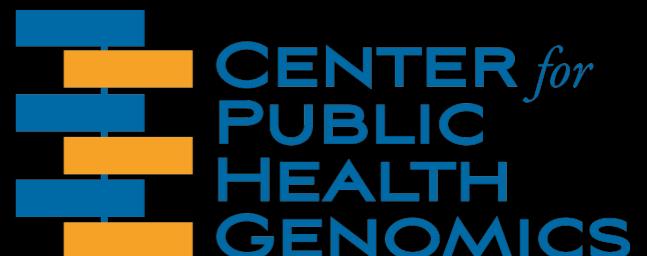


Exploring the structure and function of genomes.

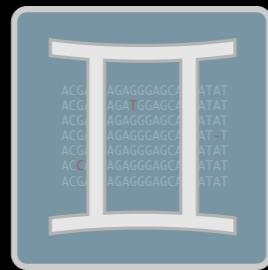
Aaron Quinlan
quinlanlab.org

University of Virginia, Charlottesville VA
Center for Public Health Genomics
Biochemistry and Molecular Genetics



What do we work on?

Software tools for genome research

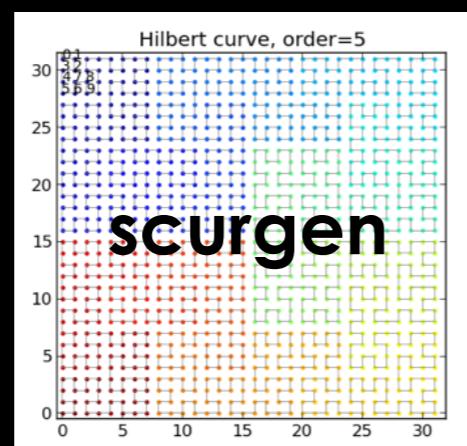


gemini

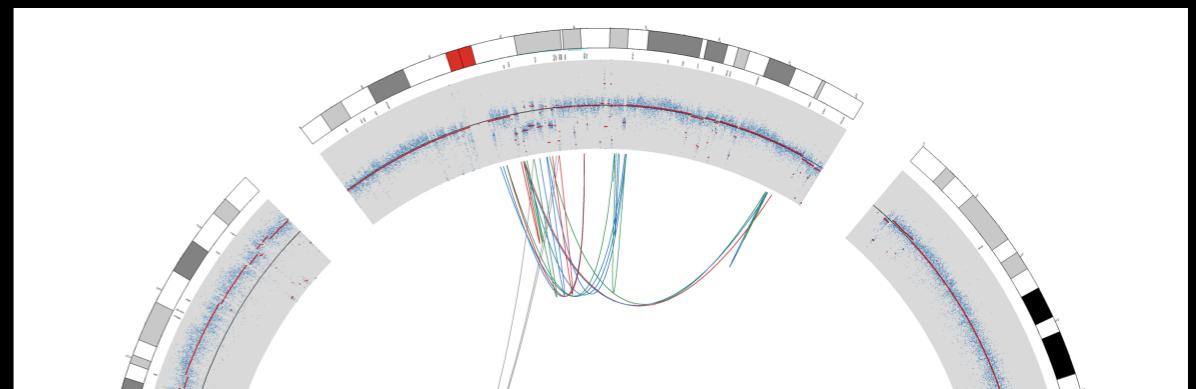
lumpy



genome query
language
(GQL)



Cancer genomics



Cellular and molecular origins of GBM

Initiating mutations underlying OV?

Tumor evolution and mutational mechanisms.

Disease genetics

Type 1 diabetes

SLE (Lupus)

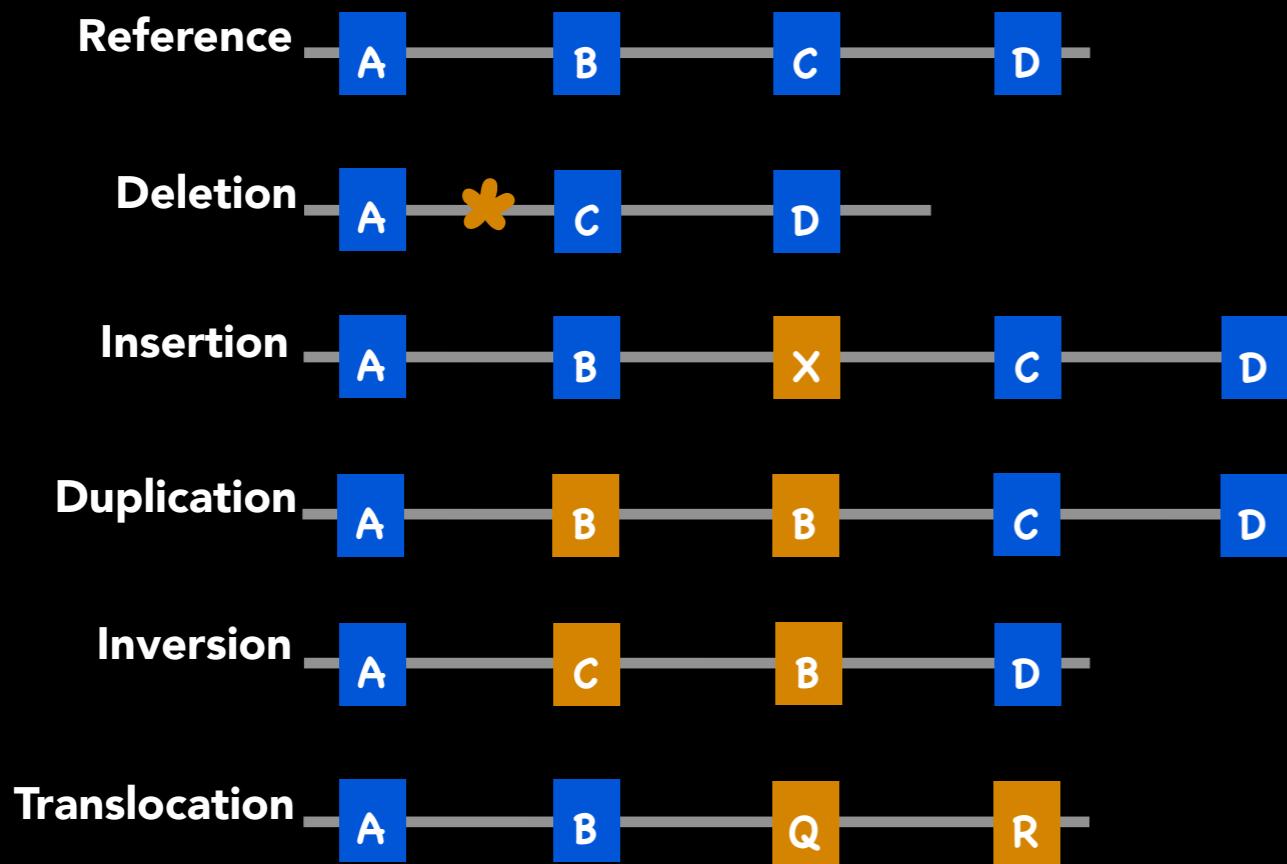
Unexplained developmental disorders

Radiation hypersensitivity

Outline

- Discovery of structural variation (SV)
- Complex SV in 64 tumor genomes
- Inferring function from complex genomic datasets

SV definitions



structural variant (SV): a difference in the copy number, orientation or location of genomic segments >100bp

genomic rearrangement: ditto

copy number variant (CNV), or alteration (CNA): an SV that alters DNA copy number

breakpoint: The junction(s) between structurally variable genomic segments

complex SV: 2 or more breakpoints that arise through a single mutational event, but cannot be explained by one DNA exchange or end-joining reaction

Evolution of our SV discovery tools

Hydra
(2010)



Paired-end mapping (PEM)

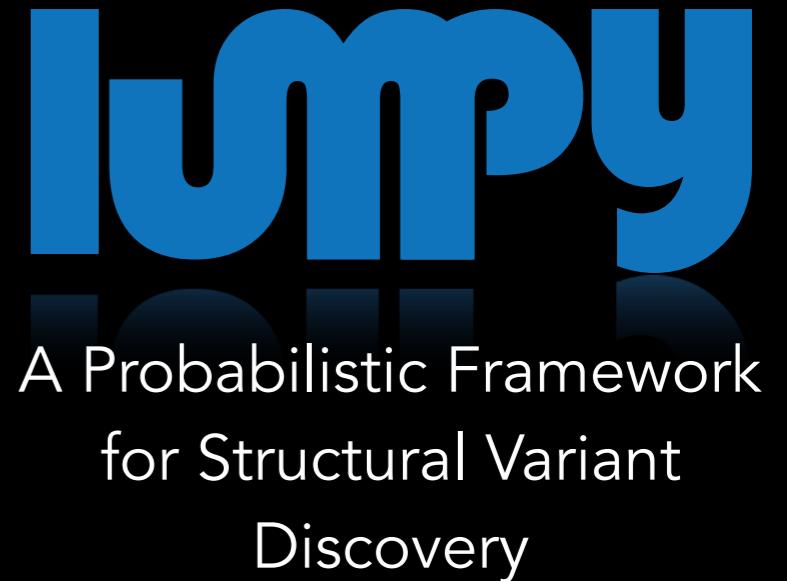
1 signal, 1 sample

Hydra_Multi
(2011)



