

Measuring significant relationships between sets of genomic features

Aaron Quinlan

quinlanlab.org

April 25, 2013

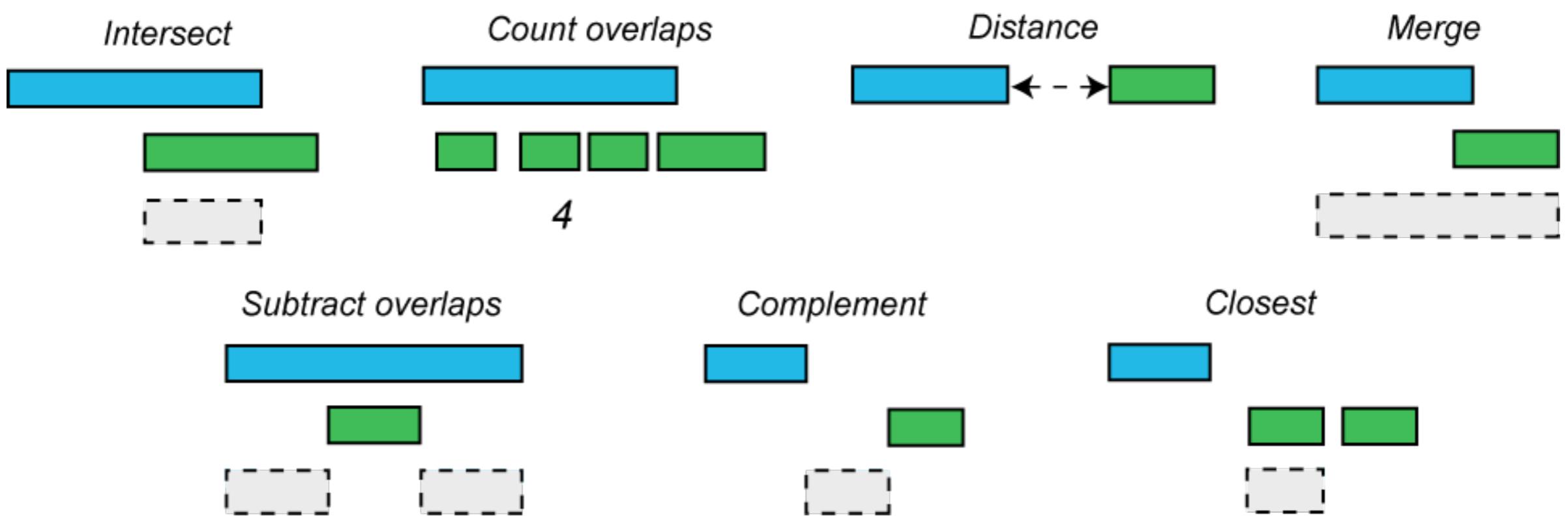
Public Health Sciences

Center for Public Health Genomics

Biochemistry and Molecular Genetics



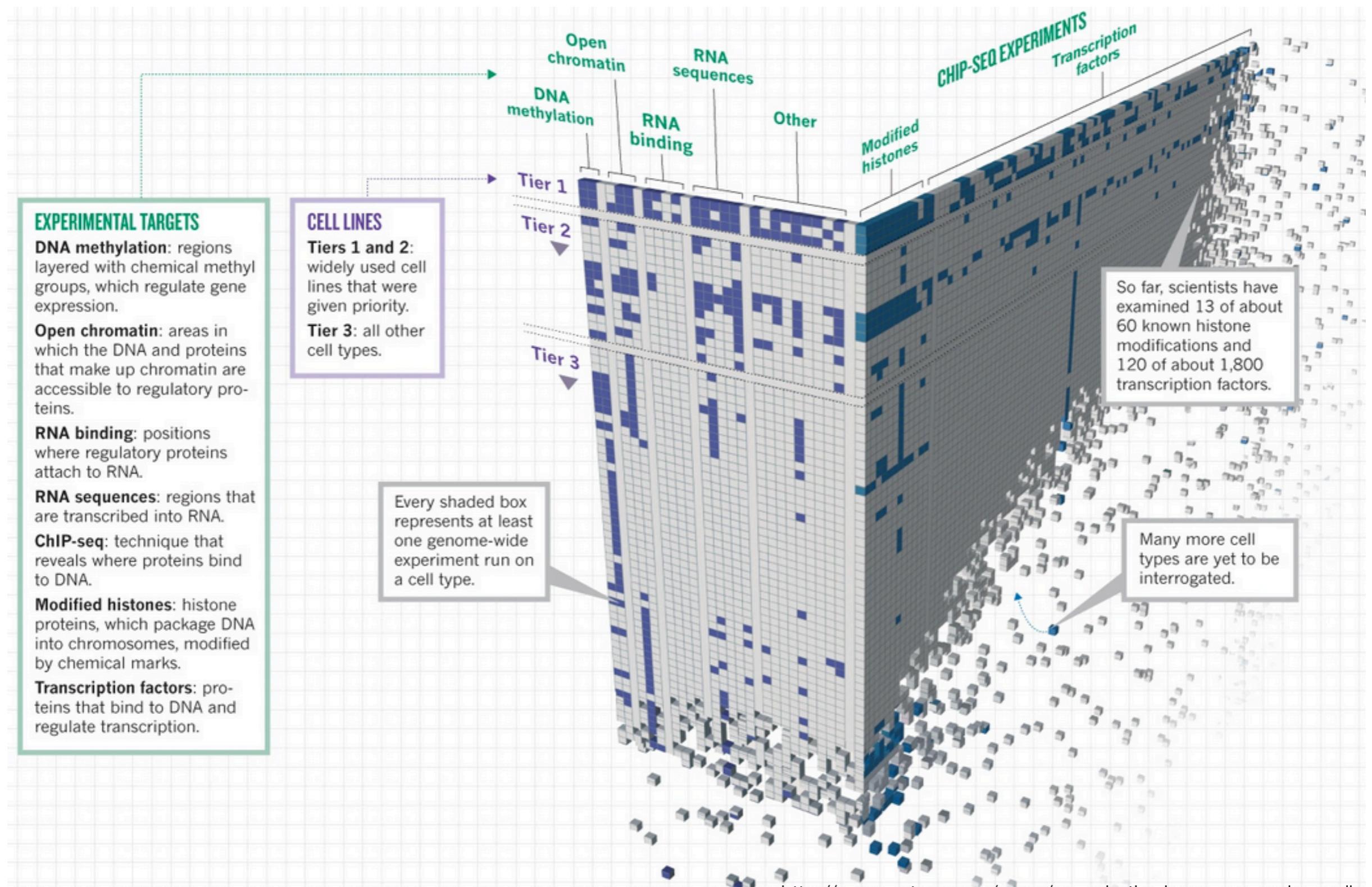
genome arithmetic



analytic limitations

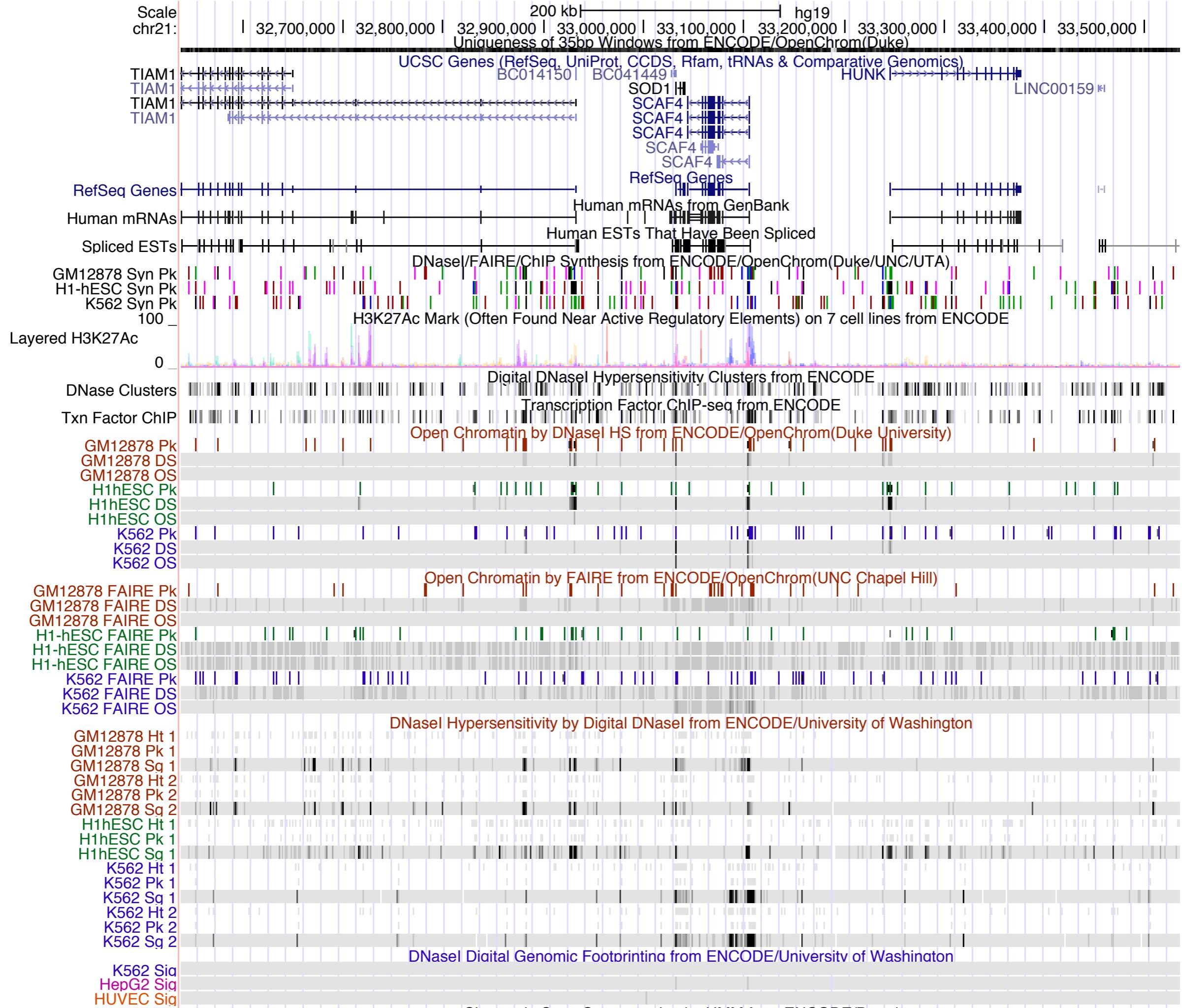
- Existing methods such as bedtools work primarily on pairs of datasets
- These methods exhaustively list the **individual details** of relationships between genome annotations and experimental datasets.
- How do we reduce the huge dimensionality to something we can understand and compare?

How do we make sense of **complex, multi-dimensional** datasets to gain insights into genome biology?



