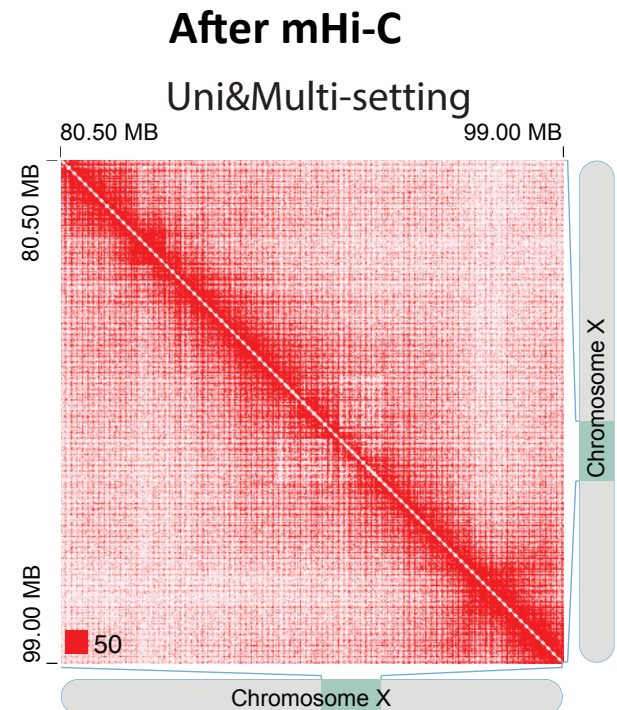# Hi-C experimental protocol

Data gets summarized as a contact count matrix.

# Just like any sequencing dataset, Hi-C analysis start with read alignment

# Signals from repetitive regions are under-represented



**Before mHi-C**

Uni-setting

# Signals from repetitive regions are under-represented



**Before mHi-C**

**After mHi-C**

Uni-setting

Uni&Multi-setting

Chromosome X

# Evaluation: 6 independent studies, with 8 datasets, and multiple replicates per dataset

**Table 1.** Hi-C Data Summary

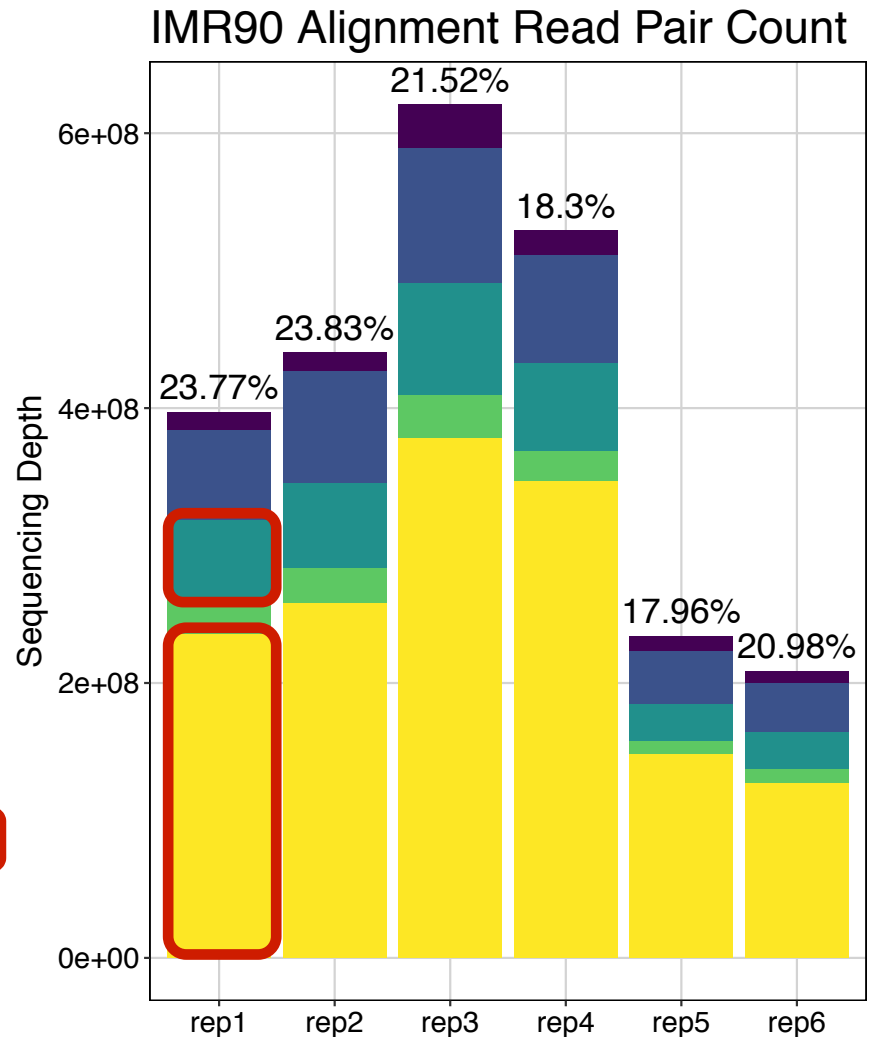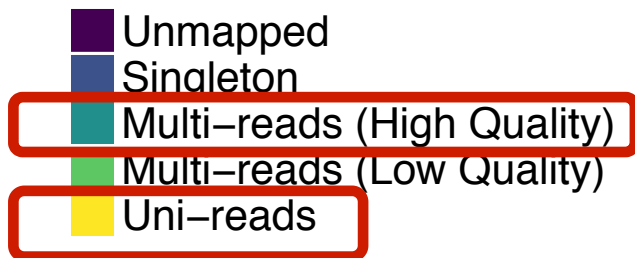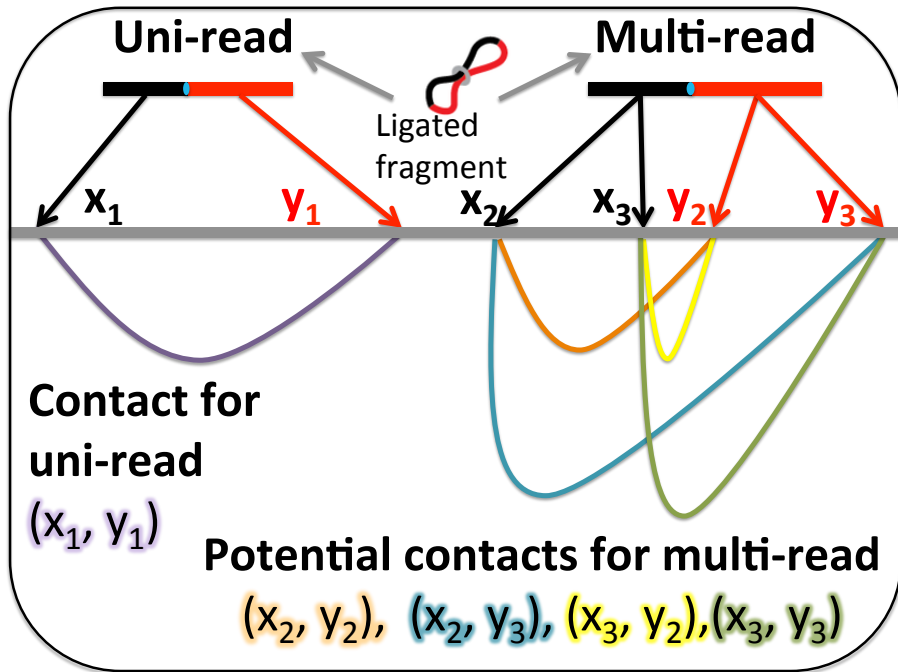| Cell line | Replicate | Read length (bp) | Restriction Enzyme | HiC Protocol | Source | Resolution (kb) |
|---|---|---|---|---|---|---|
| IMR90 | rep1-6 | 36 | HindIII | dilution | *Jin et al. (2013)* | 40 |
| GM12878 | rep2-9 | 101 | MboI | in situ | *Rao et al. (2014)* | 5, 10*, 40* |
| GM12878 | rep32, rep33 | 101 | DpnII | in situ | *Rao et al. (2014)* | 5 |
| A549 | rep1-4 | 151 | MboI | in situ | *Dixon et al. (2018)* | 10, 40 |
| ESC(2012) | rep1, rep2 | 36 | HindIII | dilution | *Dixon et al. (2012)* | 40 |
| ESC(2017) | rep1-4 | 50 | DpnII | in situ | *Bonev et al. (2017)* | 10, 40 |
| Cortex | rep1-4 | 50 | DpnII | in situ | *Bonev et al. (2017)* | 10, 40 |
| P.falciparum | 3 stages | 40 | MboI | dilution | *Ay et al. (2014b)* | 10, 40 |

