

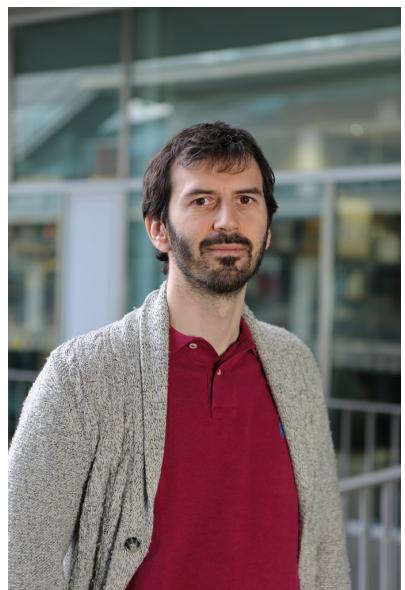
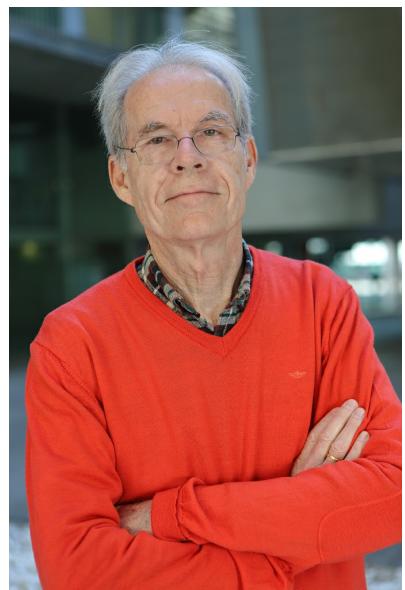
Interrogating genome structure

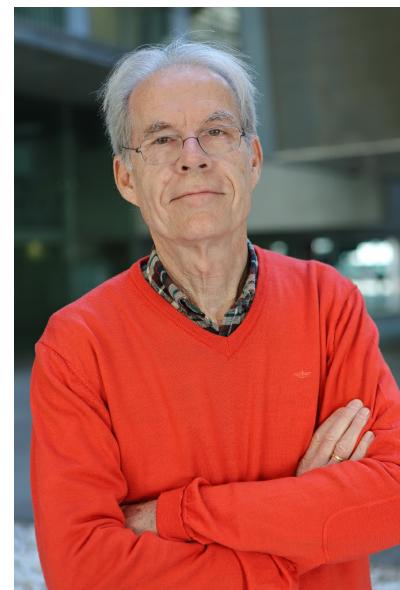
Reproducibility and biases

Münster, 12 July 2018

Enrique (Quique) Vidal
 enrique.vidal@crg.eu

 @qenvio  qenvio





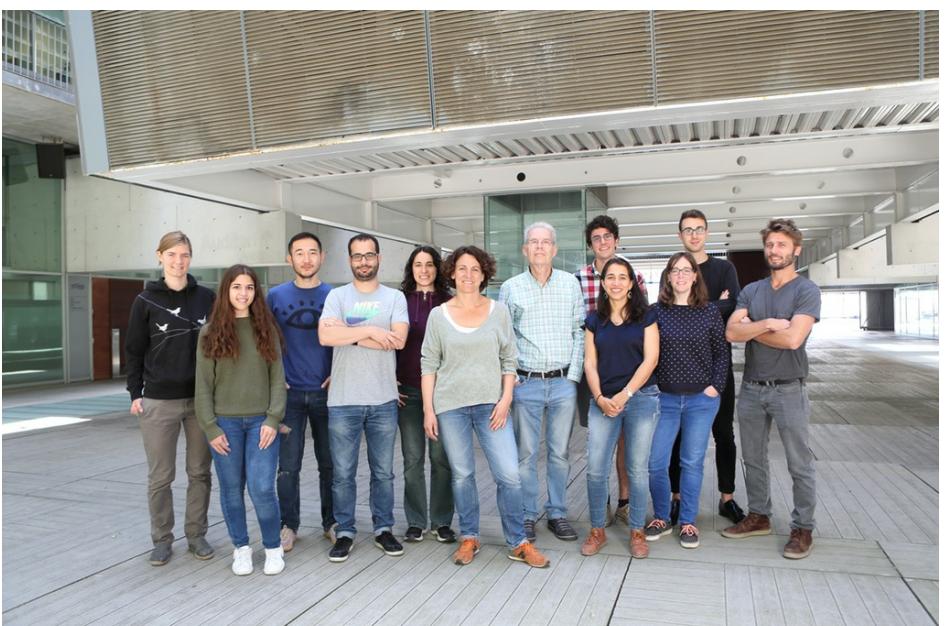
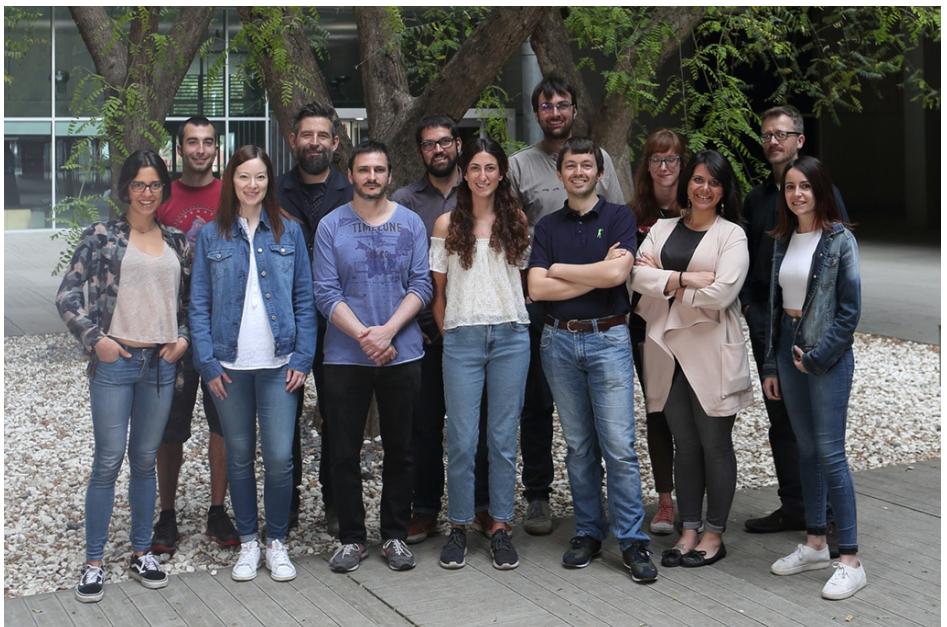
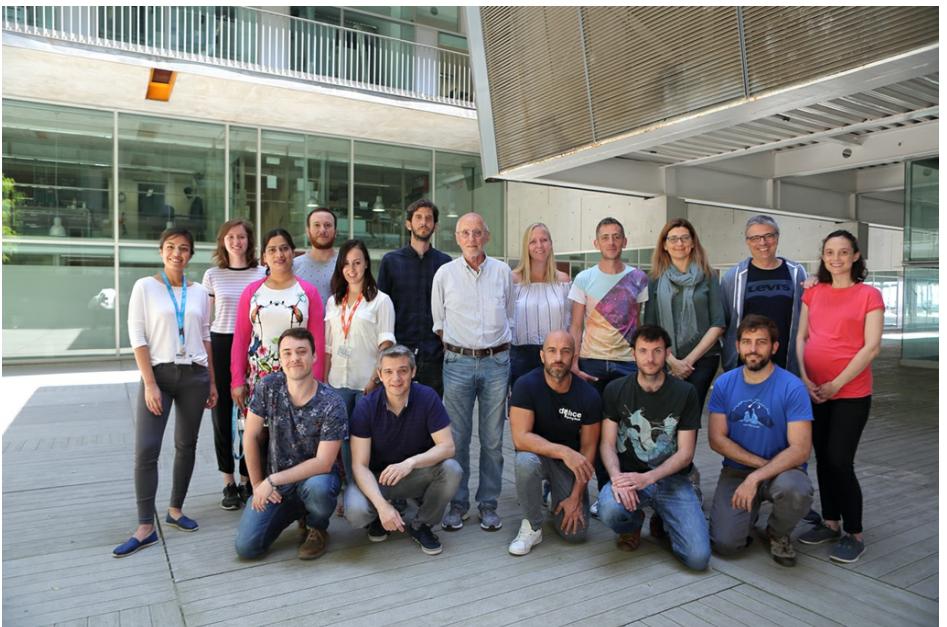
4DGenome

erc European Research Council

CRG Centre for Genomic Regulation

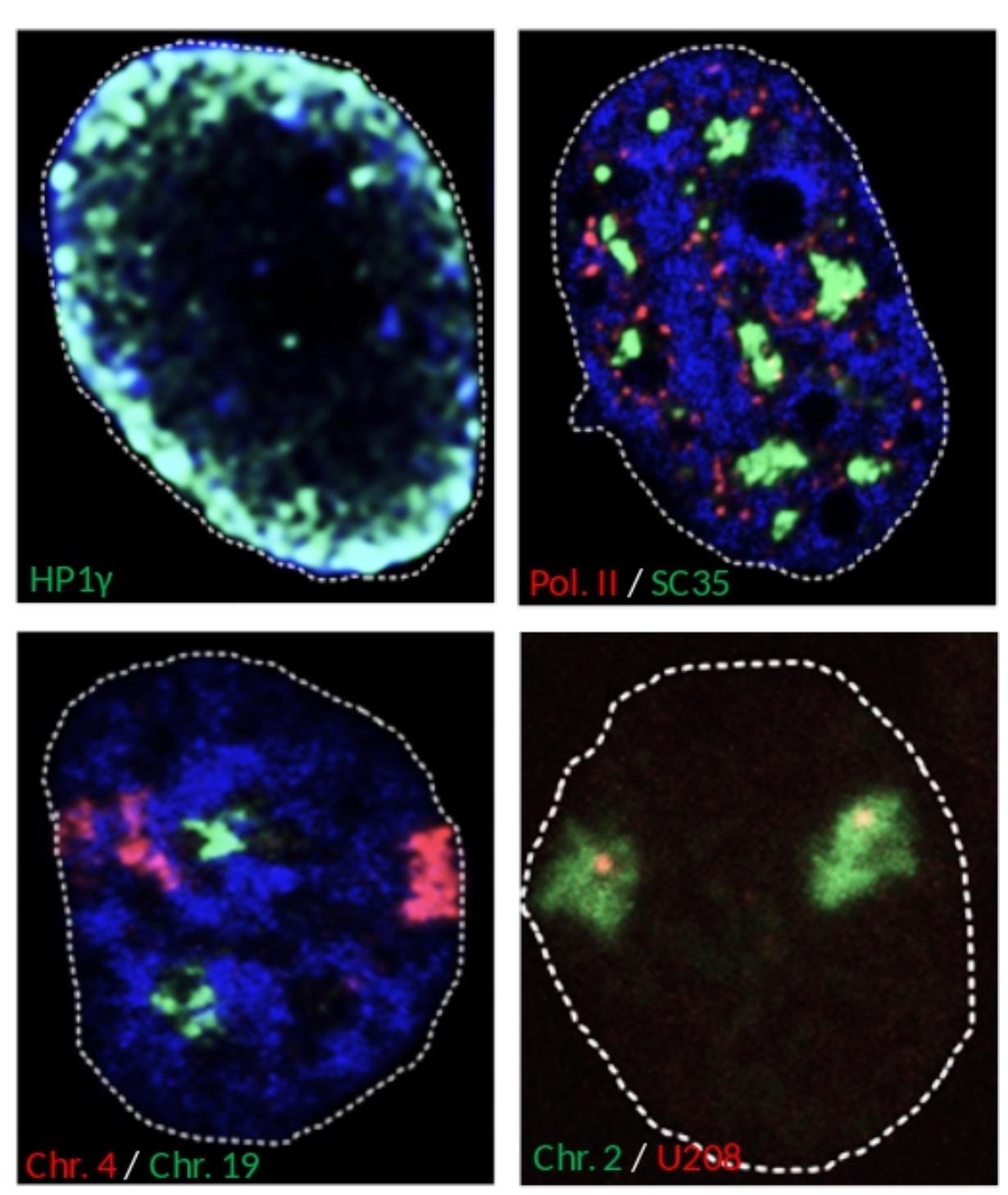
cnag centre nacional d'anàlisi genòmica
centro nacional de análisis genómico

4DGenome



Motivation

Genome structure is not random



Chromosome conformation capture

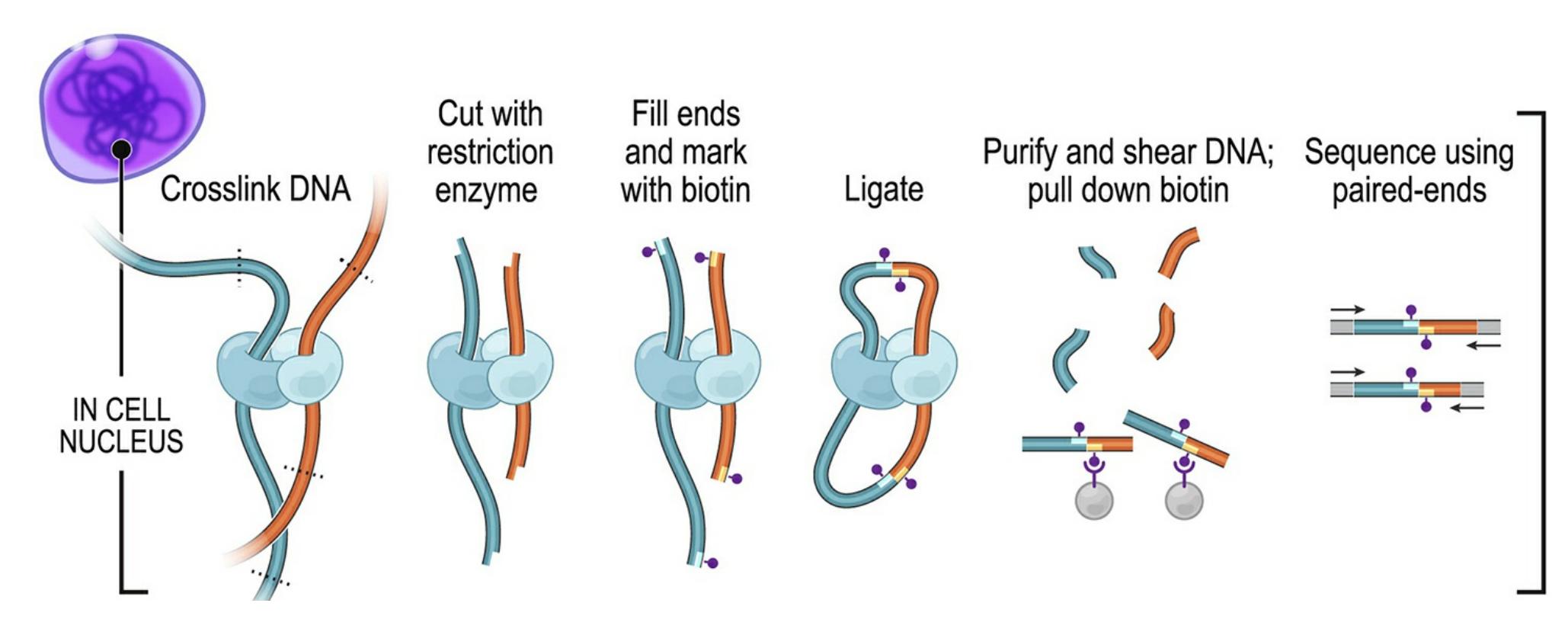
Dekker, J. et al. (2002)

Science

Chromosome conformation capture

Dekker, J. et al. (2002)

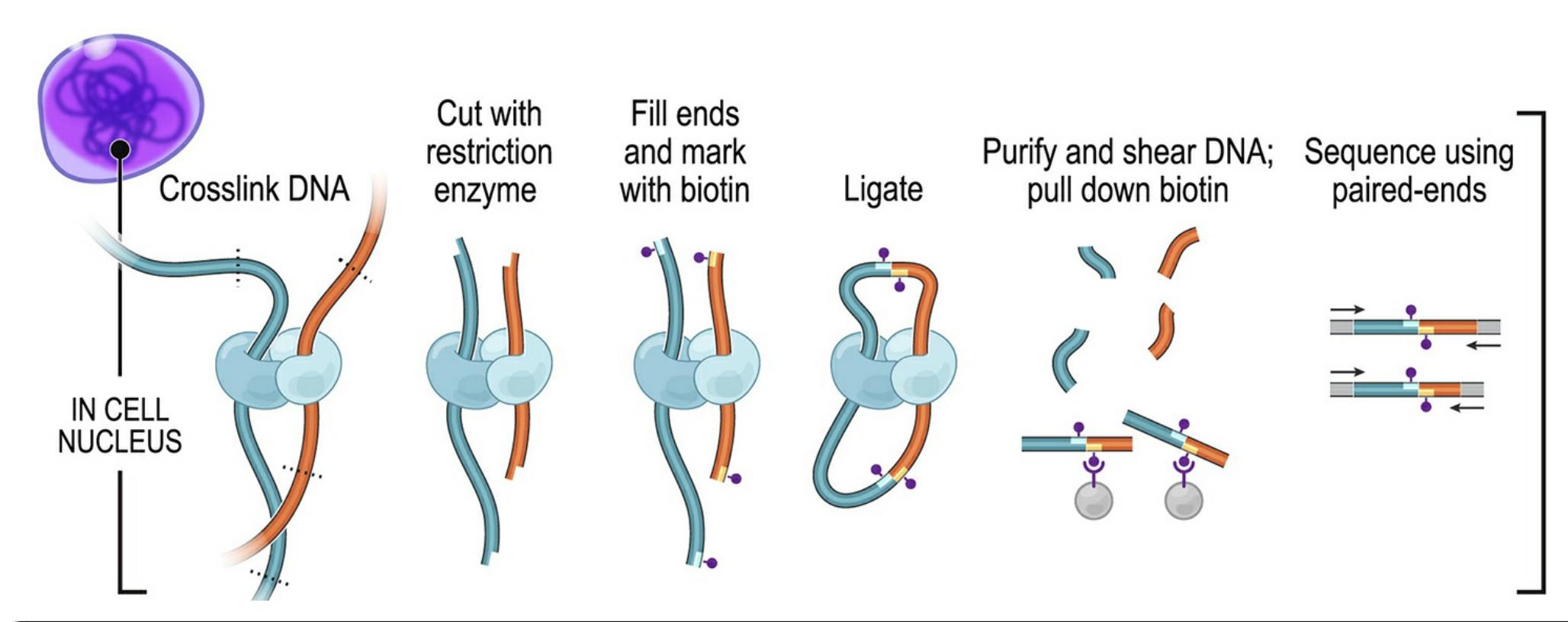
Science



Chromosome conformation capture

Dekker, J. et al. (2002)

Science

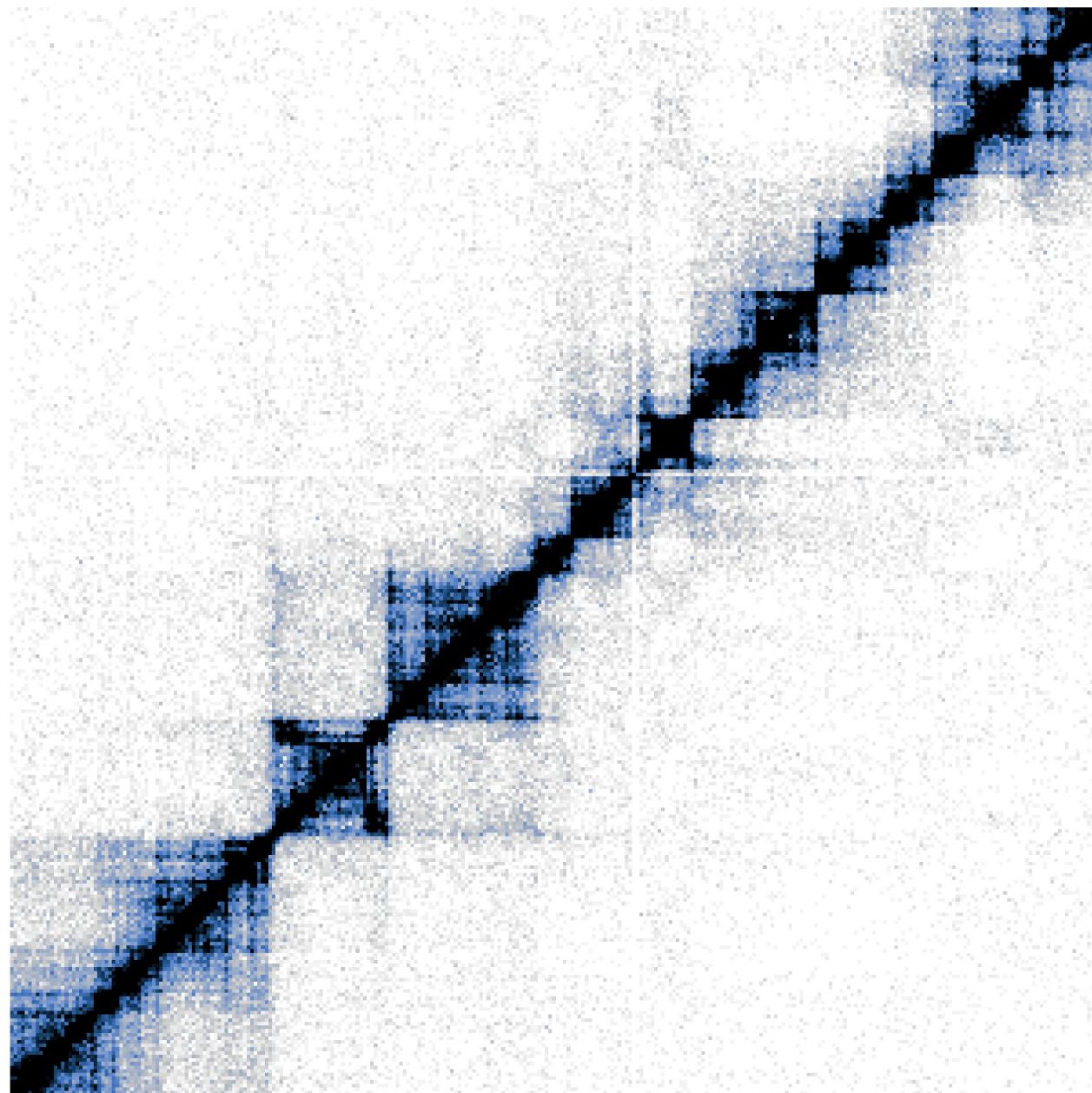


Lieberman-Aiden, E. et al. (2009)

Science

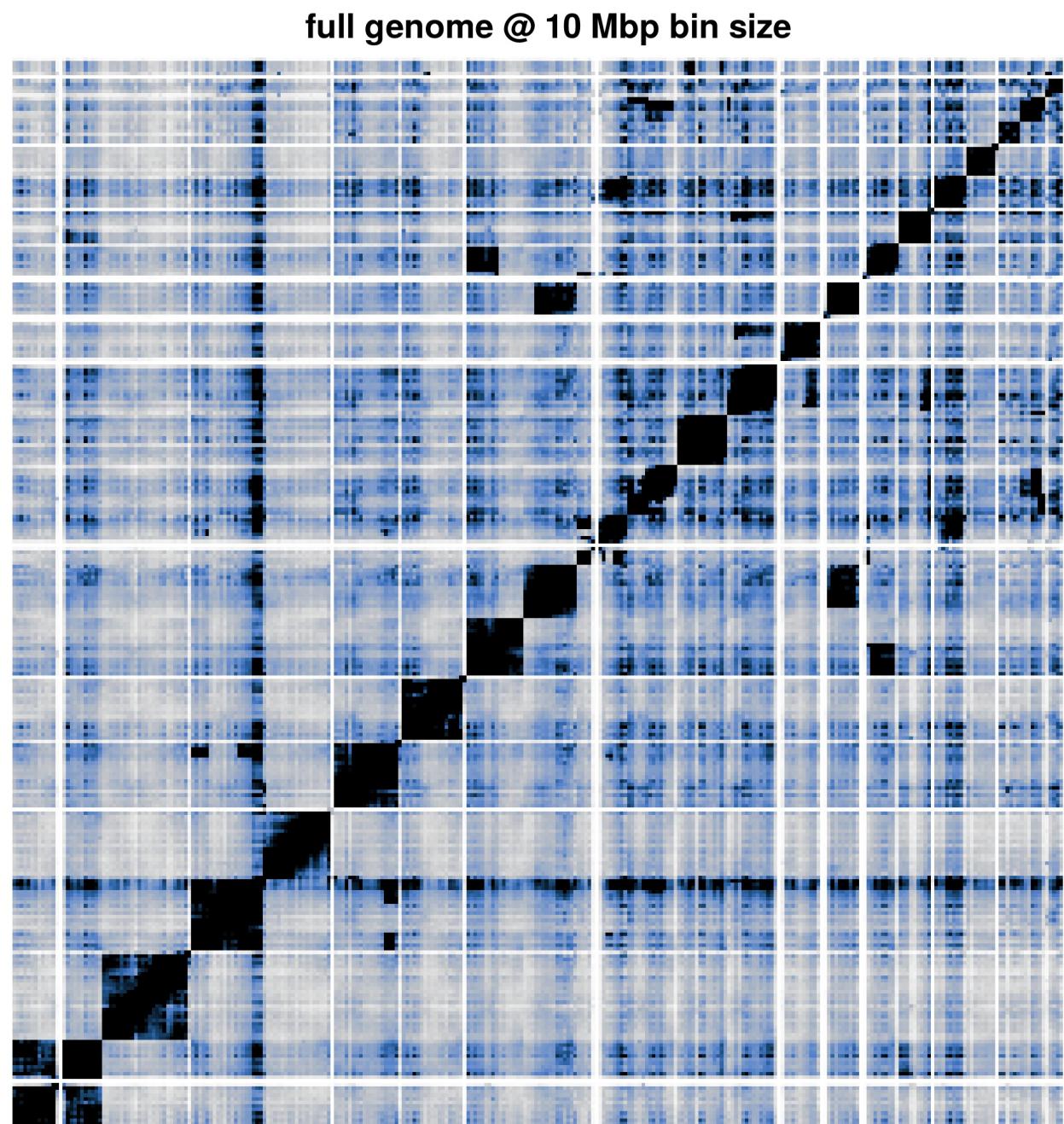
Hi-C

16.5 Mbp region @ 50 Kpb bin size



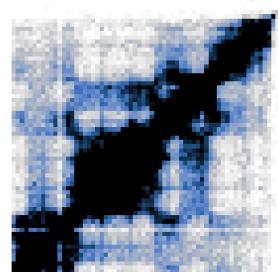
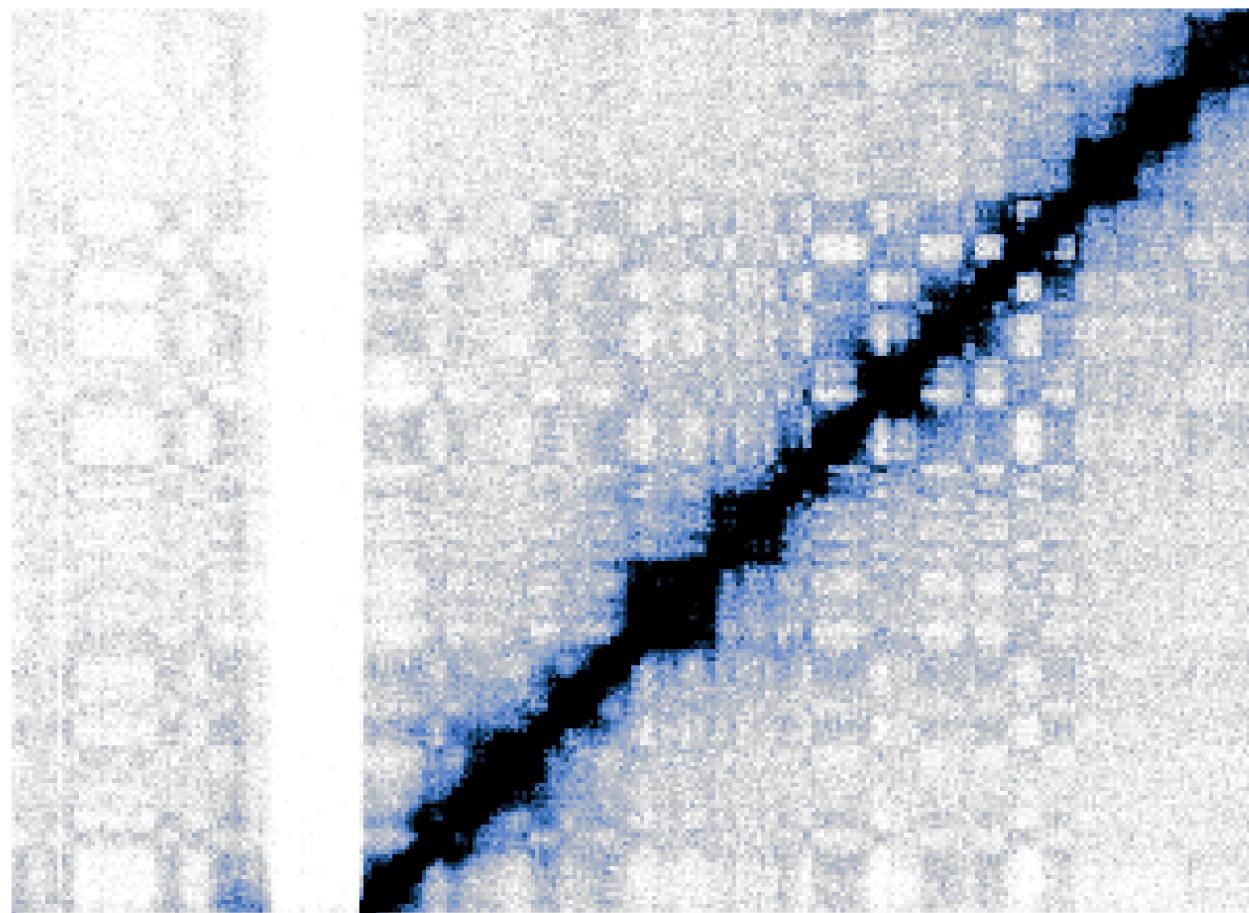
What can we get from Hi-C?