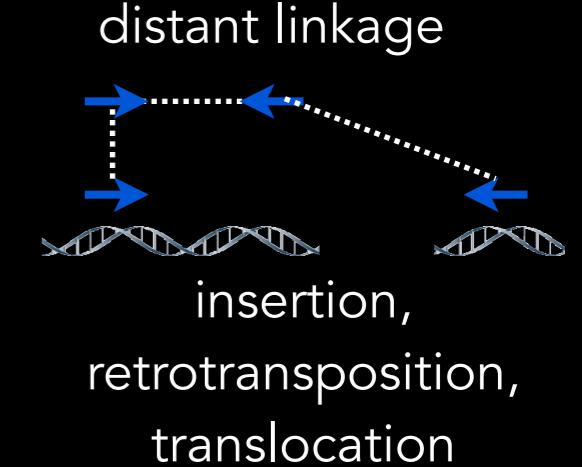
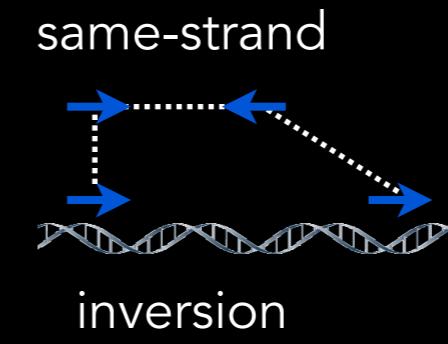
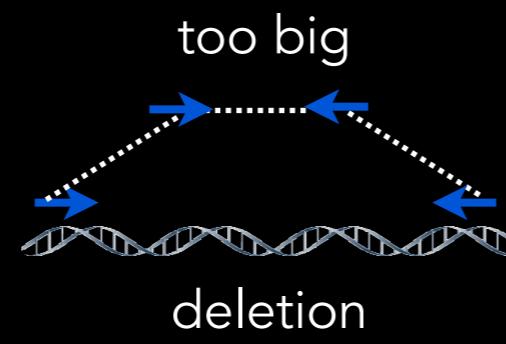
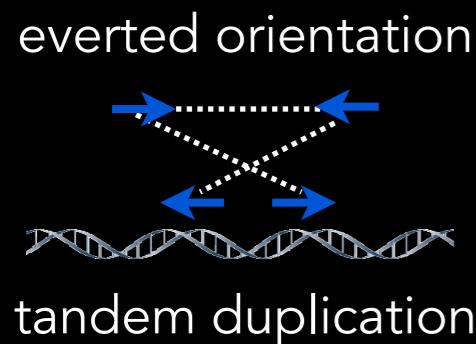
1 signal, ∞ samples

LUMPY
(2012)



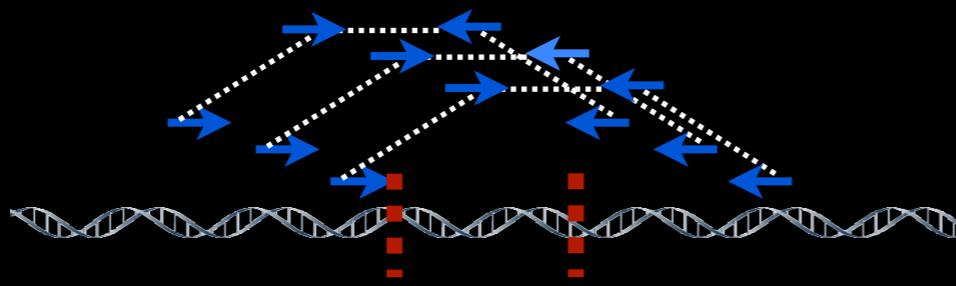
many signals,
many samples

PEM clusters ***discordant*** mappings



↓
Cluster to localize
breakpoints

ref. genome



Hydra



1 signal (PEM), 1 sample

The Hydra algorithm:

- simple & fast
- comprehensive: detects all breakpoint classes
- combinatorial: optionally uses multiple mappings to detect mobile element insertions
- Quinlan et al., 2010. *Genome Research*;
<http://code.google.com/p/hydra-sv/>

The dirty secrets of PEM

Secret #1.

Often many false positives

- Short reads + heuristic alignment + rep. genome
= systematic alignment artifacts (false calls)
- Chimeras and duplicate molecules
- Ref. genome errors (e.g., gaps, mis-assemblies)
- **ALL** SV mapping studies use strict filters for above

Secret #2.

The false negative rate is also generally very high.

- Most current datasets have low to moderate physical coverage due to small insert size (~10-20X)

Secret #2.

The false negative rate is also generally very high.

- Most current datasets have low to moderate physical coverage due to small insert size (~10-20X)
- Breakpoints are enriched in repetitive genomic regions that pose problems for sensitive read alignment

Secret #2.

The false negative rate is also generally very high.

- Most current datasets have low to moderate physical coverage due to small insert size (~10-20X)
- Breakpoints are enriched in repetitive genomic regions that pose problems for sensitive read alignment
- FILTERING!

Secret #2.

The false negative rate is also generally very high.

- Most current datasets have low to moderate physical coverage due to small insert size (~10-20X)
- Breakpoints are enriched in repetitive genomic regions that pose problems for sensitive read alignment
- FILTERING!
- The false negative rate is usually hard to measure, but is thought to be extremely high for most PEM studies (>30%)

Secret #2.

The false negative rate is also generally very high.

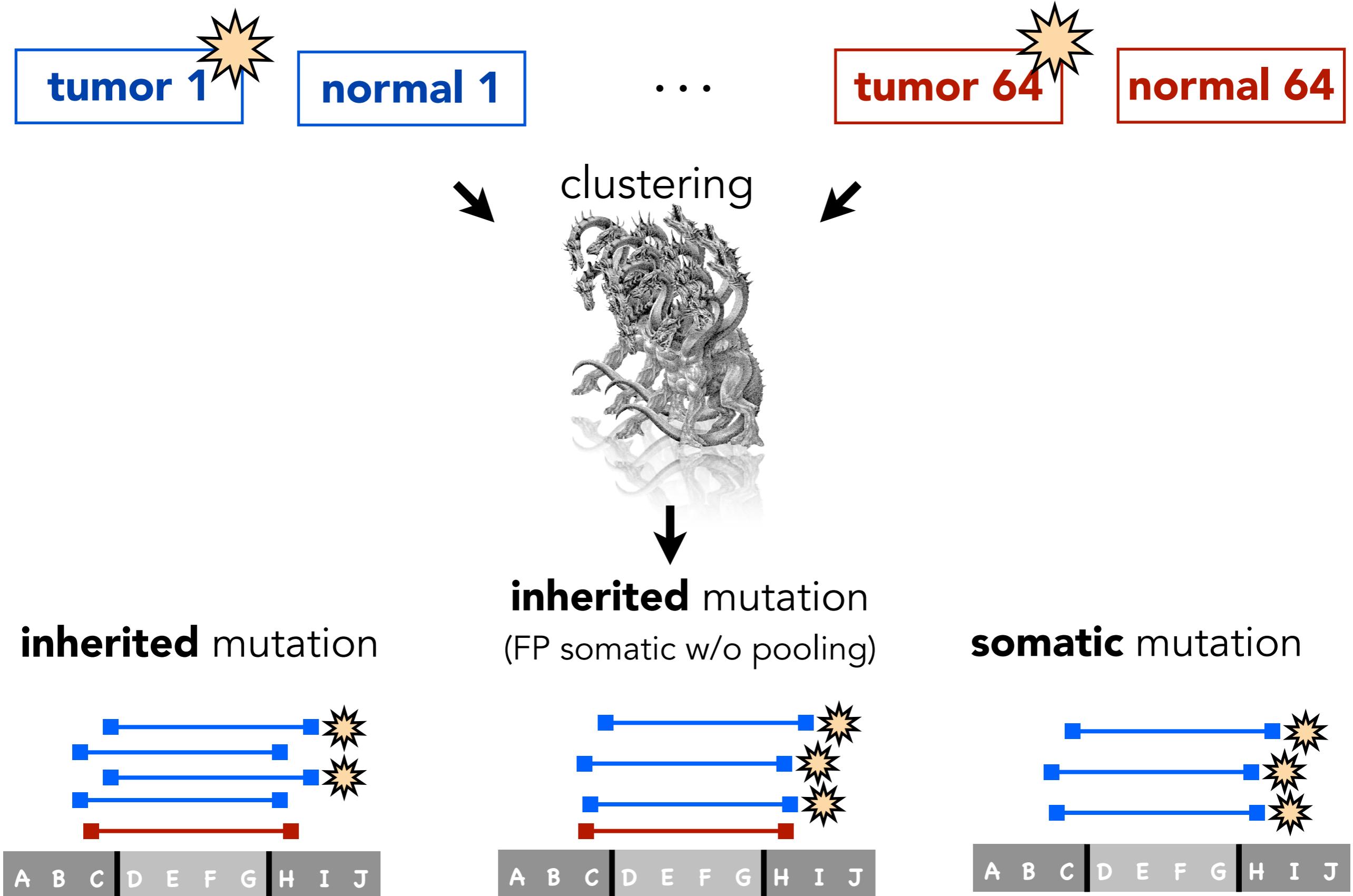
- Most current datasets have low to moderate physical coverage due to small insert size (~10-20X)
- Breakpoints are enriched in repetitive genomic regions that pose problems for sensitive read alignment
- FILTERING!
- The false negative rate is usually hard to measure, but is thought to be extremely high for most PEM studies (>30%)
- When searching for spontaneous mutations in a family or a tumor/normal comparison, a false negative call in one sample can be a false positive somatic or de novo call in another.