\* Replicates 2, 3, 4, and 6 of the GM12878 cell line datasets were process at 10kb and 40kb resolutions.
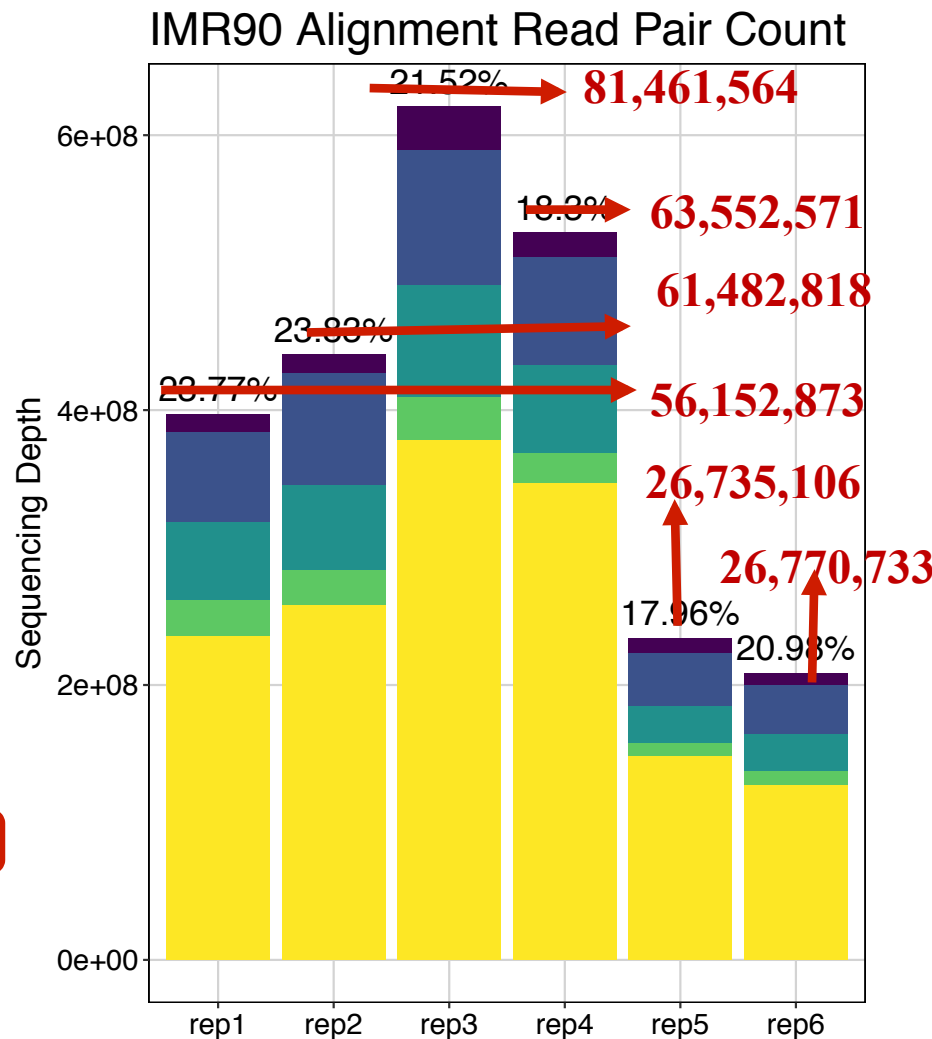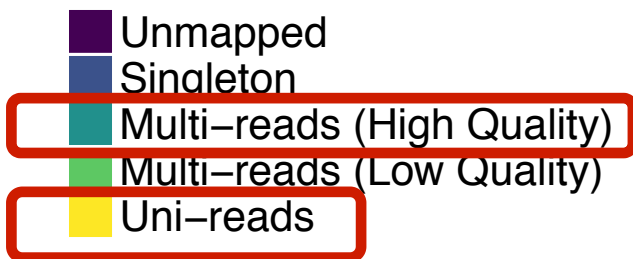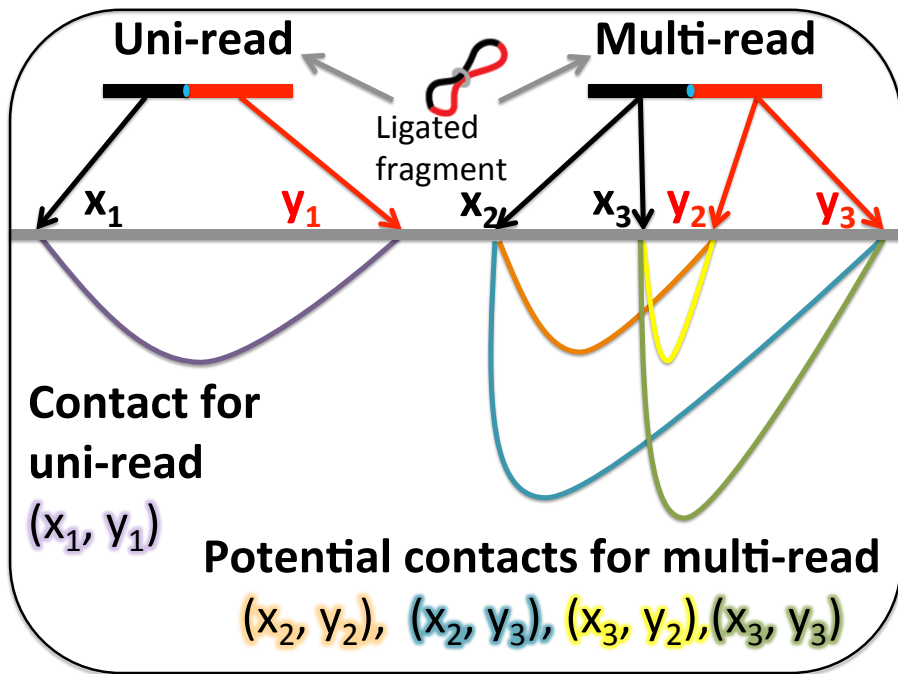
## Criteria for selection

- Genome size (large, small)
- Sequencing depth, coverage
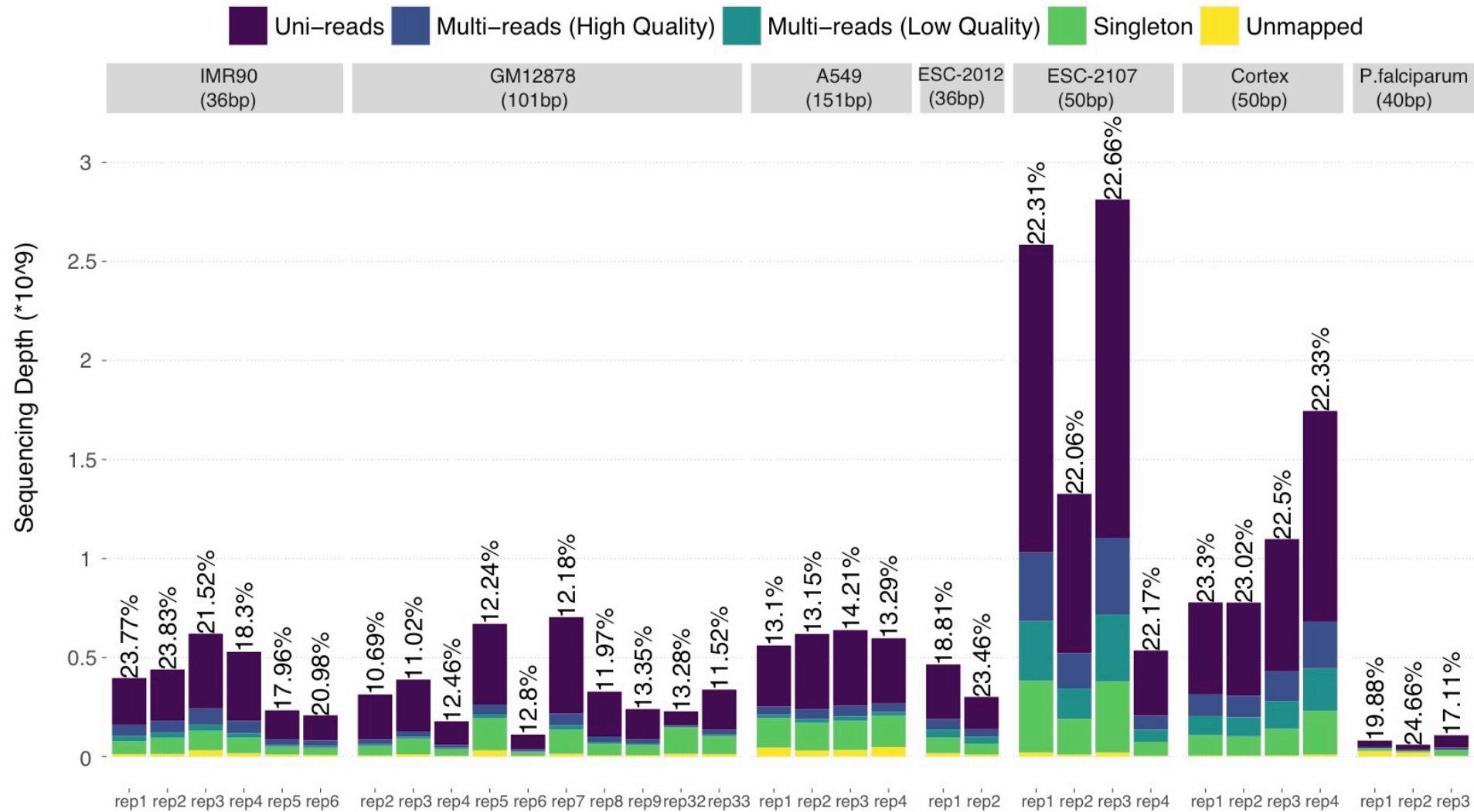- *Cis*-to-*Trans* ratio
- Proportion of mappable and valid reads
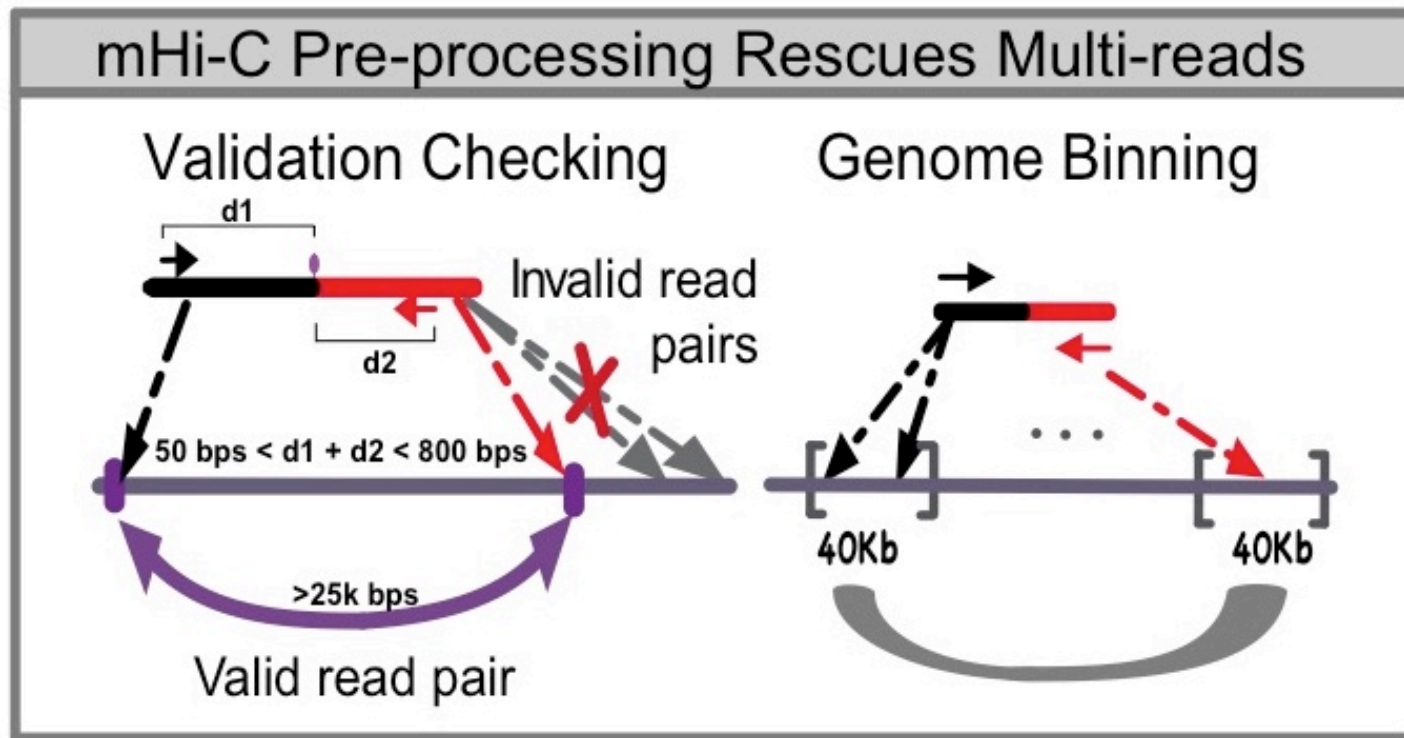
# Multi-reads are abundant