Territories



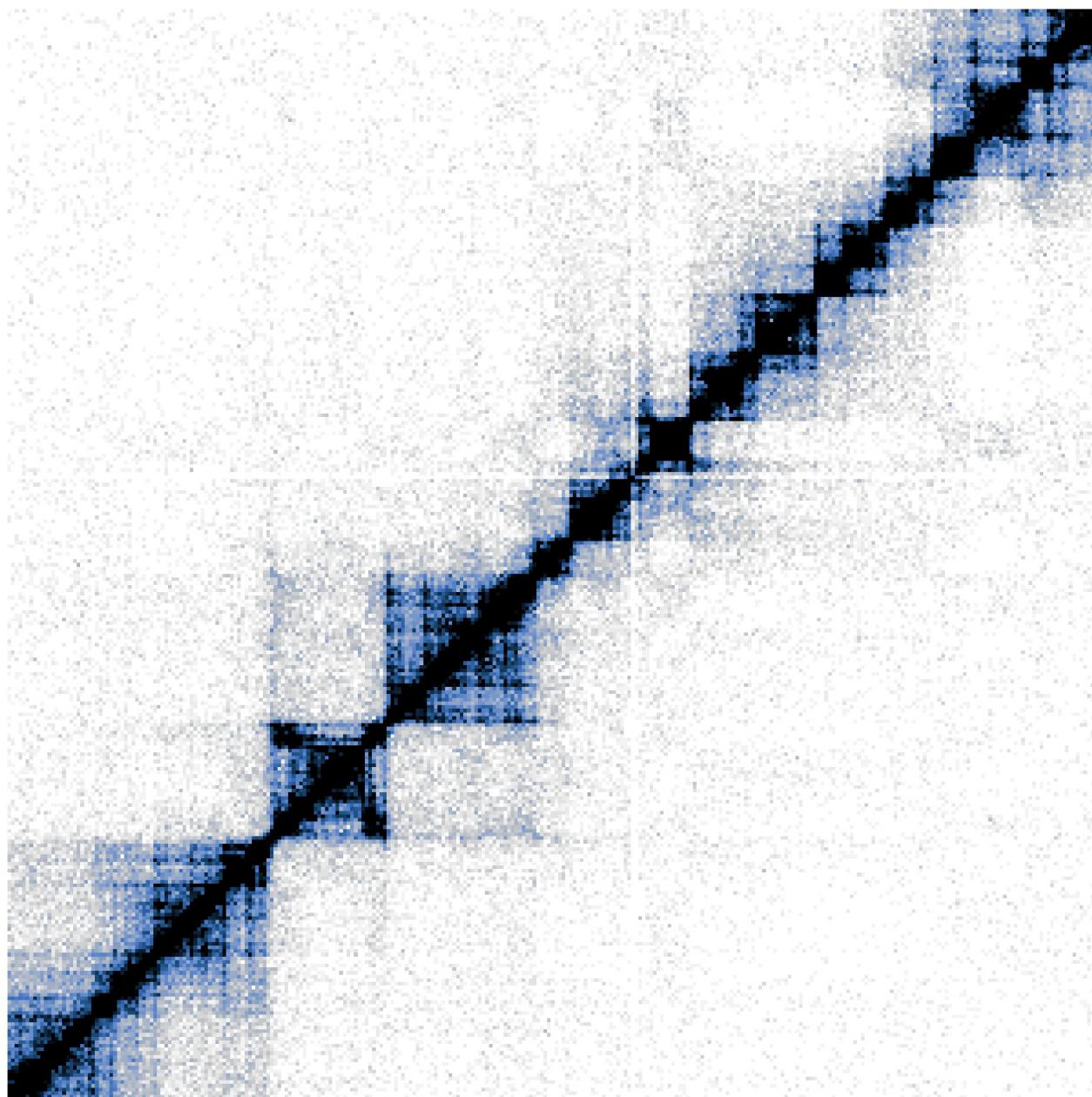
Compartments

chromosome 17 @ 250 Kbp bin size



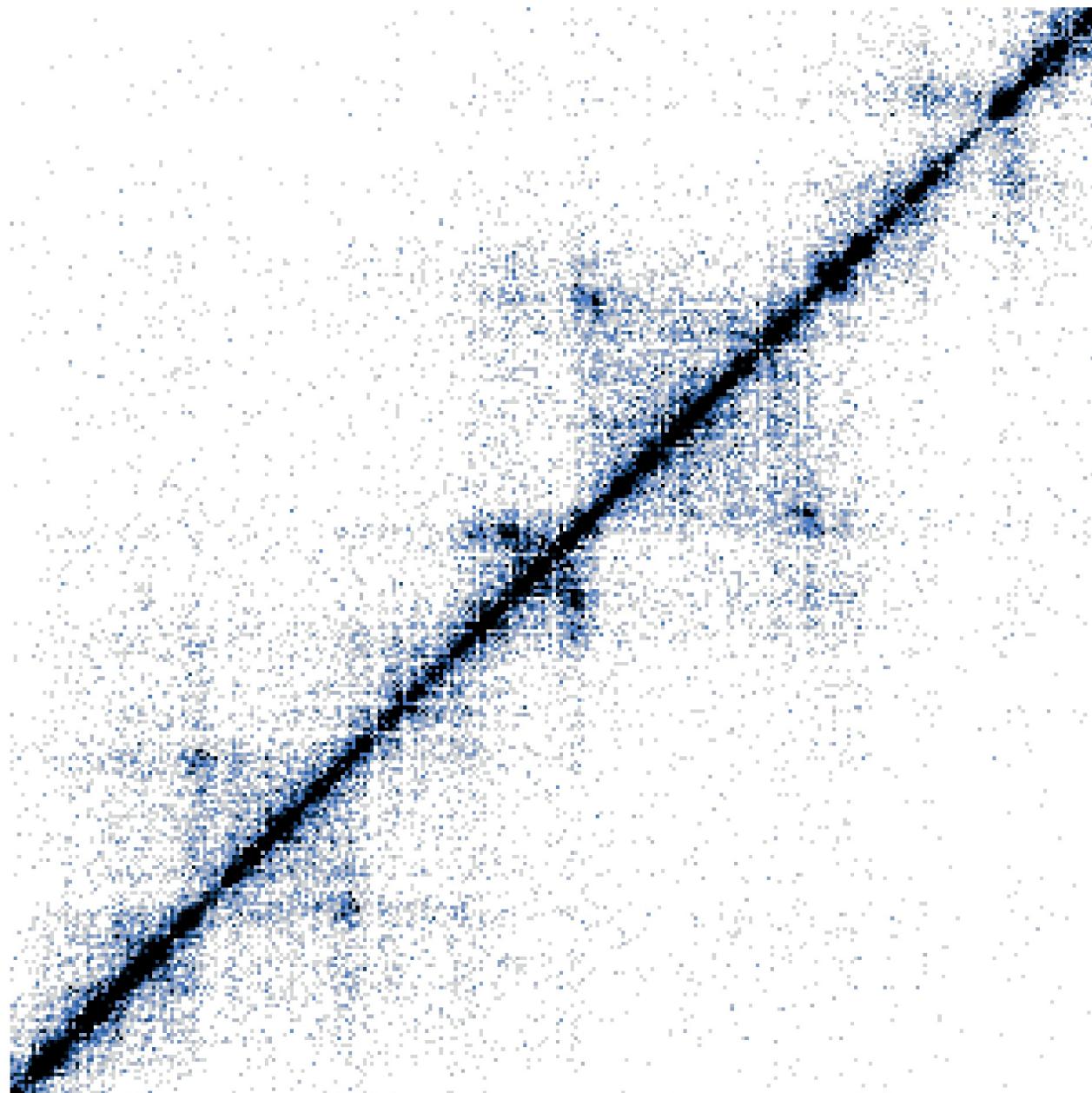
Domains (TADs)

16.5 Mbp region @ 50 Kpb bin size

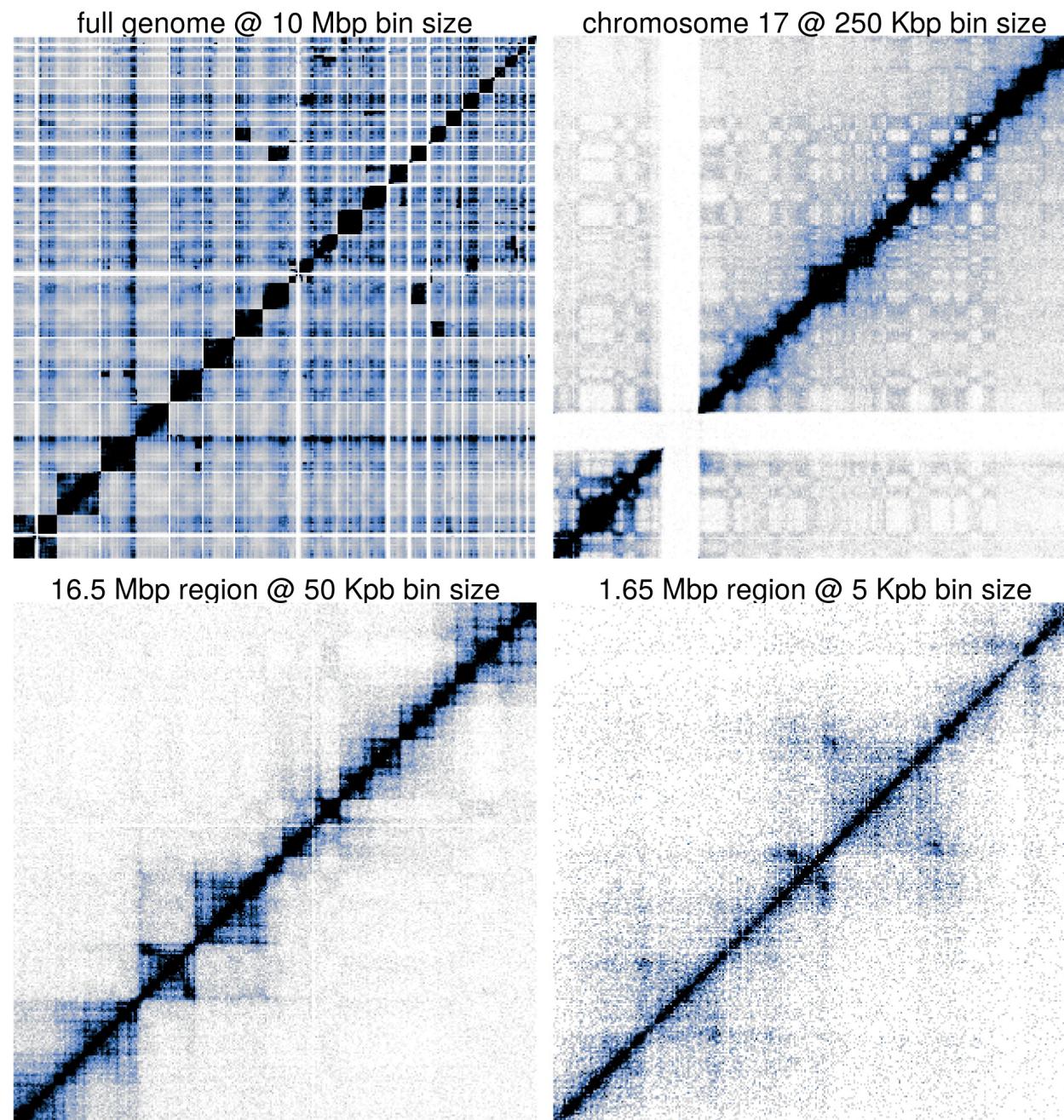


Loops

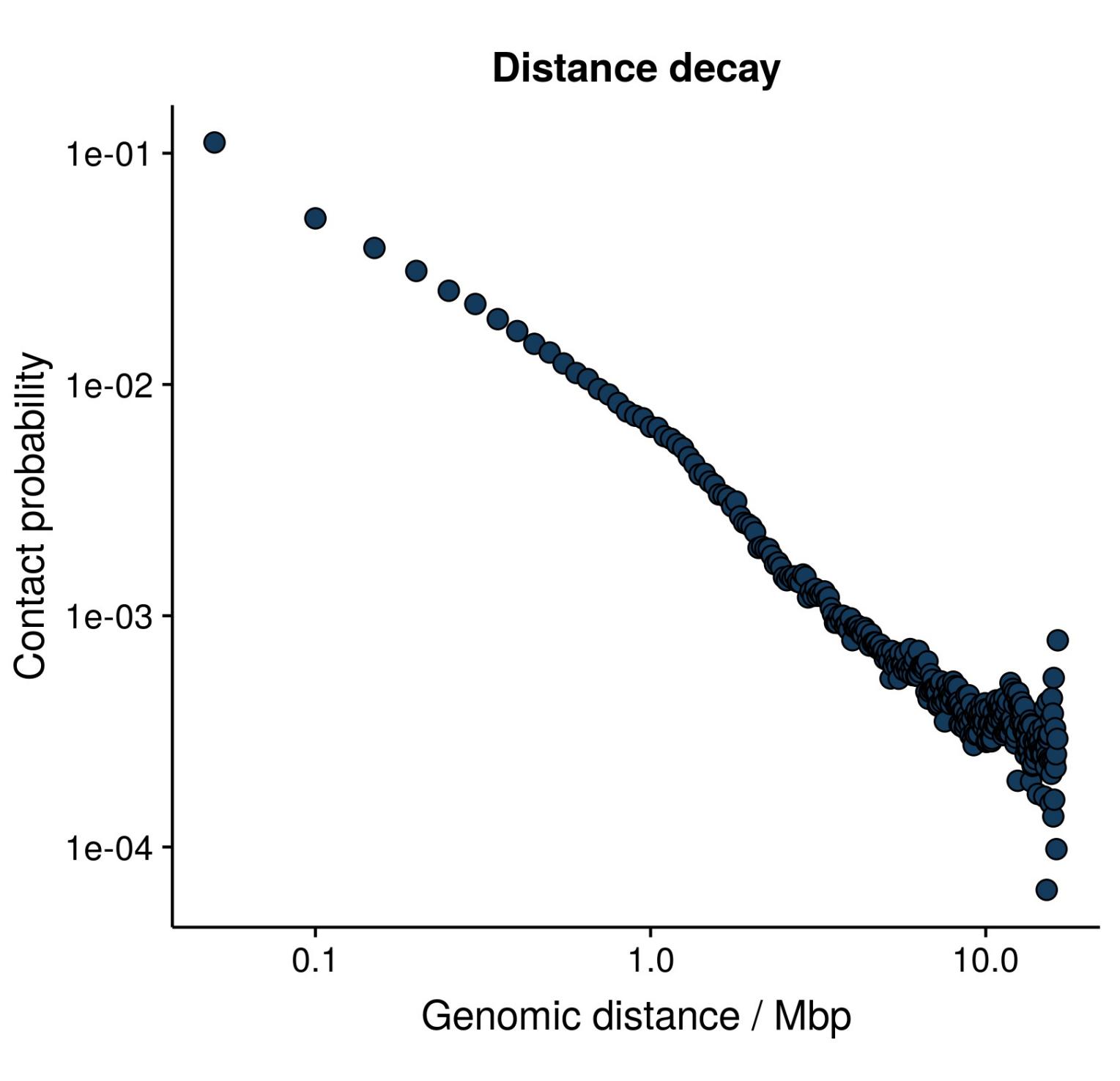
1.65 Mbp region @ 5 Kpb bin size



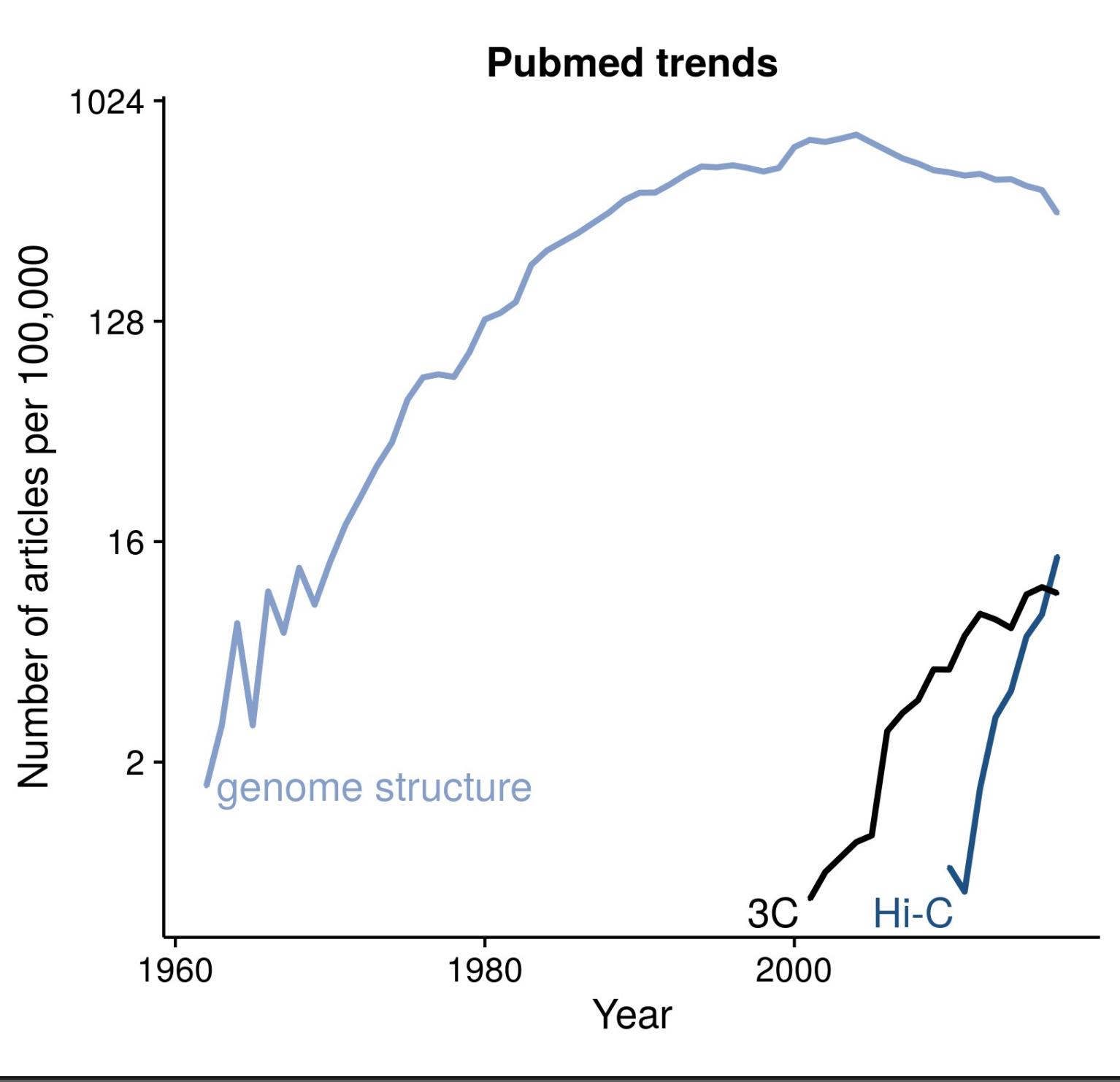
In a nutshell



Distance decay



Papers



Papers

Lieberman-Aiden, E. et al. (2009)

Science

Yaffe, E. and Tanay, A. (2011)

Nature genetics

Imakaev, M. et al. (2012)

Nature methods

Dixon, J. R. et al (2012)

Nature

Nora, E. P. et al. (2012)

Nature

Rao, S.S. et al. (2014)

Cell

Hug, B. H. (2017)

Cell

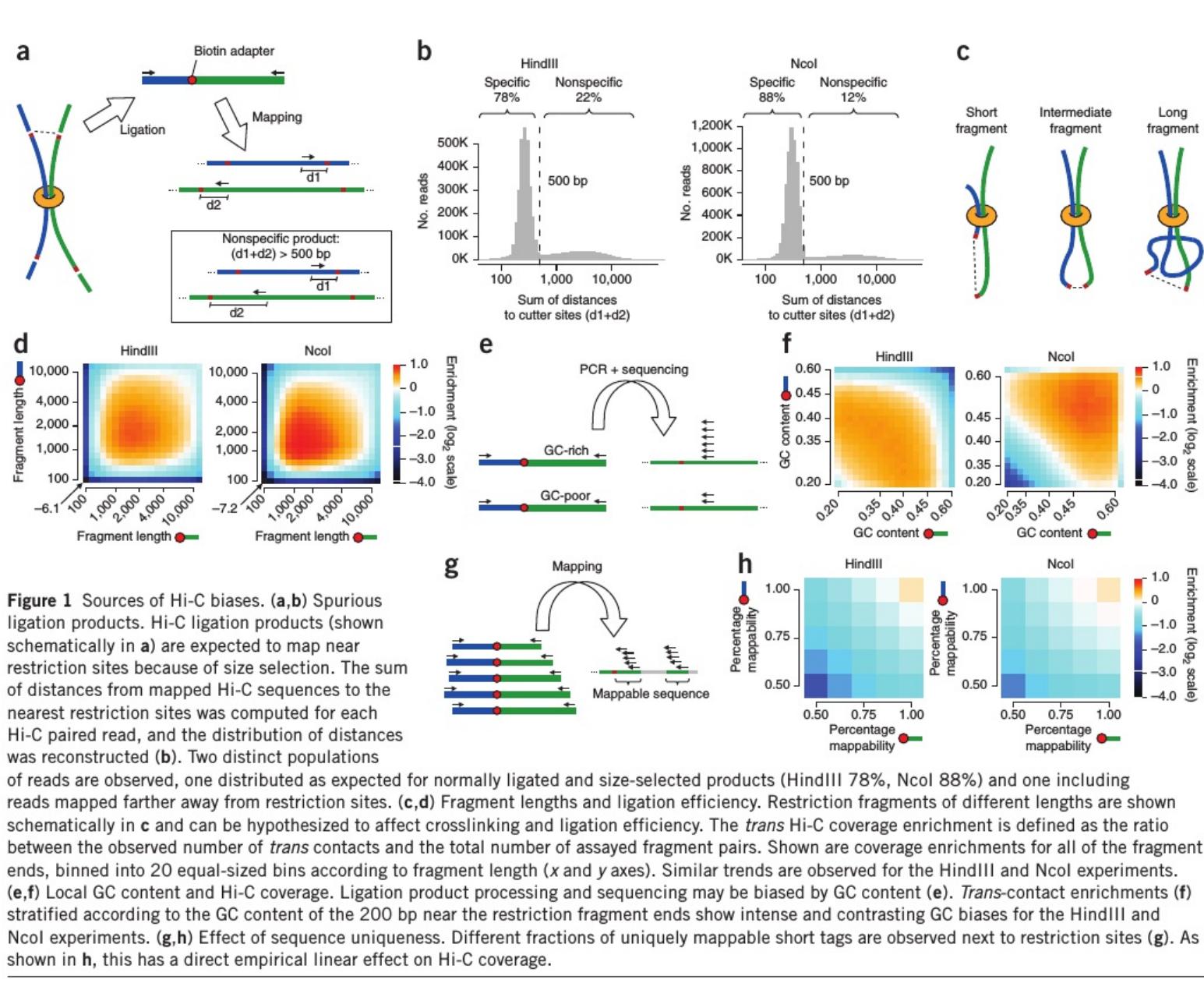
Stadhouders, R. and Vidal, E. et al (2018)

Naure genetics

Challenge 1

Systematic biases

Genomic features affect Hi-C



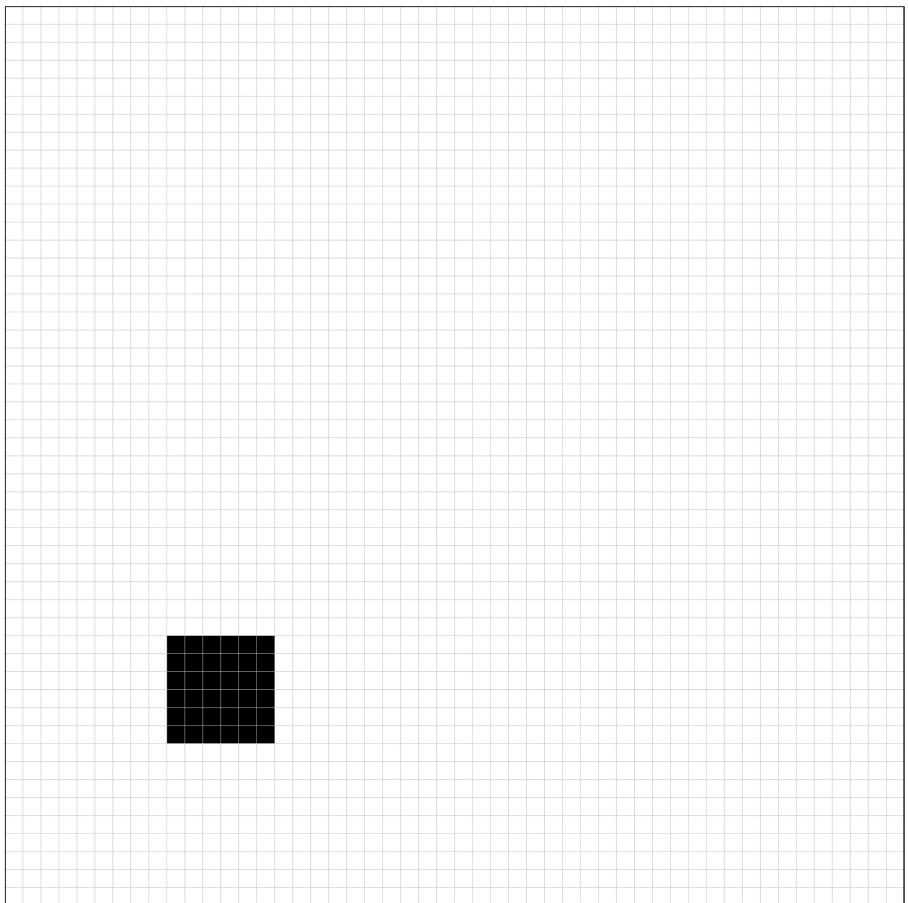
Yaffe, E. and Tanay, A. (2011)

Nature genetics

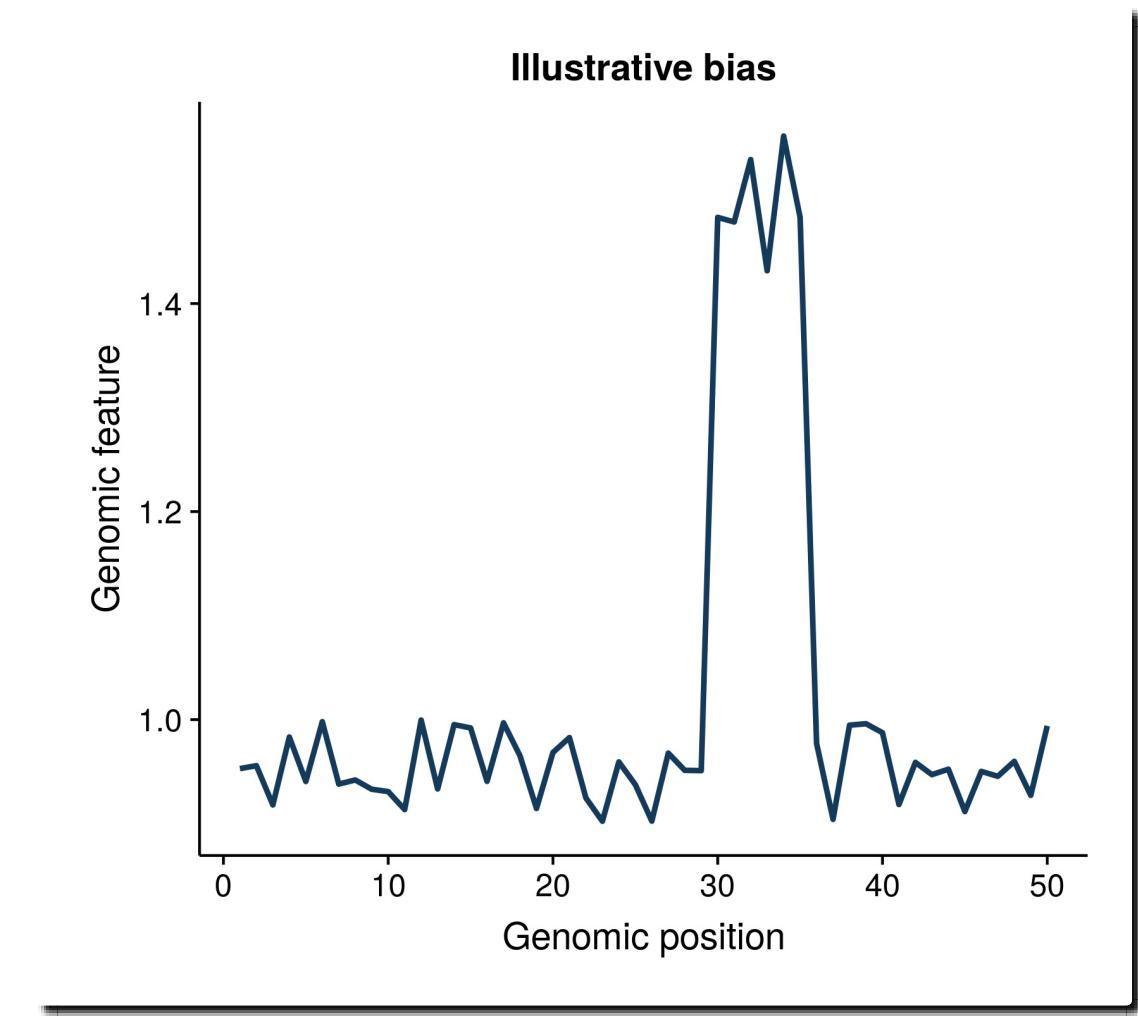
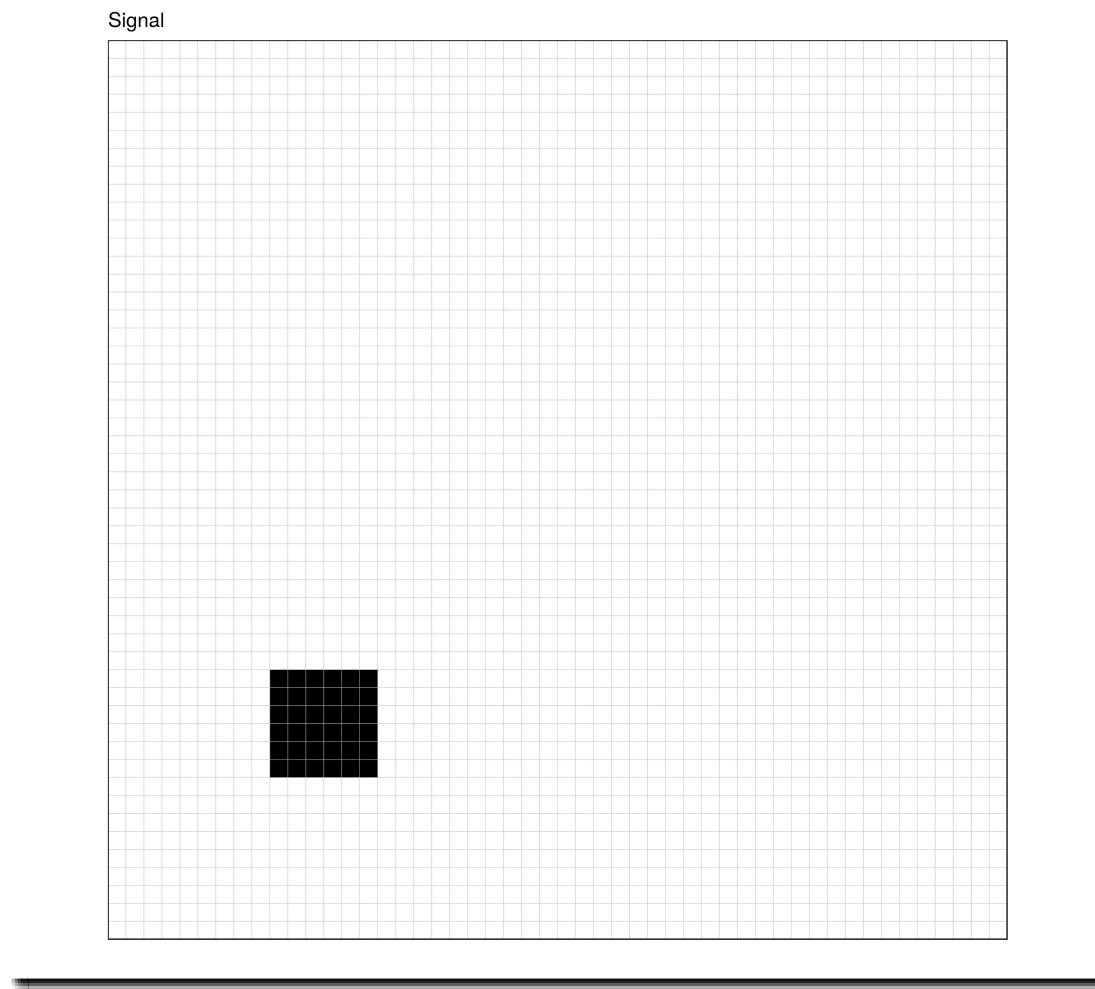
13 / 50