Solution to both FPs and FNs:

Pool data from multiple samples.

*It improves SNP and INDEL calling, so
why not SVs?*

Pooling prevents false somatic calls



The landscape of complex variation in 64 cancer genomes.

64 Tumors and 65 matched normals (1 dup.)



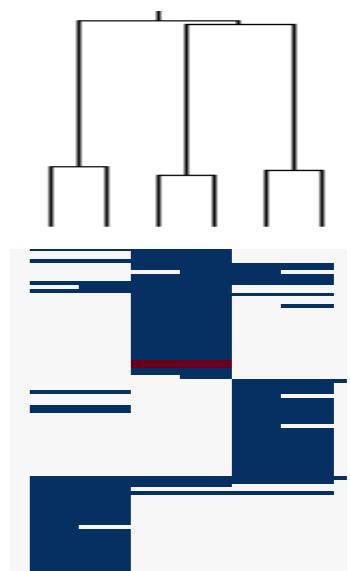
- 12 breast invasive carcinomas (BRCA)
- 3 colon adenocarcinomas (COAD)
- 18 glioblastoma multiforme (GBM)
- 6 lung adenocarcinoma (LUAD)
- 13 lung squamous cell carcinoma (LUSC)
- 11 ovarian serous cystadenocarcinoma (OV)
- 2 rectum adenocarcinoma (READ)

64 out of 64 tumor / normal pairs cluster as
nearest neighbors

12096 SV breakpoints

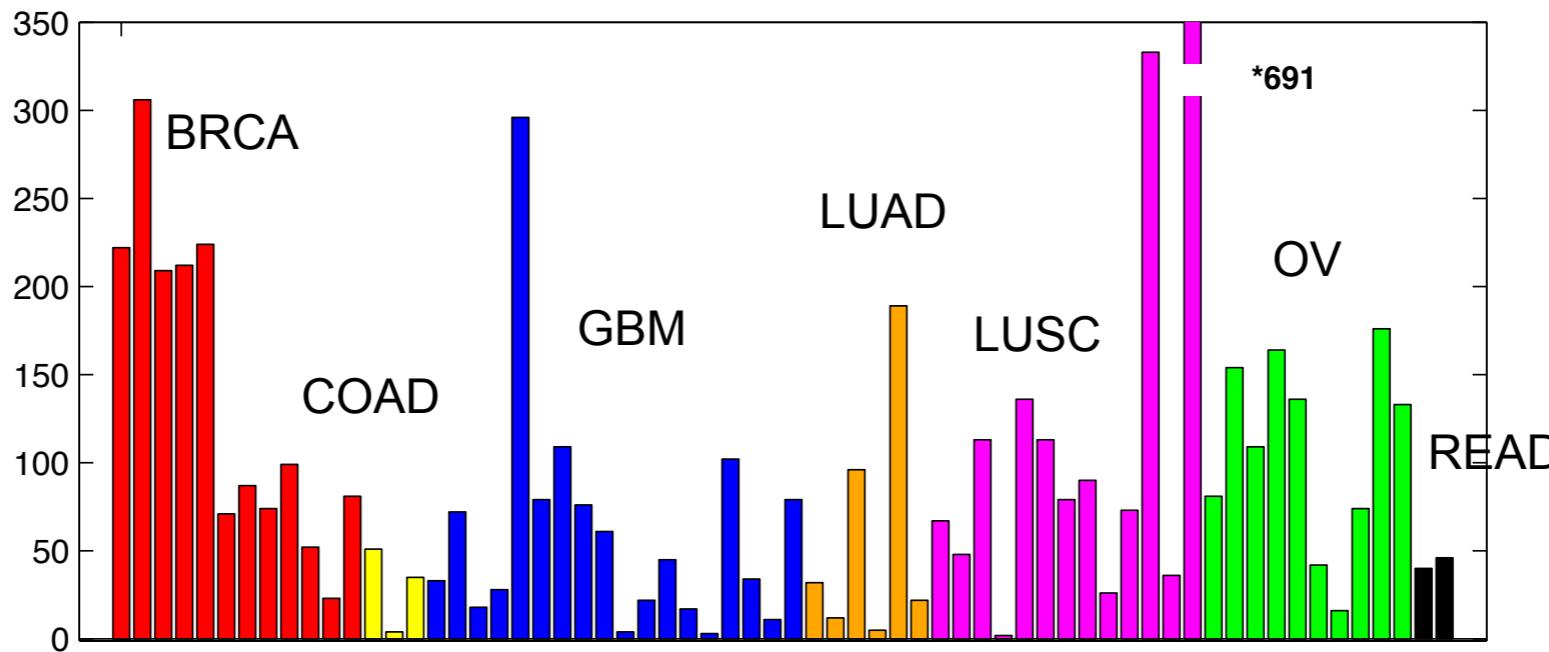


3 tumor/
normal pairs

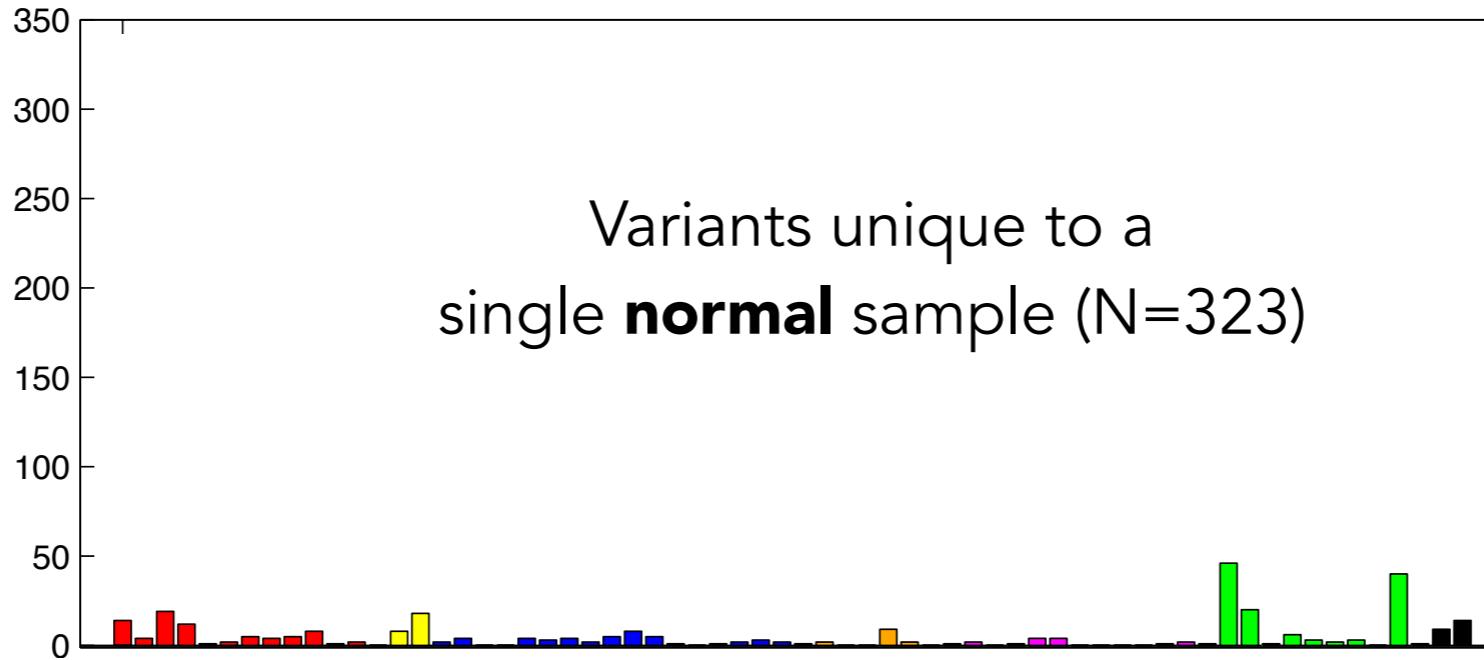


Pooling yields accurate predictions of somatically-acquired SVs in tumors.