**Uni-read**     **Multi-read**
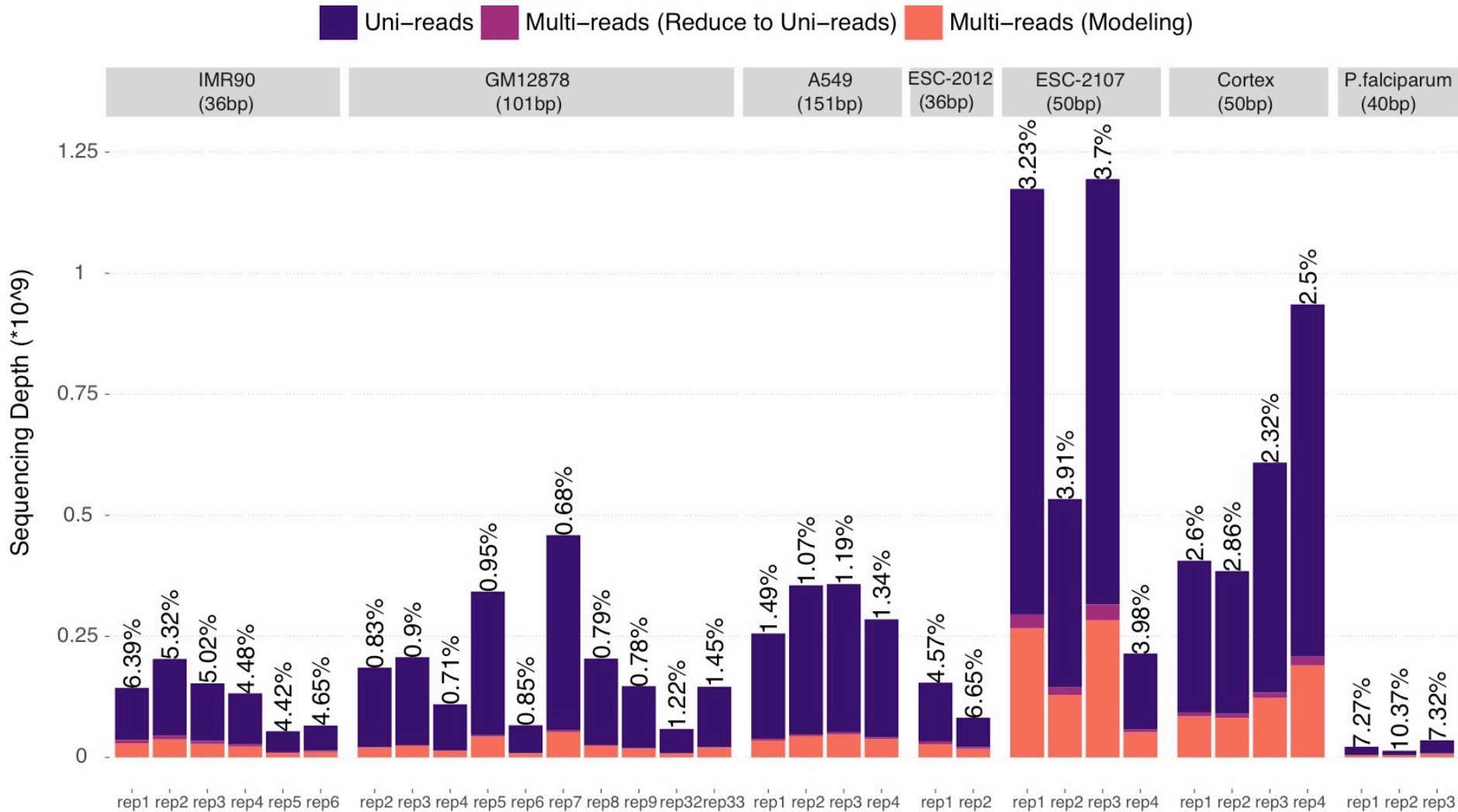
Ligated fragment

$x_1$   $y_1$   $x_2$   $x_3$   $y_2$   $y_3$

...tial contacts for multi-read

$(y_2)$, $(x_2, y_3)$, $(x_3, y_2)$, $(x_3, y_3)$

IMR90 Alignment Read Pair Count

21.52%

18.3%

23.77%    23.83%

17.96%

20.98%

Sequencing Depth

6e+08

4e+08

2e+08

0e+00

rep1   rep2   rep3   rep4   rep5   rep6

- Unmapped
- Singleton
- Multi–reads (High Quality)
- Multi–reads (Low Quality)
- Uni–reads

17.96%

20.98%

# Multi-reads are abundant



**Uni-read**  **Multi-read**

Ligated fragment

$x_1$  $y_1$  $x_2$  $x_3$  $y_2$  $y_3$

**tial contacts for multi-read**

$(x_2, y_2)$, $(x_2, y_3)$, $(x_3, y_2)$, $(x_3, y_3)$

- Unmapped
- Singleton
- Multi–reads (High Quality)
- Multi–reads (Low Quality)
- Uni–reads

17.96%

20.98%

## IMR90 Alignment Read Pair Count

21.52% → **81,461,564**

18.3% → **63,552,571**

**61,482,818**

23.83%

23.77% → **56,152,873**

**26,735,106**

**26,770,733**

17.96%

20.98%

Sequencing Depth

6e+08

4e+08

2e+08

0e+00

rep1  rep2  rep3  rep4  rep5  rep6

# Results across eight studies

# Sometimes, there is free lunch

# Sometimes, there is free lunch

# Sometimes, there is free lunch

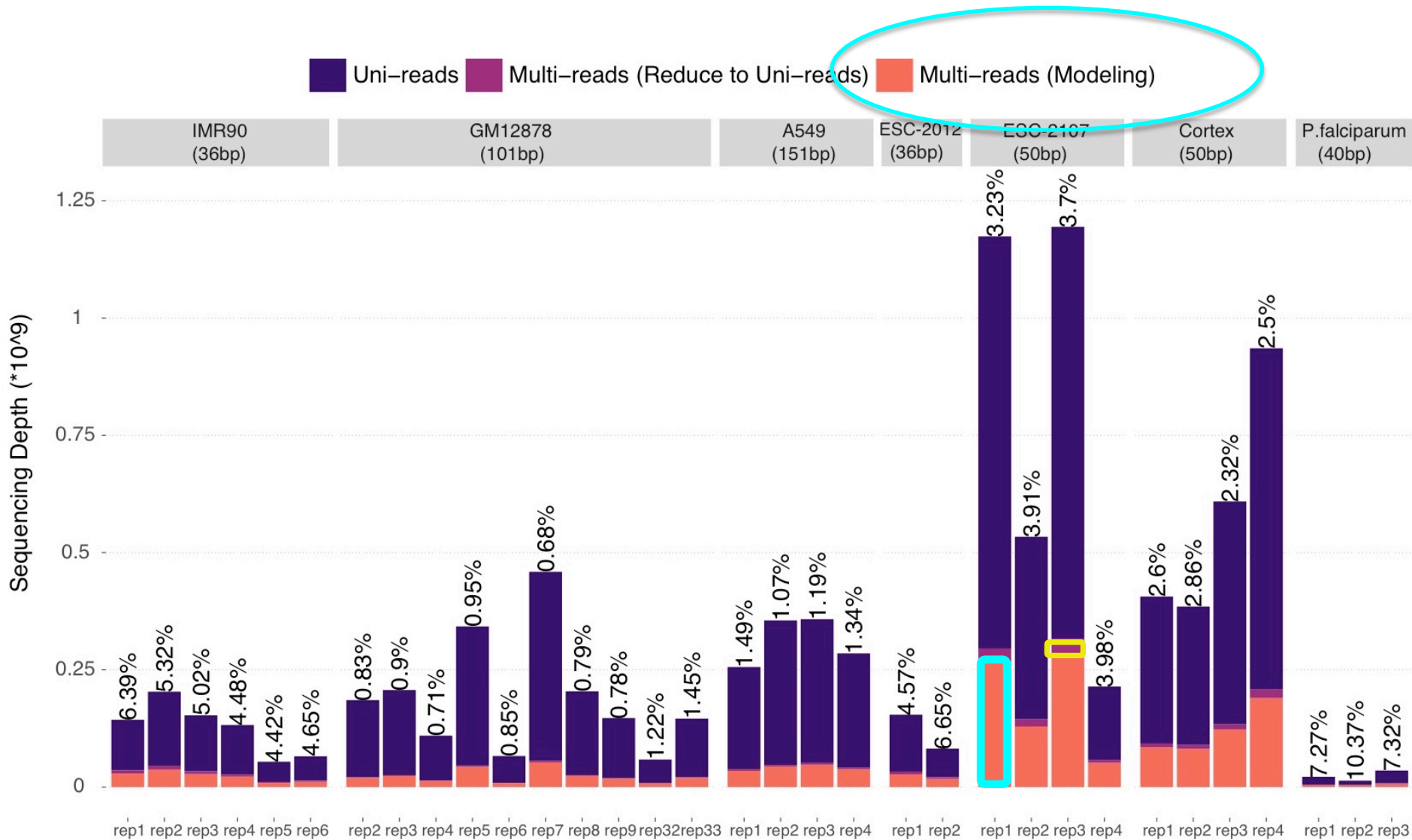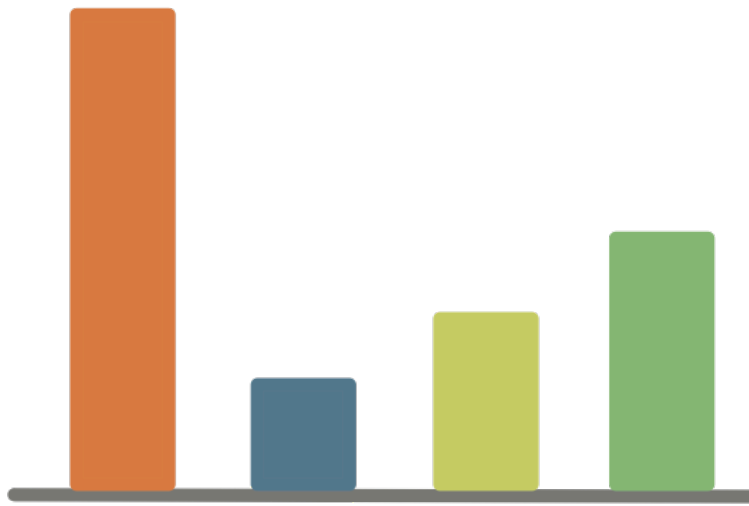No-cost multi-reads: add ~5%