Illustrating biases

Signal

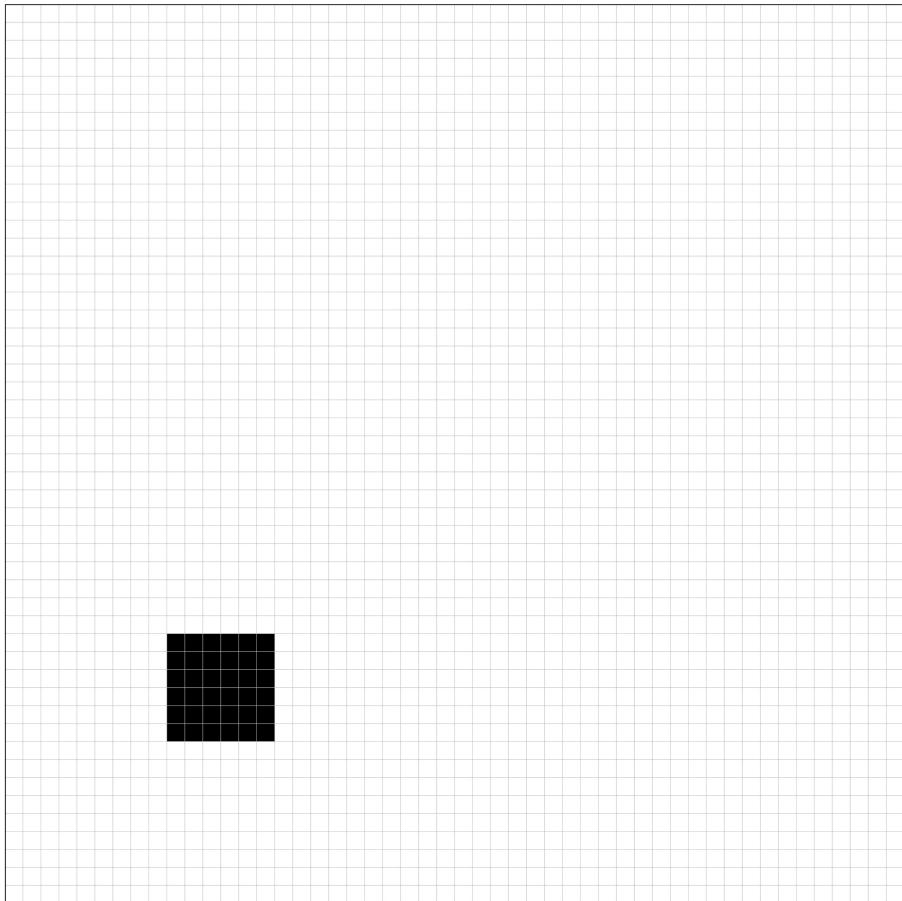


Illustrating biases

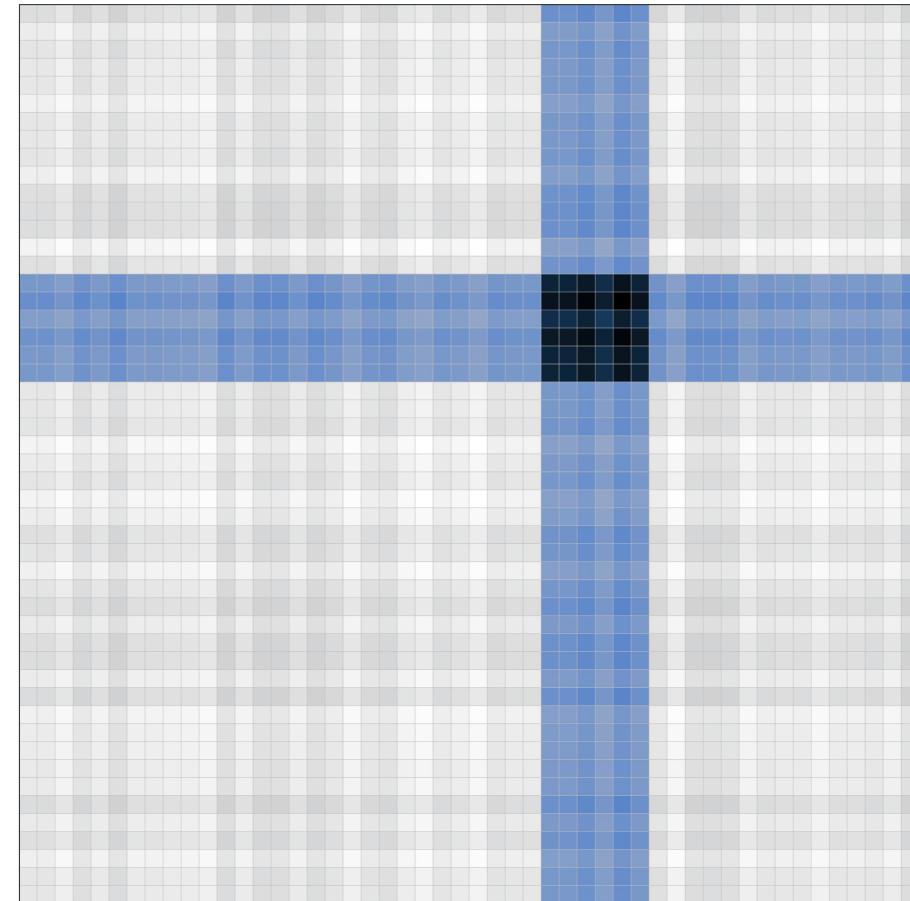


Illustrating biases

Signal

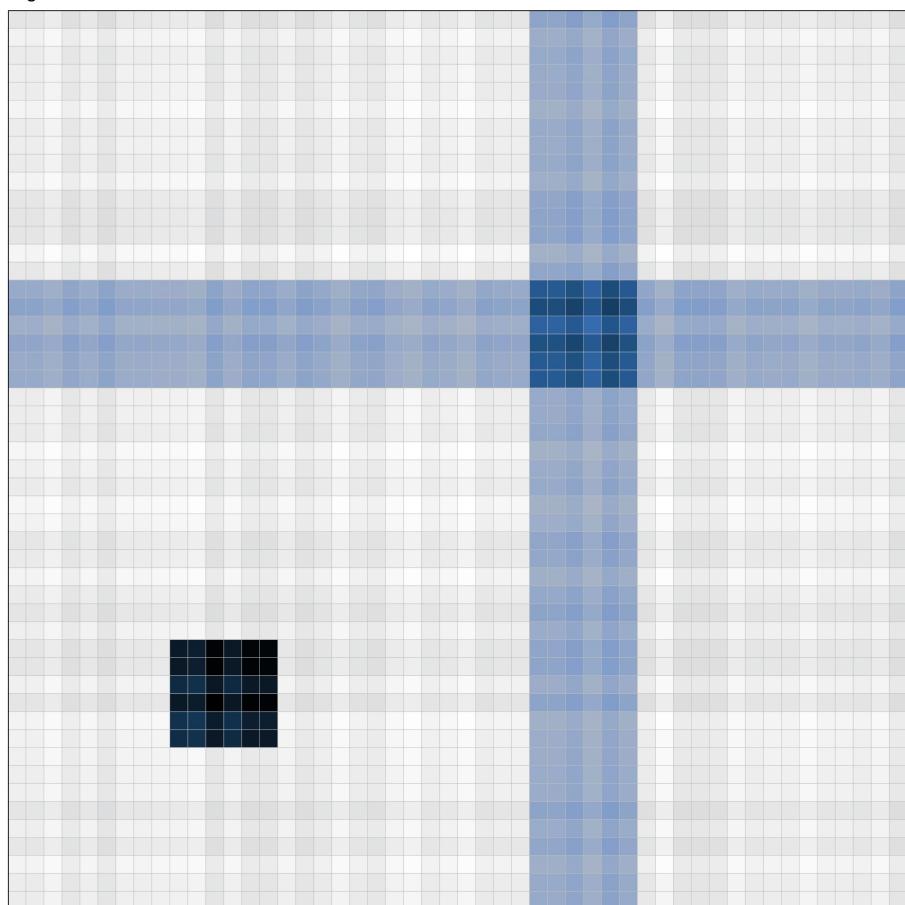


Bias

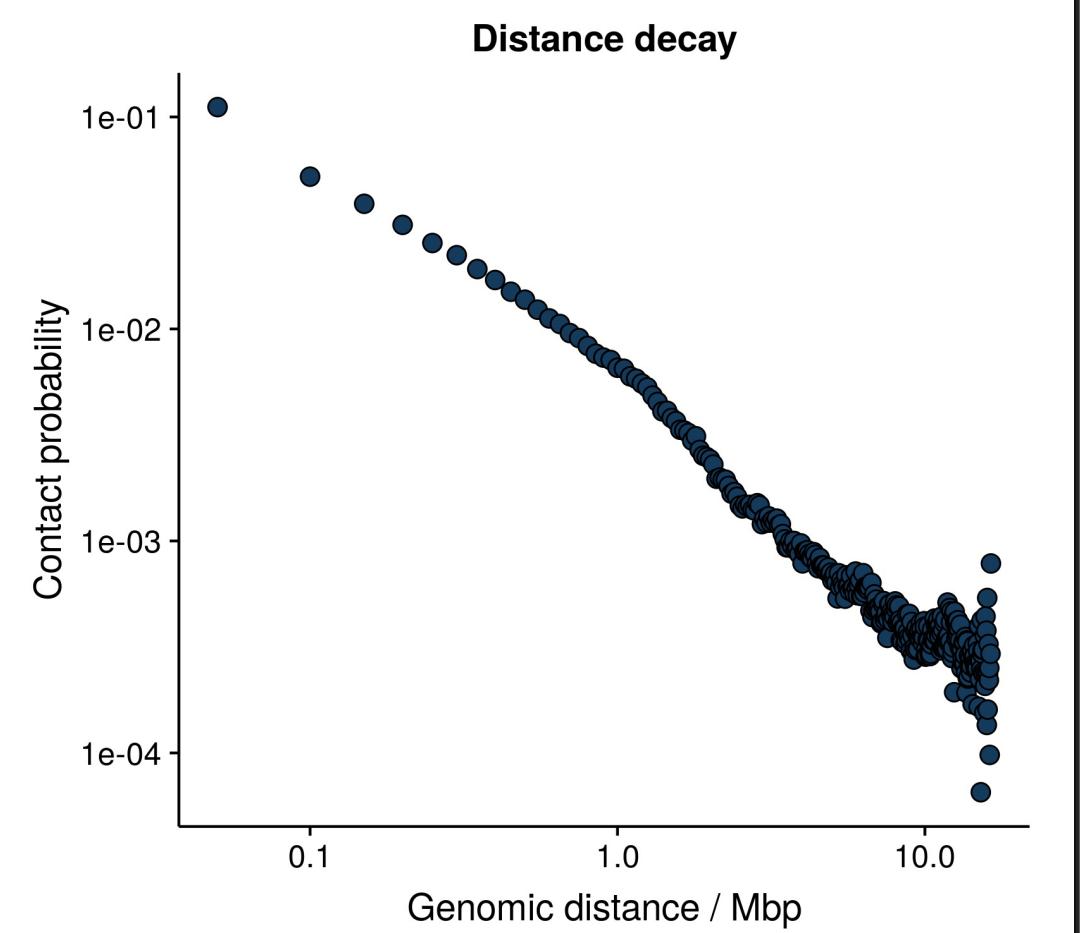
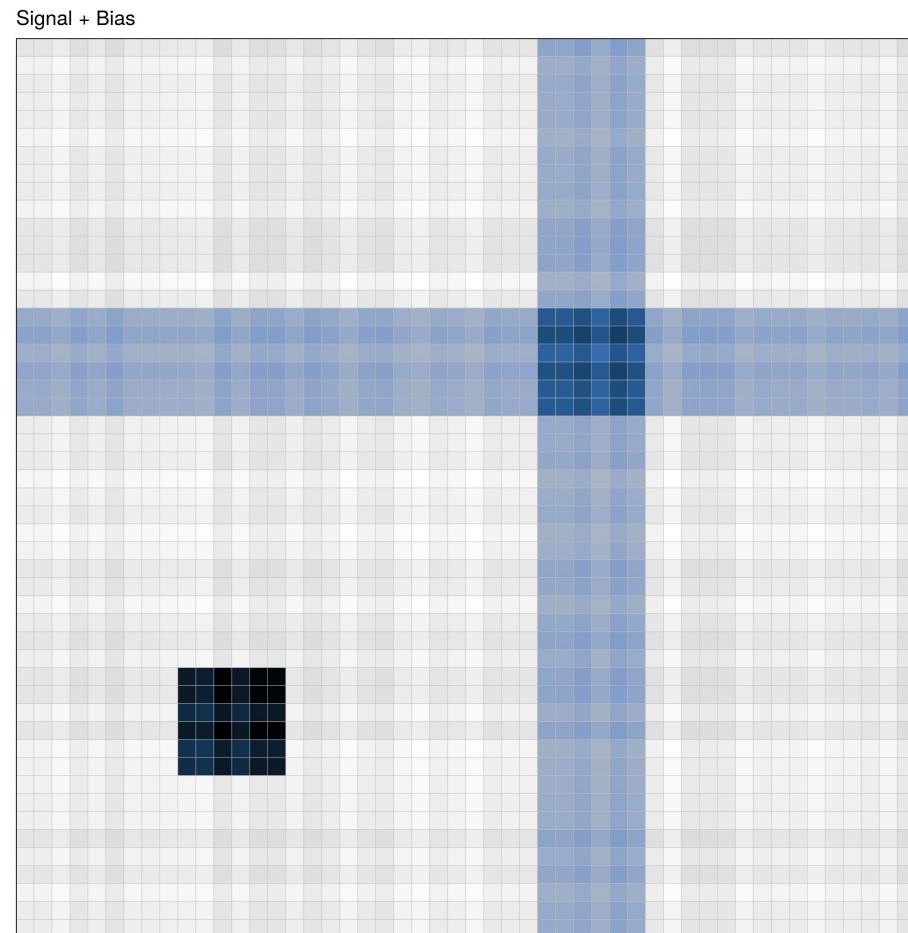