Variants unique to a single **cancer** sample (N=6,179)



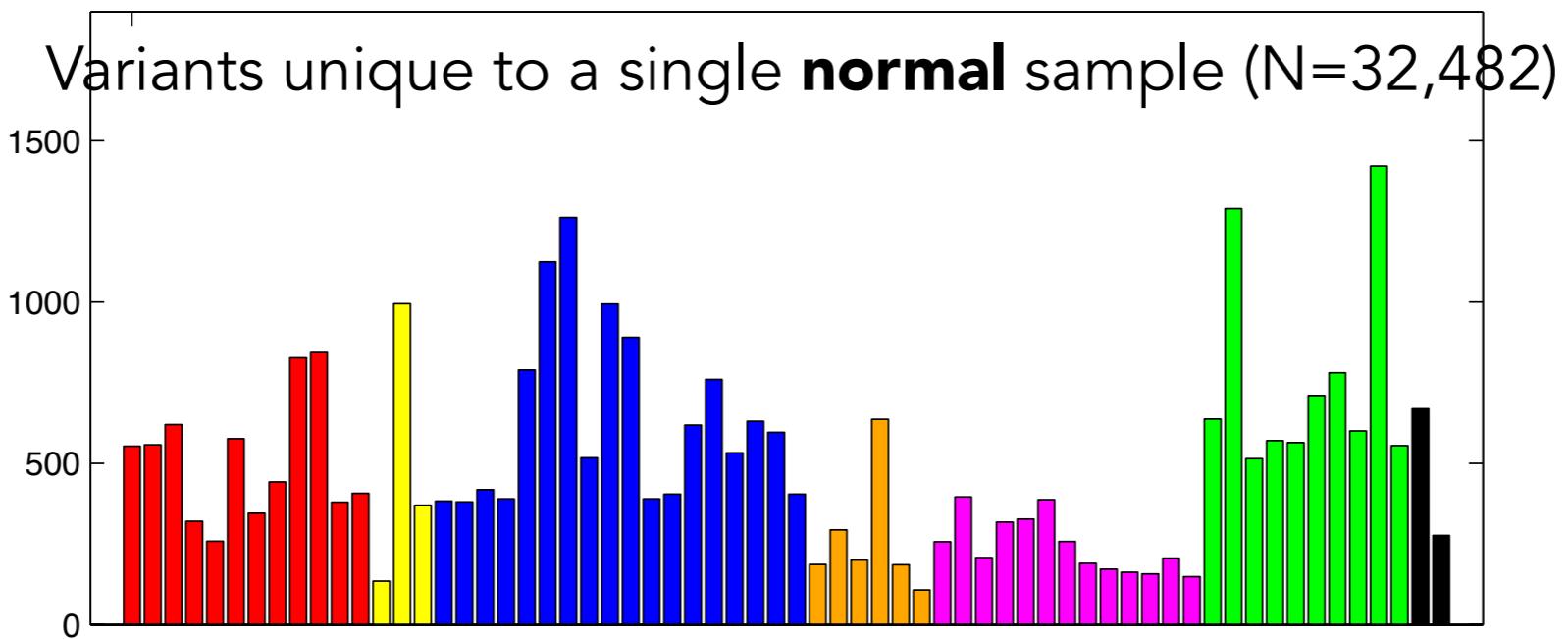
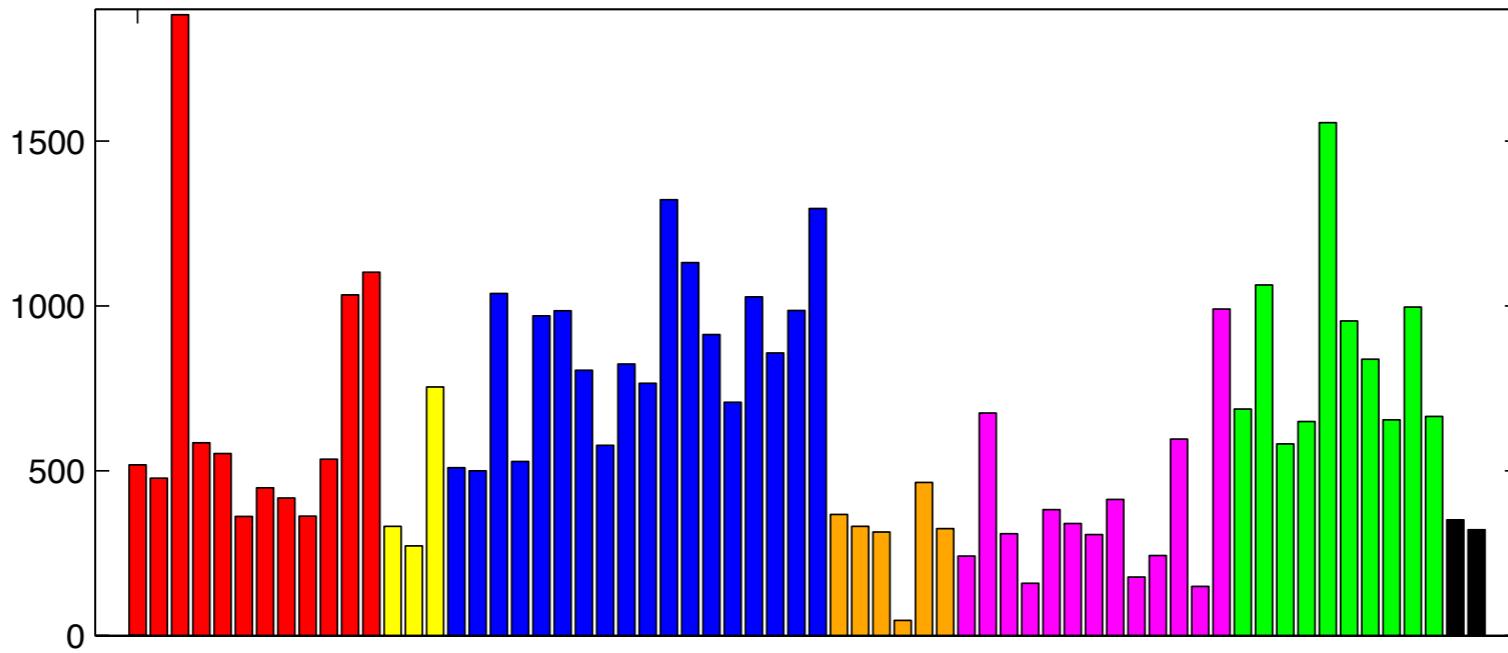
Variants unique to a single **normal** sample (N=323)



Assuming all normal-only calls are false, suggests 5% somatic prediction error rate.
Likelihood of LOH suggests it is actually lower.

Much worse if we just did a simple tumor/ normal comparison (the standard)

Variants unique to a single **cancer** sample (N=41,510)



Somatic misclassification rate jumps from 5% with pooling to 86%!

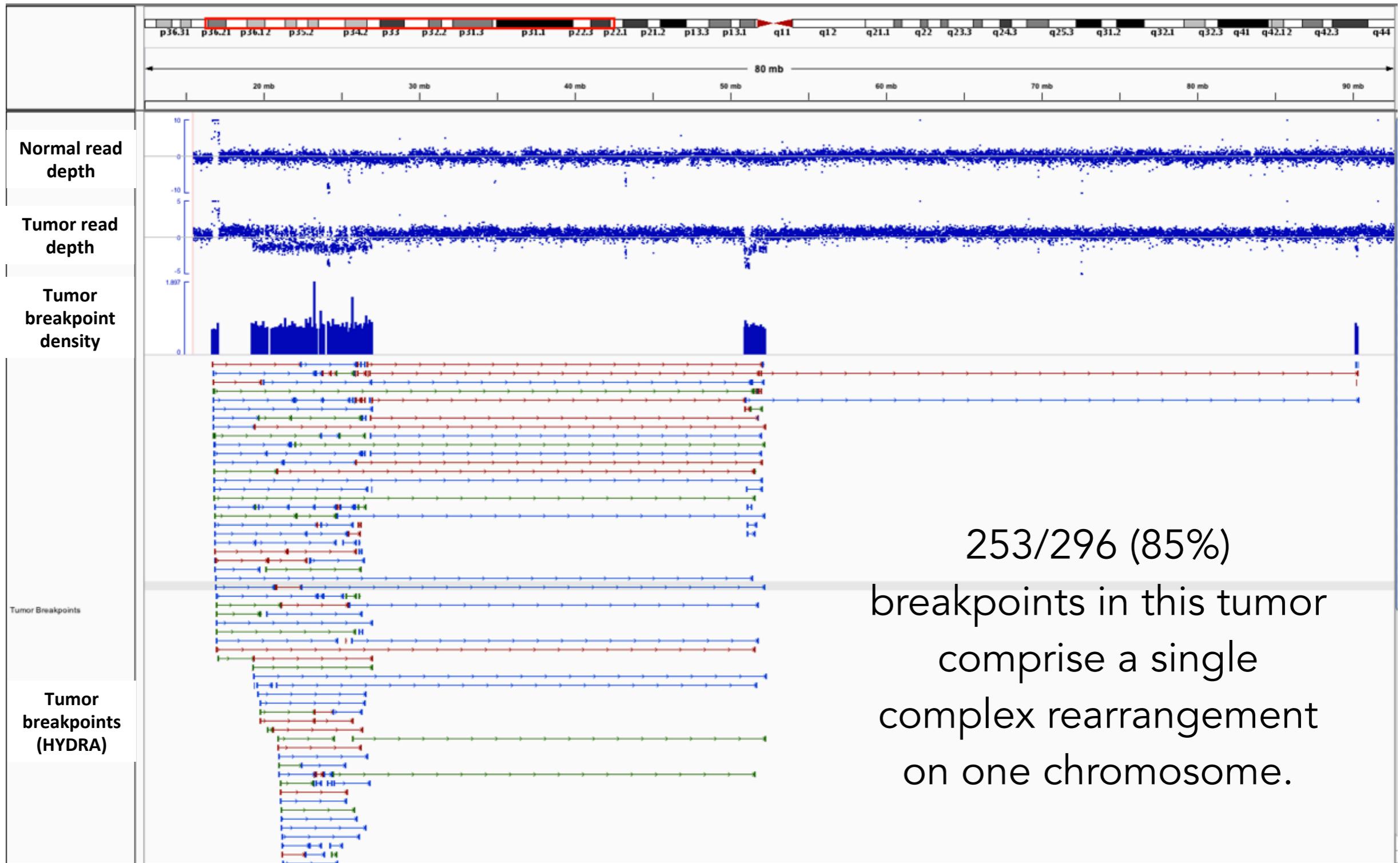
We have a high-quality set of somatic rearrangements from multiple tumors.