“Low dimensional representations and
human understanding are synonymous.”

Ewan Birney



This is a hard yet important problem.

- Understand every base pair's function (or lack thereof) in different cell types and contexts.
- Challenges (among many):
 - Basic exploratory data analysis: slicing and dicing very large, heterogeneous datasets.
 - Visualization: unbiased exploration; let the data tell its story.
 - **Testing for significant spatial relationships**

Motivation

- How do we decide whether two observations (e.g., ChIP-seq peaks from two different assays) are correlated with one another?
- They overlap, but is it more than one would expect by chance?

How do we detect genomic co-association?

*That is, do two sets of genomic features co-occur
(overlap or have spatial consistency) more than
expected?*

Genome annotation
(e.g. conserved
sequences)



Experimental
dataset

| Any overlap? | Yes | Yes | No | Yes | Yes |
|---------------|-----|-----|----|-----|------|
| Bases overlap | 32% | 40% | 0% | 50% | 100% |

How do we detect genomic co-association?

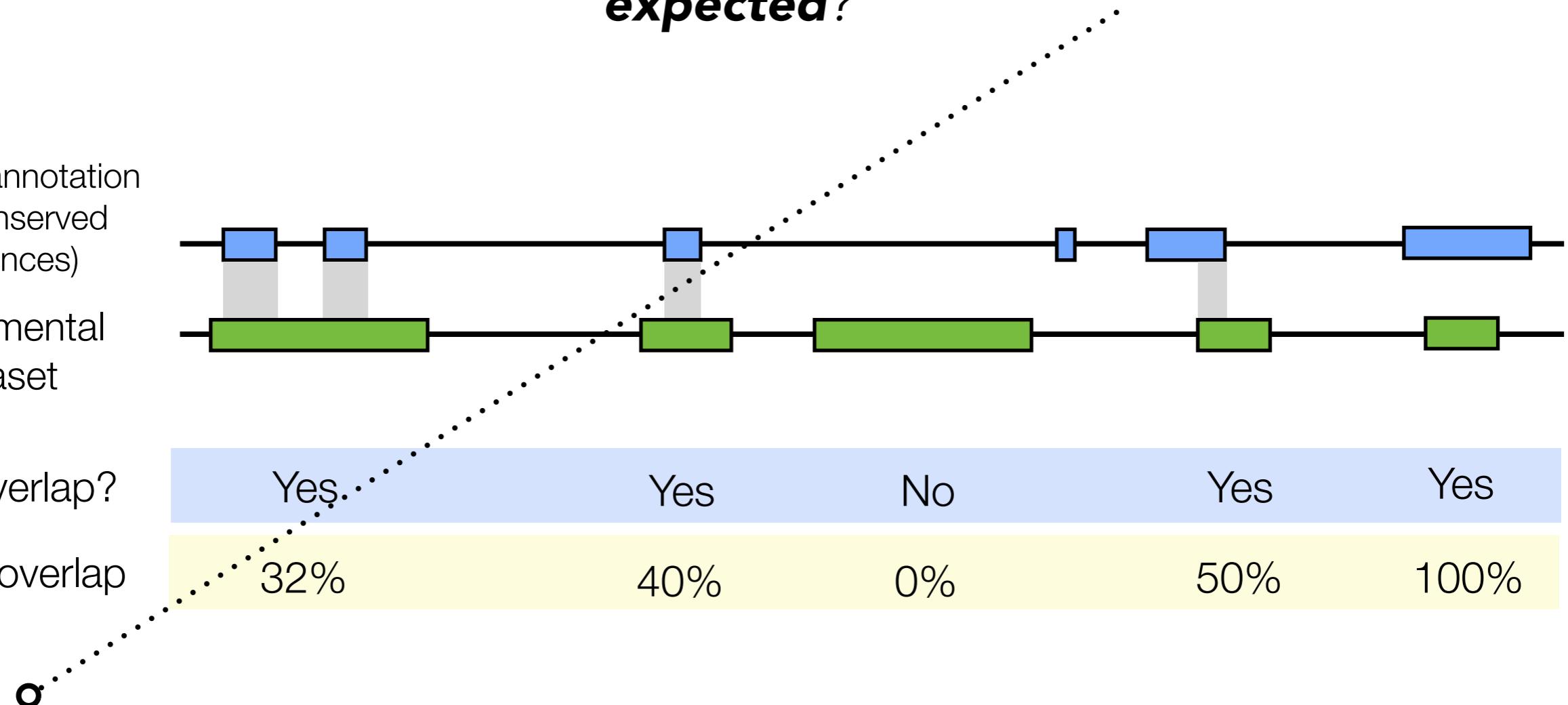
*That is, do two sets of genomic features co-occur
(overlap or have spatial consistency) more than
expected?*

Genome annotation
(e.g. conserved
sequences)

Experimental
dataset

Any overlap?

Bases overlap



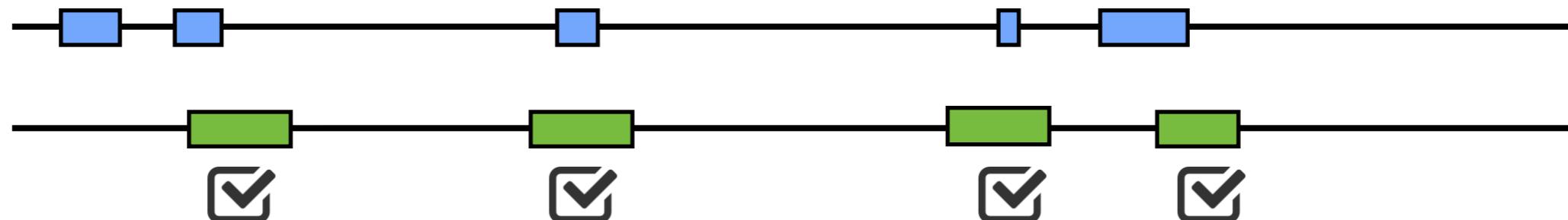
How do we develop a proper null expectation
in order to reduce the dimensionality of the
data to an informative statistic?