Illustrating biases

Signal + Bias

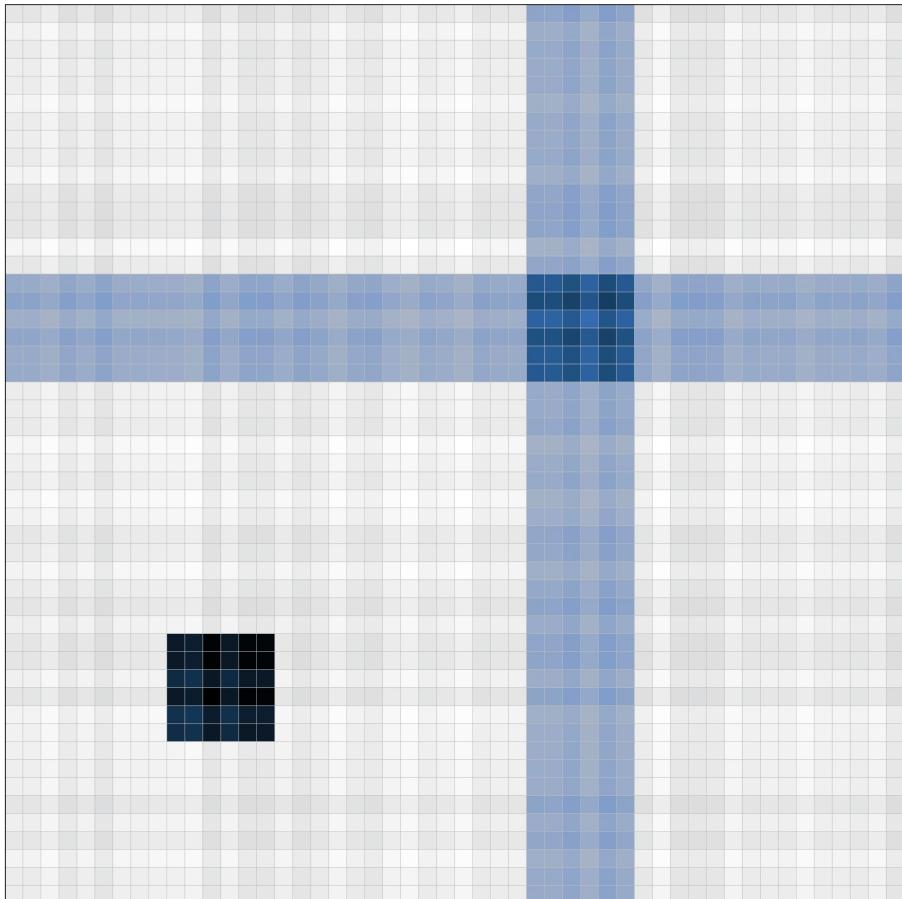


Illustrating biases

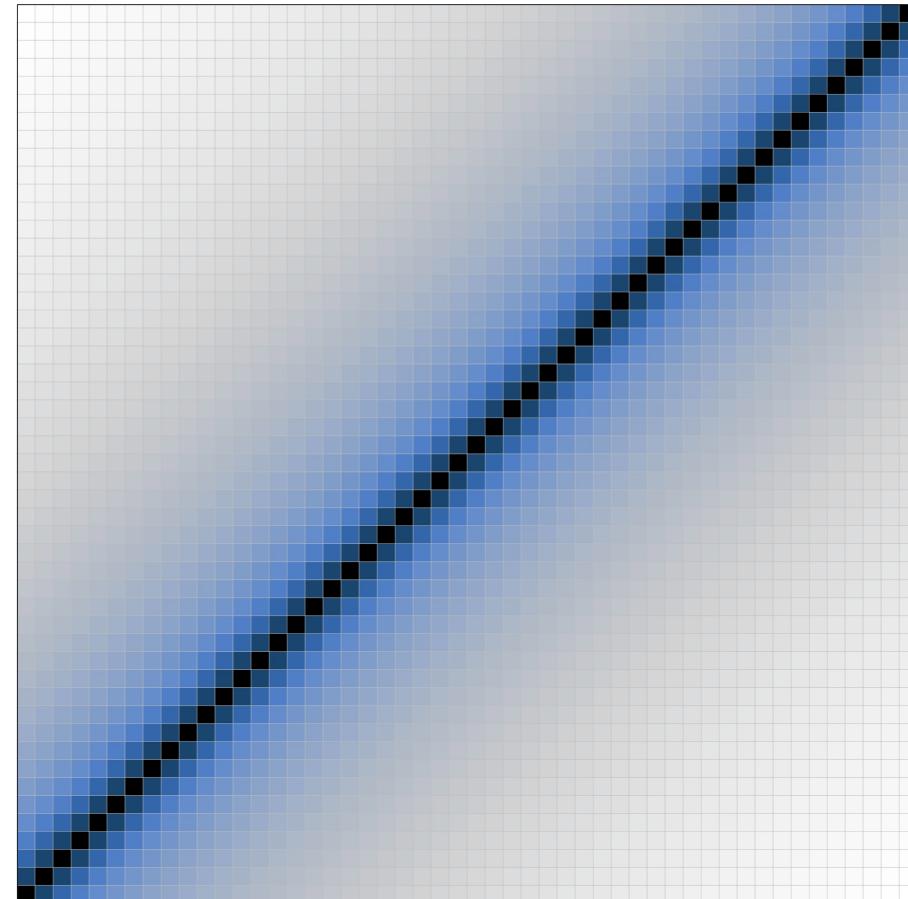


Illustrating biases

Signal + Bias

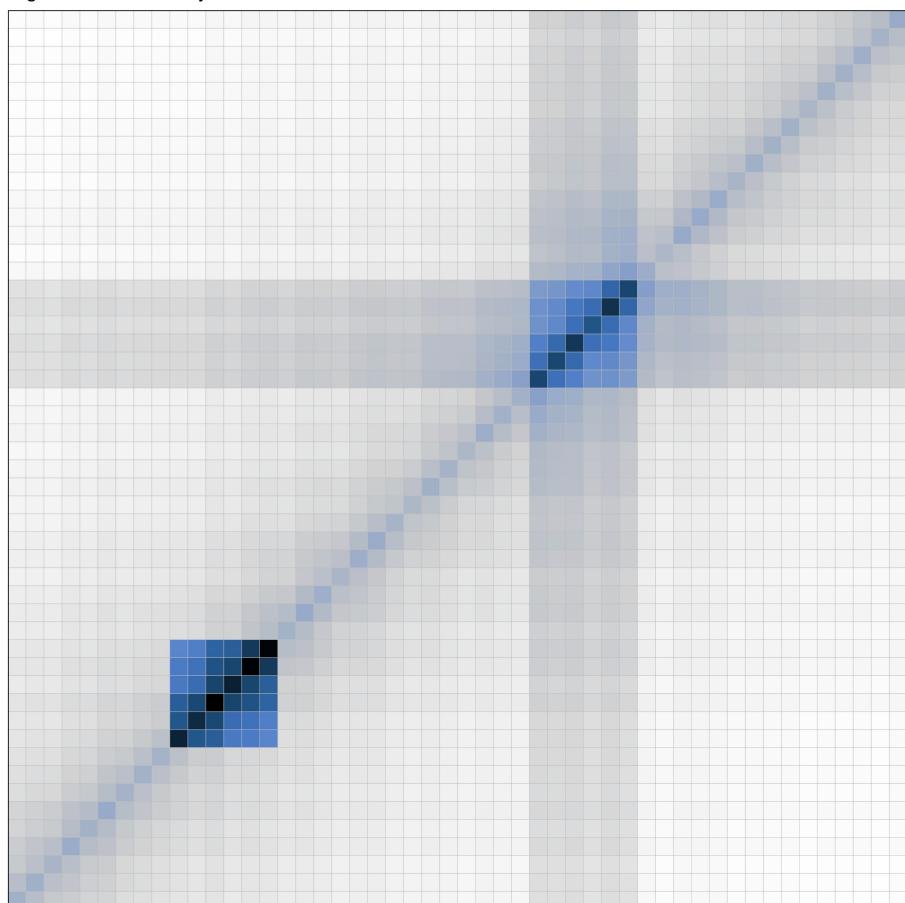


Decay



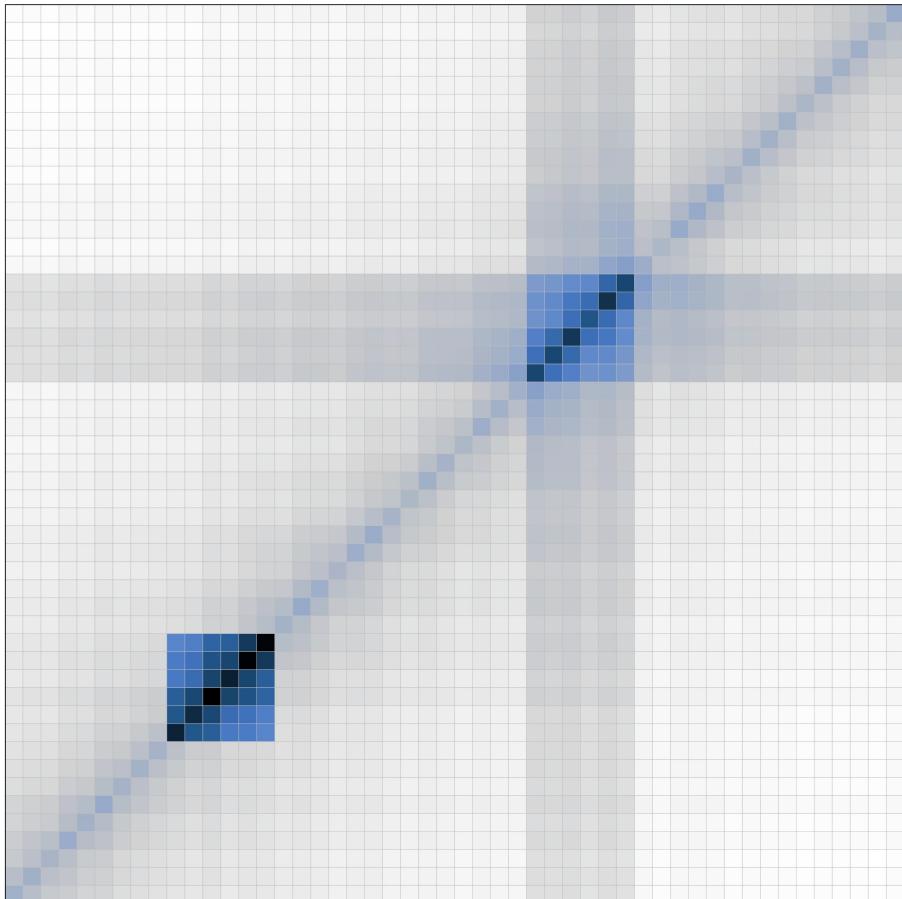
Illustrating biases

Signal + Bias + Decay

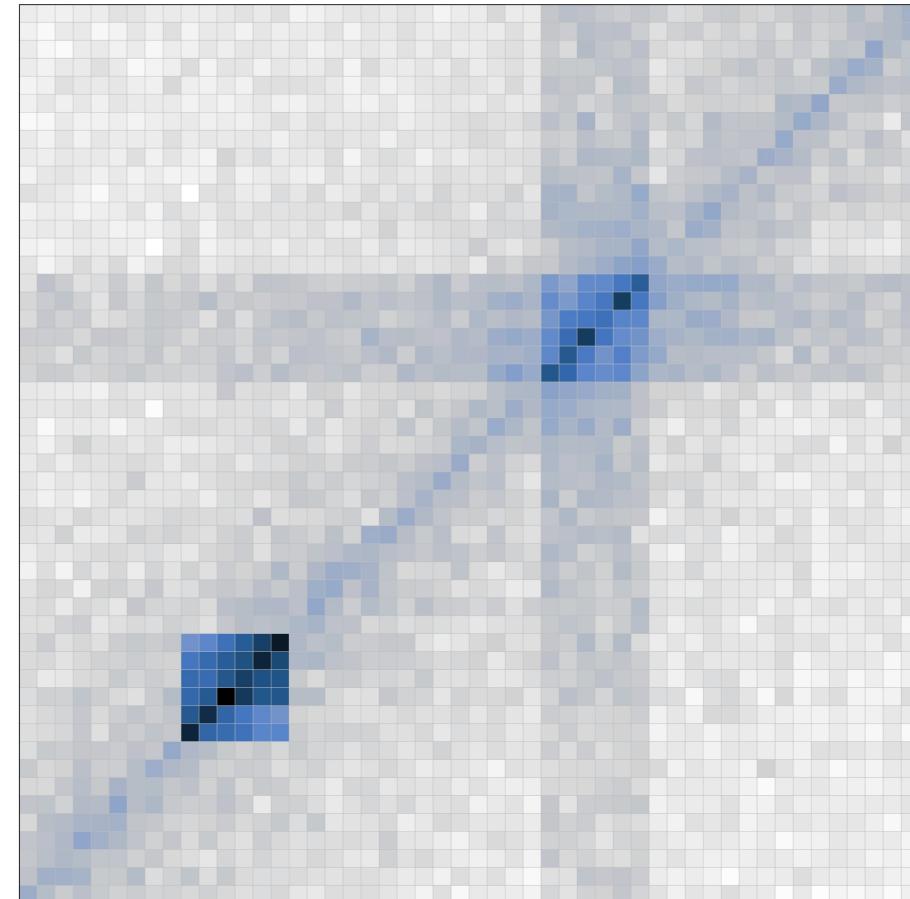


Illustrating biases

Signal + Bias + Decay

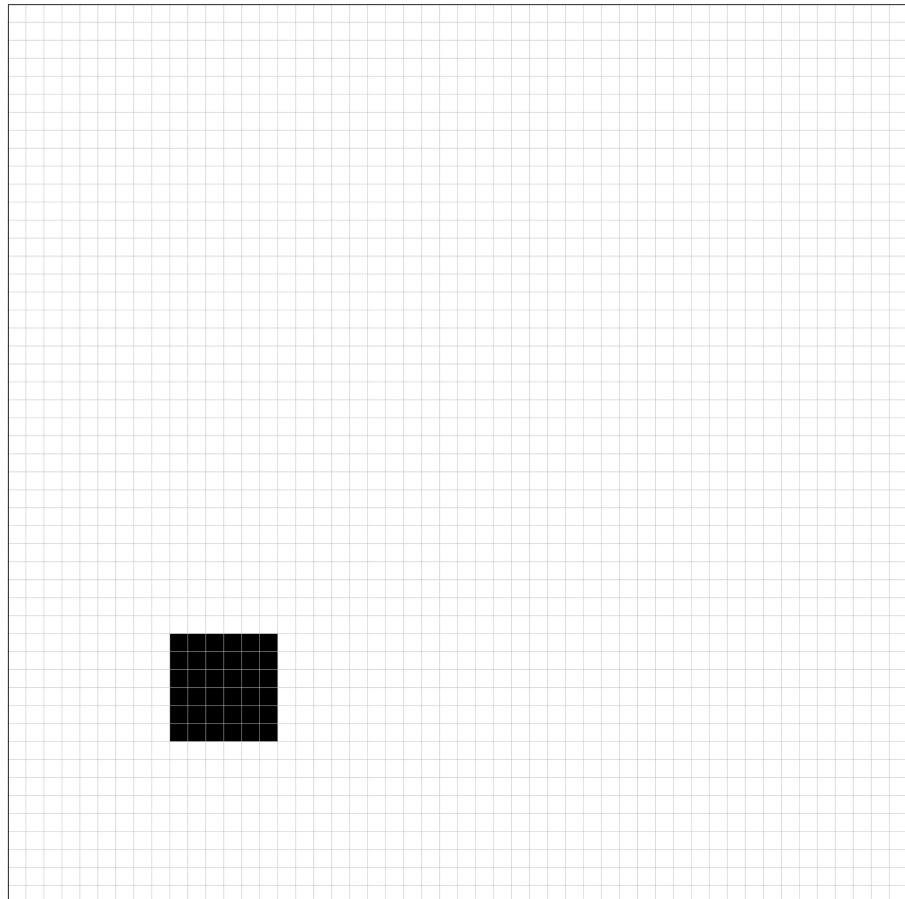


Signal + Bias + Decay + Noise

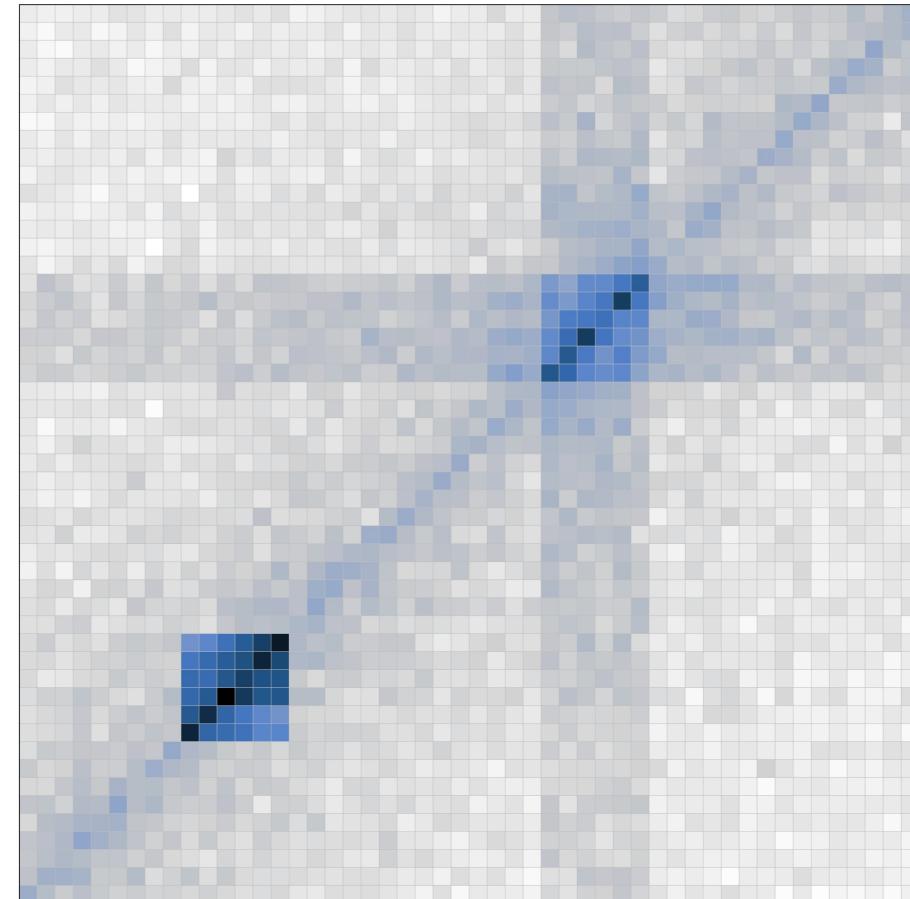


Illustrating biases

Signal



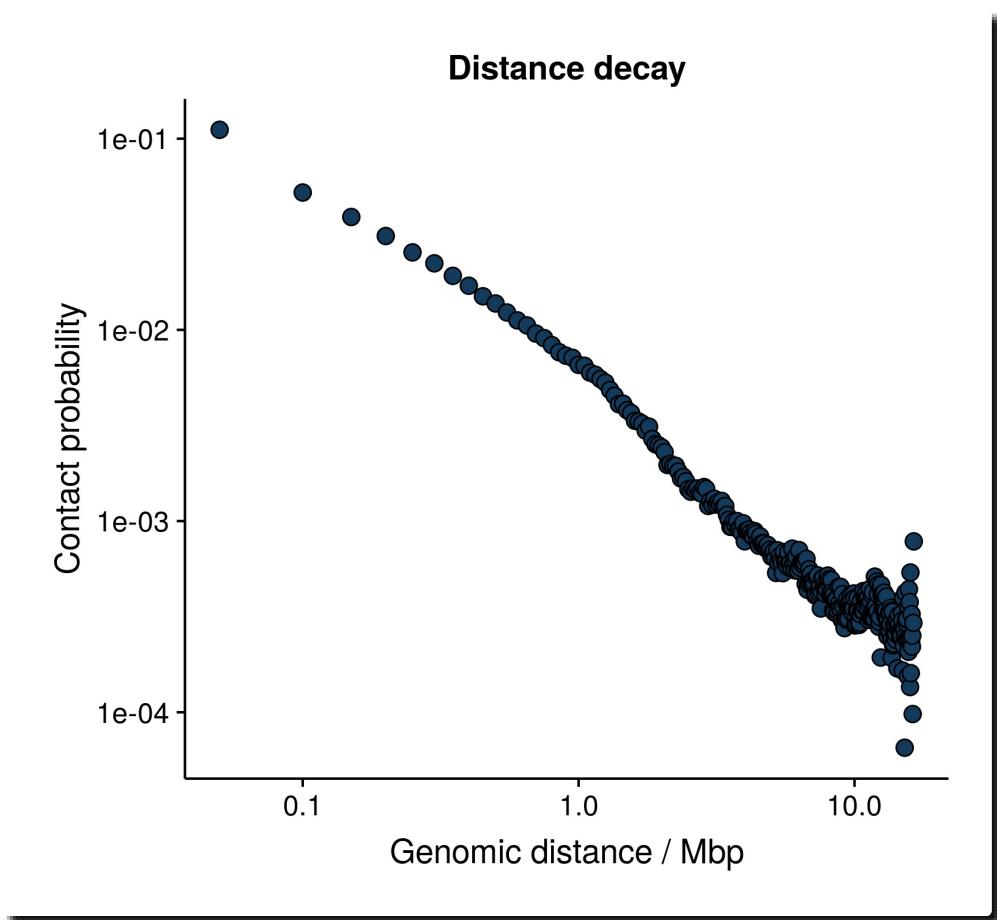
Signal + Bias + Decay + Noise



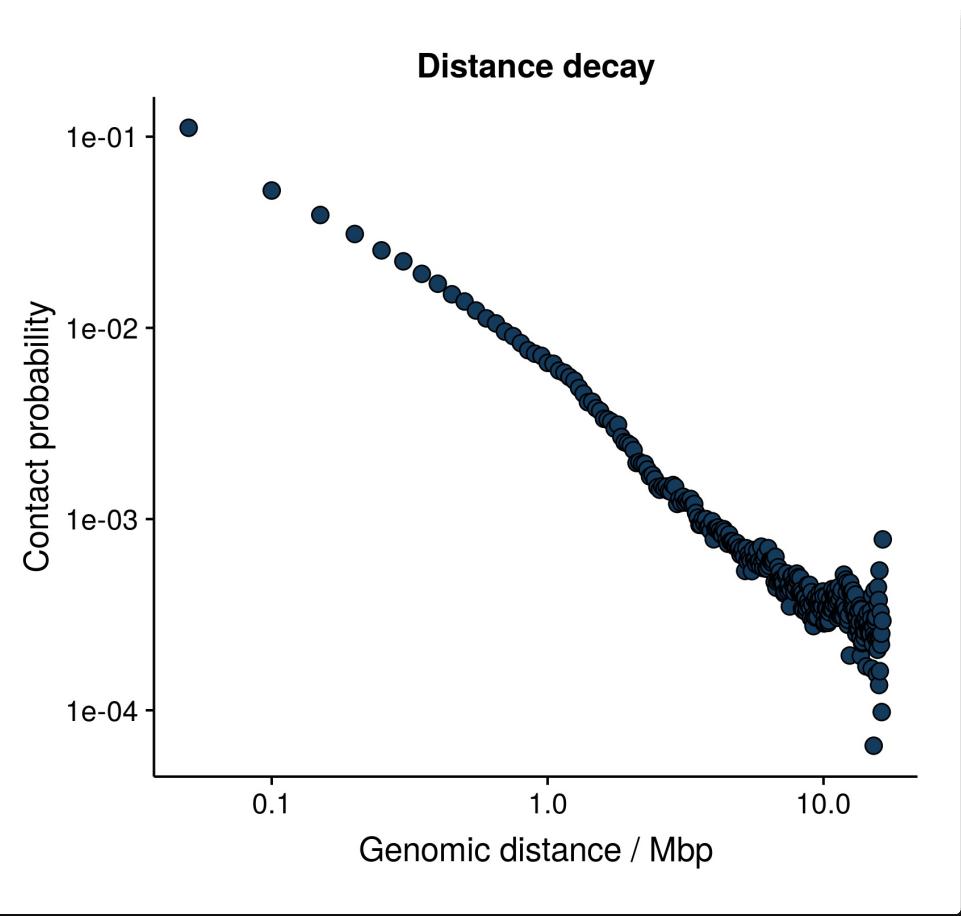
Challenge 2

Reproducibility

Decay drives Hi-C matrix comparison



Decay drives Hi-C matrix comparison



- Correct by expected decay **O / E**

- Compute by distance **SCC**

Yang, T. et al. (2017)

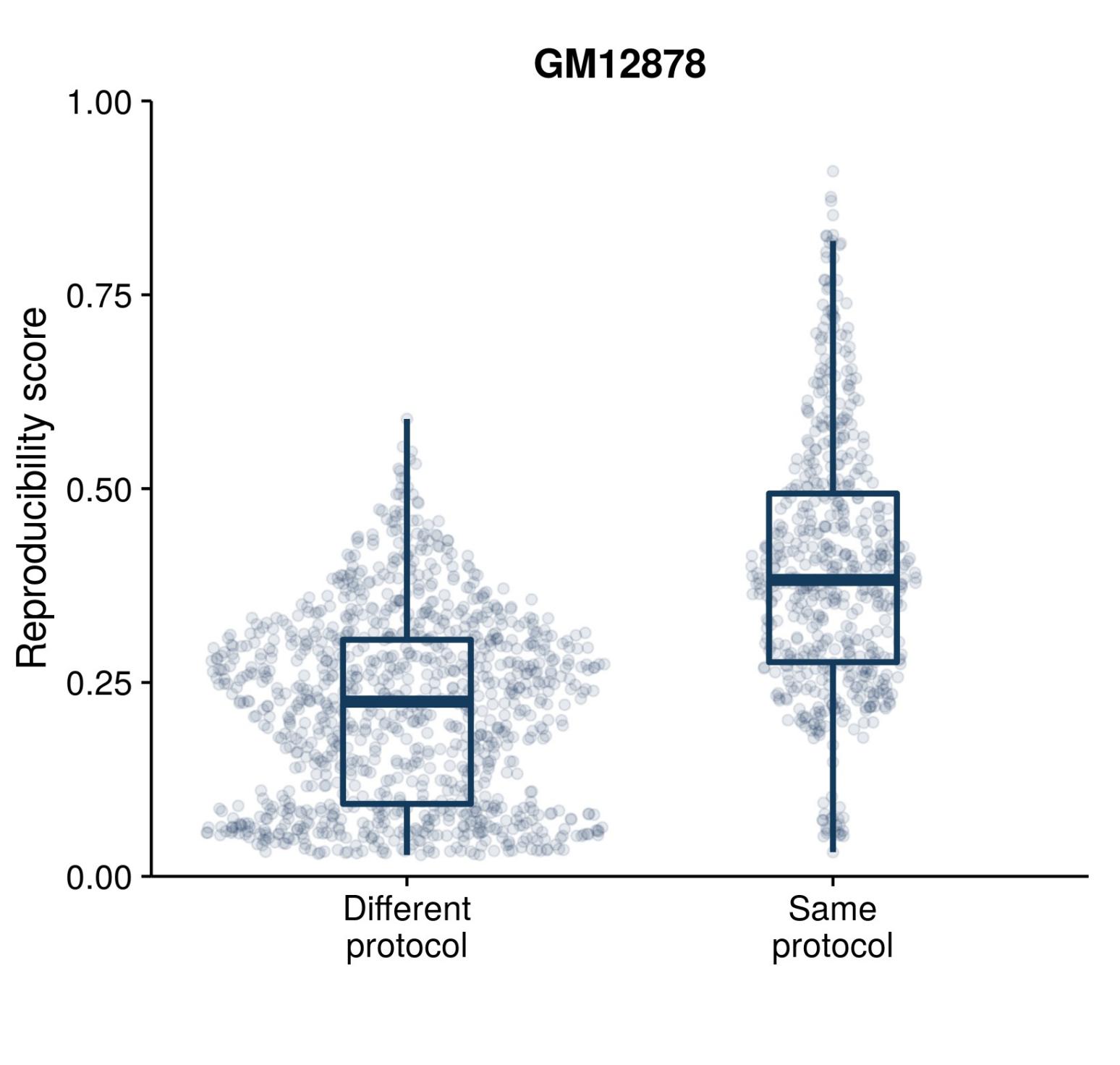
Genome research

- Spectral decomposition
Reproducibility score

Yan, K. K. et al. (2017)

Bioinformatics

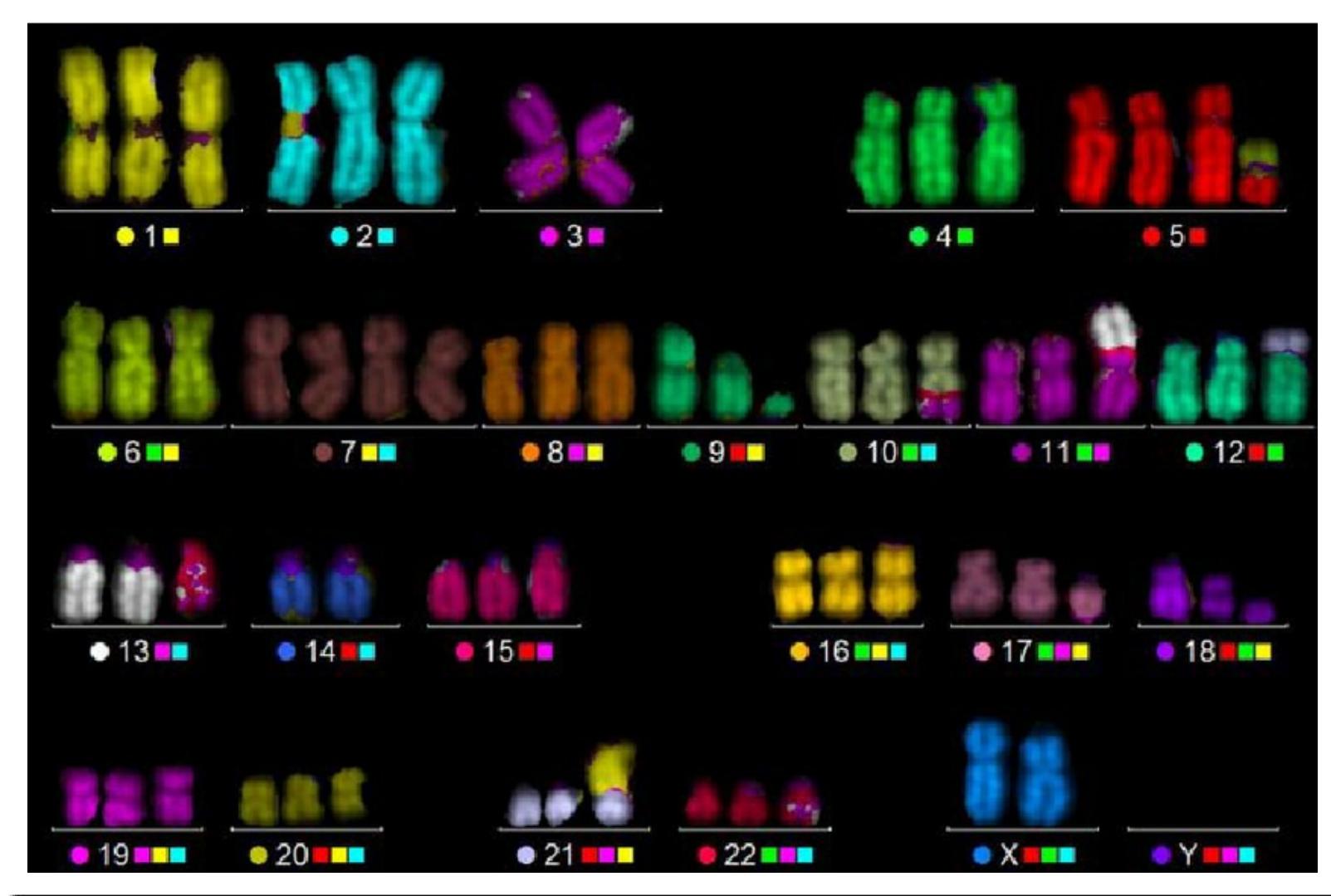
Reproducibility



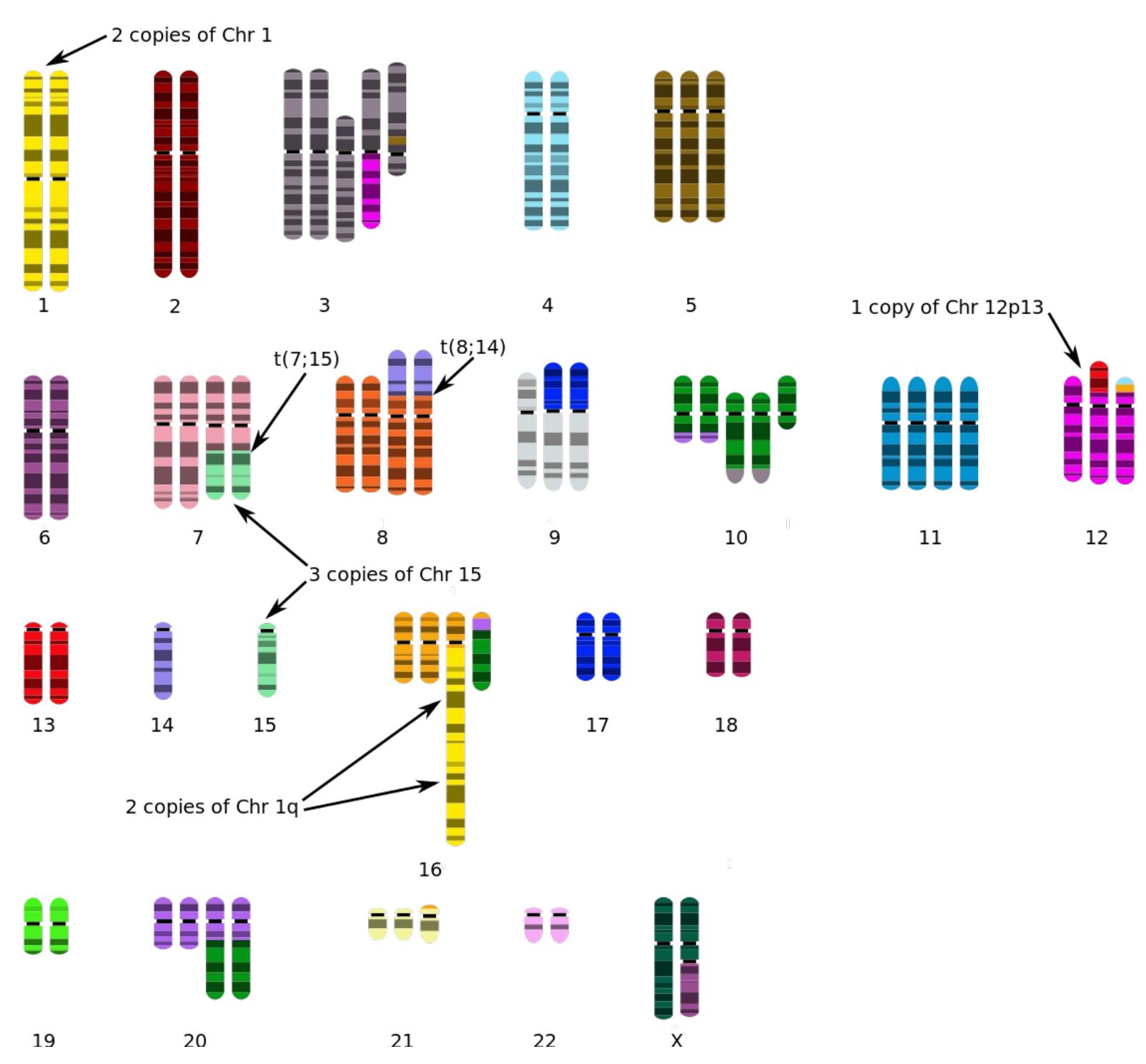
Challenge 3

Aberrant karyotypes

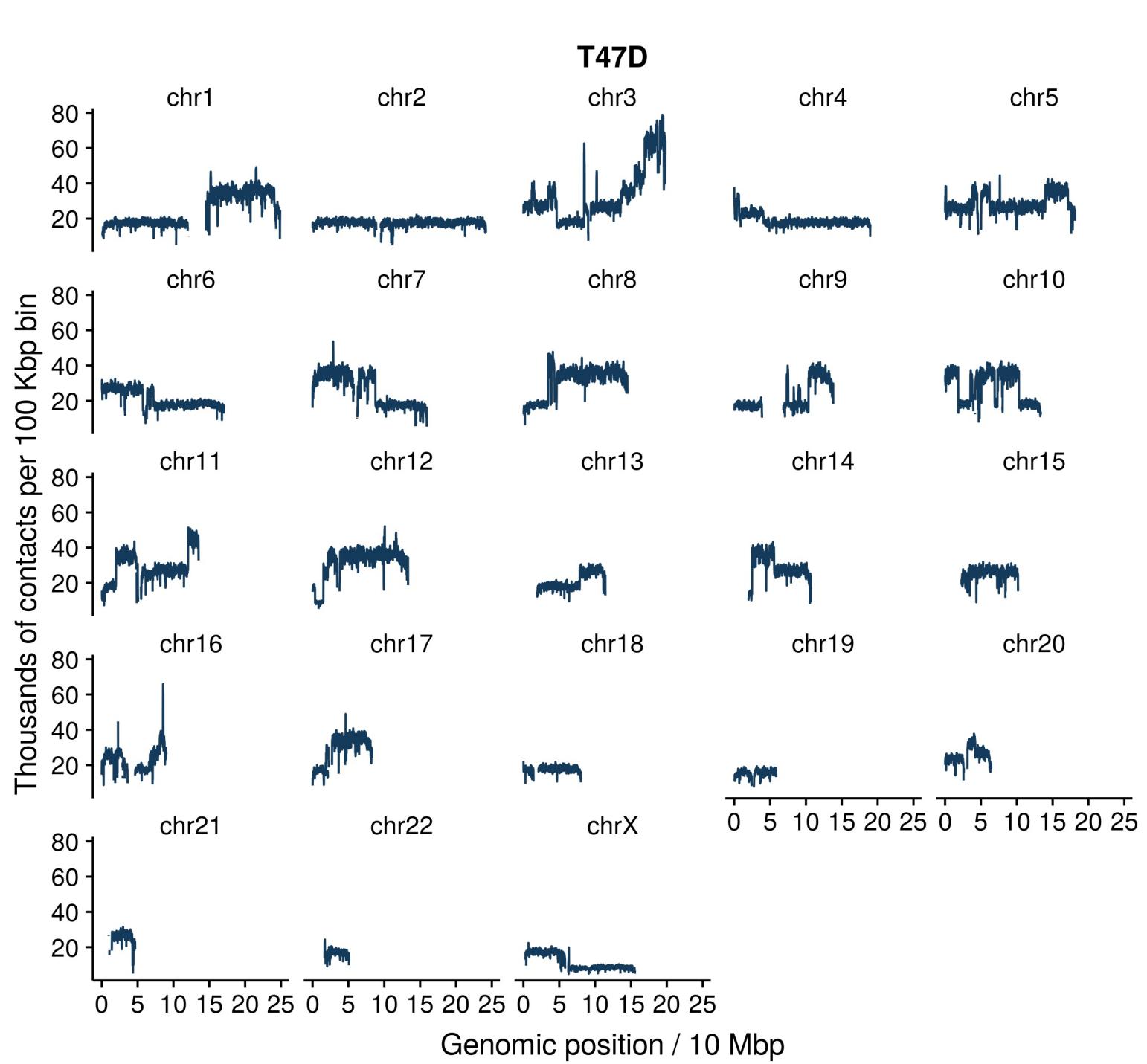
Cancer cell line: K562



Cancer cell line: T47D



Cancer cell line: T47D



Previous approaches

HiCNorm

Explicit model of biases

Regression model

All matrix entries

Hu, M. et al. (2012)

Bioinformatics

Yaffe, E. and Tanay, A. (2011)

Nature genetics

HiCNorm

ICE

Explicit model of biases

Regression model

All matrix entries

Implicit correction

Matrix balancing

Equal visibility

Hu, M. et al. (2012)

Bioinformatics

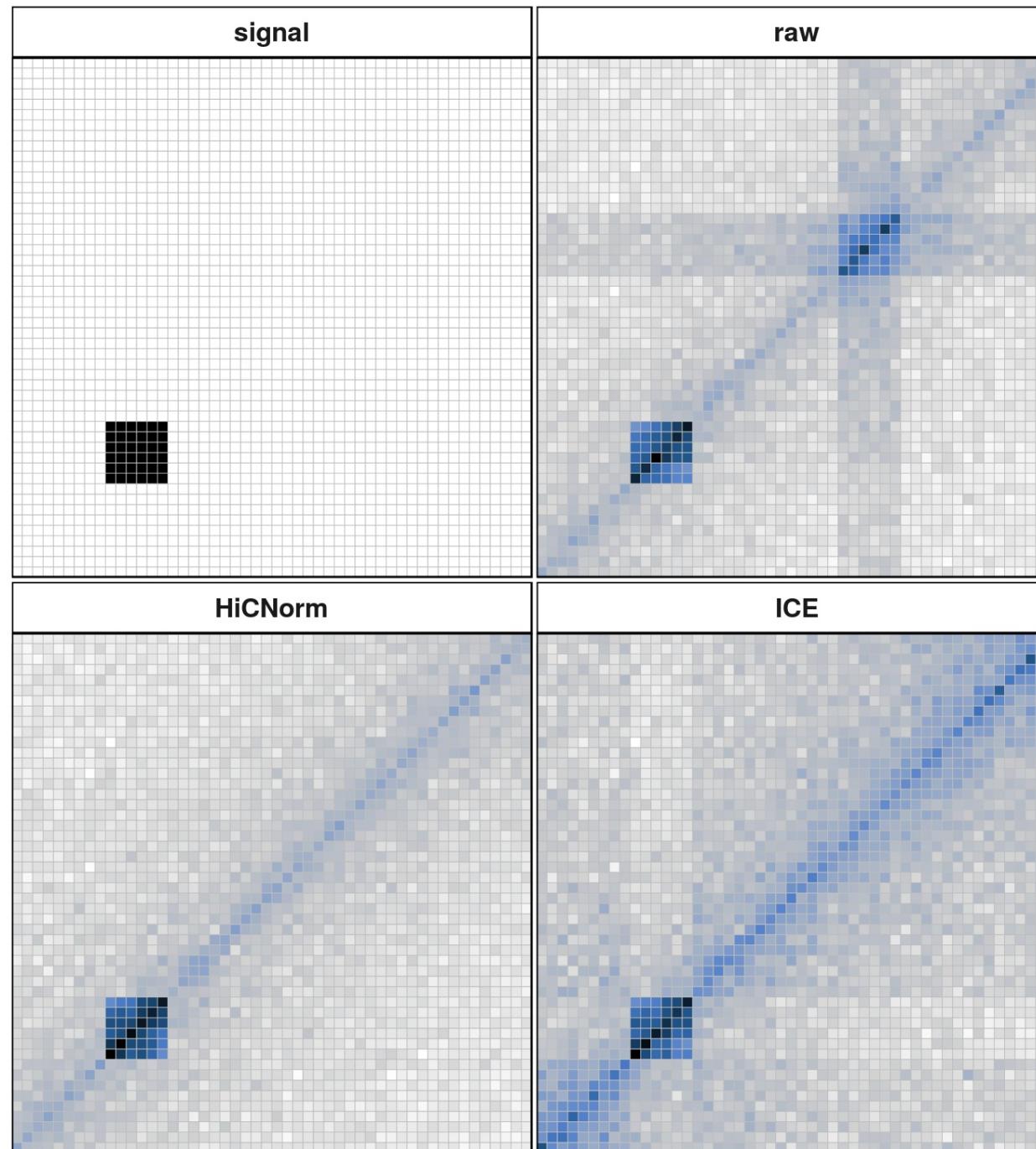
Yaffe, E. and Tanay, A. (2011)