# Sometimes, there is free lunch

# mHi-C: multi-read allocation for Hi-C
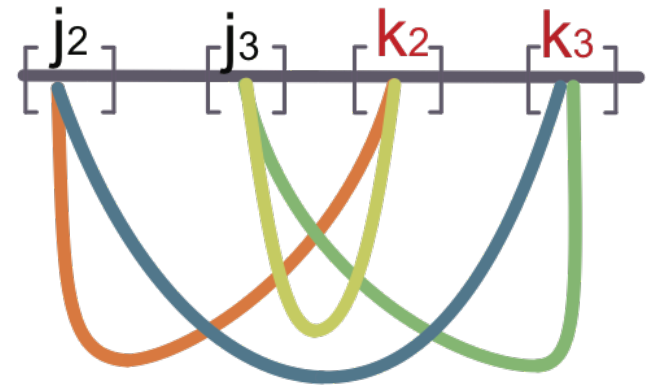


Local Bin-pair Contact Counts
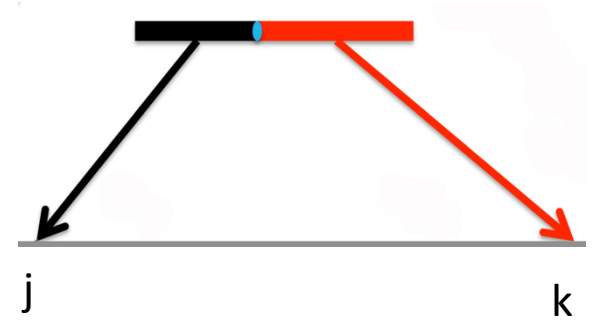
# mHi-C model

Observed:  $Y_{i,(j,k)} = 1.$

Valid read pair $i$ aligned

to contact unit $(j, k)$.

# mHi-C model

Observed: $Y_{i,(j,k)} = 1.$
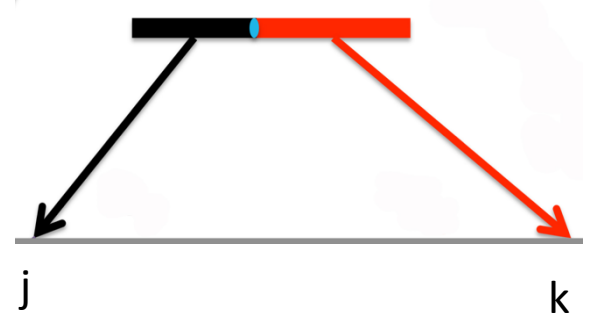
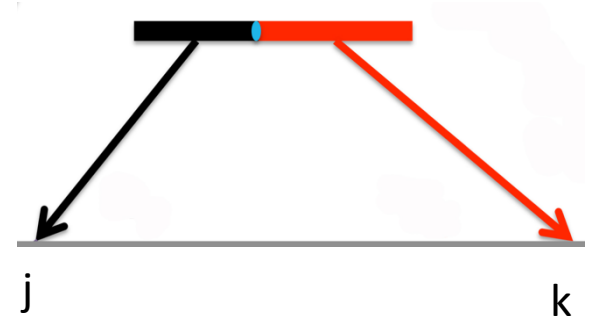Valid read pair $i$ aligned

to contact unit $(j, k)$.

Uni

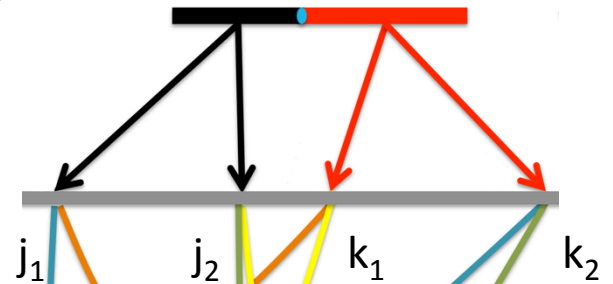$$\sum_{j,k} Y_{i,(j,k)} = 1, \quad i = 1, \cdots, N$$

# mHi-C model

Observed: $Y_{i,(j,k)} = 1.$

Valid read pair $i$ aligned

to contact unit $(j,k)$.

Uni



j                    k

$$\sum_{j,k} Y_{i,(j,k)} = 1, \quad i = 1, \cdots, N$$

Multi



$j_1$        $j_2$       $k_1$            $k_2$

e.g.

$$\sum_{j,k} Y_{i,(j,k)} = 4$$

$$\sum_{j,k} Y_{i,(j,k)} \geq 1 \quad , i = 1, \cdots, N$$

# mHi-C model

**Observed:** $Y_{i,(j,k)} = 1.$

Valid read pair $i$ aligned

to contact unit $(j, k)$.



Uni

**Hidden:** $Z_{i,(j,k)} = 1,$

Valid read pair $i$ originated

from contact unit $(j, k)$.

Multi

# mHi-C model

Observed: $Y_{i,(j,k)} = 1.$

Valid read pair $i$ aligned

to contact unit $(j, k)$.

Hidden: $Z_{i,(j,k)} = 1,$

Valid read pair $i$ originated
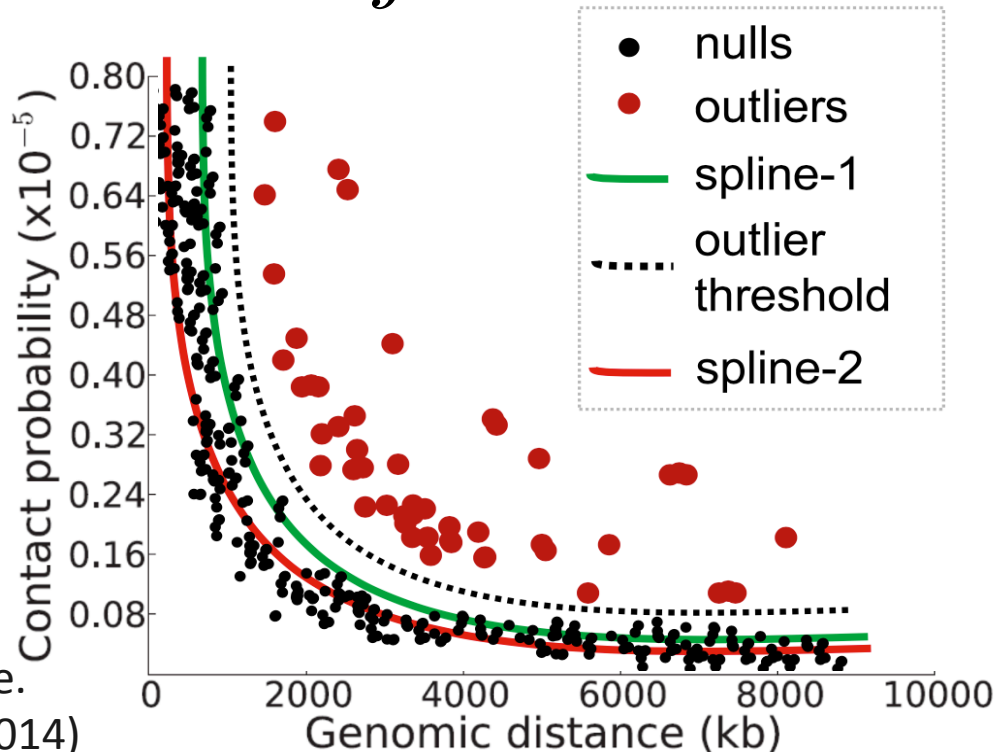from contact unit $(j, k)$.

Uni

$$\sum_{j,k} Z_{i,(j,k)} = 1$$

Multi

$$\sum_{j,k} Z_{i,(j,k)} = 1$$

# mHi-C model

$$Z_i \sim \mathrm{Multinomial}(\pi_{(1,2)}, \pi_{(j,k)}, \cdots, \pi_{(M,M-1)})$$

$$\pi \sim \mathrm{Dirichlet}(\gamma_{(1,2)}, \cdots, \gamma_{(j,k)}, \cdots, \gamma_{(M,M-1)})$$

$\gamma_{(j,k)}$ is modeled as a function of the

distance between contact units $j$ and $k$

$\gamma_{(j,k)}$ play the role of
**pseudo-counts** in the Dirichlet-
Multinomial framework.