Monte-Carlo simulation (slow)

Are the observed feature overlaps more than expected by chance?

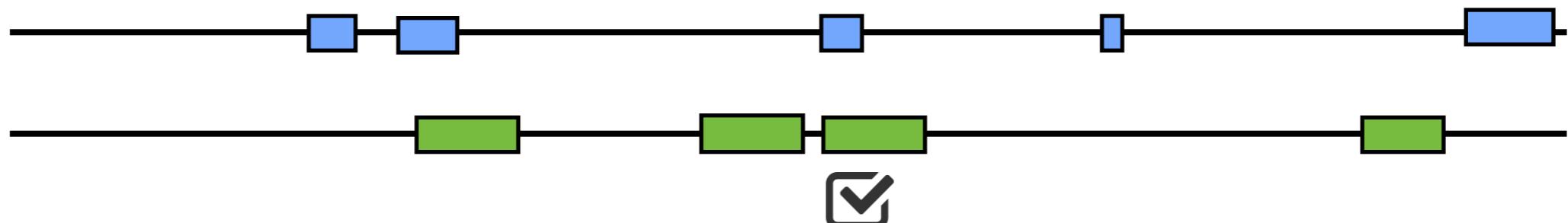
Observed

(4)



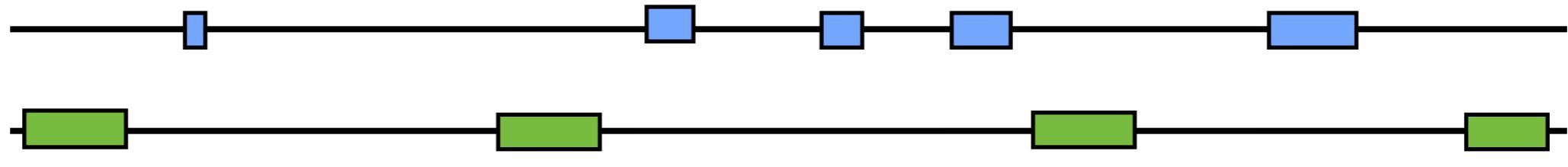
Simulation 1

(1)



Simulation 2

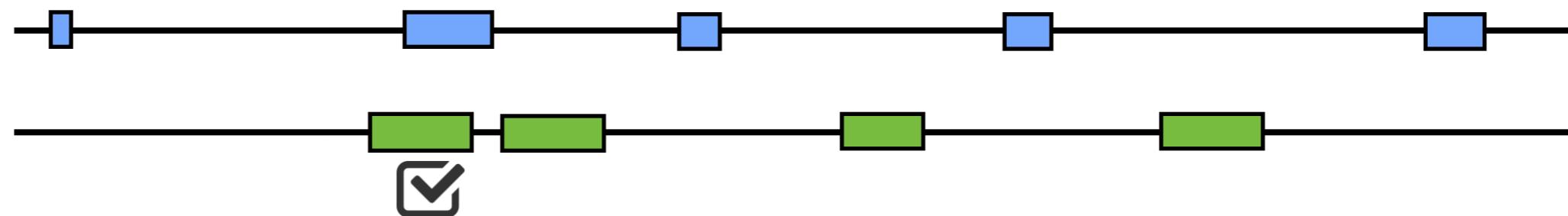
(0)



Simulation

10000

(1)



How do we speed this up?

We invent a scalable new algorithm, of course!

**Binary Interval Search (BITS):
A Scalable Algorithm for Counting Interval Intersections**

Ryan M. Layer¹, Kevin Skadron¹, Gabriel Robins¹, Ira M. Hall², and Aaron R. Quinlan^{3*}

¹Department of Computer Science, University of Virginia, Charlottesville, VA

²Department of Biochemistry and Molecular Genetics, University of Virginia, Charlottesville, VA

³Department of Public Health Sciences and Center for Public Health Genomics, University of Virginia, Charlottesville, VA