Nature genetics

Imakaev, M. et al. (2012)

Nature methods

In action



HiCNorm

✓ fewer artifacts

✗ very slow

✗ no high resolution

ICE

✗ more artifacts

✓ fast

✓ high resolution

HiCNorm

✓ fewer artifacts

✗ very slow

✗ no high resolution

ICE

✗ more artifacts

✓ fast

✓ high resolution

✗ aberrant karyotypes

Wish list

Wish list

♥ Suitable aberrant karyotypes

Wish list

♥ Suitable aberrant karyotypes

♥ As good as existing

Wish list

- ♥ Suitable aberrant karyotypes
- ♥ As good as existing
- ♥ Not worse on normal karyotypes

Wish list

- ♥ Suitable aberrant karyotypes
- ♥ As good as existing
- ♥ Not worse on normal karyotypes
- ♥ Fast

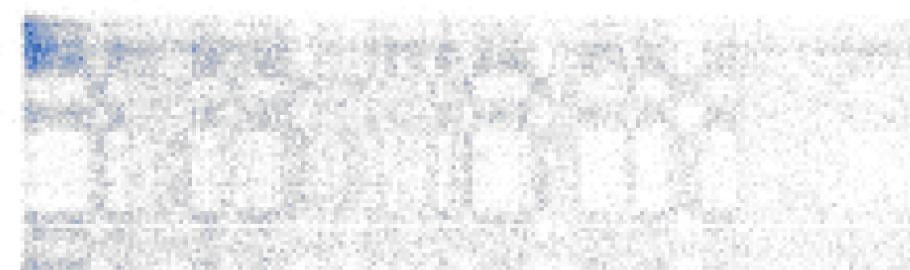
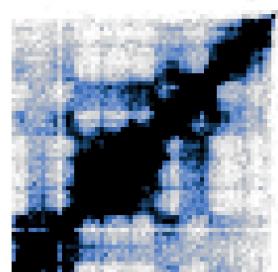
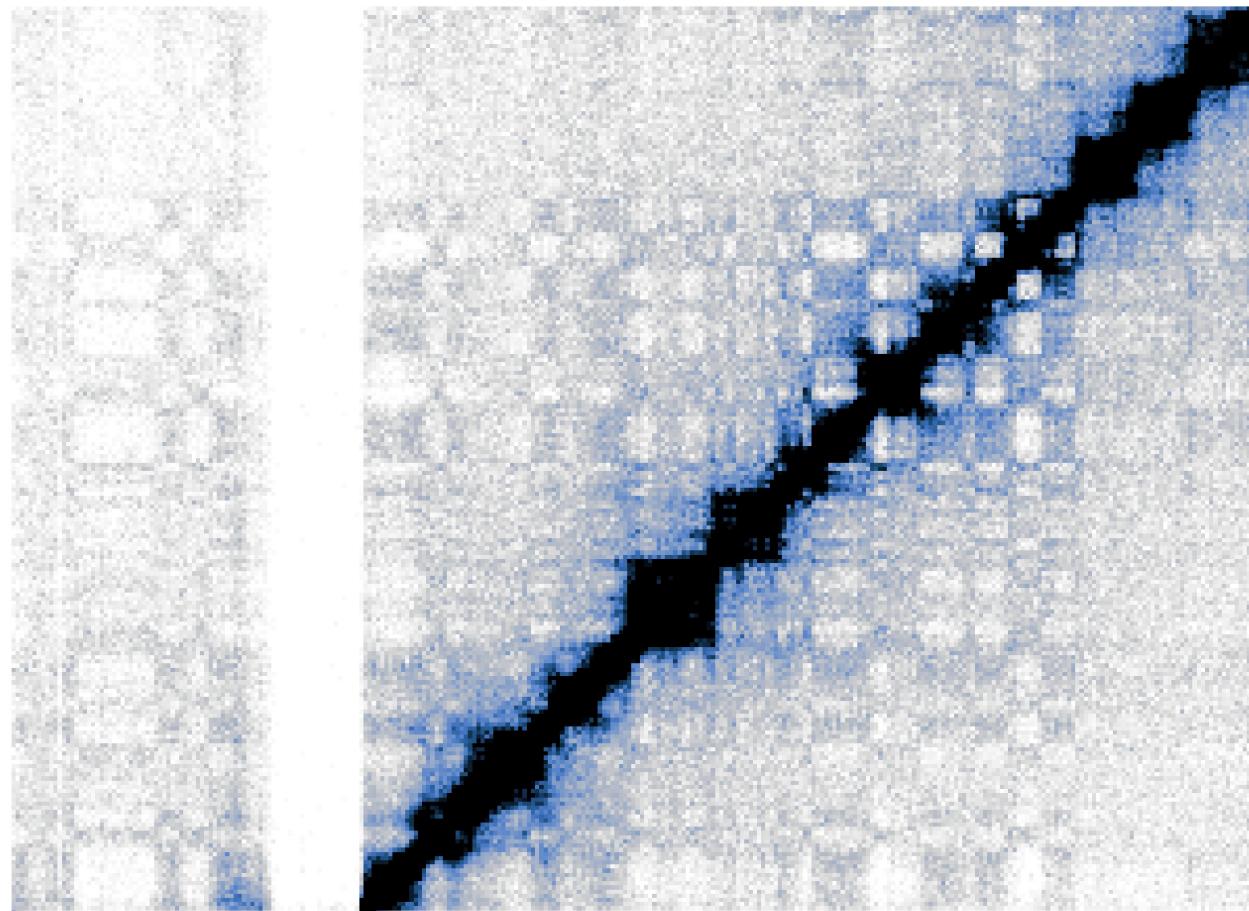
Wish list

- ♥ Suitable aberrant karyotypes
- ♥ As good as existing
- ♥ Not worse on normal karyotypes
- ♥ Fast
- ♥ Usable at high resolution

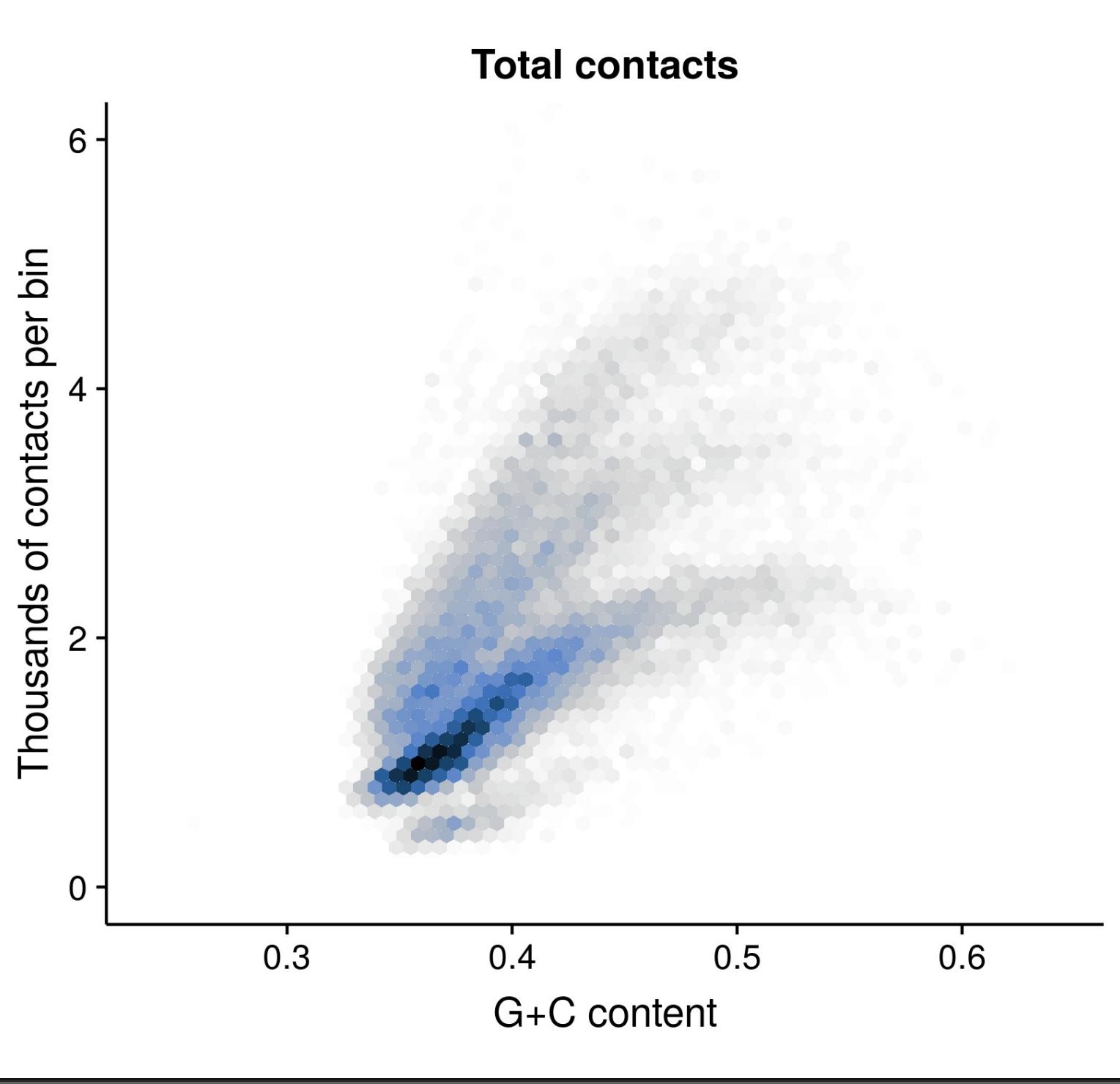
2D → 1D

2D → 1D

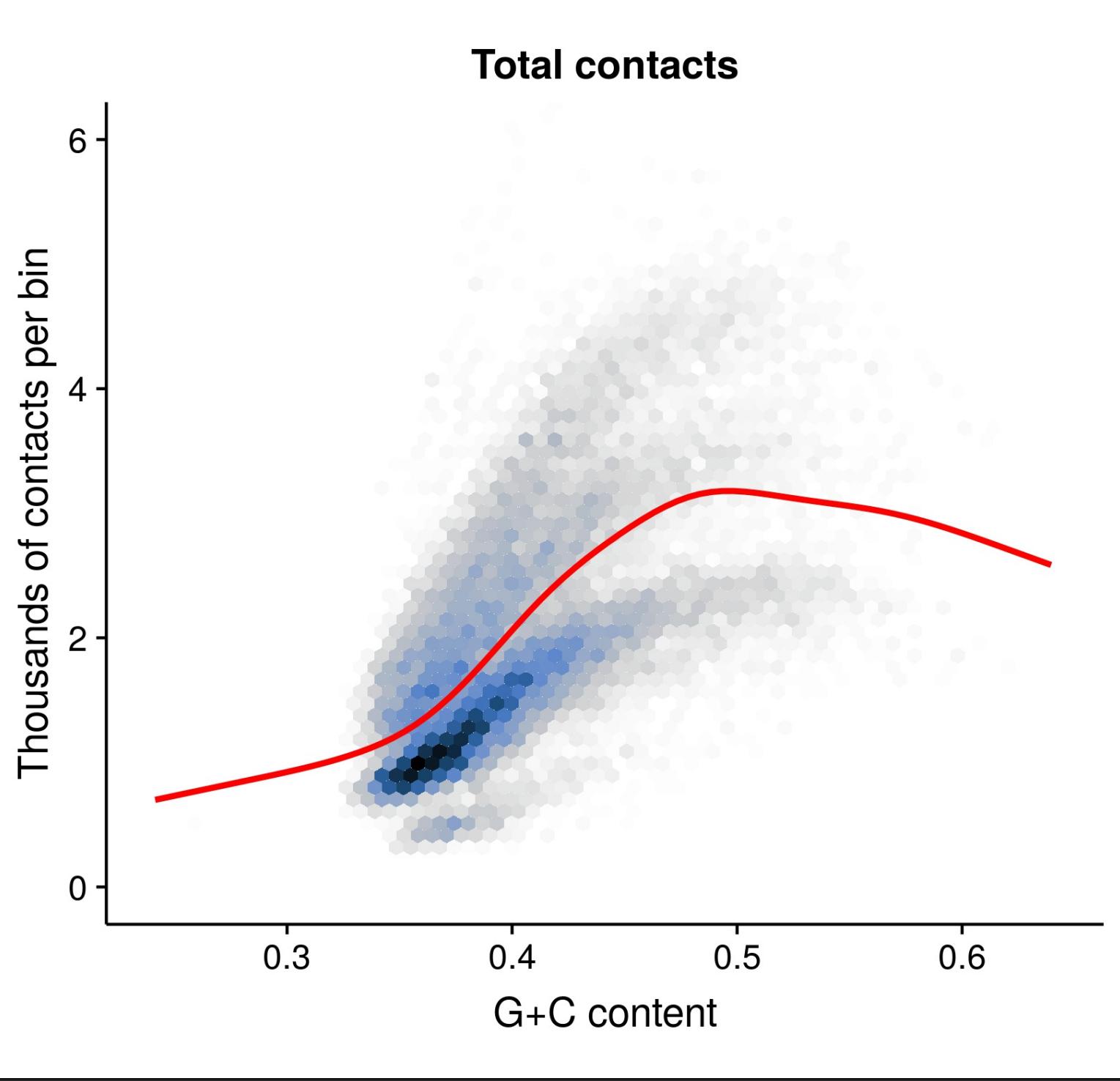
chromosome 17 @ 250 Kbp bin size



Totals vs. Genomic features

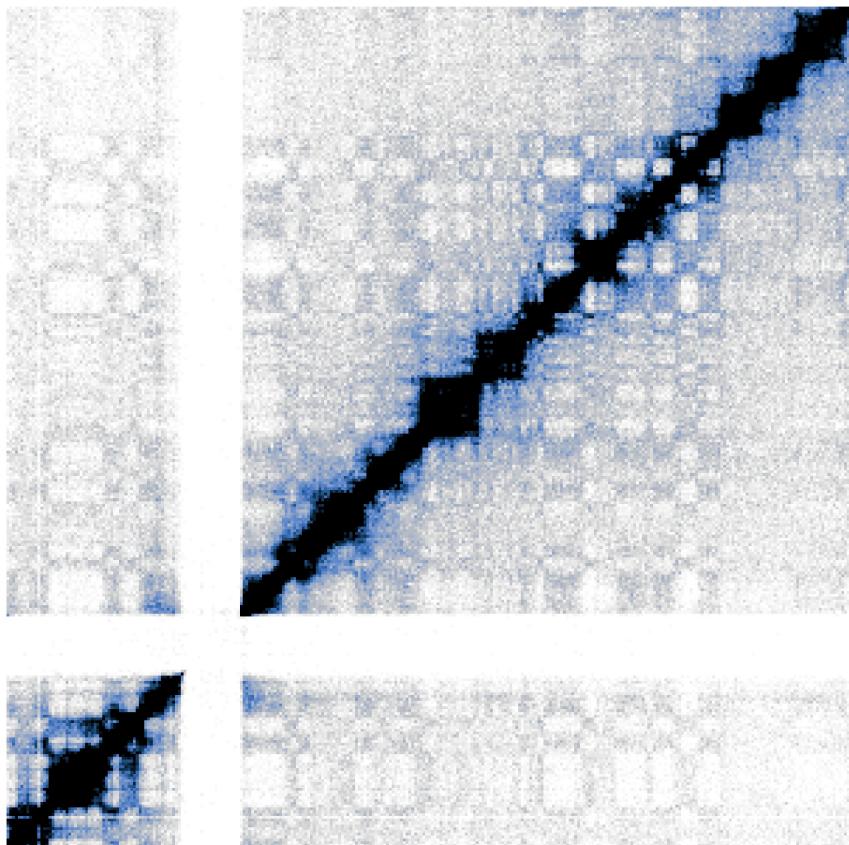


Non-linear relationship

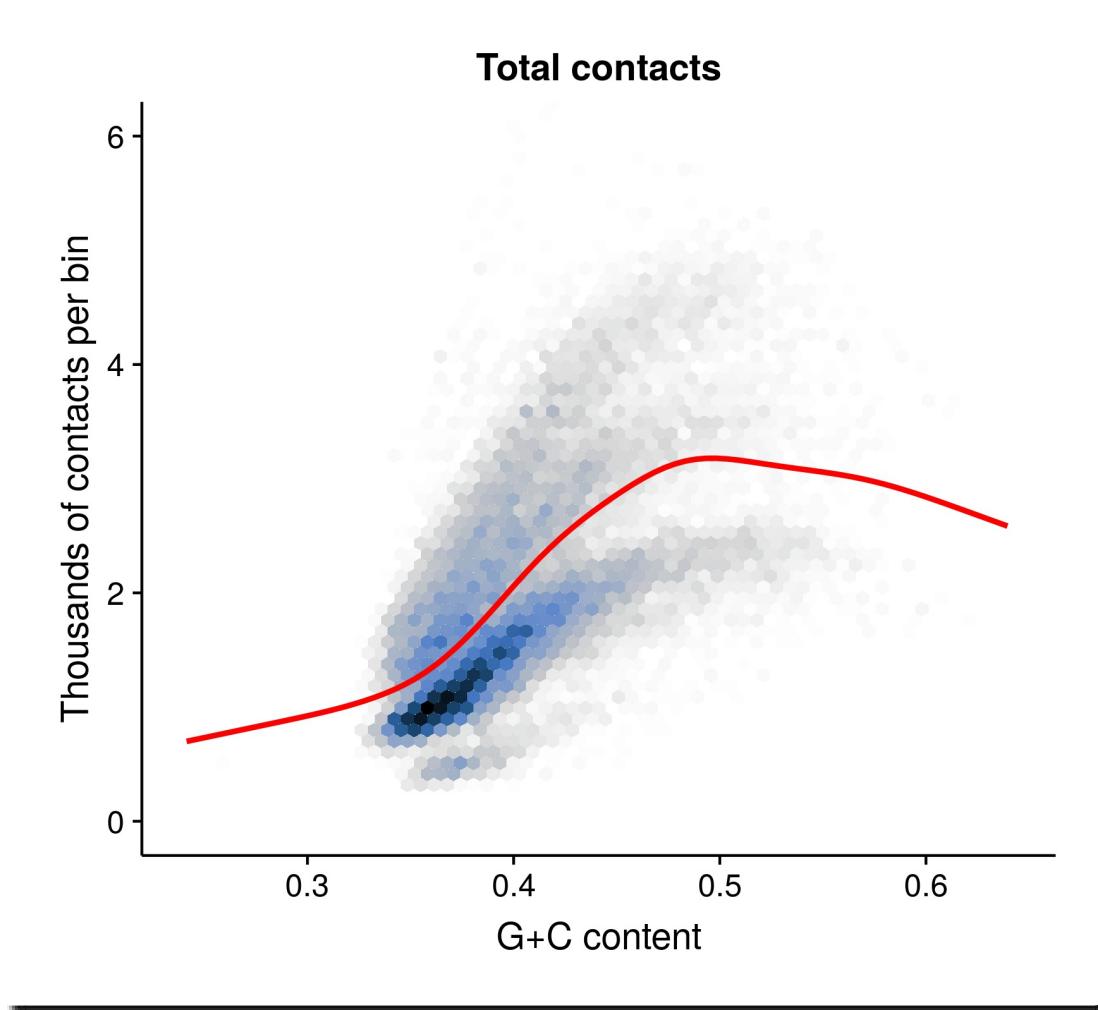


OneD model

chromosome 17 @ 250 Kbp bin size

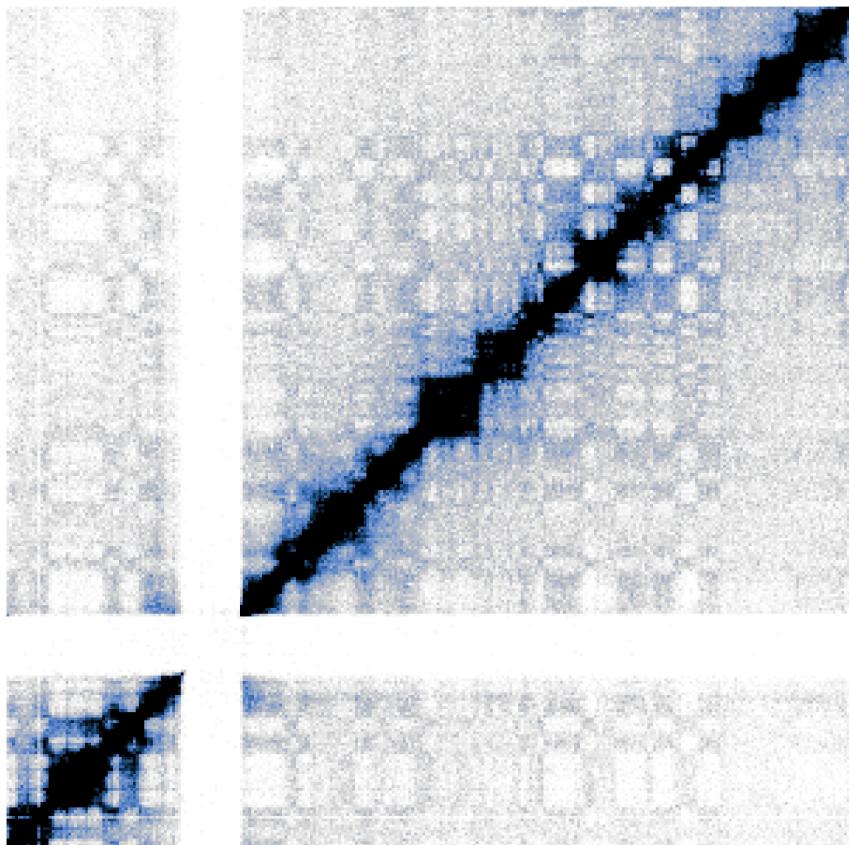


OneD model

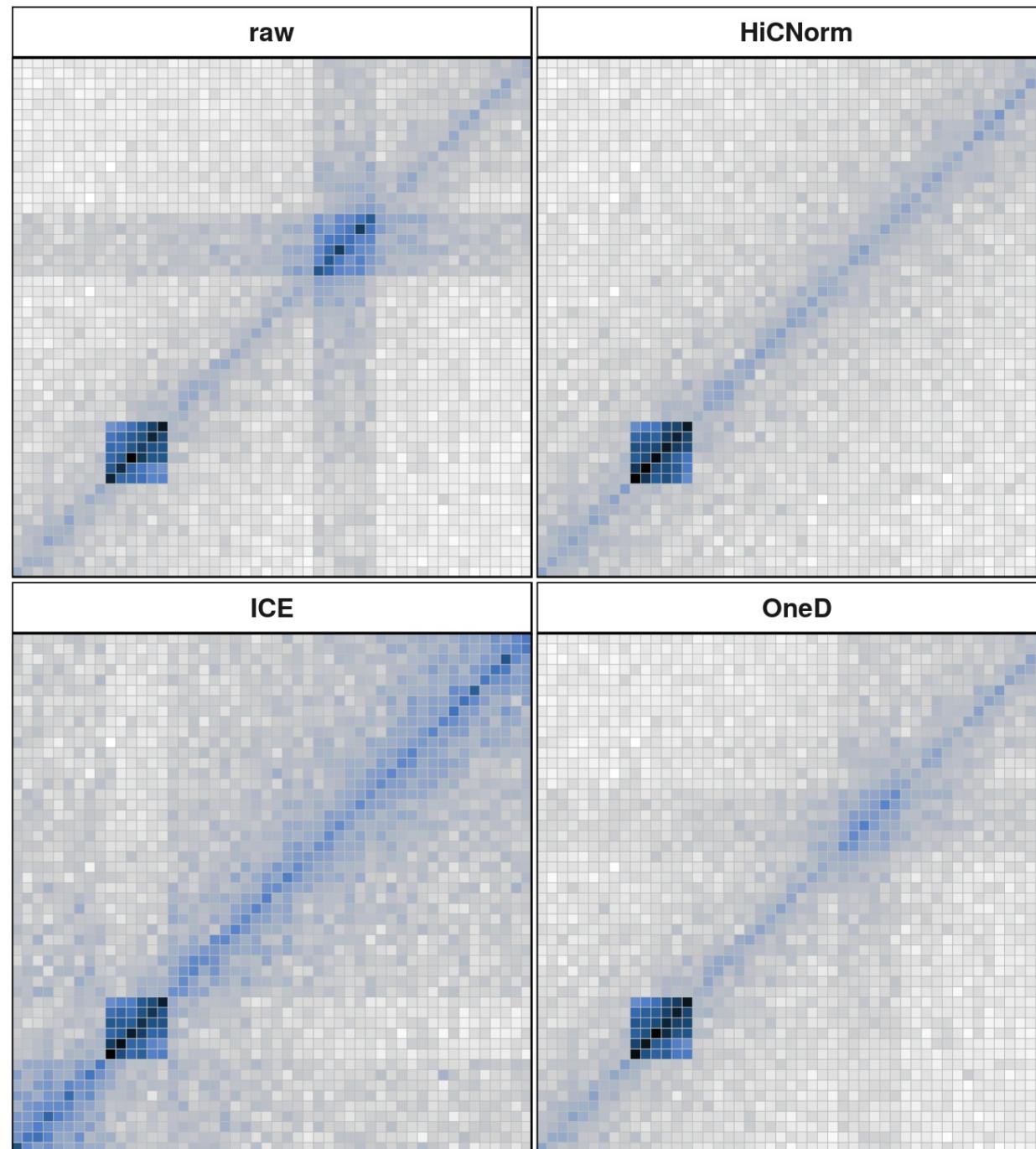


OneD model

chromosome 17 @ 250 Kbp bin size



In action



OneD performance

Benchmark strategy

Different methods

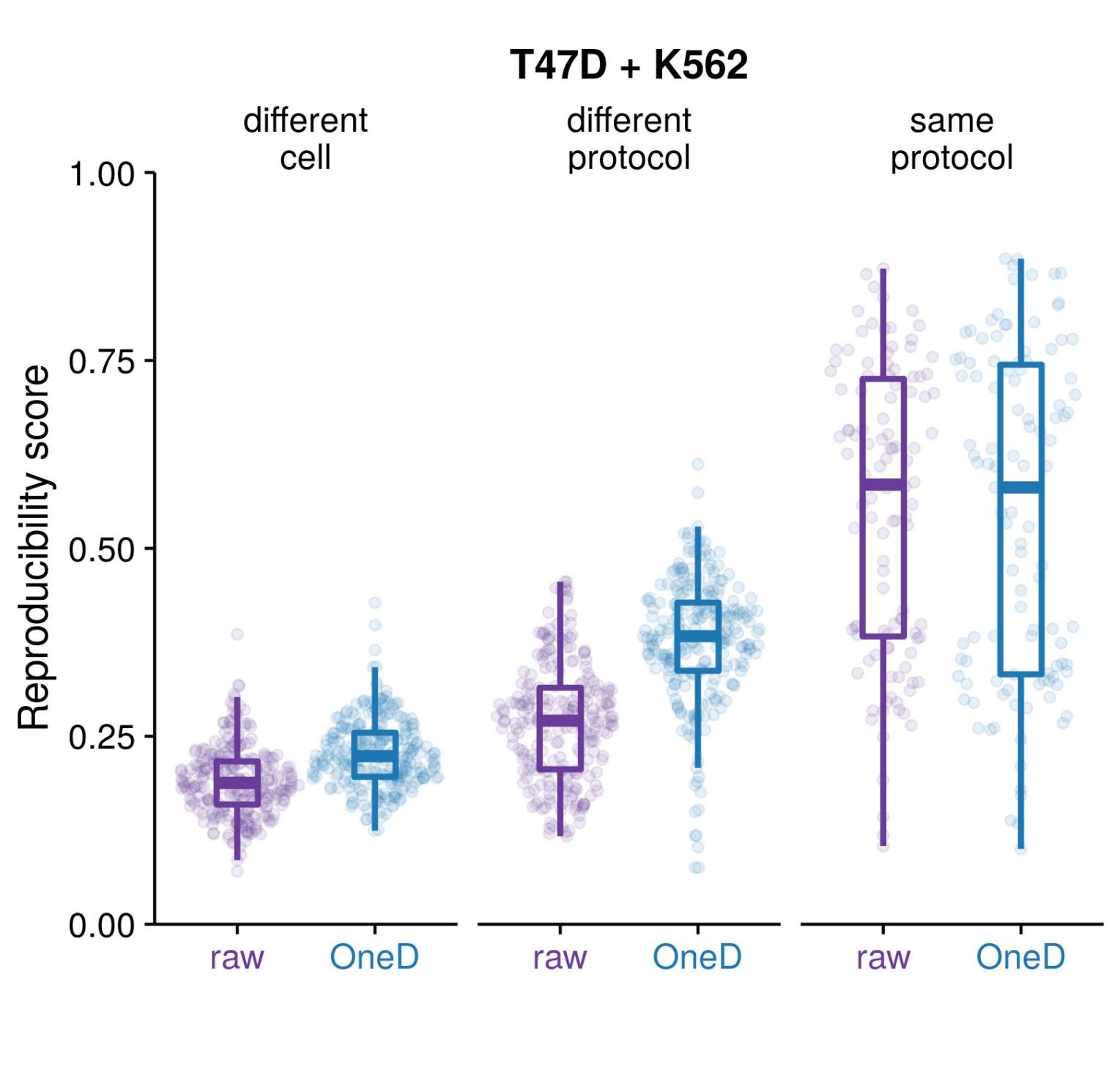
Including doing nothing (raw)

Different protocols

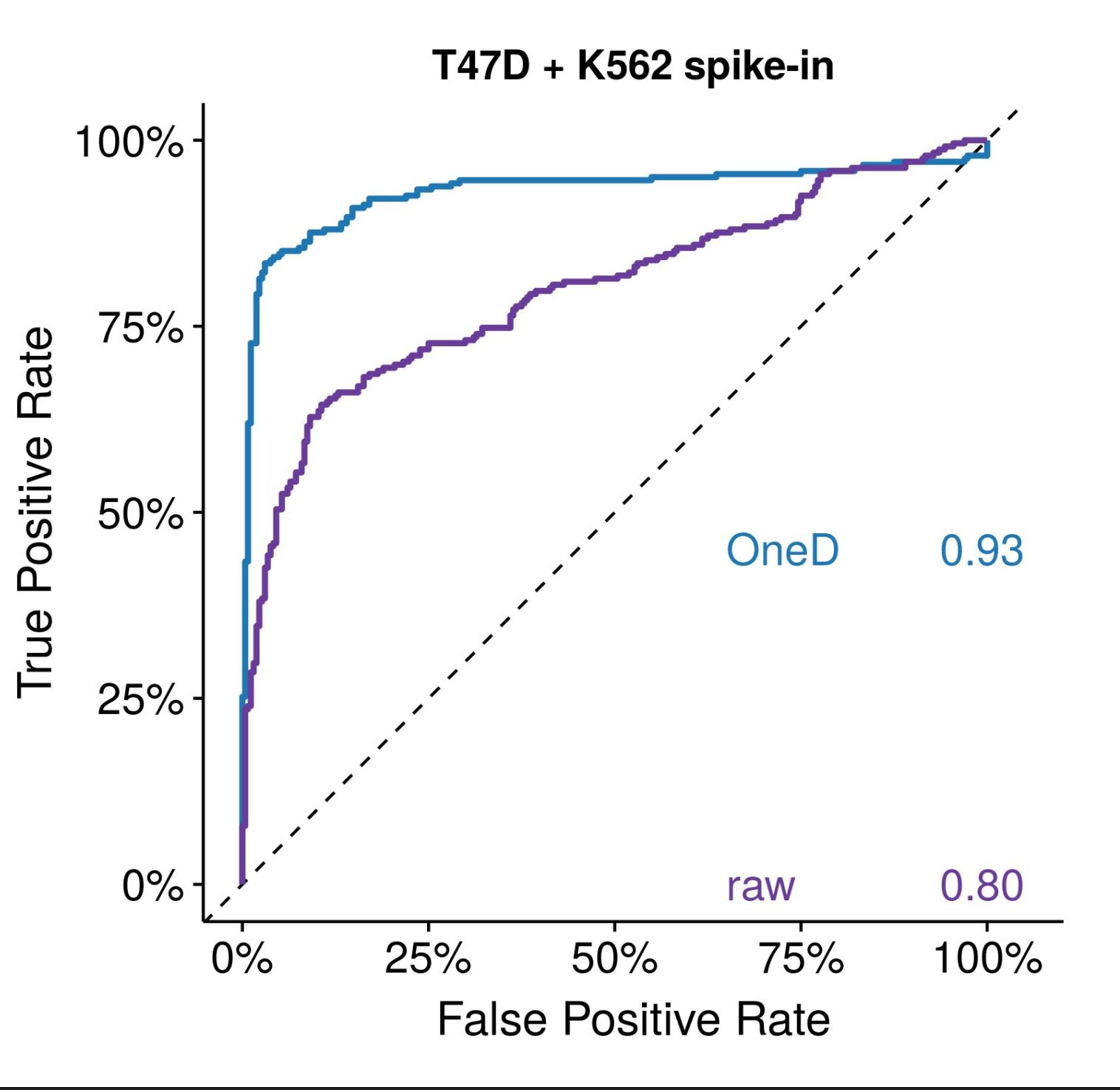
Restriction enzyme, in-situ / diluted, lab, etc ...