Ay, Bailey, and Noble.
*Genome research* (2014)

# mHi-C

$$P(Z_{i,(j,k)} = 1 \mid Y_{i,(j',k')}, \forall j', k')$$

Threshold posterior probabilities to use resulting alignments with existing significant contact identification methods (e.g., fit-HiC).

# mHi-C: from read-pairs to significant contacts

Process reads to get valid read pairs

Partition genome into non-overlapping intervals
(5-300Kb or 10 RE sized units)

Generate raw contact map

mHiC makes these steps multi-read aware
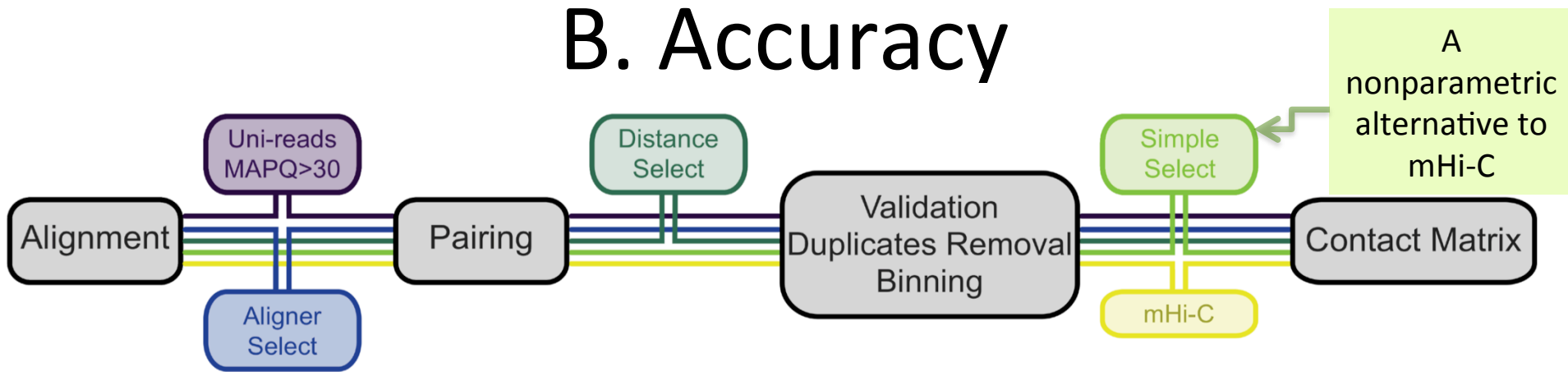
Normalize contact map

Identify significant contacts

# Evaluation

| | |
|---|---|
| A. Sequencing depth | ✔ |
| B. Accuracy of multi-read assignment by trimming experiments | |
| C. Impact on coverage | |
| D. Reproducibility across replicates: both raw contact count matrix and also identified contacts | |
| E. Biological impact: Novel promoter-enhancer interactions | |
| F. Biological impact: TAD inference | |

# B. Alternative read rescue

# B. Accuracy

A nonparametric alternative to mHi-C

Alignment

Uni-reads MAPQ>30

Aligner Select

Pairing

Distance Select

Validation
Duplicates Removal
Binning

Simple Select

mHi-C

Contact Matrix

**Trimming experiments:**

Start with long read datasets (e.g., ≥100 bp).

Align and get uni-reads (long uni-reads).

Trim the long uni-reads to generate short reads.

Align trimmed reads, some of which are now multi-reads.

Evaluate them against their true alignment positions from the longer uni-read set.

# B. Accuracy

# B. Accuracy

# B. Accuracy

# B. Recovering the full length contact matrix
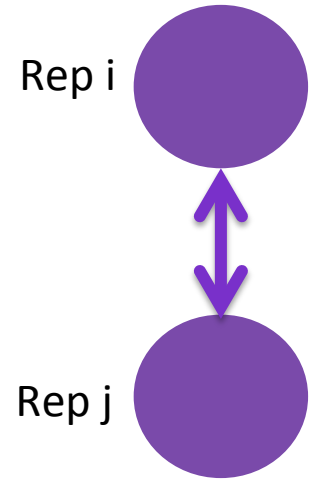
# C. Major improvement in coverage
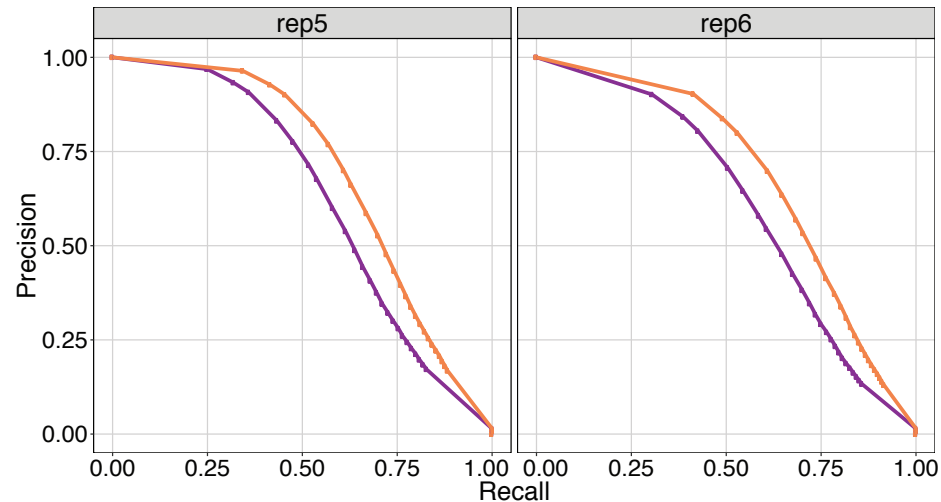
# D. Reproducibility of the contact matrix