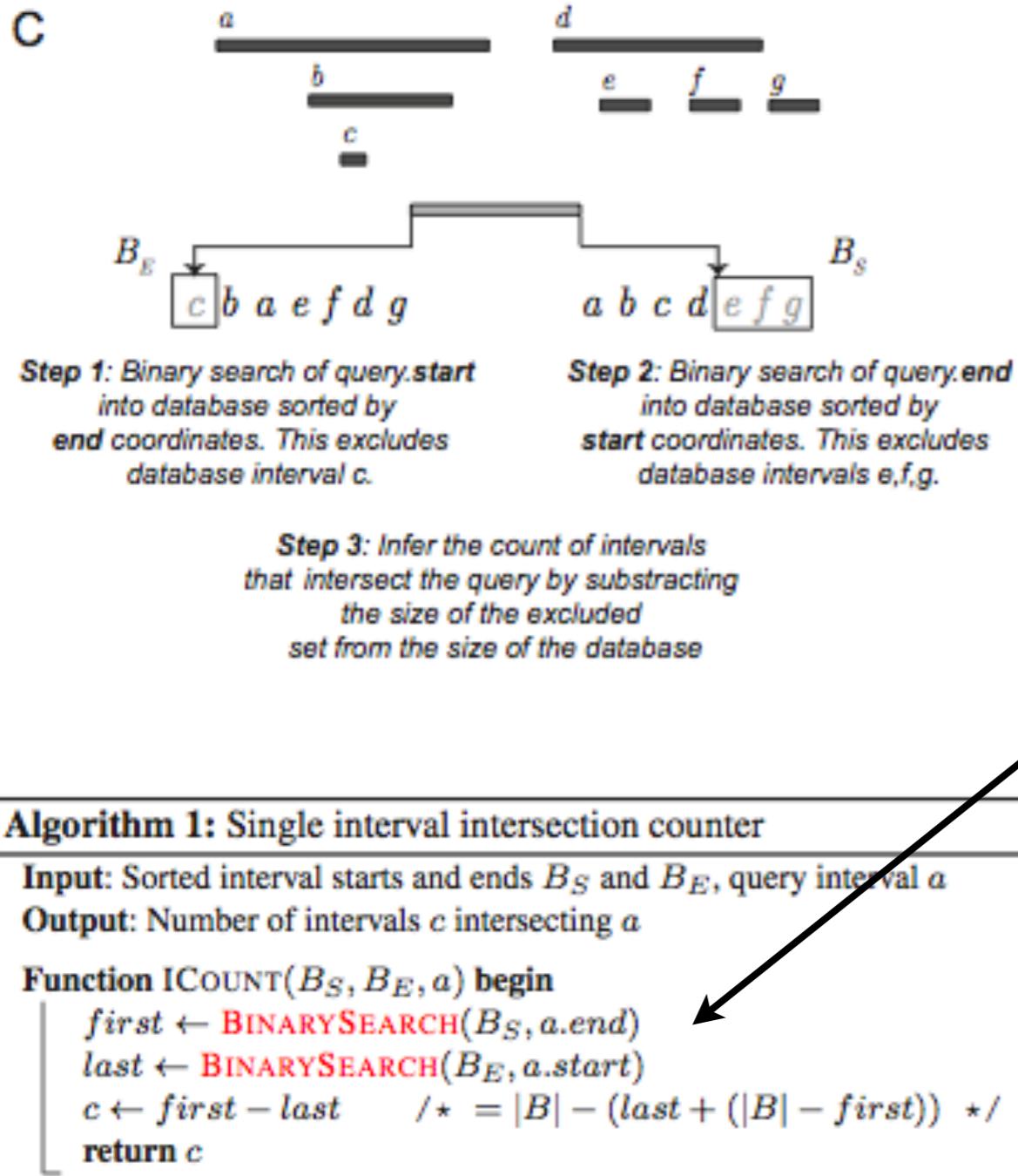


Ryan Layer

<https://github.com/arq5x/bits>

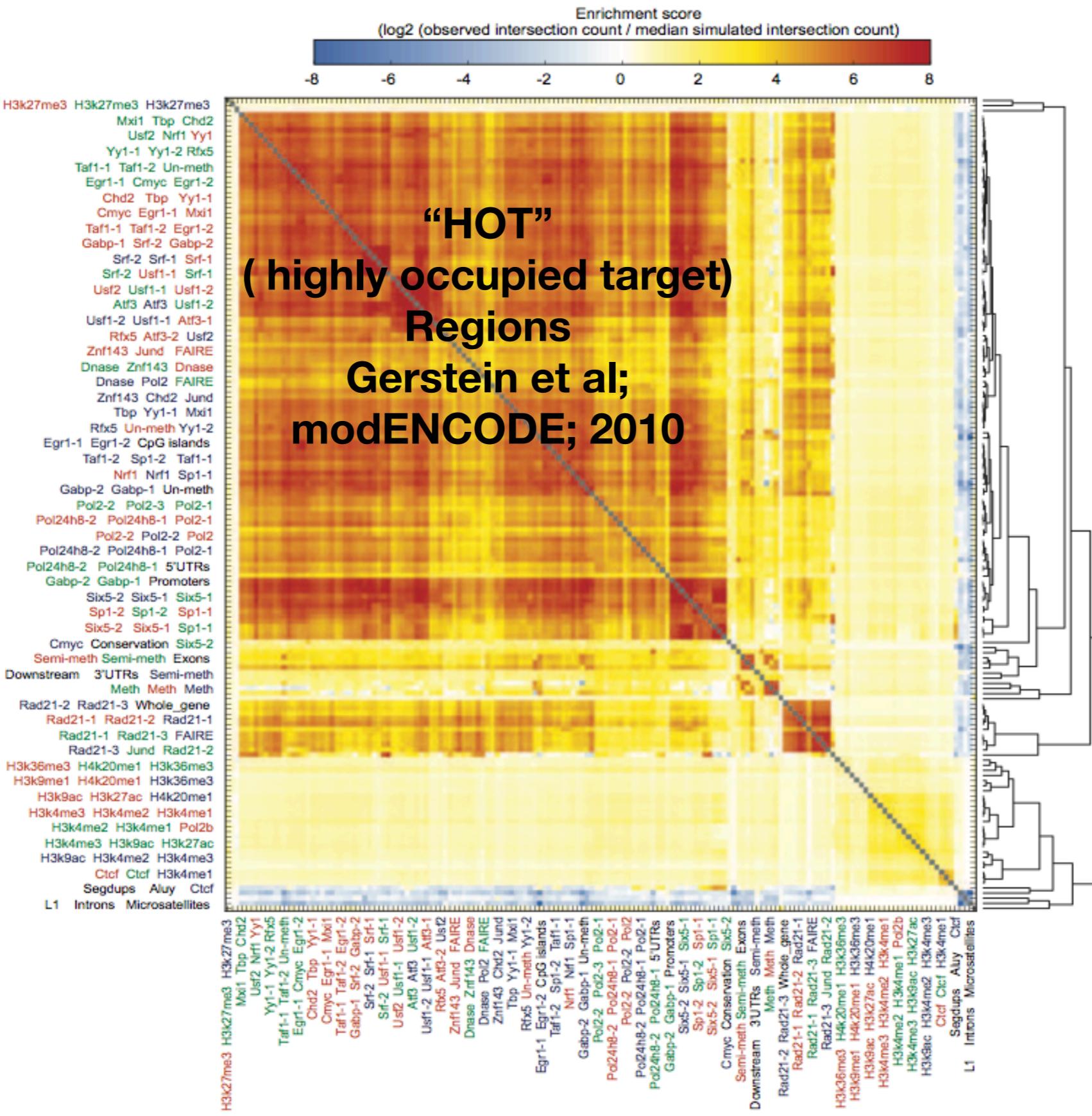
Advance access at Bioinformatics

Binary InTerval Search BITS



- Entirely novel algorithm for detecting genomic interval intersections.
- Clever aspect: unlike any other algorithm, it can deduce the **count** of overlaps without having to **enumerate** each individual intersection.
- Uses two binary searches. Very fast. **Spaceballs fast on a GPU.**