Different cell types

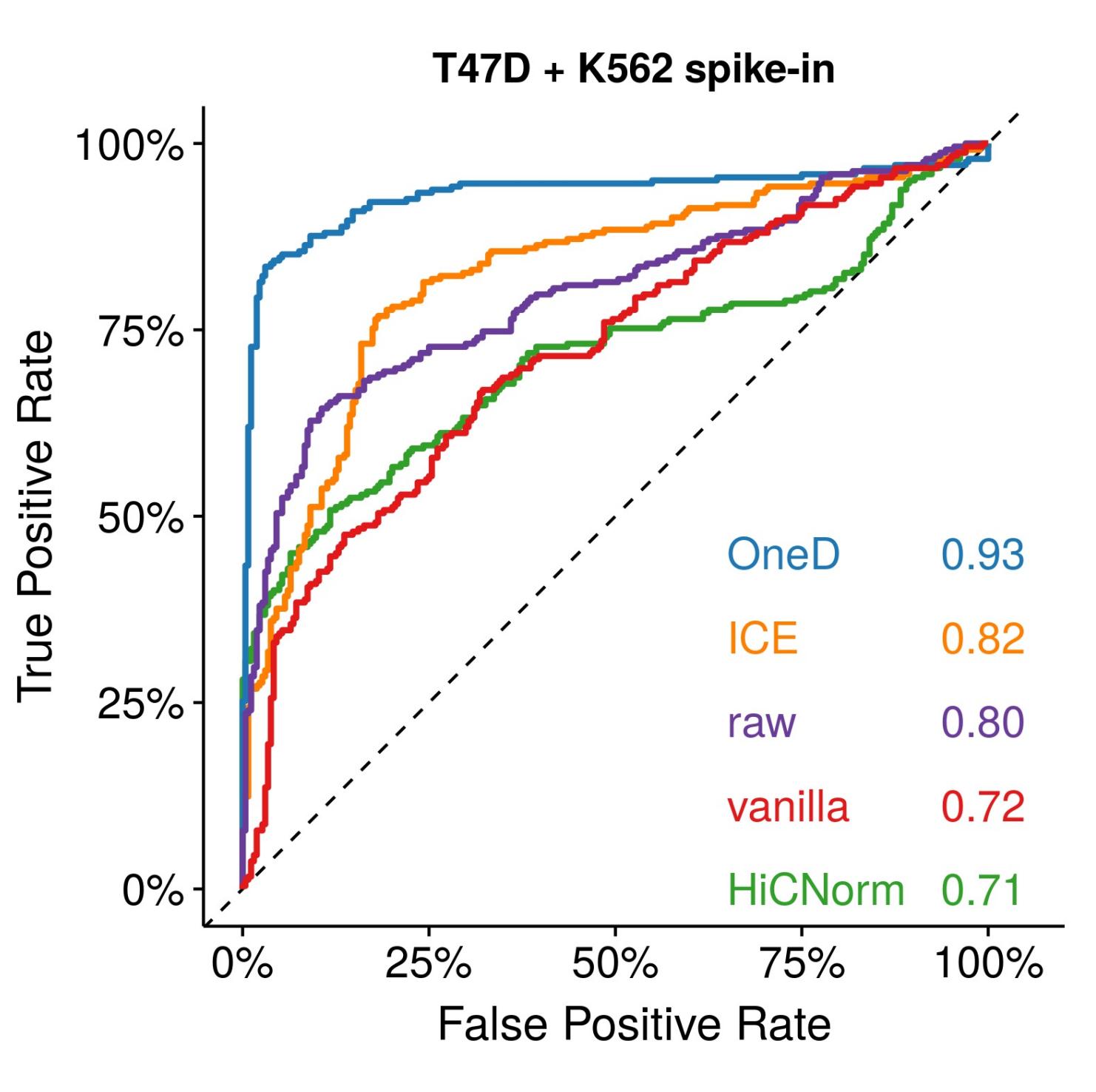
Pair-wise comparison



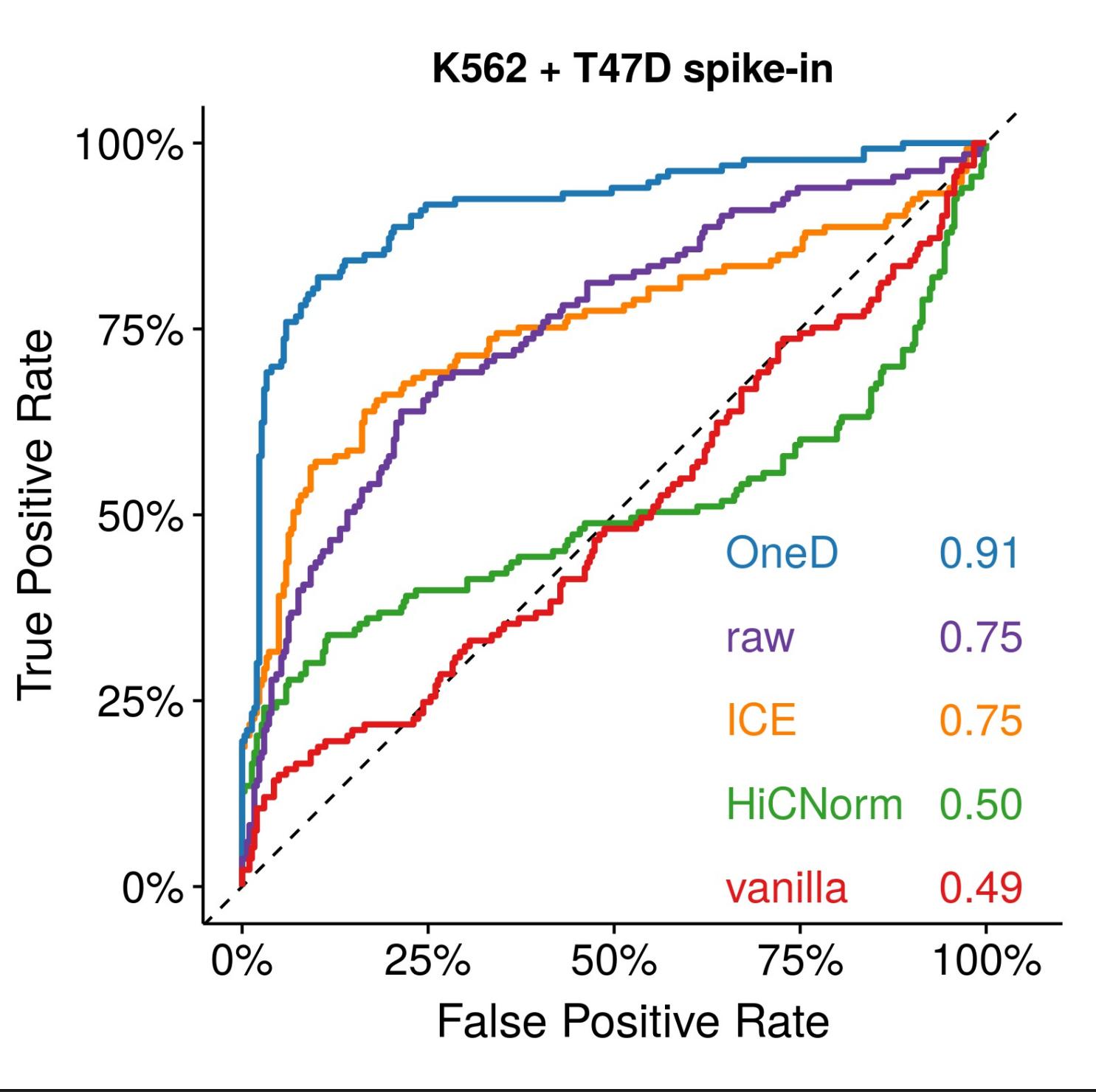
ROC and AUC



ROC and AUC



ROC and AUC



Wish list

♥ Suitable aberrant karyotypes ✓

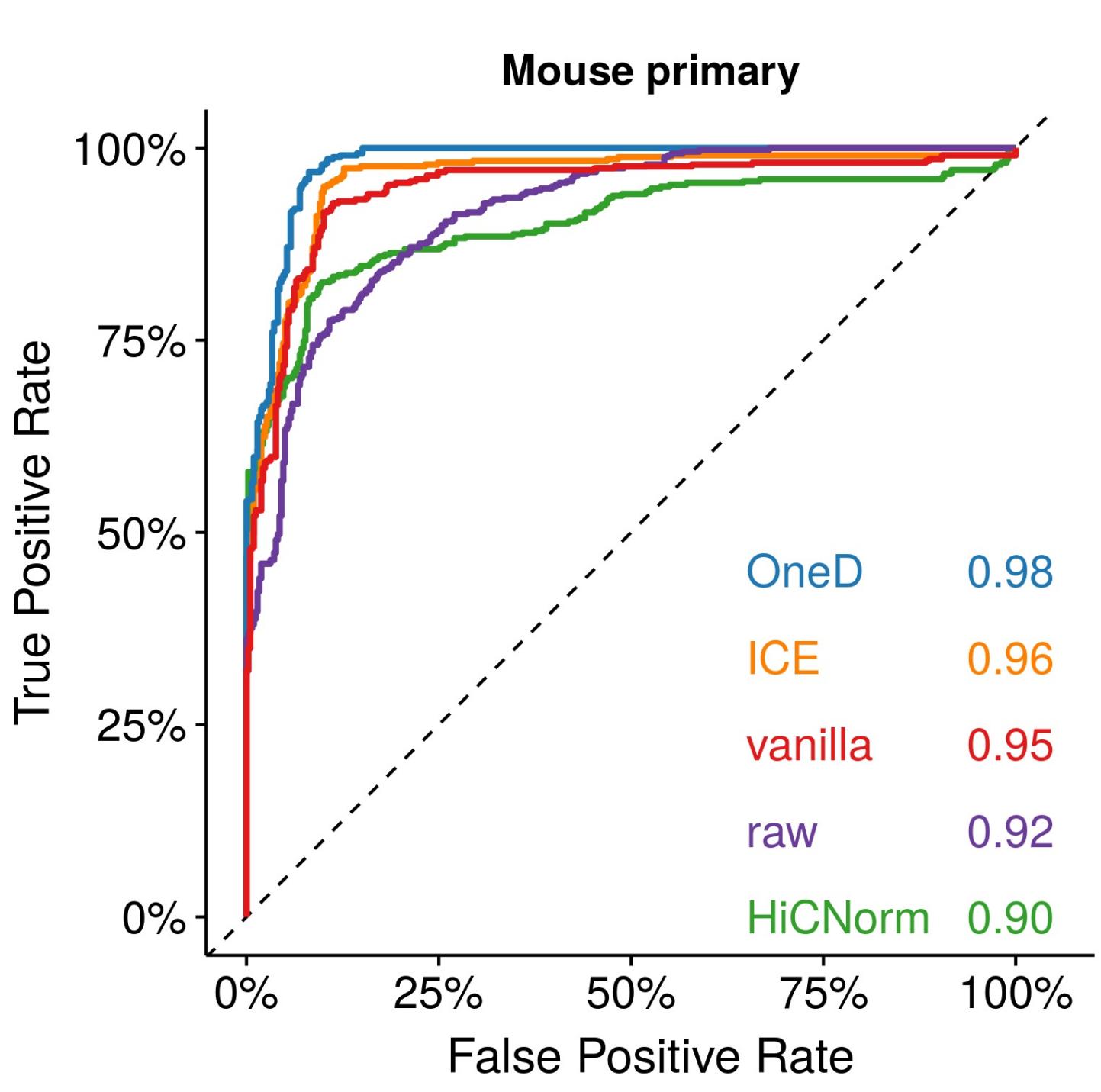
♥ As good as Better than existing ✓

♥ Not worse on normal karyotypes

♥ Fast

♥ Usable at high resolution

ROC and AUC (diploid)



Wish list

♥ Suitable aberrant karyotypes ✓

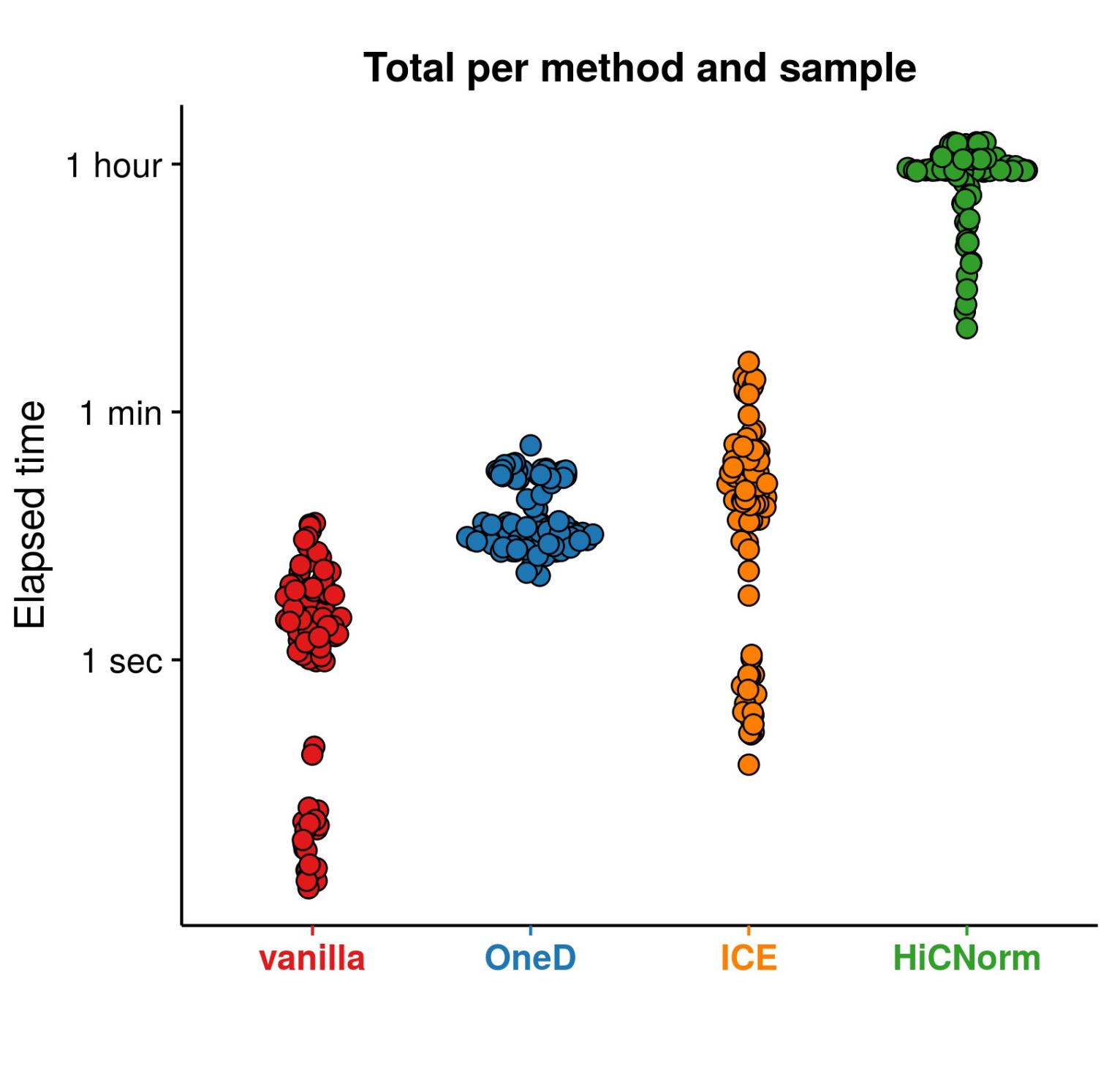
♥ Better than existing ✓

♥ Not worse on normal karyotypes ✓

♥ Fast

♥ Usable at high resolution

Speed



Wish list

♥ Suitable aberrant karyotypes ✓

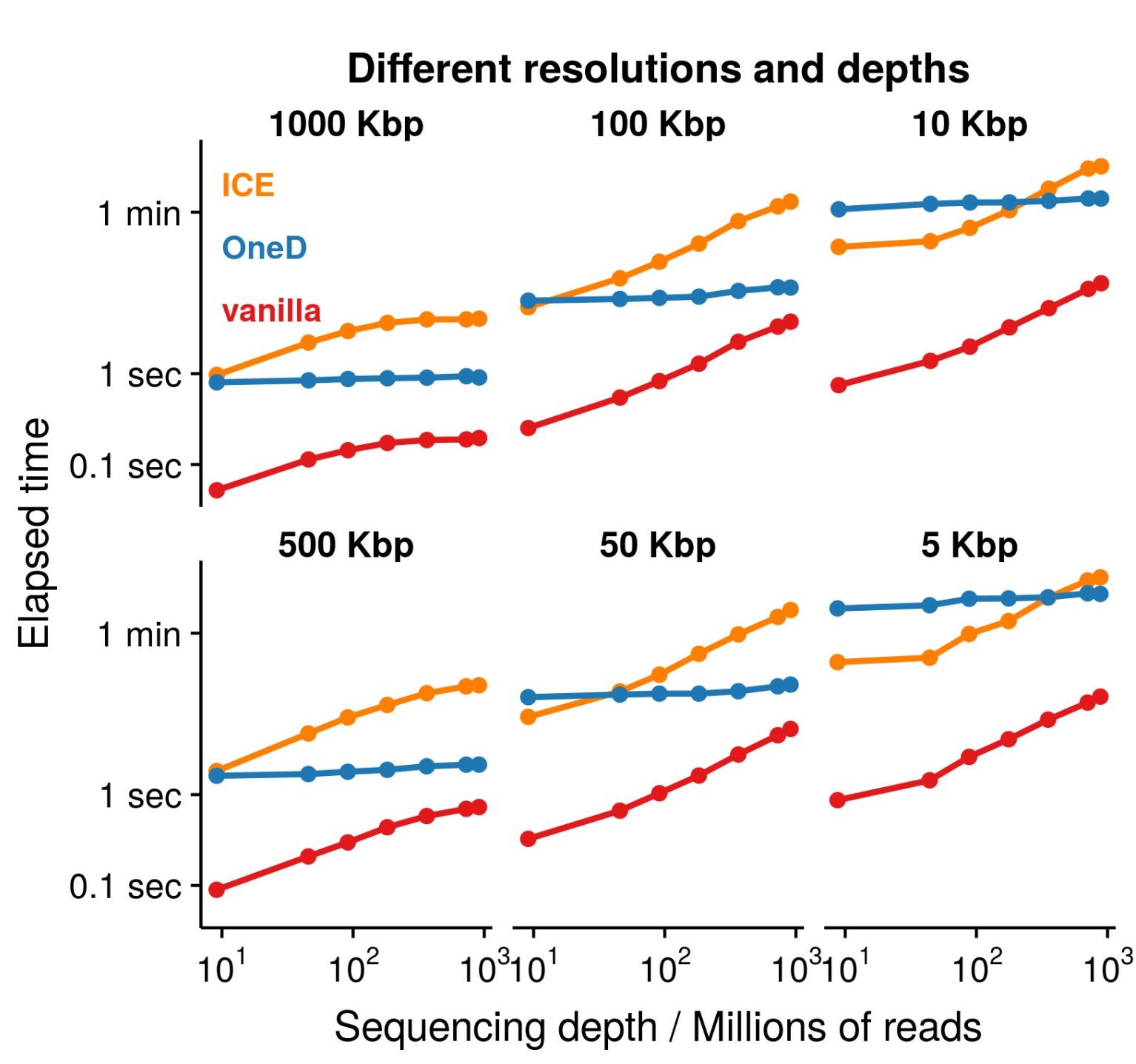
♥ Better than existing ✓

♥ Not worse on normal karyotypes ✓

♥ Fast ✓

♥ Usable at high resolution

Resolution



Wish list

♥ Suitable aberrant karyotypes ✓

♥ Better than existing ✓

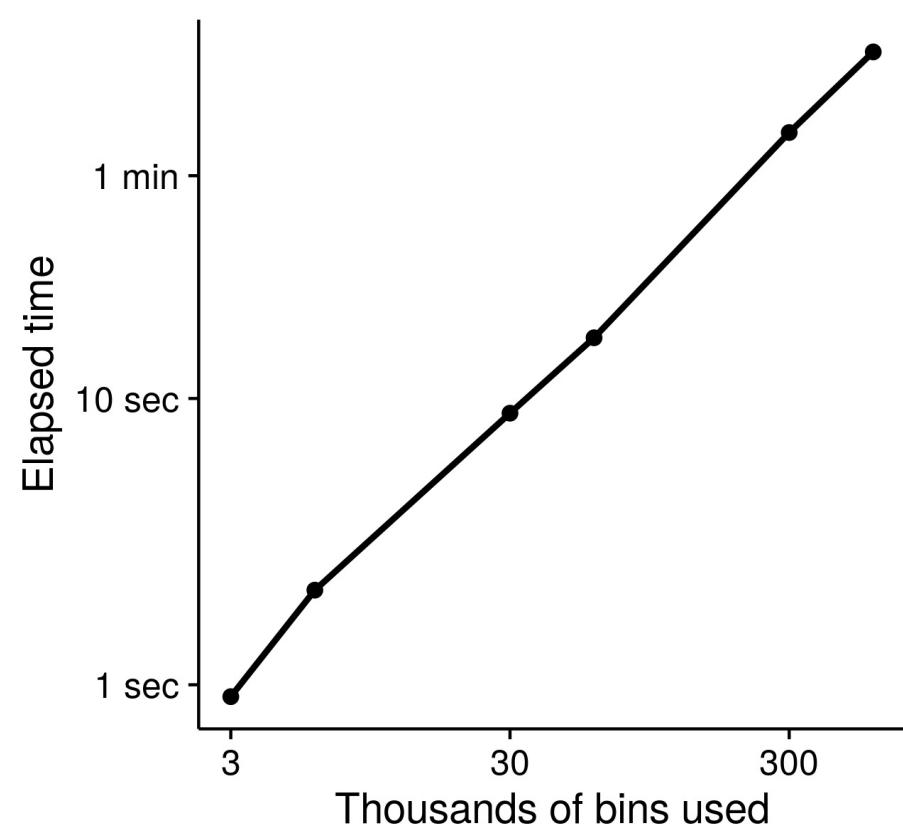
♥ Not worse on normal karyotypes ✓

♥ Fast ✓

♥ Usable at high resolution ✓

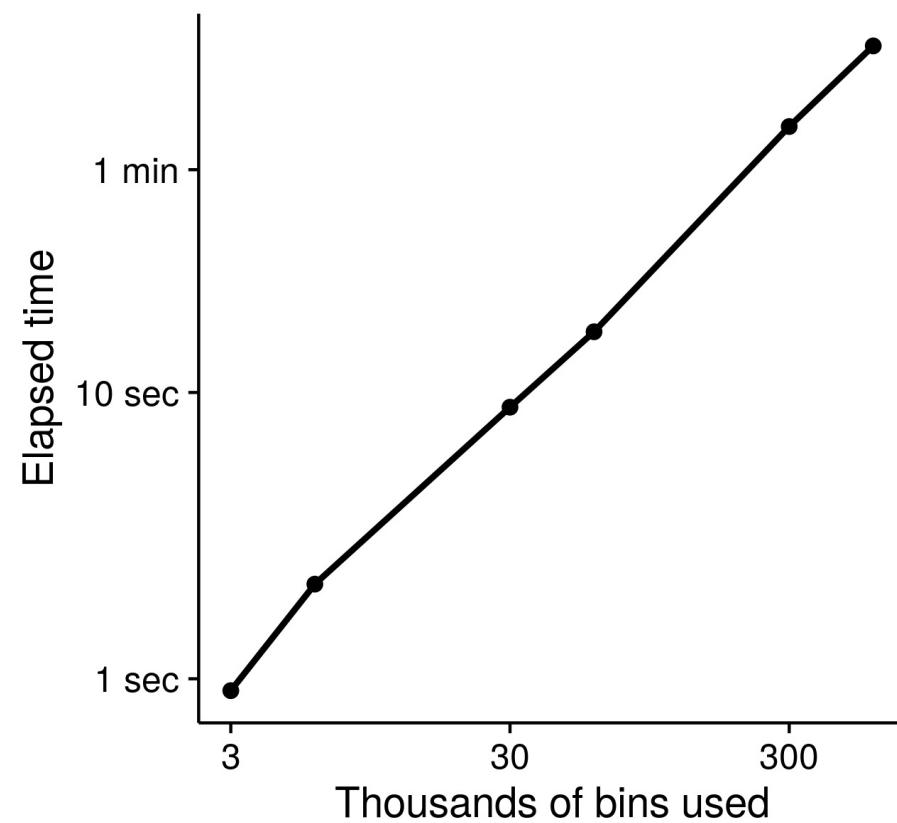
Resolution (with some tricks)

Resolution (with some tricks)

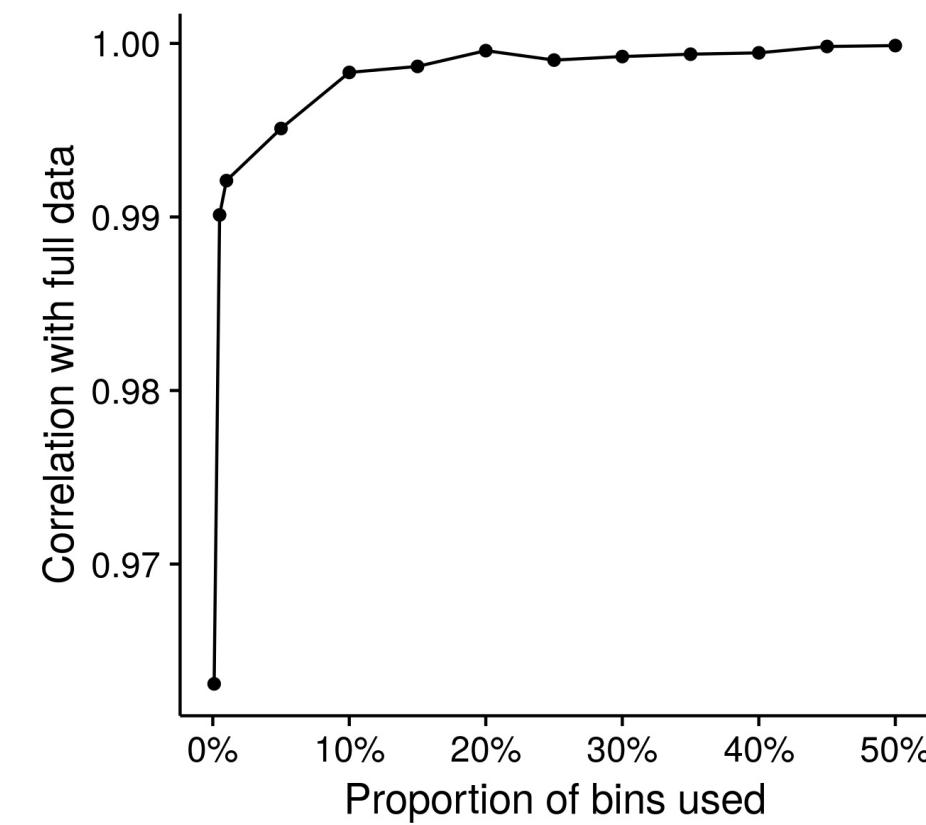


Less data, less time

Resolution (with some tricks)



Less data, less time



Less data, same performance

Summary

Summary

Hi-C can interrogate genome structure

Hi-C experiments have biases (including copy number)

Summary

Hi-C can interrogate genome structure

Hi-C experiments have biases (including copy number)

OneD removes biases increasing reproducibility

OneD is fast and suitable at high resolutions

OneD availability

R package

<https://github.com/qenvio/dryhic>

Paper

Vidal, E. et al. (2018)