IMR90

GM12878

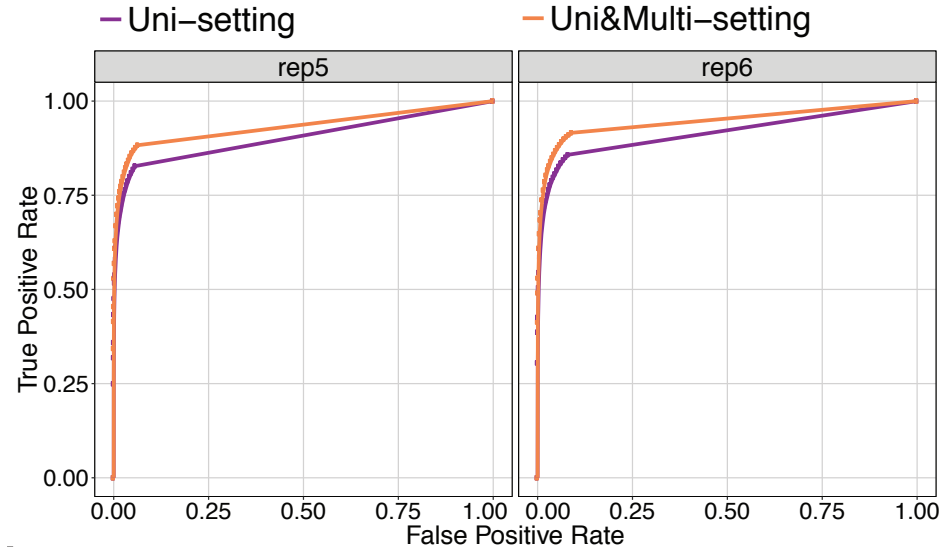Reproducibility of Contact Matrices

# D. Reproducibility of the significant interactions

Uni-reads

Uni-&Multi-reads

Rep i

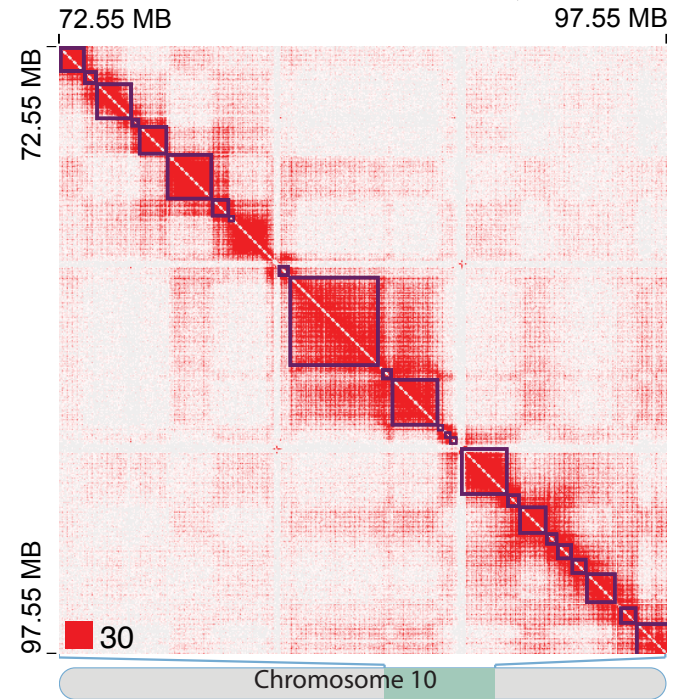Rep j

# D. ROC- and PR-based on replicate gold standard
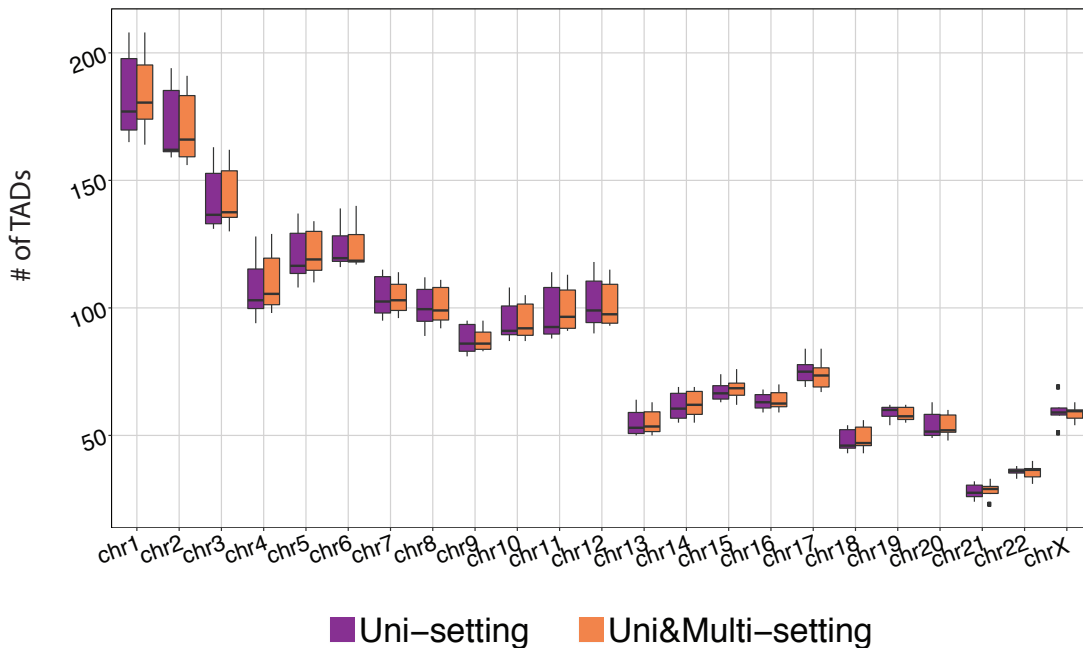


High depth replicates are used to define "true" positives and negatives.
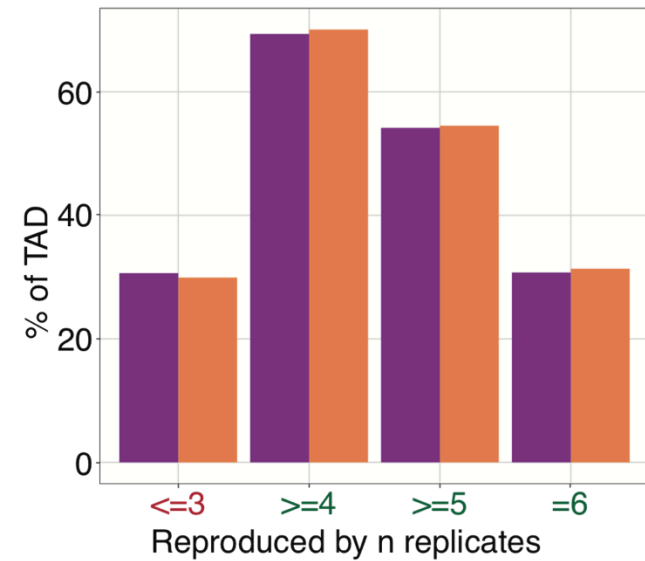
# E. Impact on TAD inference

TAD: Topologically associated domain

# of TADs detected do not change significantly.



Uni–setting   Uni&Multi–setting



Chromosome 10

30

11.18

# E. Impact on TAD inference

# of TADs detected do not change significantly.



Uni−setting    Uni&Multi−setting



# of reproducible TADs increases by 2.01%.

# of irreproducible TADs decreases by 2.36%.
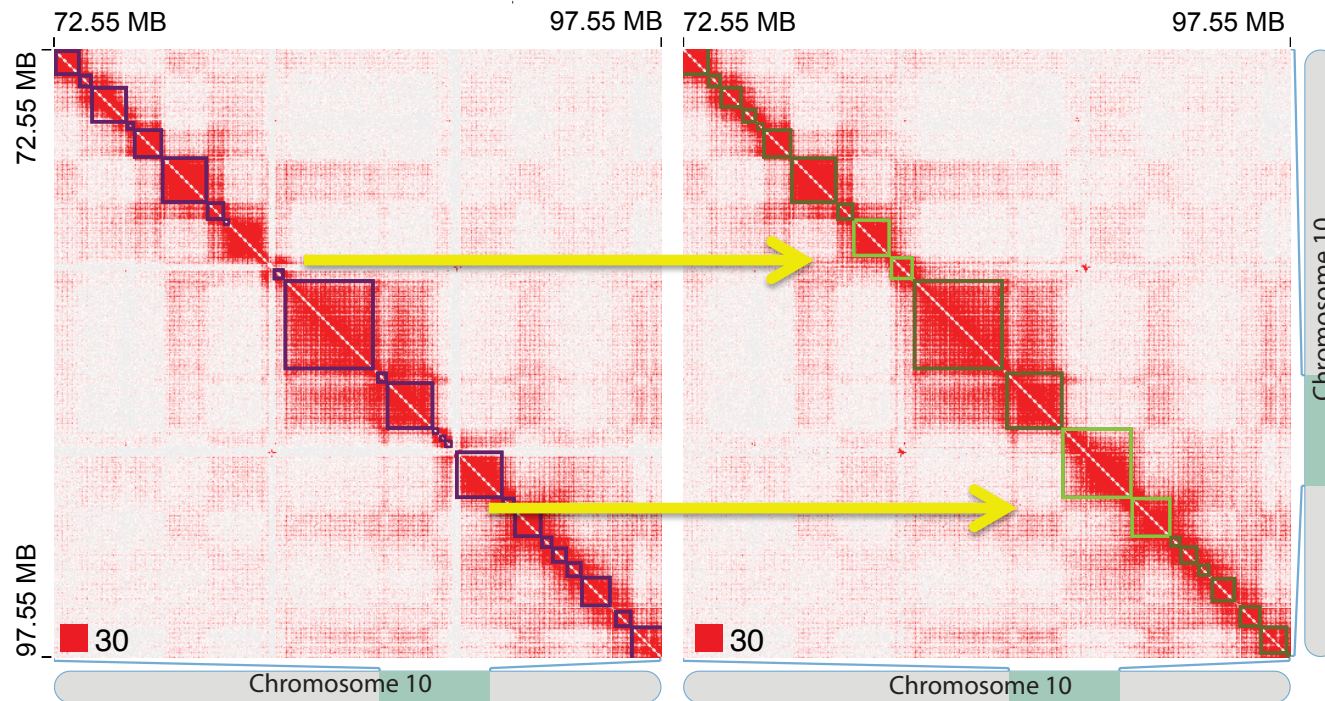
# E. Impact on TAD inference



Uni-setting

Uni&Multi-setting
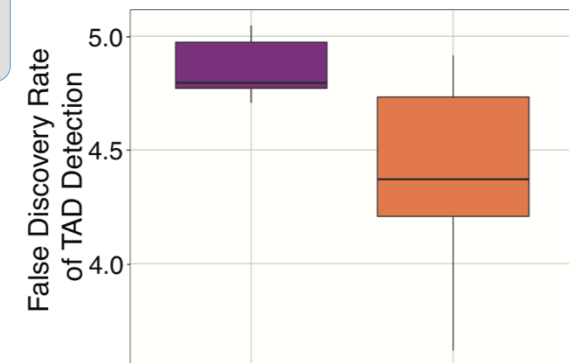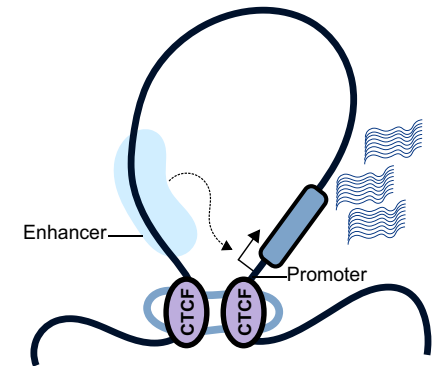
11.18

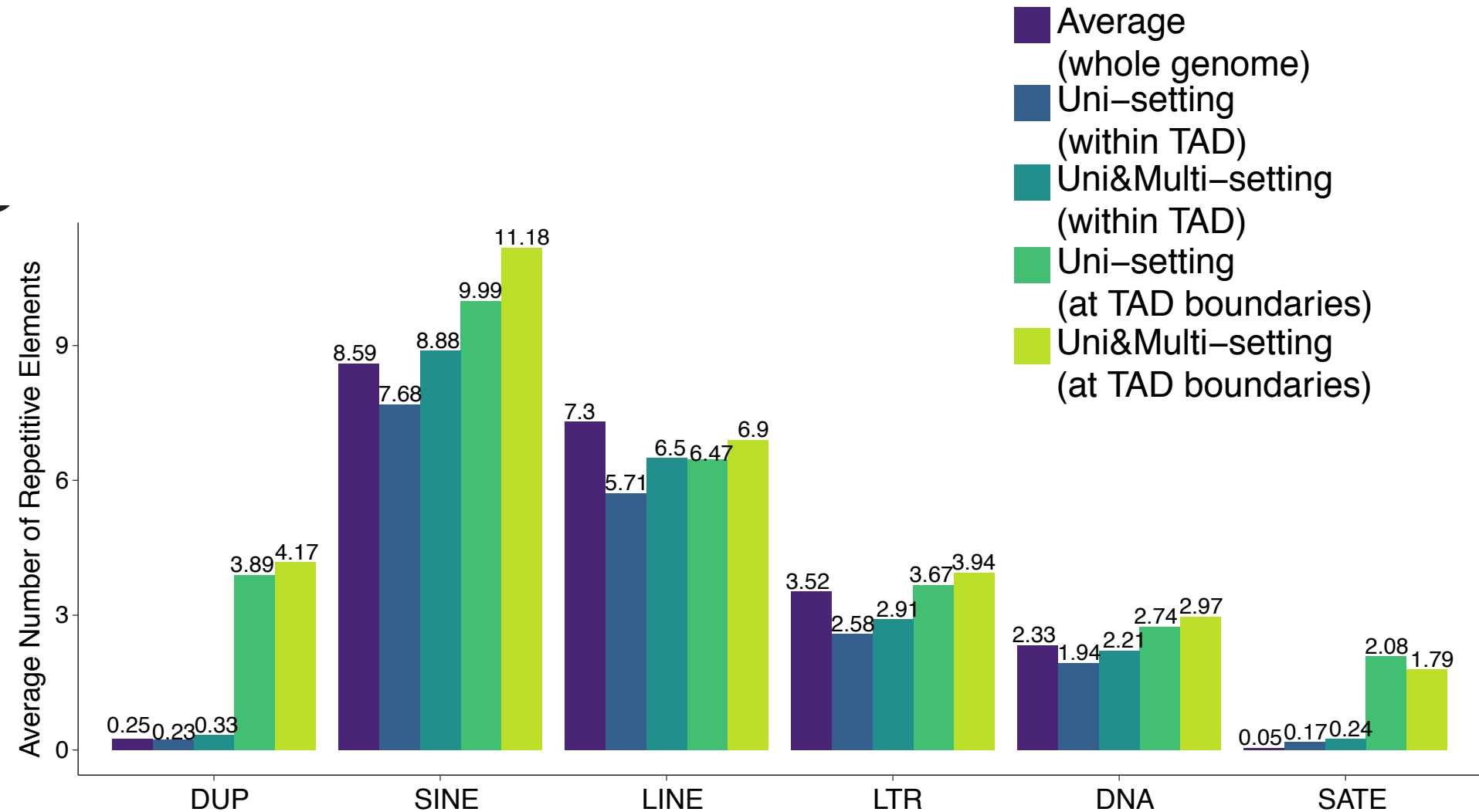# E. Impact on TAD inference



Uni-setting

Uni&Multi-setting

Arzate-Mejia *et al.*, 2018

72.55 MB     97.55 MB     72.55 MB     97.55 MB

Chromosome 10

Enhancer — CTCF CTCF — Promoter

False Discovery Rate of TAD Detection

■ Uni–setting    ■ Uni&Multi–setting

TADs that are not reproducible and lack CTCF peaks at the TAD boundaries are labeled as false positives