- If you haven't heard, faster is better.

159 ENCODE ChIPseq, RNA-seq,
from 3 cell types (GM12878, hESC, k562). UCSC tracks as well



159 ENCODE ChIPseq, RNA-seq,
from 3 cell types (GM12878, hESC, k562). UCSC tracks as well

25,281 (159^2) pairwise dataset comparisons

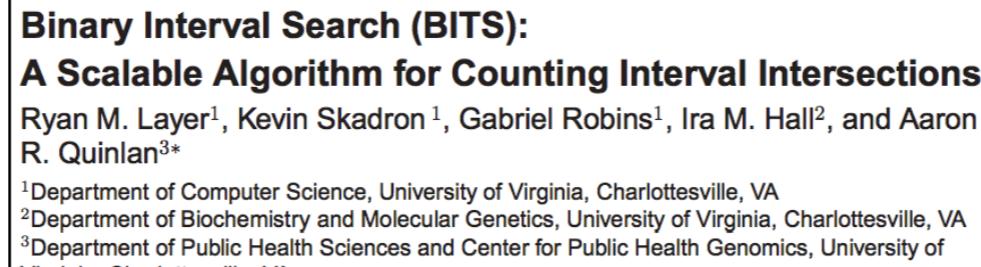
10,000 simulations per 25,281 comparisons

44 trillion intersection measurements

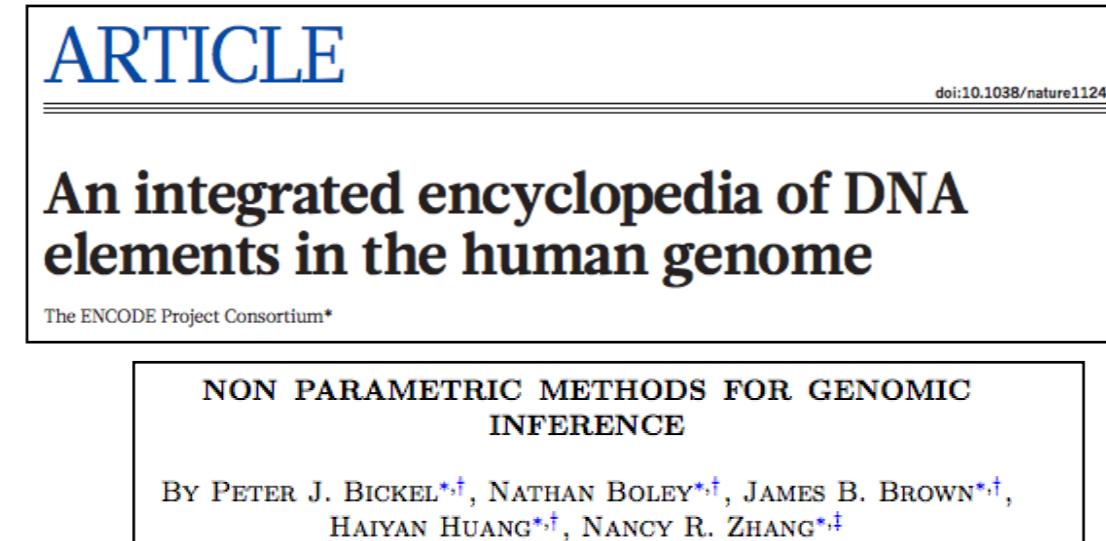
6 days is better than 138!

How do we develop a null expectation?