Nucleic acids research

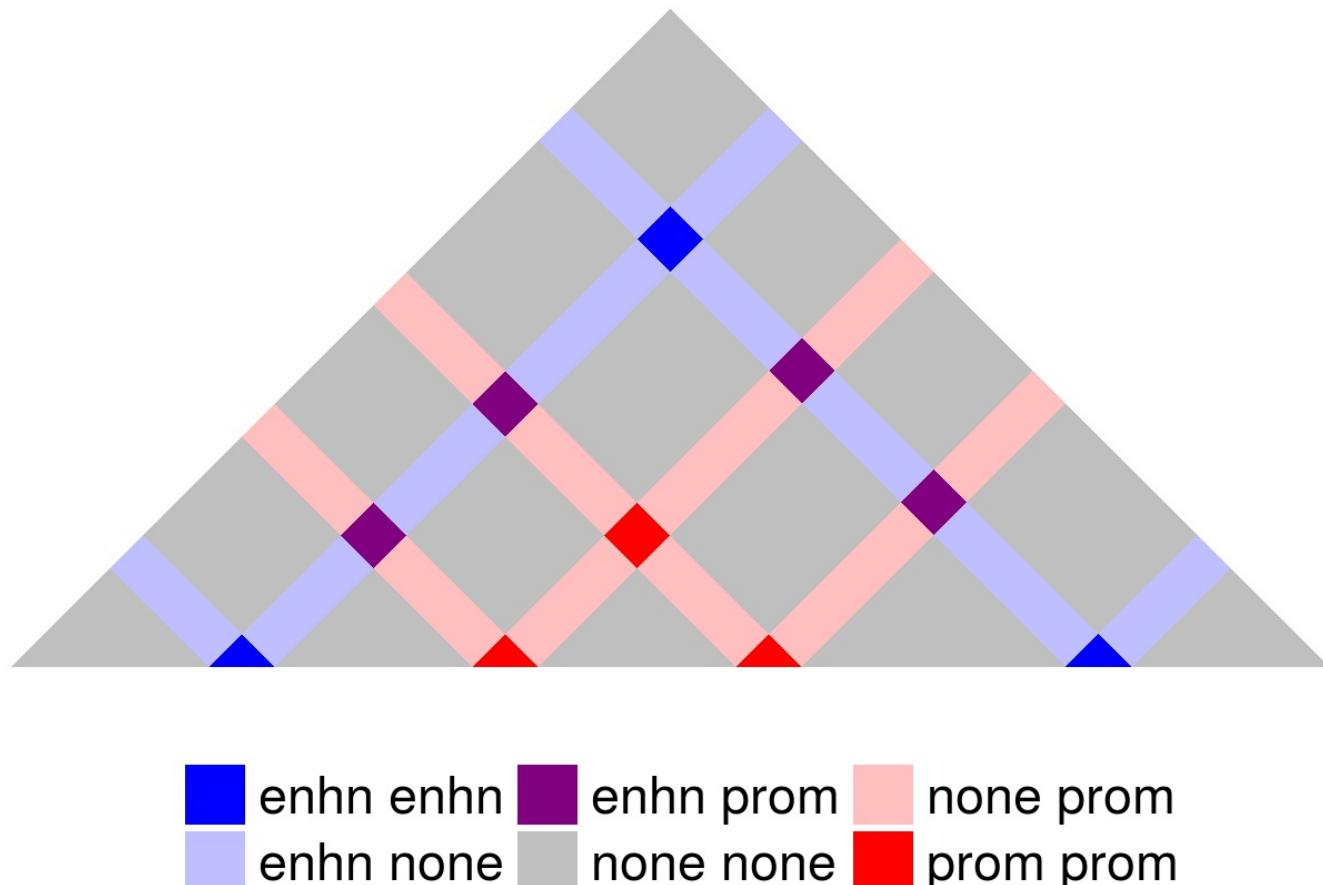
<https://doi.org/10.1093/nar/gky064>

The inner life of TADs

Genome-wide enrichment of contacts between regions

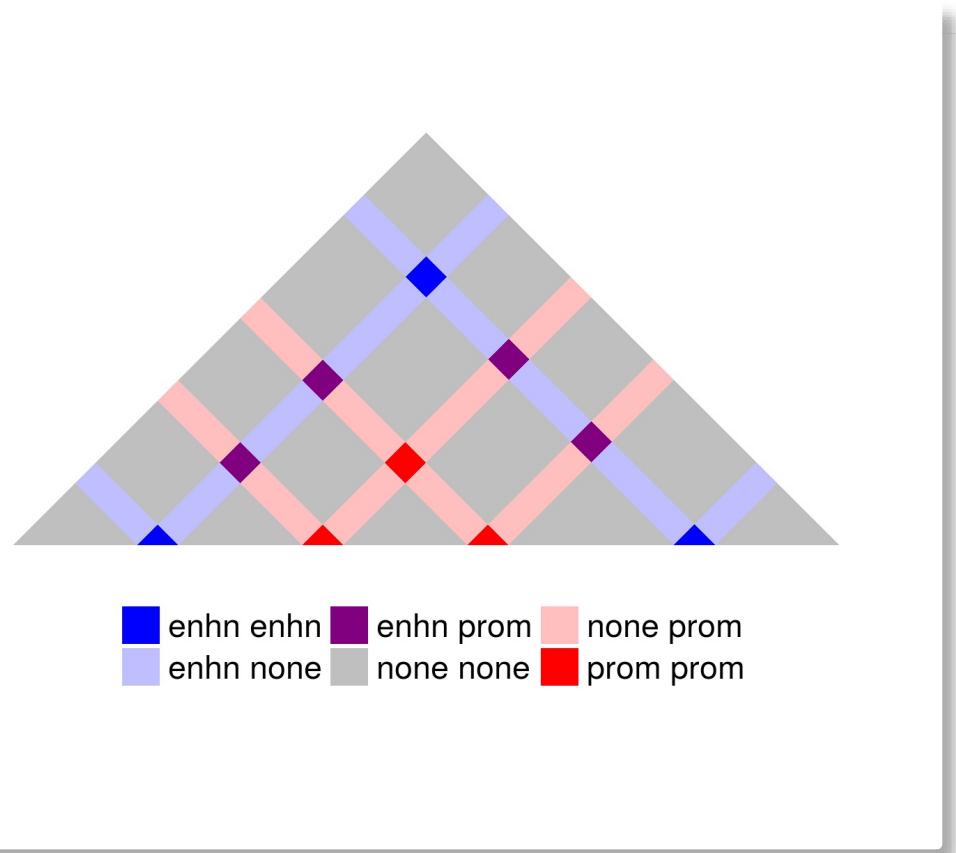
(i.e. promoters and enhancers)

Types of bins inside a TAD

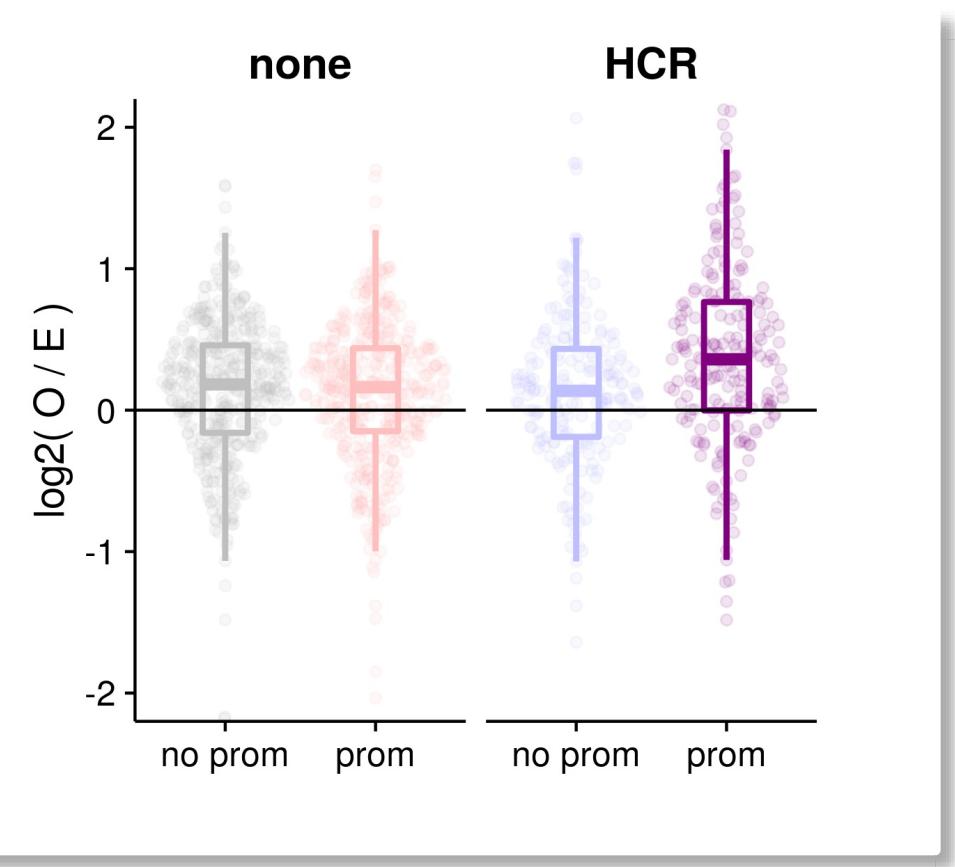
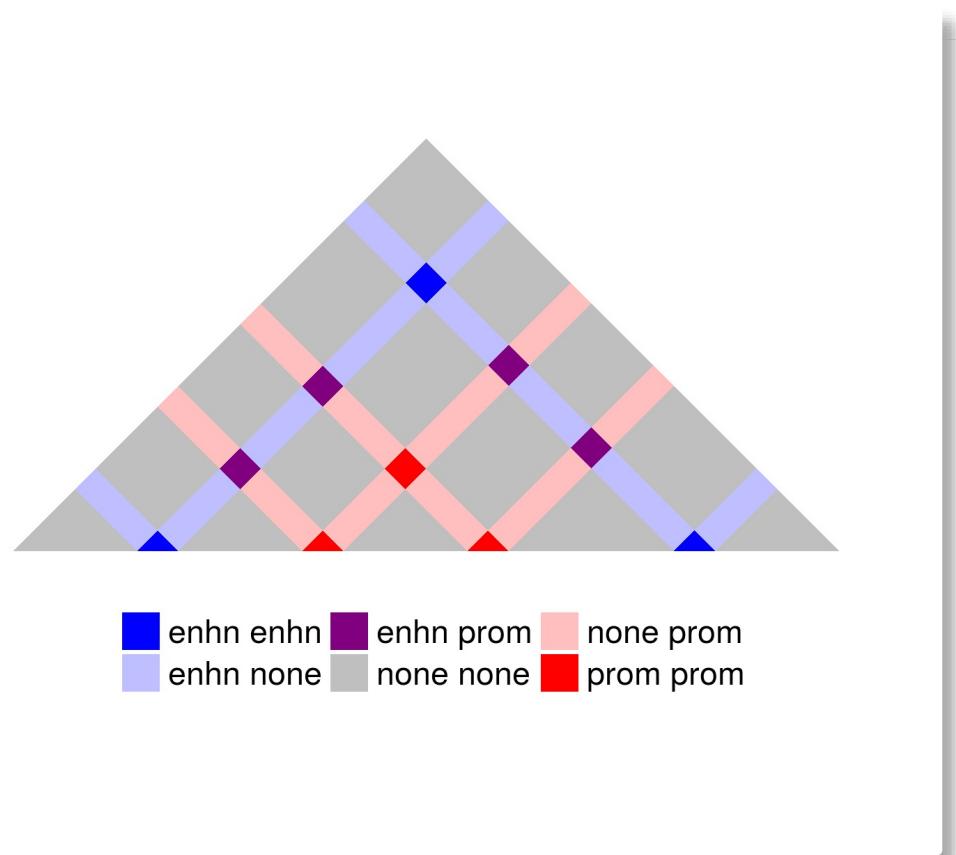


Example: HCR and H3K4me3

Example: HCR and H3K4me3



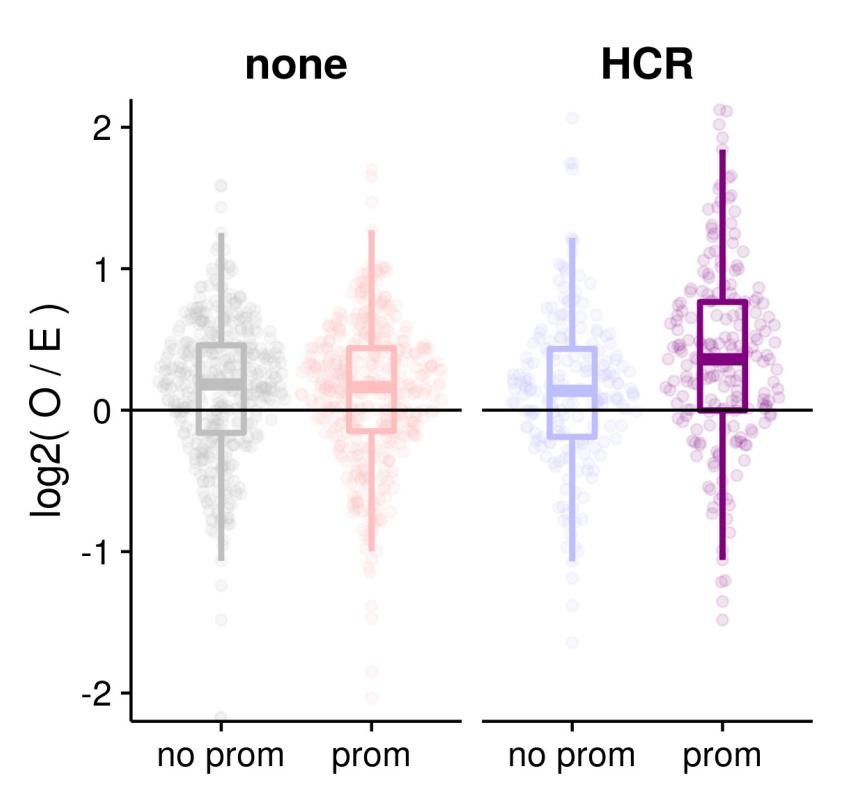
Example: HCR and H3K4me3



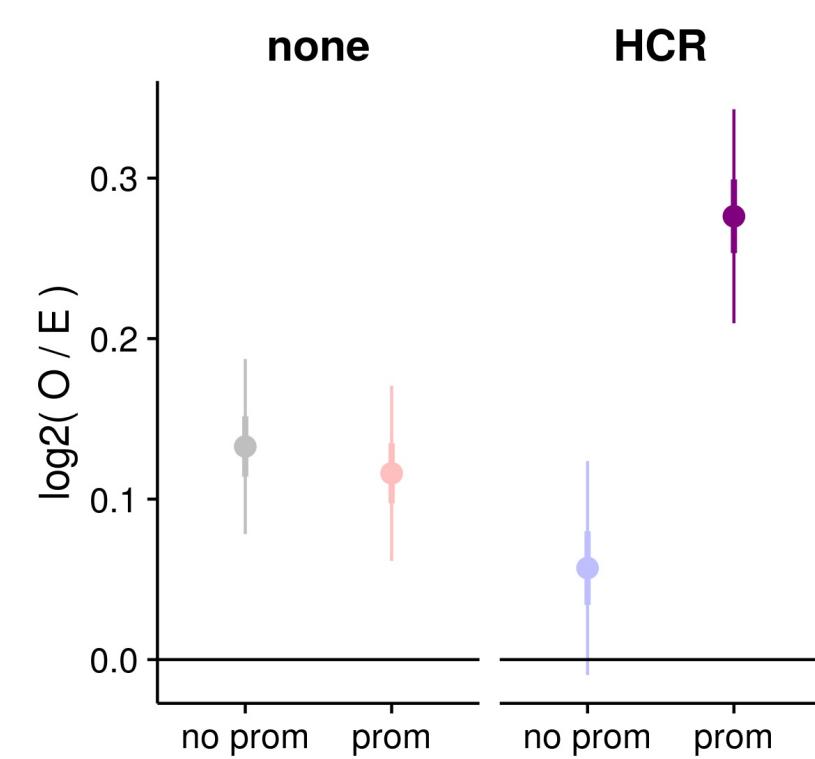
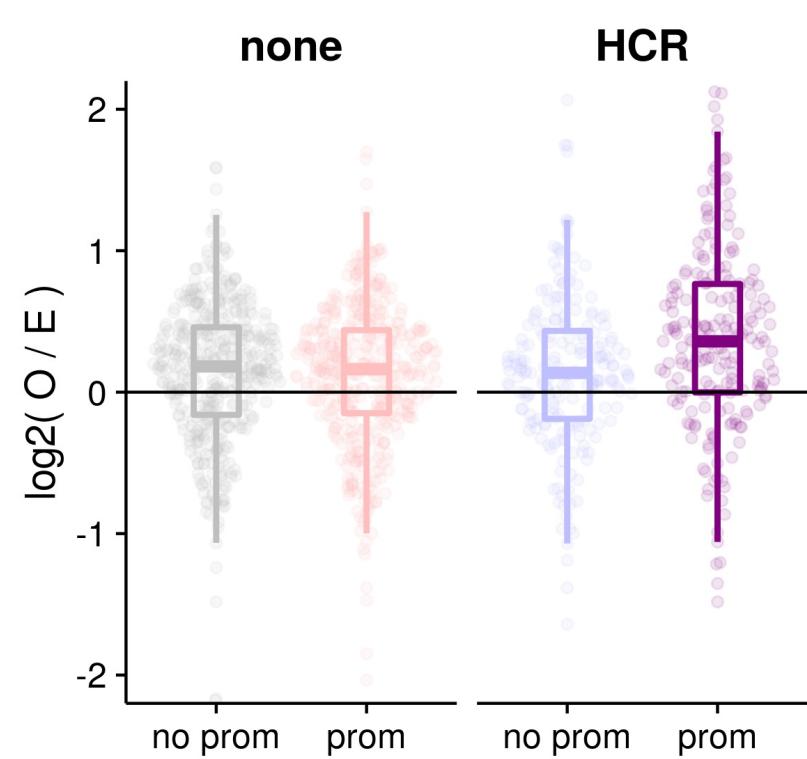
Le Dily, F. and Vidal, E. et al. (2017)

bioRxiv
<https://doi.org/10.1101/233874>

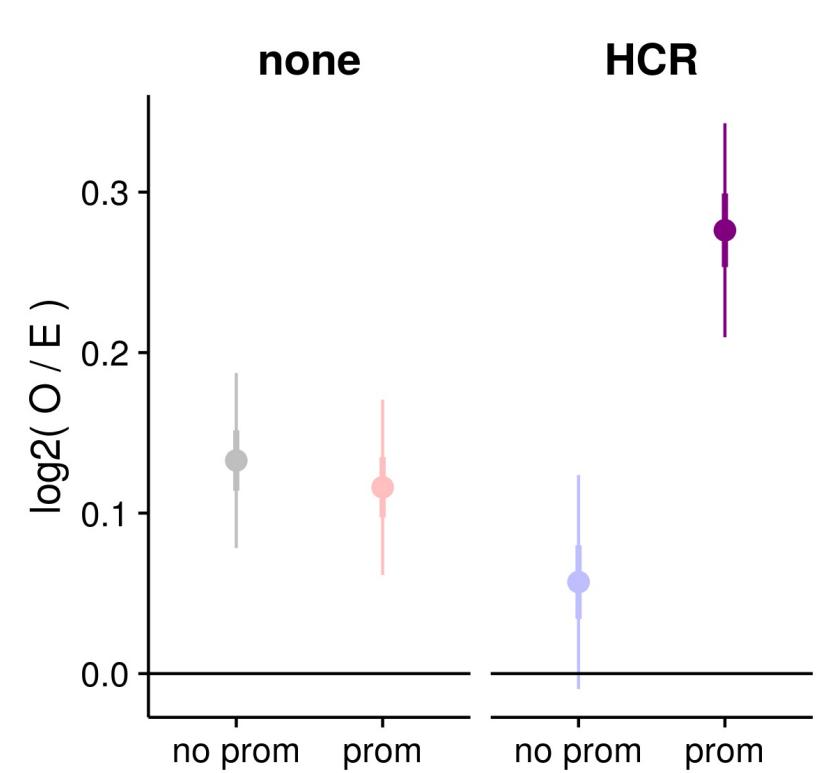
Summarizing via linear mixed model



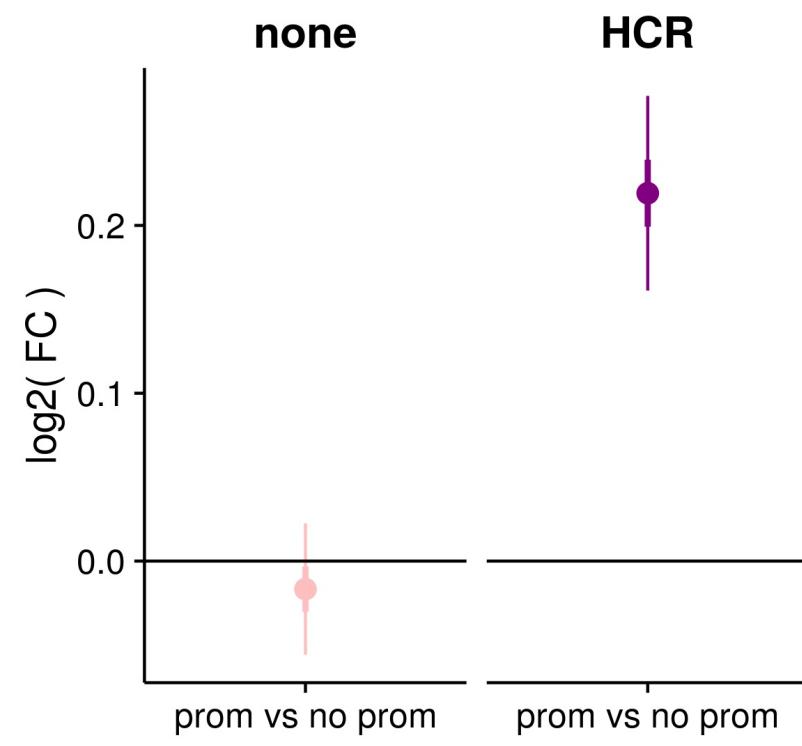
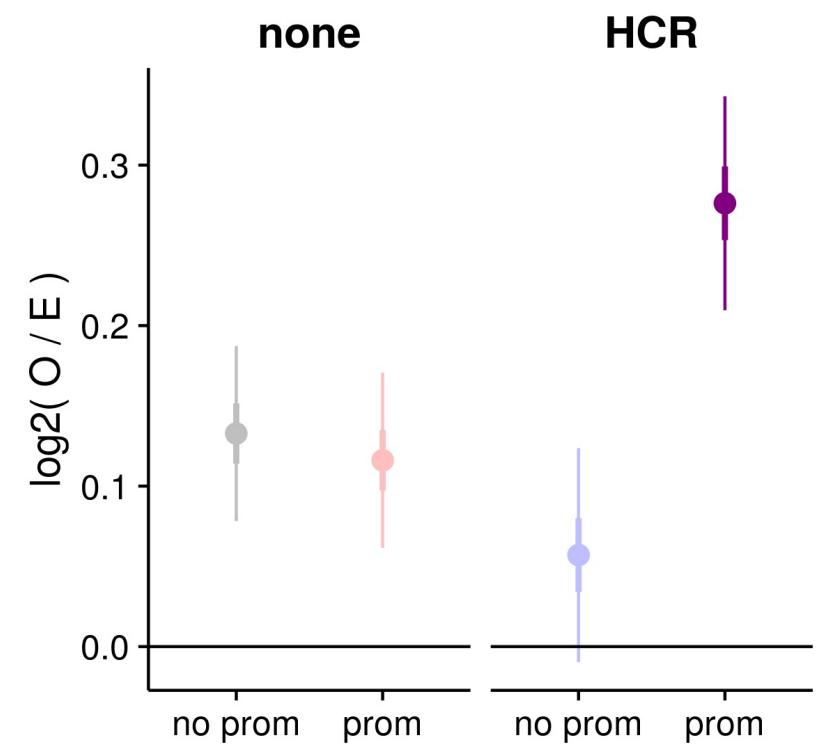
Summarizing via linear mixed model