**What do they tell us about
chromosome evolution in cancers?**

Observation 1.

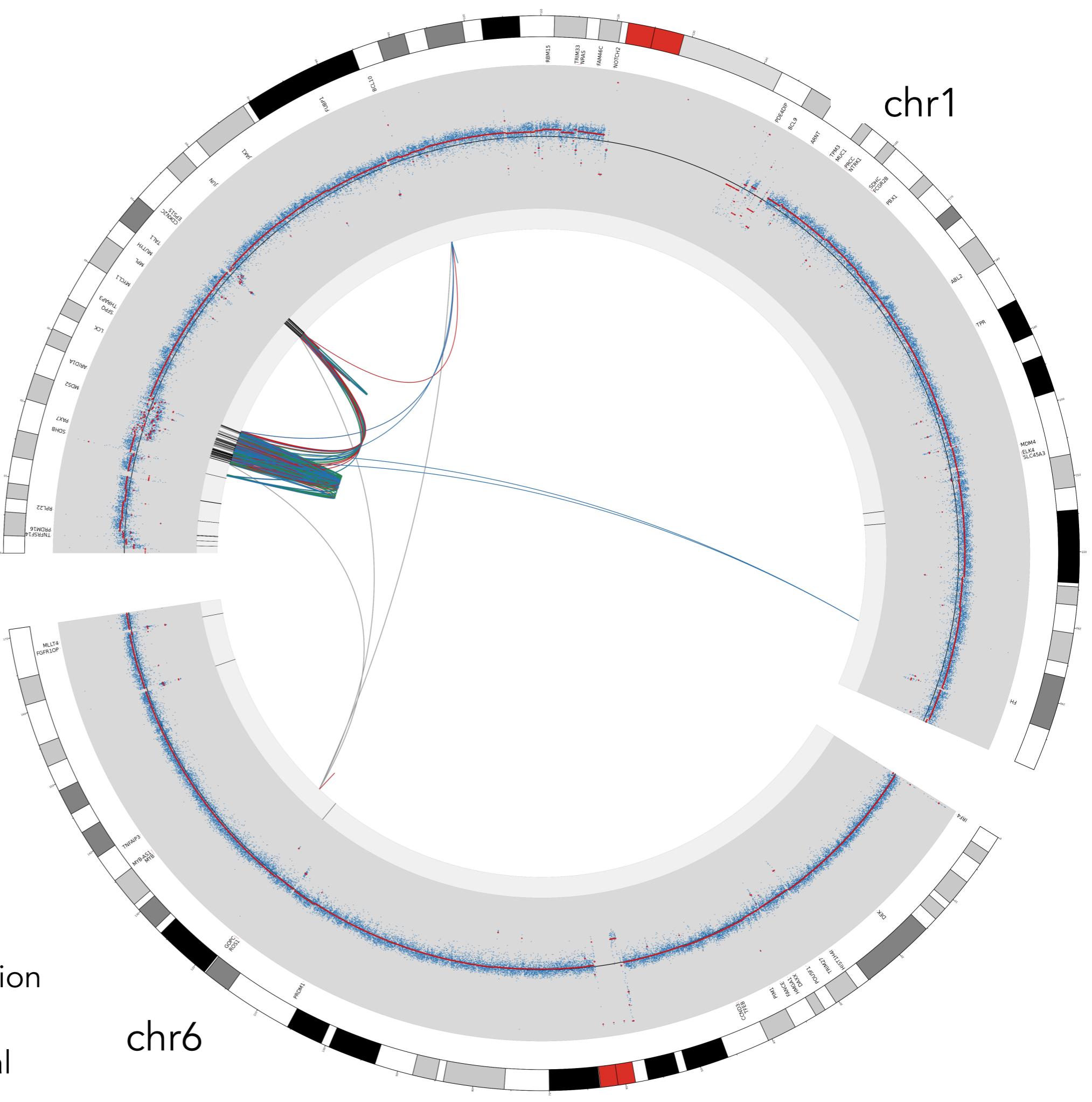
We immediately noticed a few staggeringly complex rearrangements (CRs).

A mangled tumor chromosome in IGV



This has been attributed to new mechanisms: chromothripsis (Stephens et al. 2011). The mechanism is not known.

CIRCOS plot of same event.



Chromothripis

Cell

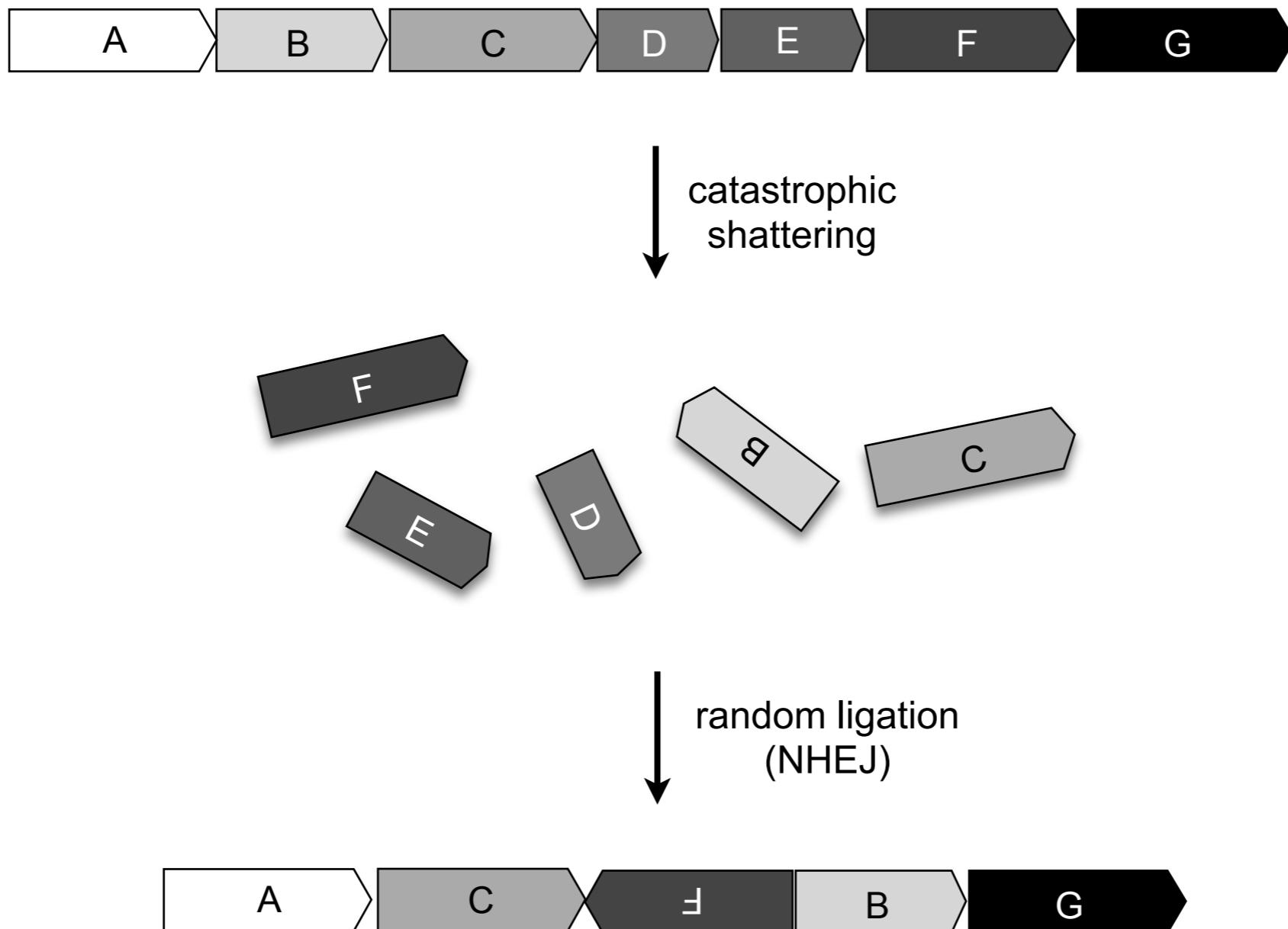
Massive Genomic Rearrangement Acquired in a Single Catastrophic Event during Cancer Development