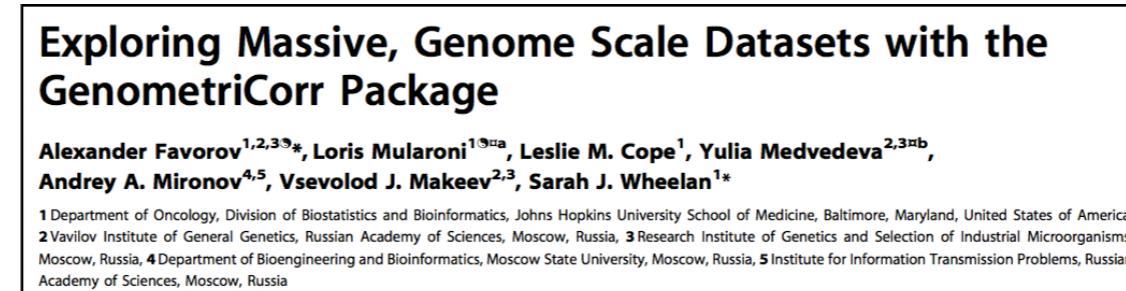
1. Monte-Carlo simulation



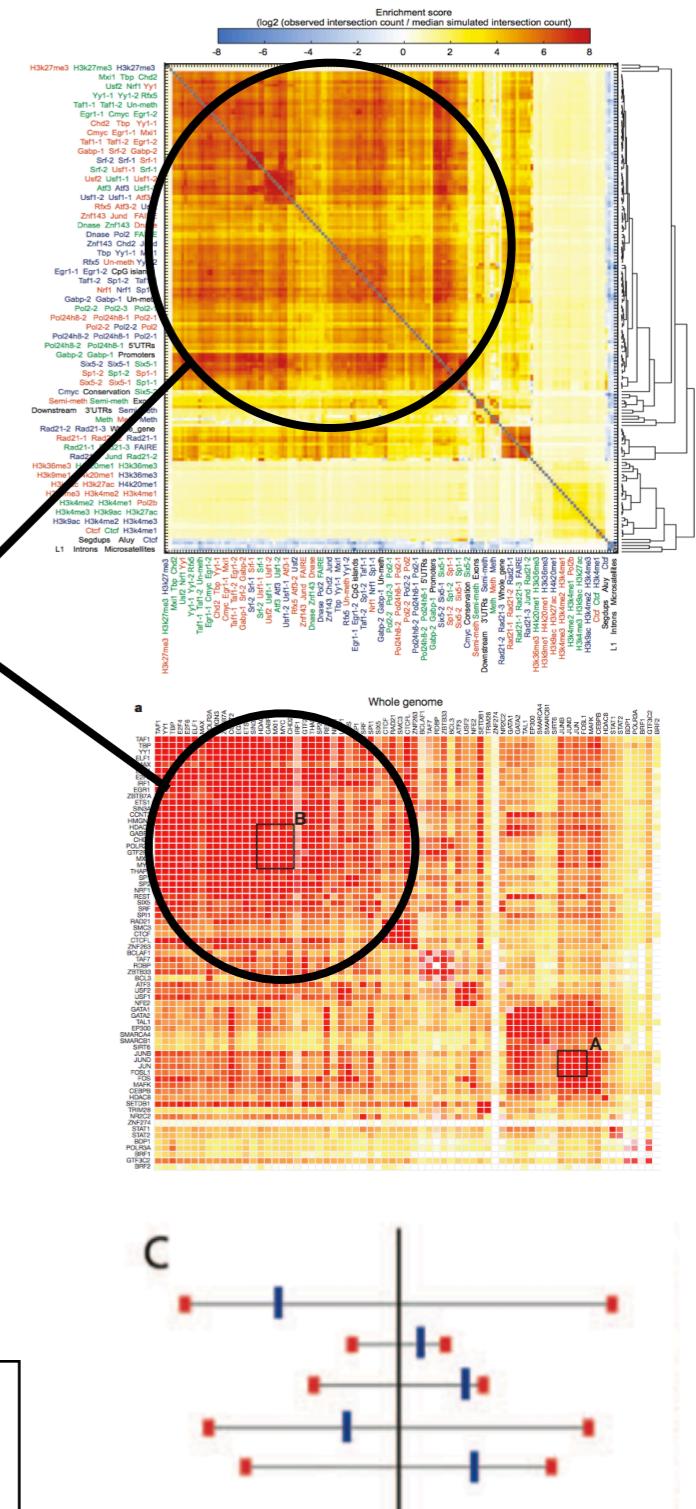
2. Block bootstrap sampling (Peter Bickel)



3. Spatial methods (Sarah Wheelan)



Integrating these and other, novel association statistics into **bedtools2**



query is randomly distributed with respect to the reference

“Block bootstrap” method

- Sampling a la Monte Carlo
- However, don't shuffle randomly.
- Instead, preserve spacial relationships.
- <http://encodestatistics.org/>
- *Example on whiteboard.*

Alternative (better?) methods?

PLOS COMPUTATIONAL BIOLOGY

Browse For Authors About Us Search advanced search

OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE | FEATURED IN PLOS COLLECTIONS

5,064 4 38 17

VIEWS CITATIONS ACADEMIC BOOKMARKS SOCIAL SHARES

Exploring Massive, Genome Scale Datasets with the GenometriCorr Package

Alexander Favorov  , Loris Mularoni , Leslie M. Cope, Yulia Medvedeva, Andrey A. Mironov, Vsevolod J. Makeev, Sarah J. Wheelan 

| Download | Print | Share | | |
|----------|-------------------|---------|----------|-----------------|
| Article | About the Authors | Metrics | Comments | Related Content |
| | | | | |