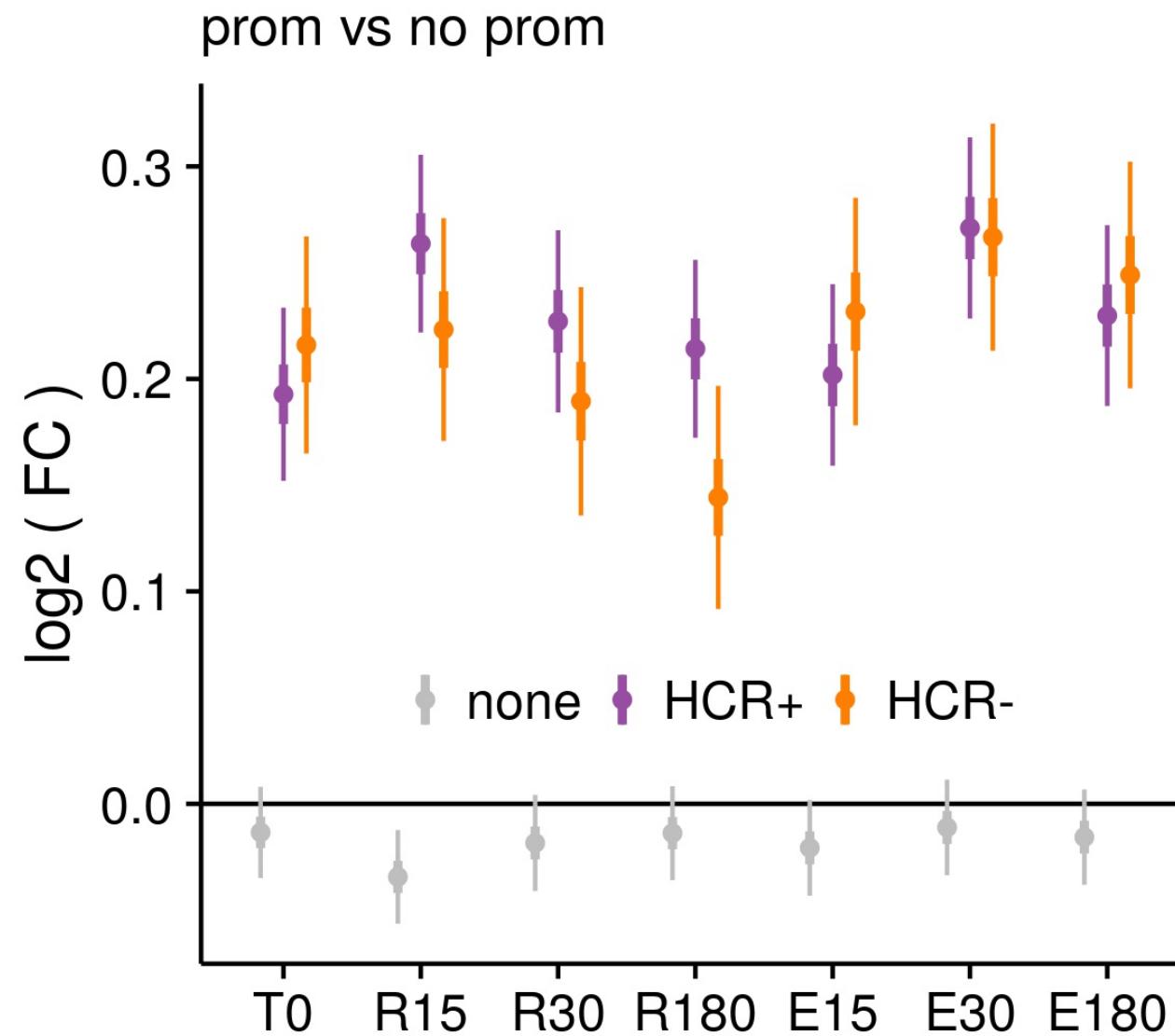
Simple comparison



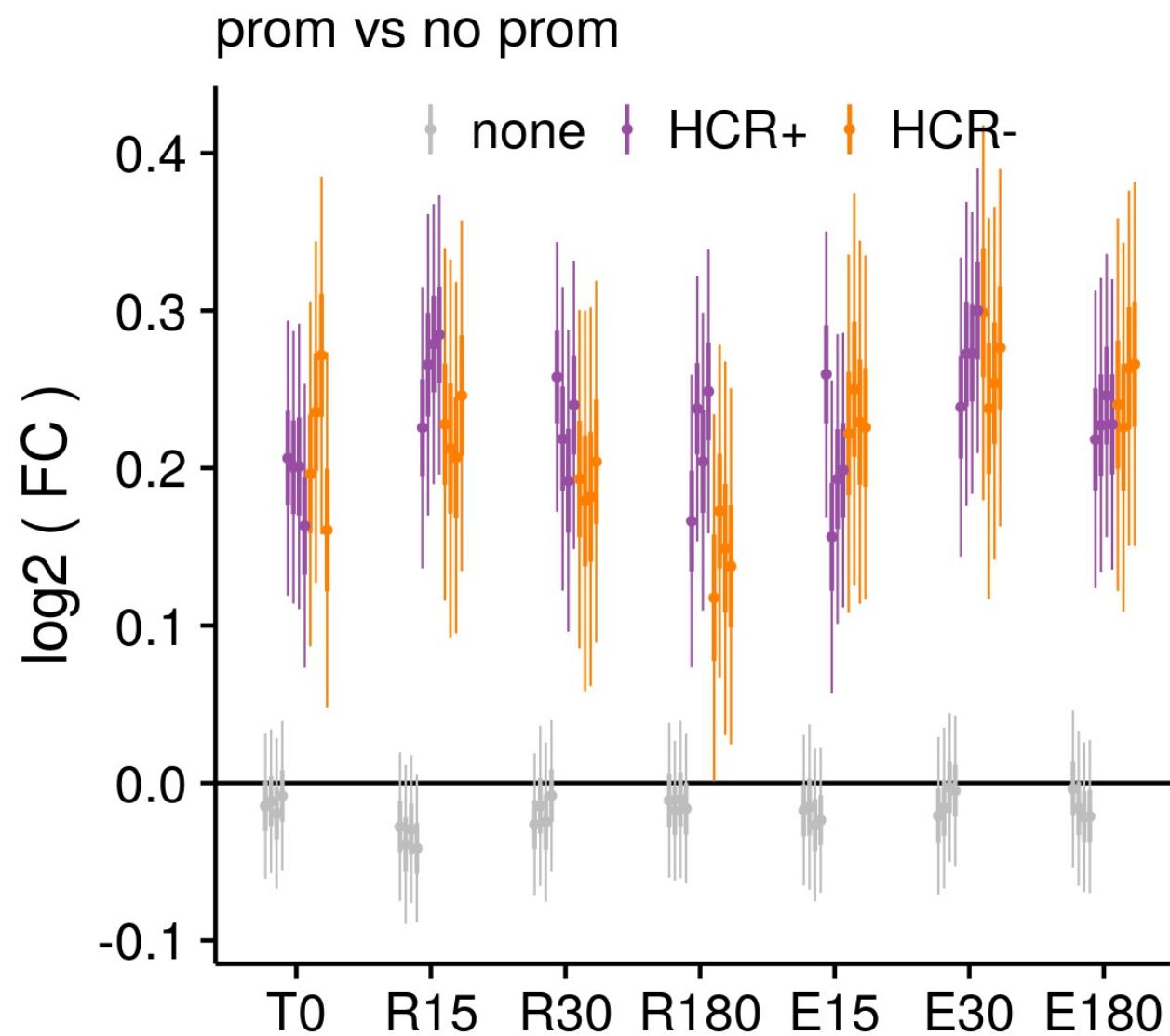
Simple comparison



More samples

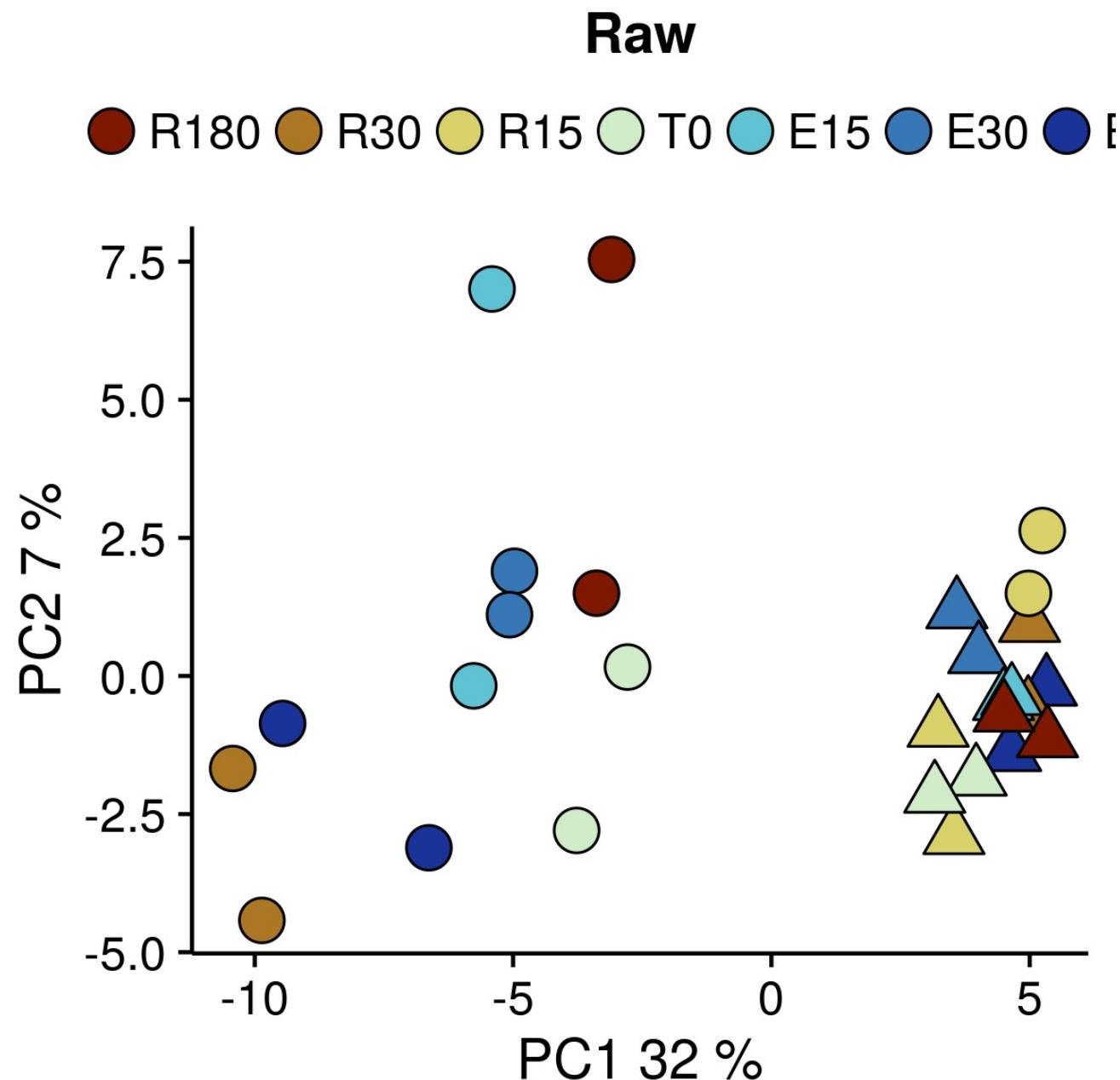


Even more samples (replicates!)

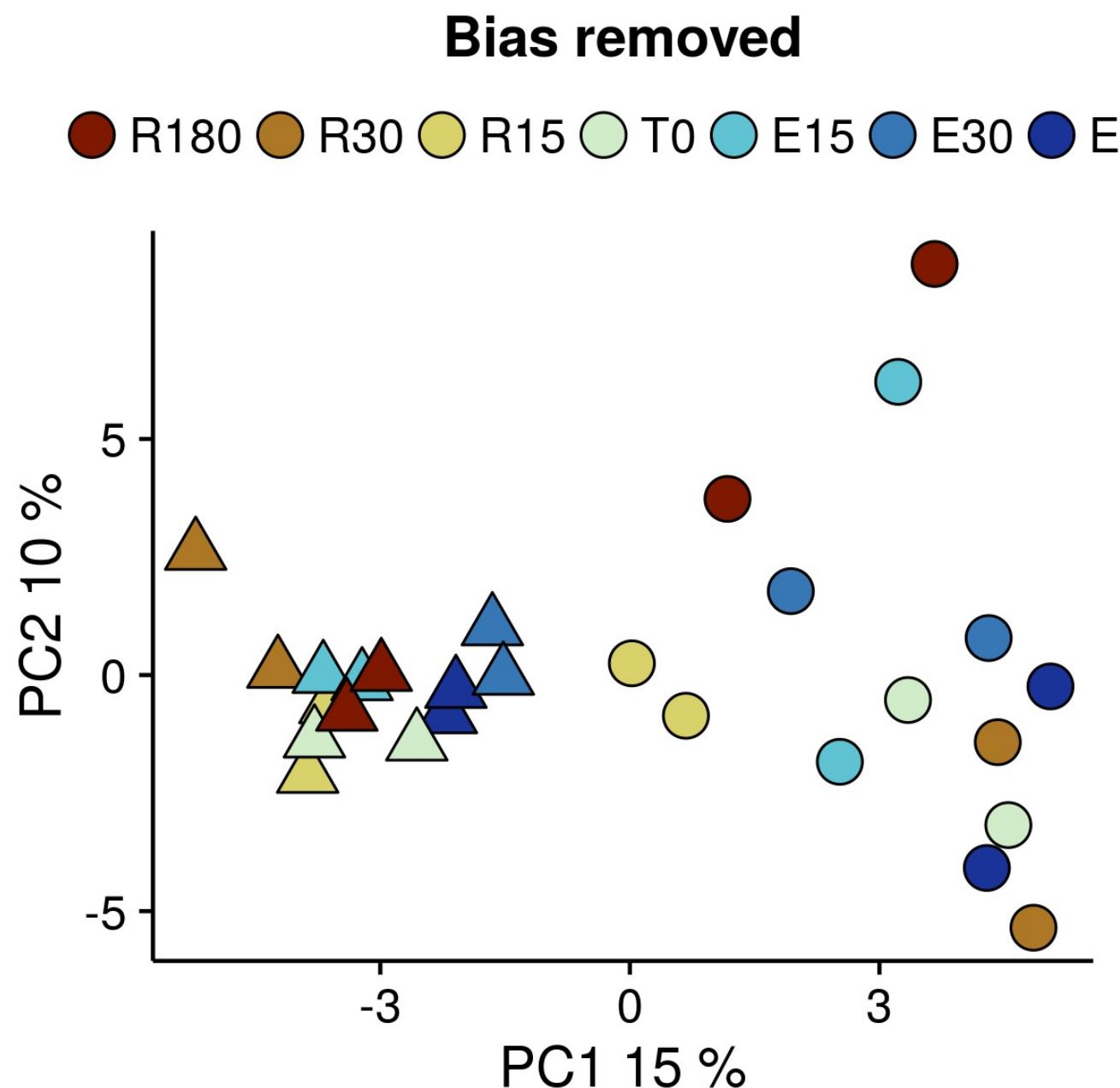


Do we have batch effect?

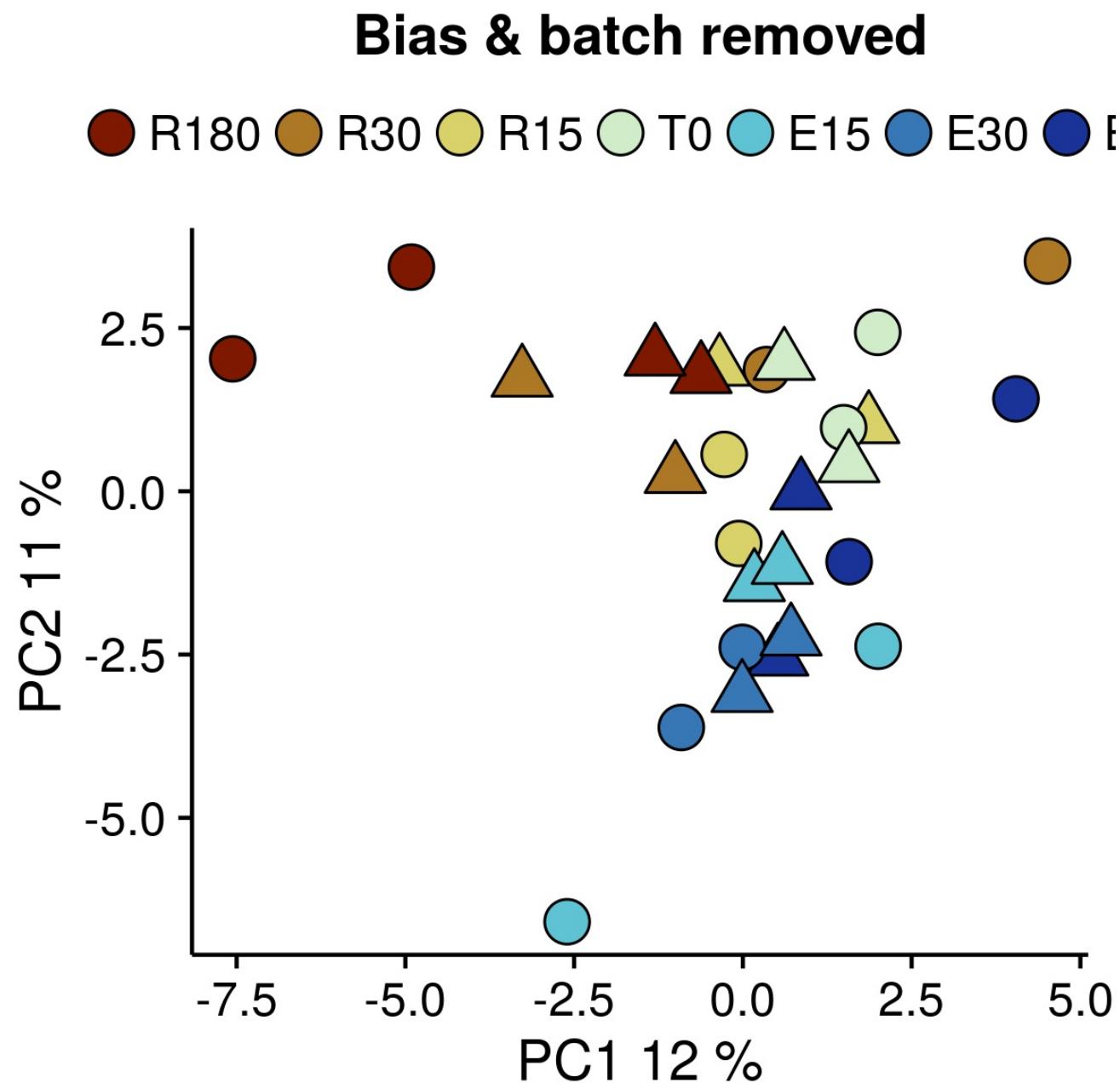
Do we have batch effect?



Removing biases helps ...

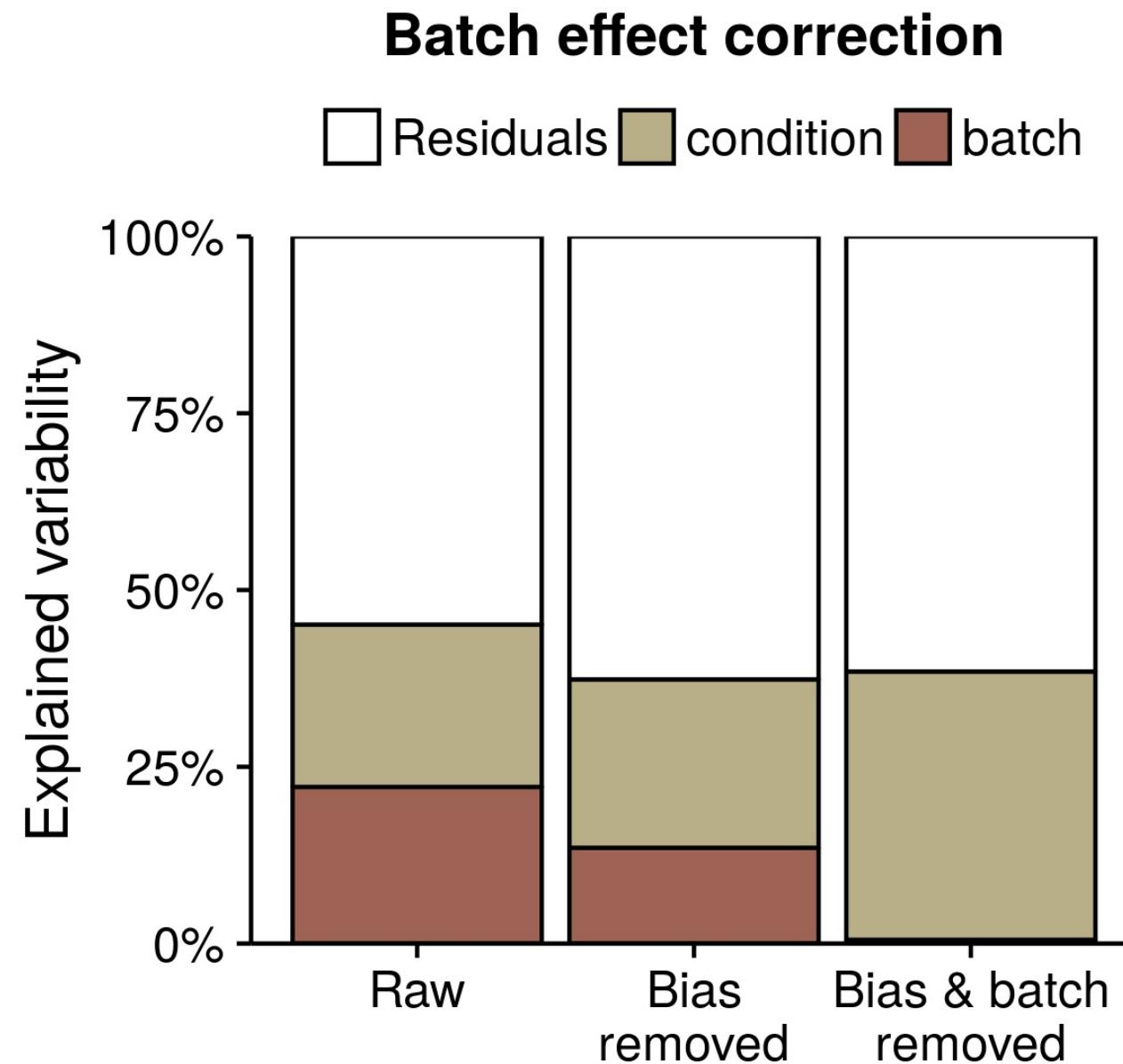


... and removing batch effect is even better



Total variance explained

Total variance explained



Grand summary

Grand summary

Genome structure plays a role and can be measured

Grand summary

Genome structure plays a role and can be measured

Increasing production of Hi-C data

Grand summary

Genome structure plays a role and can be measured

Increasing production of Hi-C data

Methodological challenges to analyze and integrate Hi-C data

Grand summary

Genome structure plays a role and can be measured

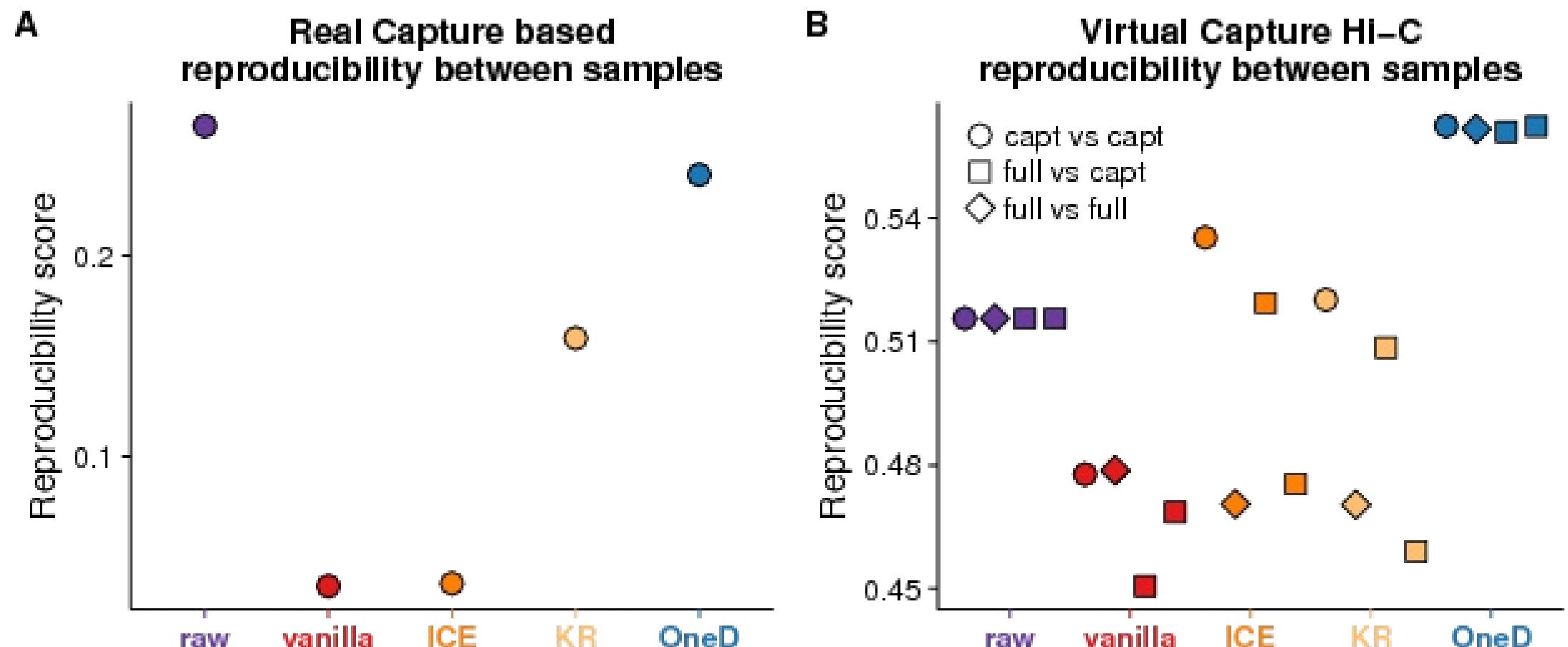
Increasing production of Hi-C data

Methodological challenges to analyze and integrate Hi-C data

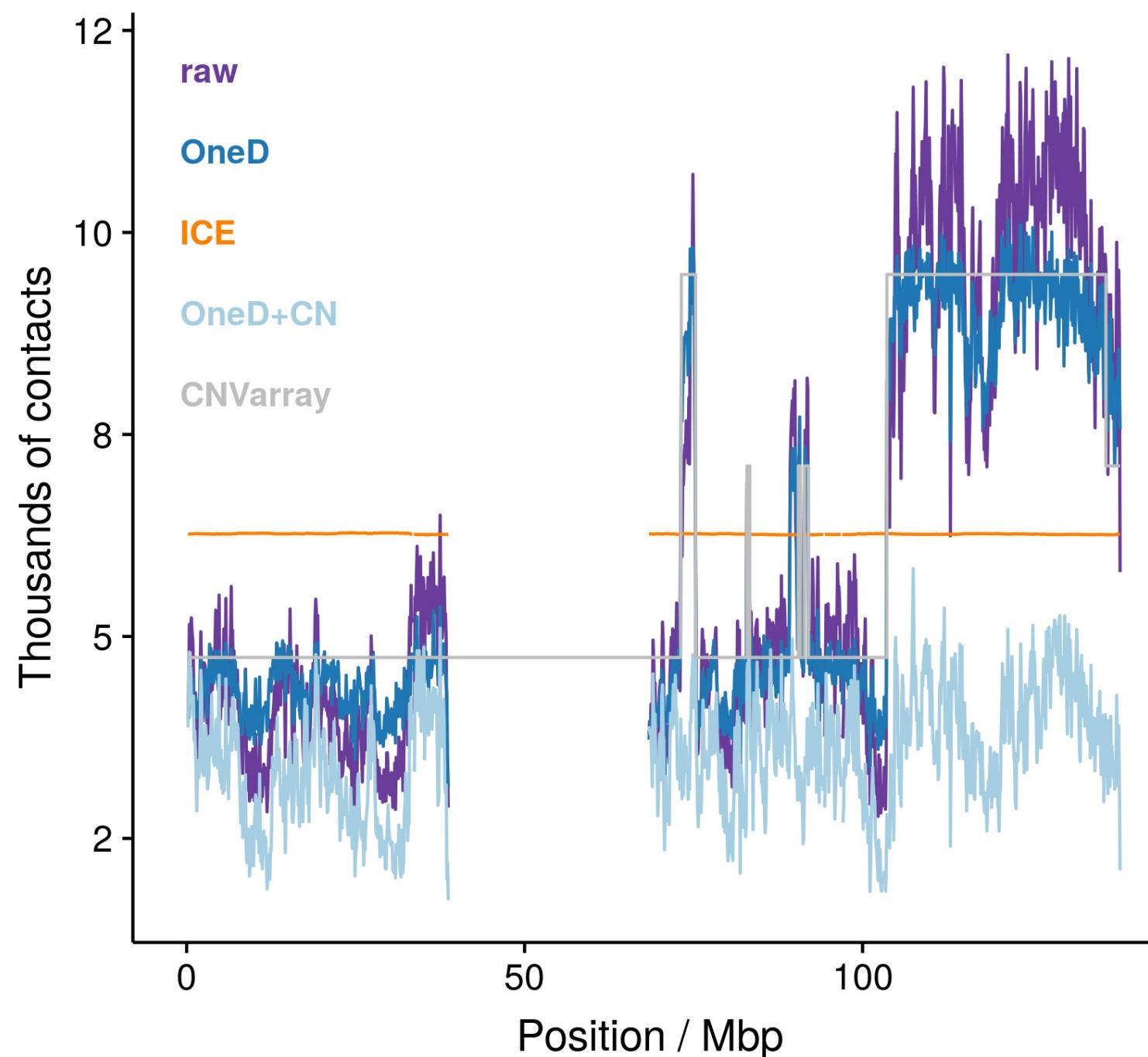
OneD:  [qenvio/dryhic](#)

OneD extras

Incomplete designs

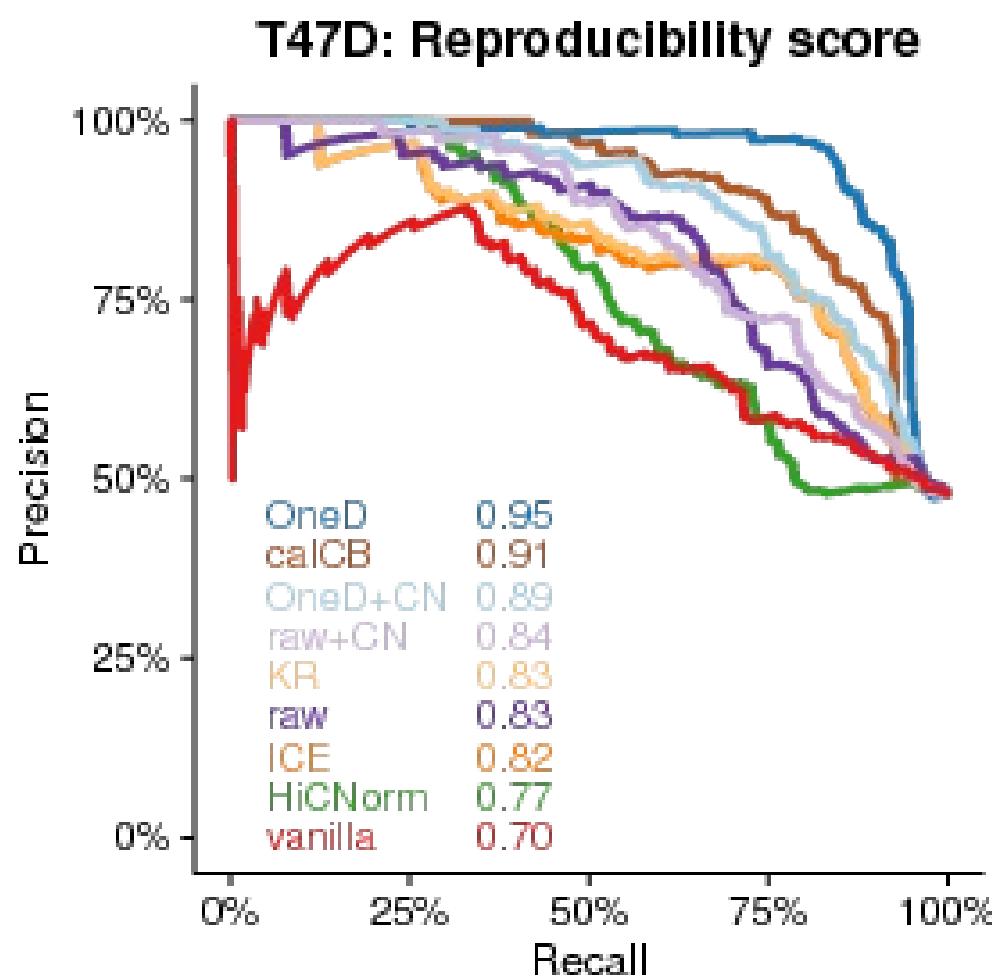


Copy number estimation

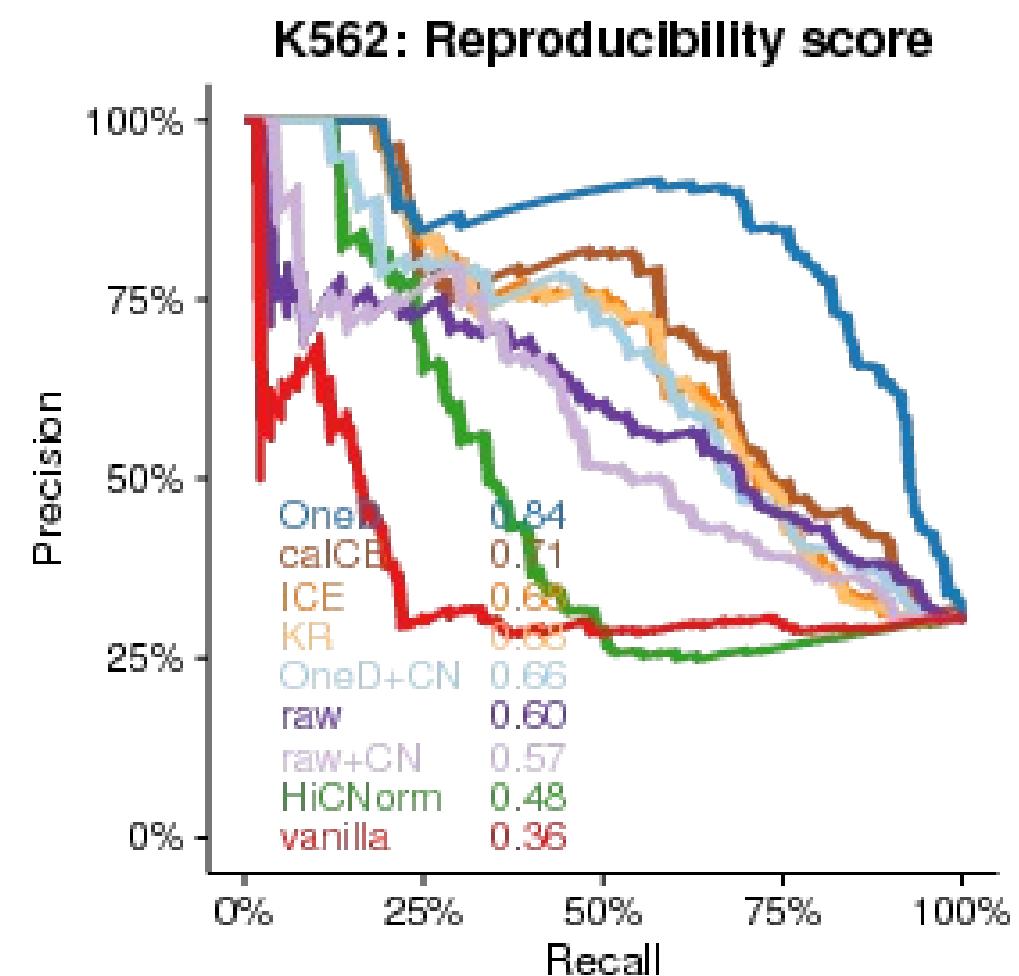


Precision-Recall

C



D



Precision-Recall