Philip J. Stephens,¹ Chris D. Greenman,¹ Beiyuan Fu,¹ Fengtang Yang,¹ Graham R. Bignell,¹ Laura J. Mudie,¹ Erin D. Pleasance,¹ King Wai Lau,¹ David Beare,¹ Lucy A. Stebbings,¹ Stuart McLaren,¹ Meng-Lay Lin,¹ David J. McBride,¹ Ignacio Varela,¹ Serena Nik-Zainal,¹ Catherine Leroy,¹ Mingming Jia,¹ Andrew Menzies,¹ Adam P. Butler,¹ Jon W. Teague,¹ Michael A. Quail,¹ John Burton,¹ Harold Swerdlow,¹ Nigel P. Carter,¹ Laura A. Morsberger,² Christine Iacobuzio-Donahue,² George A. Follows,³ Anthony R. Green,^{3,4} Adrienne M. Flanagan,^{5,6} Michael R. Stratton,^{1,7} P. Andrew Futreal,¹ and Peter J. Campbell^{1,3,4,*}

Chromosome shattering in a single, catastrophic event.

A new, punctate form of chromosome evolution.

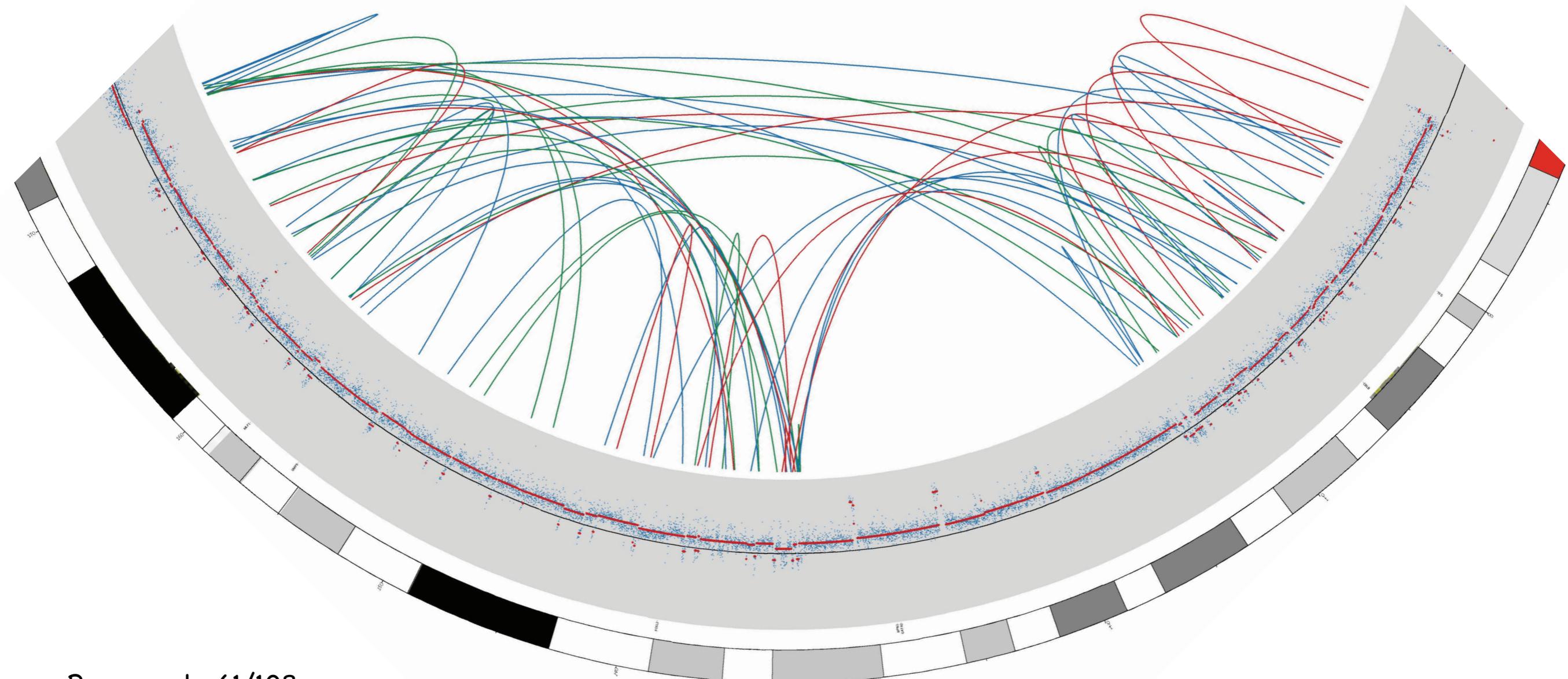
Chromothripis



Observation 2.

Complex rearrangements (CRs) are very common in tumor genomes.

Glioblastoma sample #6



Represents 61/108
breakpoints in this tumor;

5mb

GBM_8_T ; chainID=9 ; numBreaks=61

1542 of 6179 breakpoints are part of CRs.
Not random.

Experimental Data						Simulations (mean of 100 trials)			
class (num. breaks)	Tumor-specific mutations (n=6179)		Normal-specific mutations (n=323)		Tumor-specific (random shuffle)		1000 Genomes (random sample)		
	num. CGRs	total breaks	num. CGRs	total breaks	num. CGRs	total breaks	num. CGRs	total breaks	
mild (3-4)	90	298	3	9	4.7	14.3	2.2	6.7	
moderate (5-9)	32	204	0	0	0.2	1.1	0.1	0.7	
extreme (>9)	33	1045	0	0	0	0	0	0	
total	154	1542	3	9	4.9	15.4	2.3	7.4	

Observation 3.

Complex rearrangements are very common in glioblastoma.

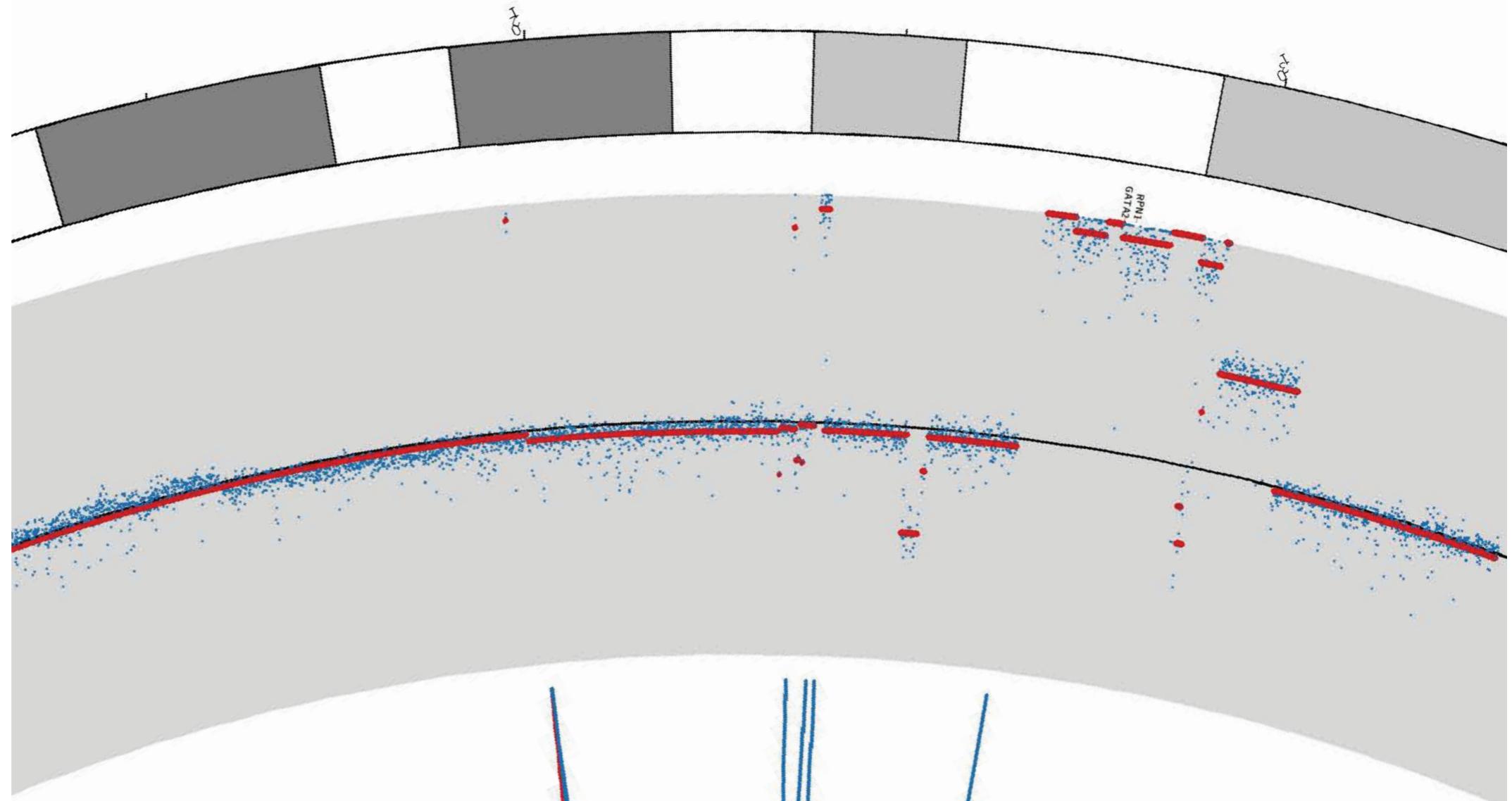
1542 of 6179 breakpoints are part of CRs.
Not random.