Included in the Following Collection

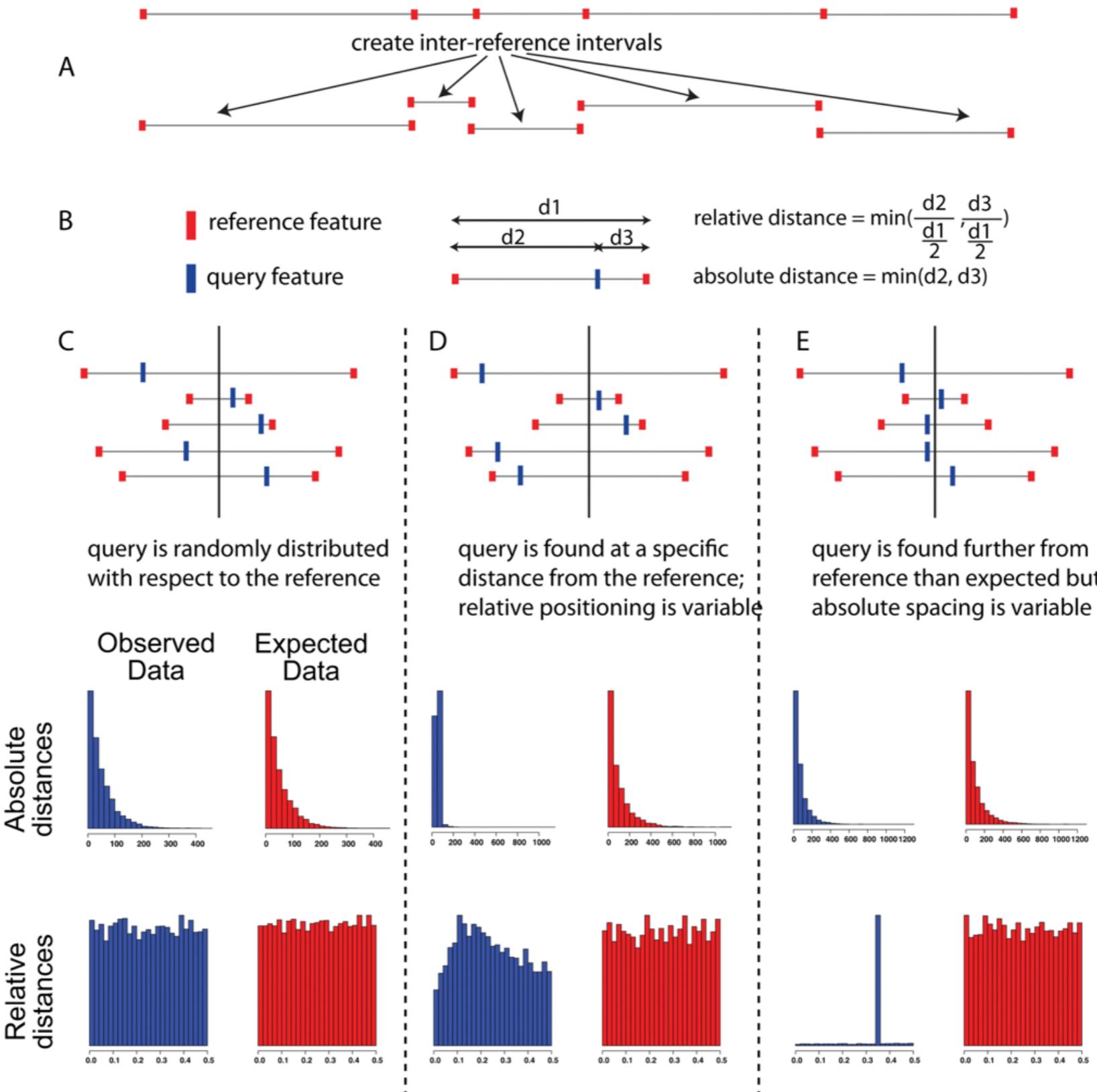
PLOS Computational Biology: Software

“Distance measures”

What if features don't *overlap*
yet are most always nearby?

For example: TF binding peaks and TSS

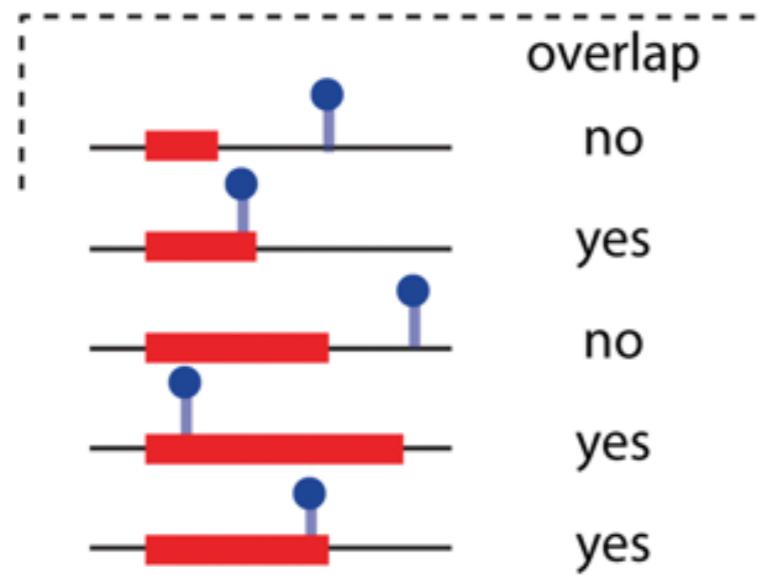
“Distance measures”



Other measures.

F

Projection test: overlap of query
with all genomic reference features



G

Jaccard test: union vs intersection for each reference feature
in the genome (union in light grey, intersection in black)