11.18

# F. Repetitive elements at the boundaries of reproducible TADs

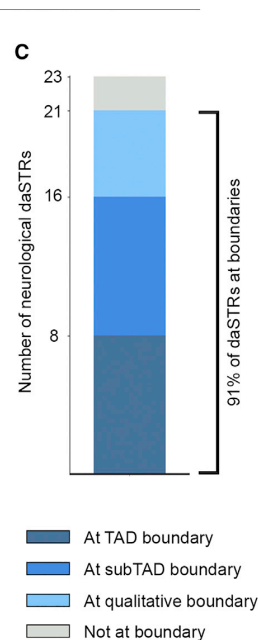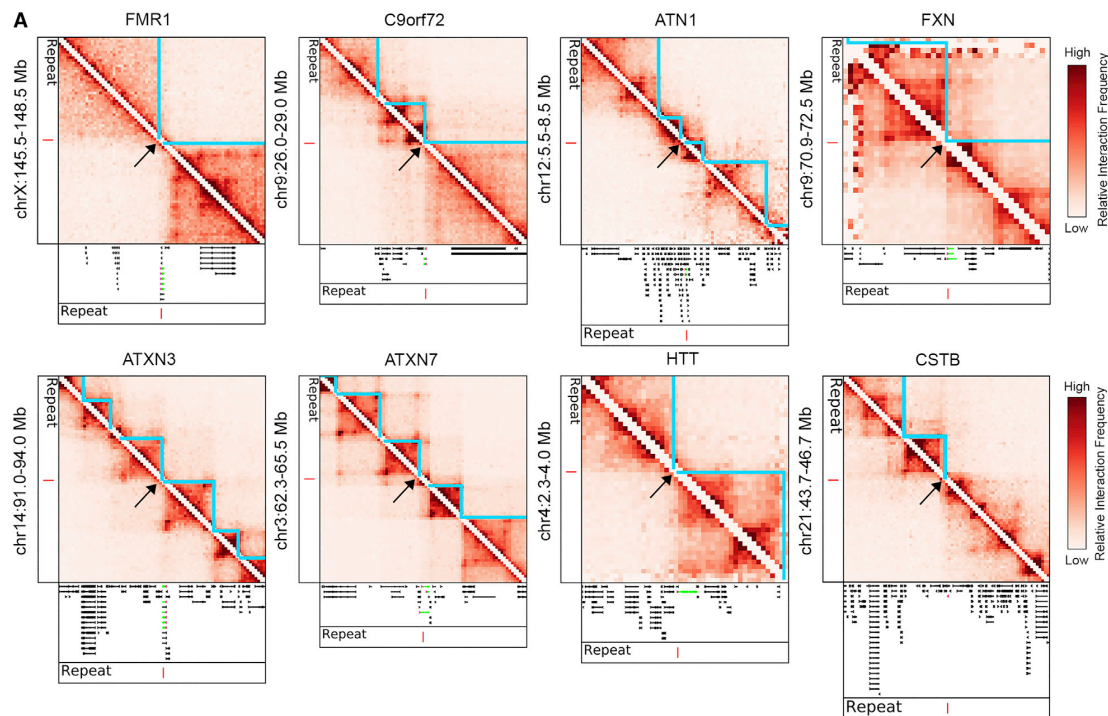# F. Disease-Associated short tandem repeats co-localize with domain boundaries

**Disease-Associated Short Tandem Repeats Co-localize with Chromatin Domain Boundaries**
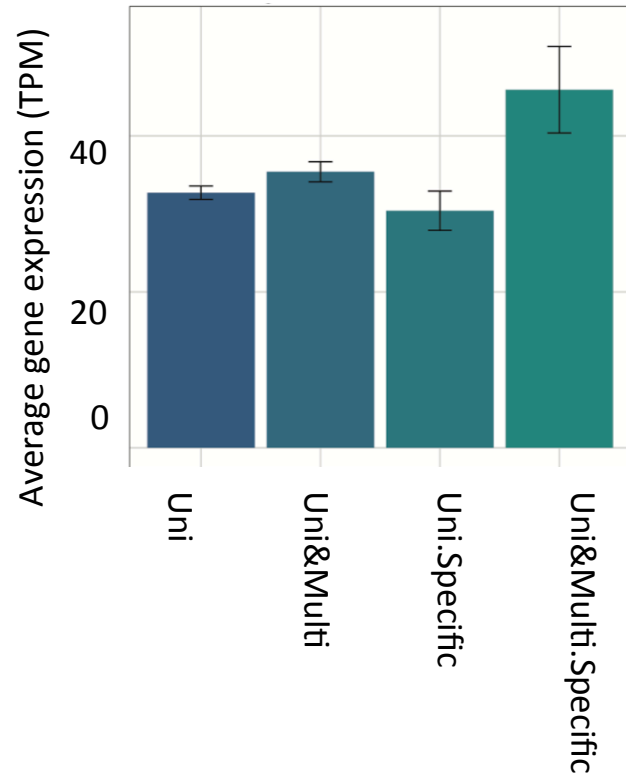
Graphical Abstract

Authors
James H. Sun, Linda Zhou,
Daniel J. Emerson, ...,
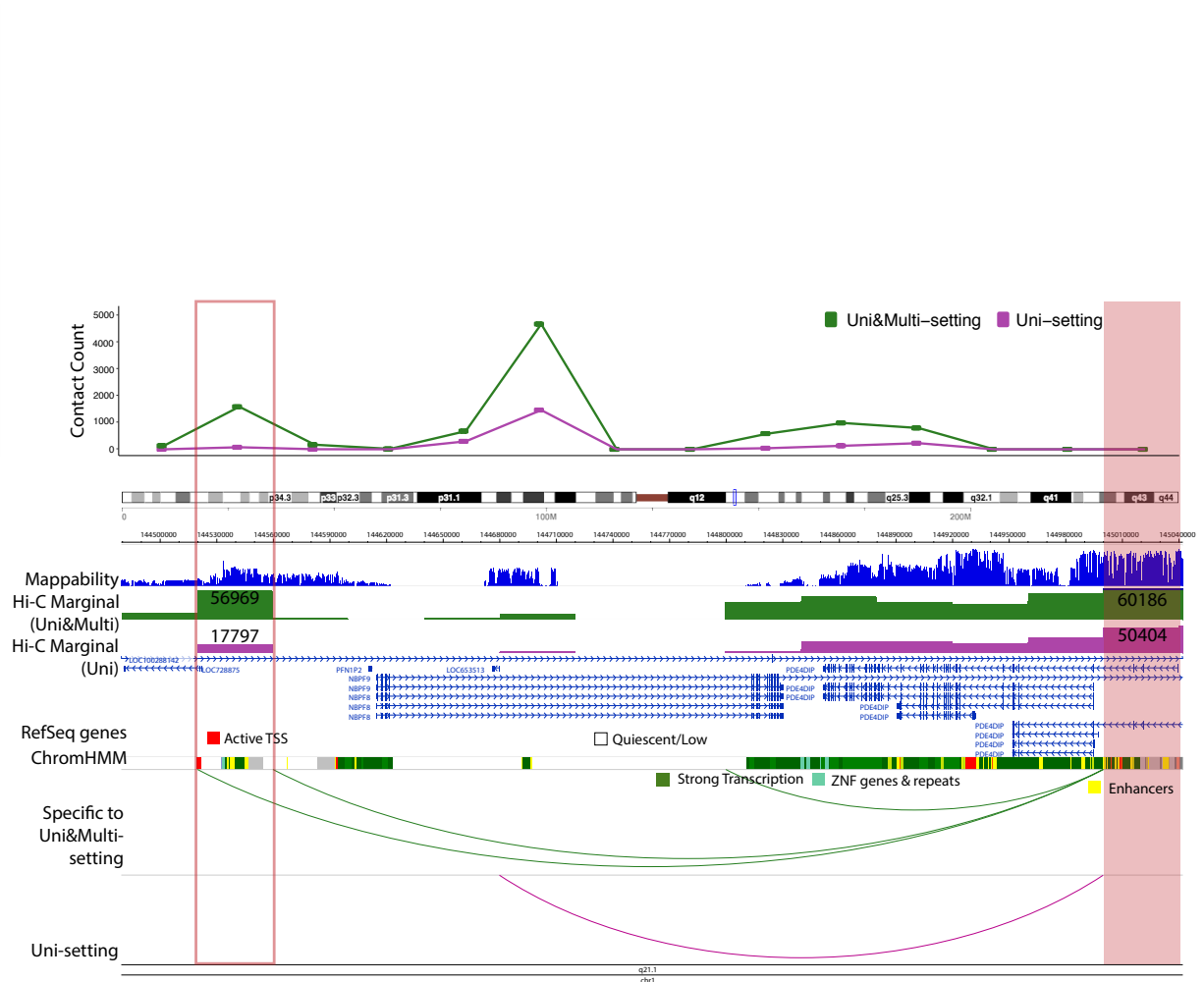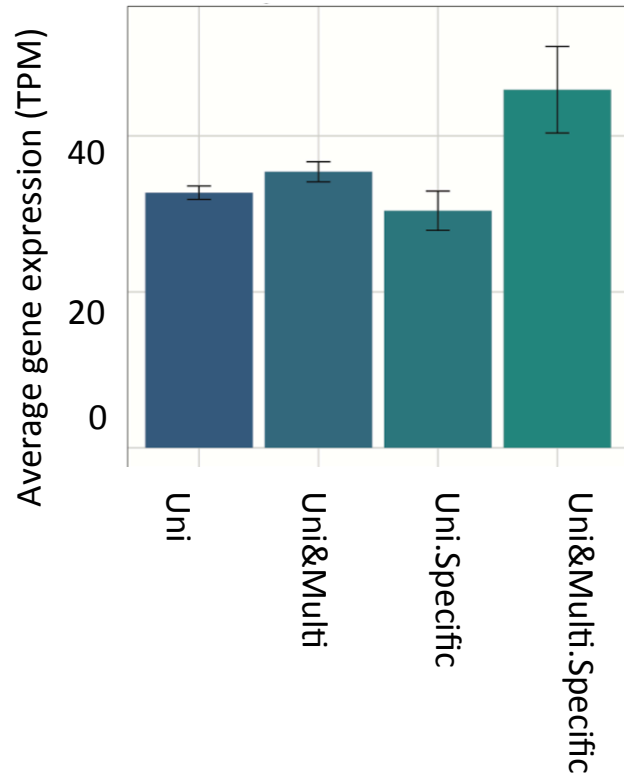Beverly L. Davidson, Flora Tassone,
Jennifer E. Phillips-Cremins