	total breaks (mean)	% complex
BRCA (n=12)	1657 (138)	4.2%
COAD (n=3)	90 (30)	10%
GBM (n=18)	1088 (60)	70%
LUAD (n=6)	356 (59)	23%
LUSC (n=13)	1806 (139)	26.7%
OV (n=11)	1096 (100)	11.6%
READ (n=2)	86 (43)	11.6%
total	6179	25%



Observation 4.

There is vast architectural diversity
in non-chromothripsis
complex variants

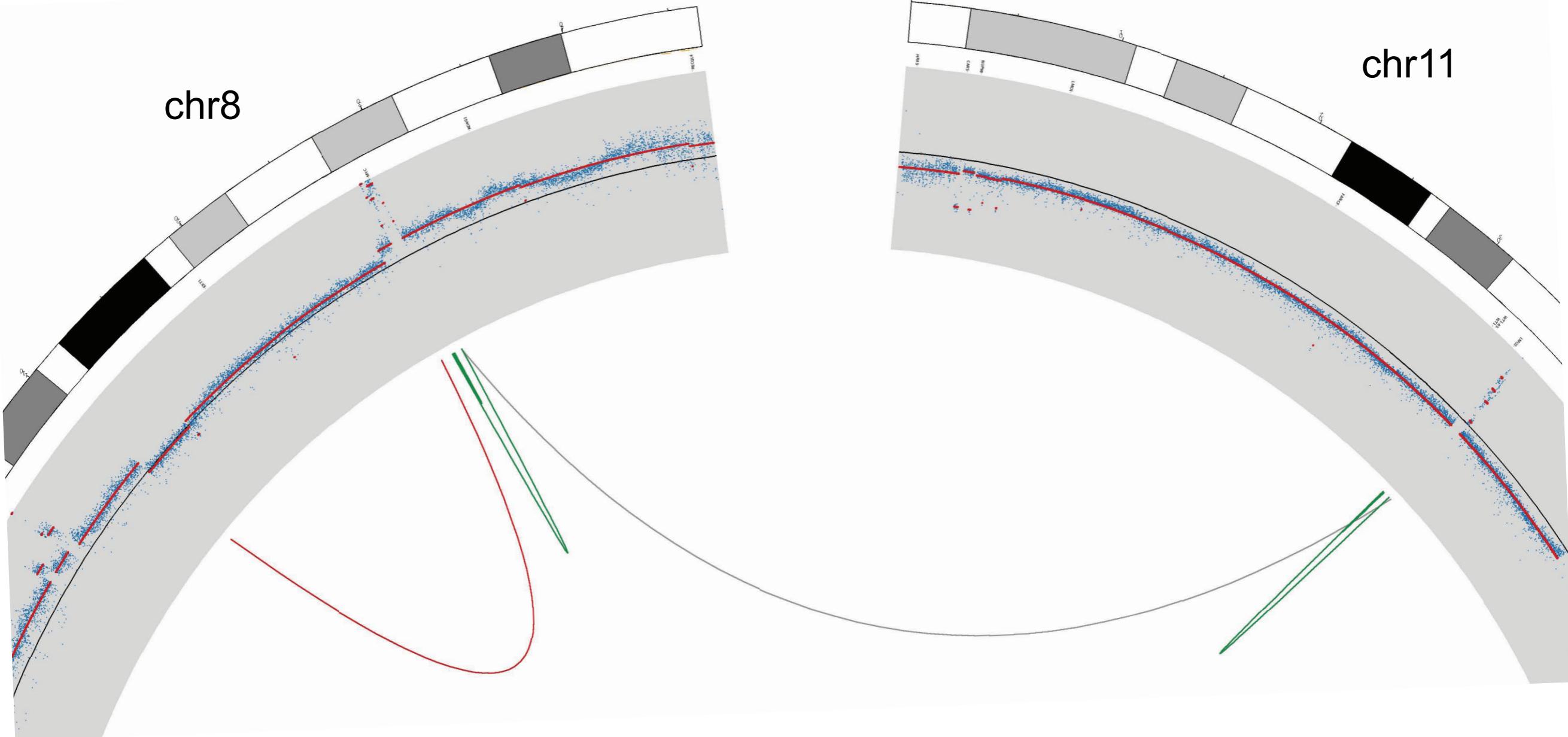


intra-chromosomal