Healthy Individuals

# F. Novel promoter-enhancer interactions

**15.8% more promoter-enhancer** interactions that are reproducible in at least 2 replicates.

# F. Novel promoter-enhancer interactions

**15.8% more promoter-enhancer** interactions that are reproducible in at least 2 replicates.

# Summary

- Software

  https://github.com/keleslab/mhic

- Paper

  https://www.biorxiv.org/content/early/2018/10/03/301705

- More results on chimeric reads, impact on differential Hi-C analysis are available in the manuscript.

# Acknowledgements

**Keleș Group**

Ye Zheng



**Collaborators**

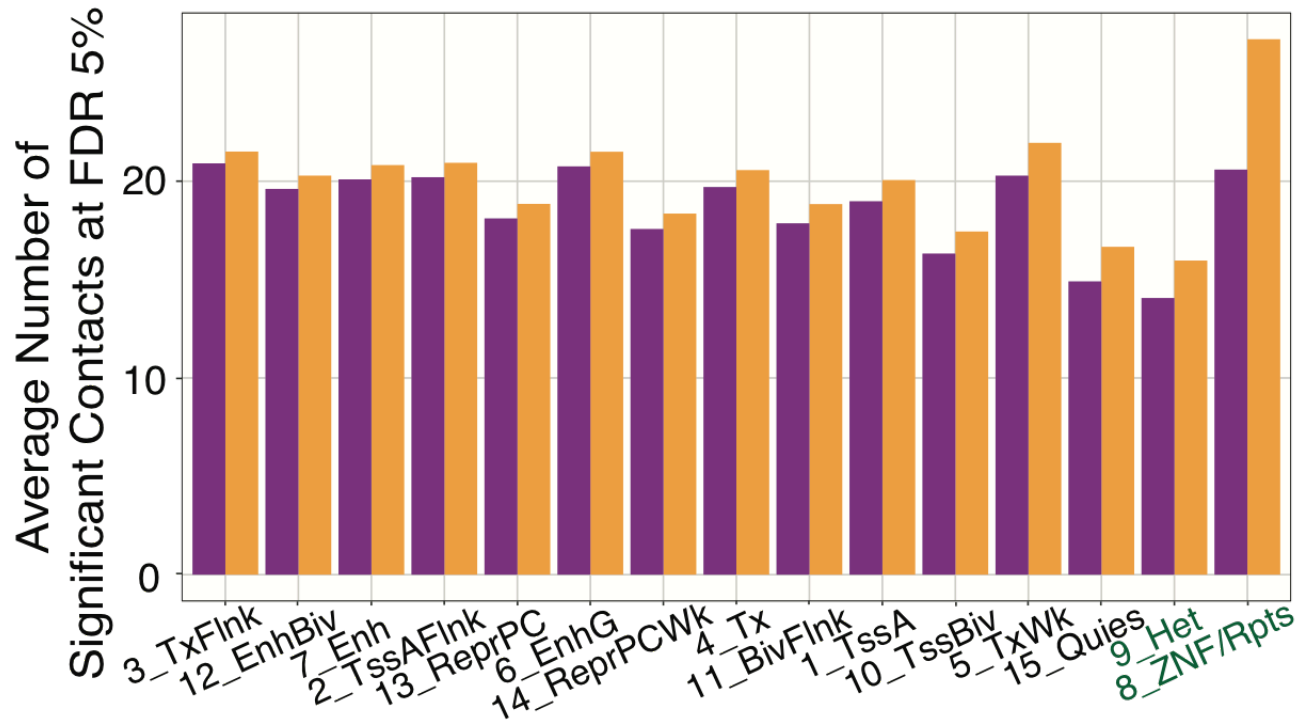Ferhat Ay (La Jolla Institute for Allergy & Immunology)

# Available positions



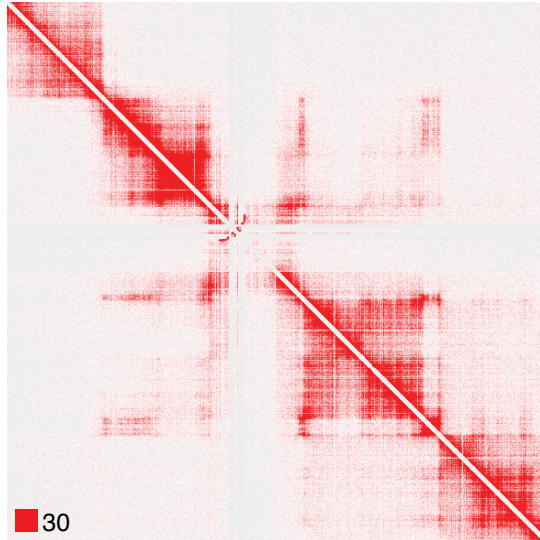1-2 postdoctoral researcher positions in statistical genomics.
If interested. send CV to keles@stat.wisc.edu
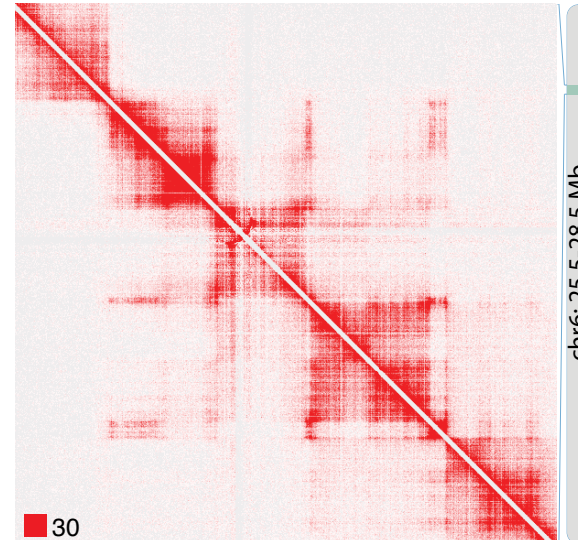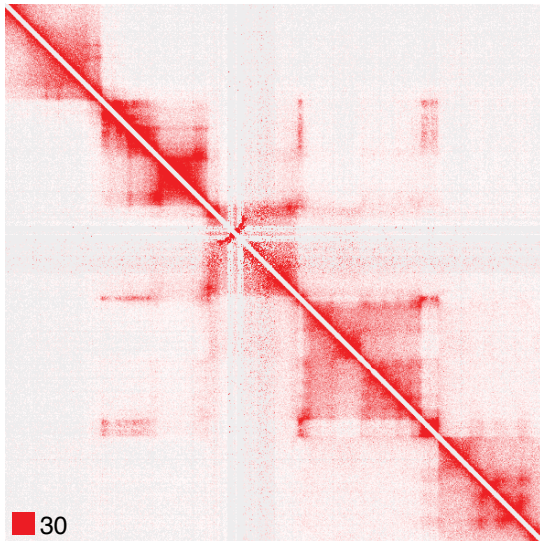
# F. Genomic characteristics
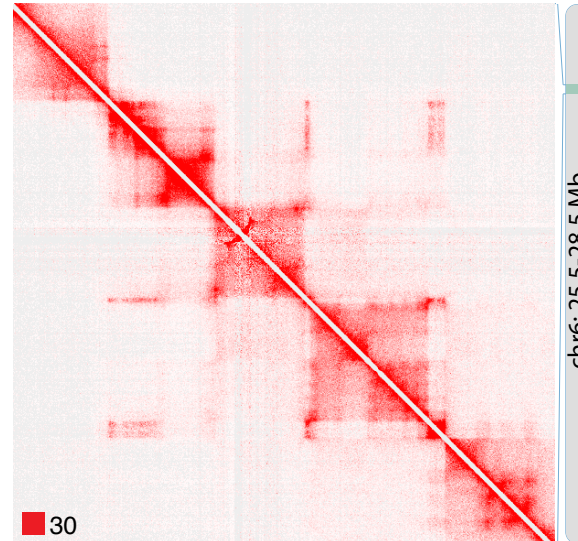
# C. Count matrices



Uni-setting (Raw Counts) — chr6: 25.5-28.5 Mb

Uni&Multi-setting (Raw Counts) — chr6: 25.5-28.5 Mb

Uni-setting (Normalized Counts)

Uni&Multi-setting (Normalized Counts) — chr6: 25.5-28.5 Mb

# POSTER
# SNPs in high LD: a formidable challenge



Labeling by
Massively parallel reporter assays (MPRA)