rearrangement between
co-amplified segments

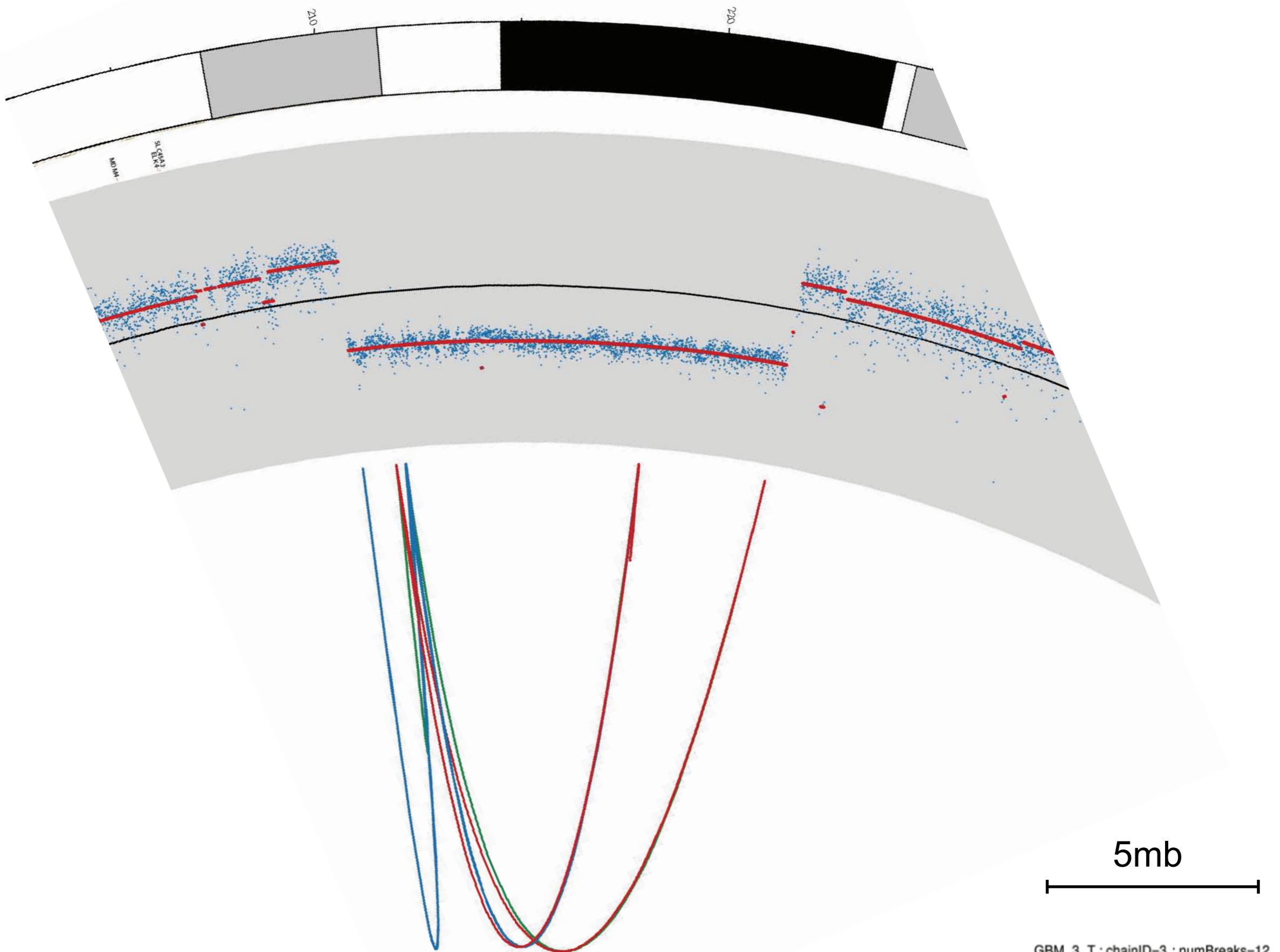
5mb

BRCA_12_T ; chainID=38 ; numBreaks=4

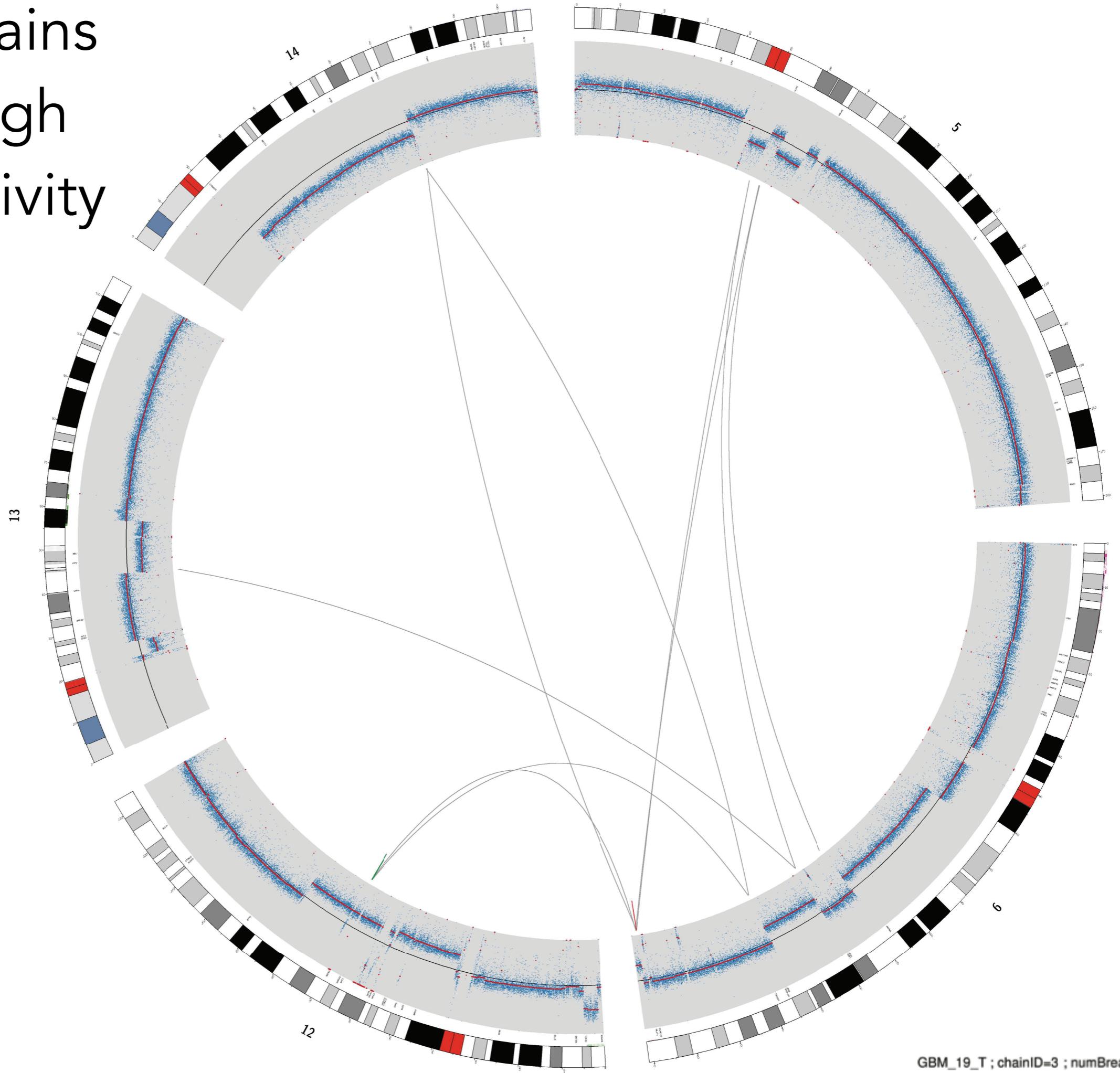
inter-chromosomal rearrangement between co-amplified segments



cryptic rearrangements

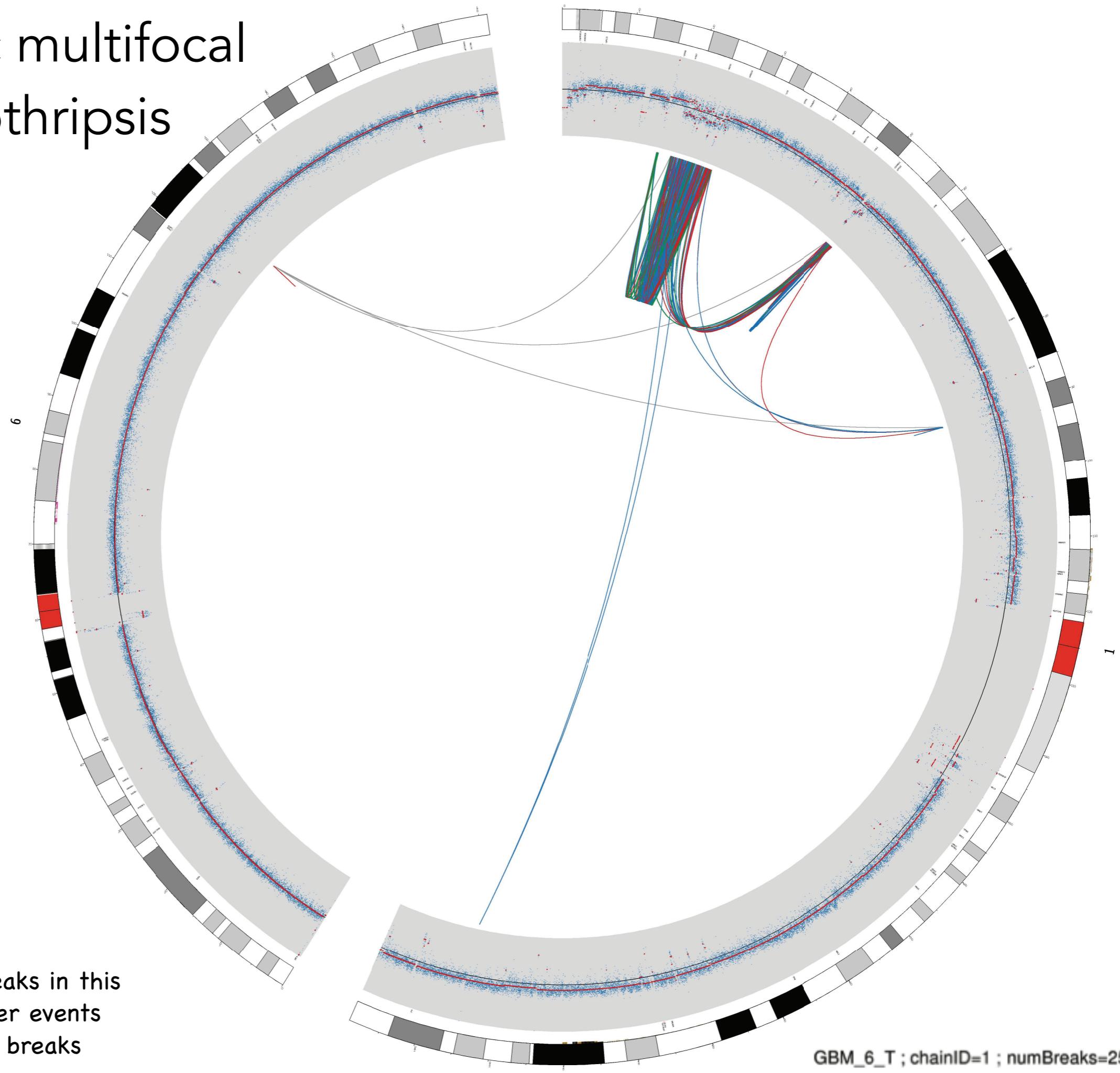


long chains with high connectivity



Many chromothripsis examples are
horrifically complex.

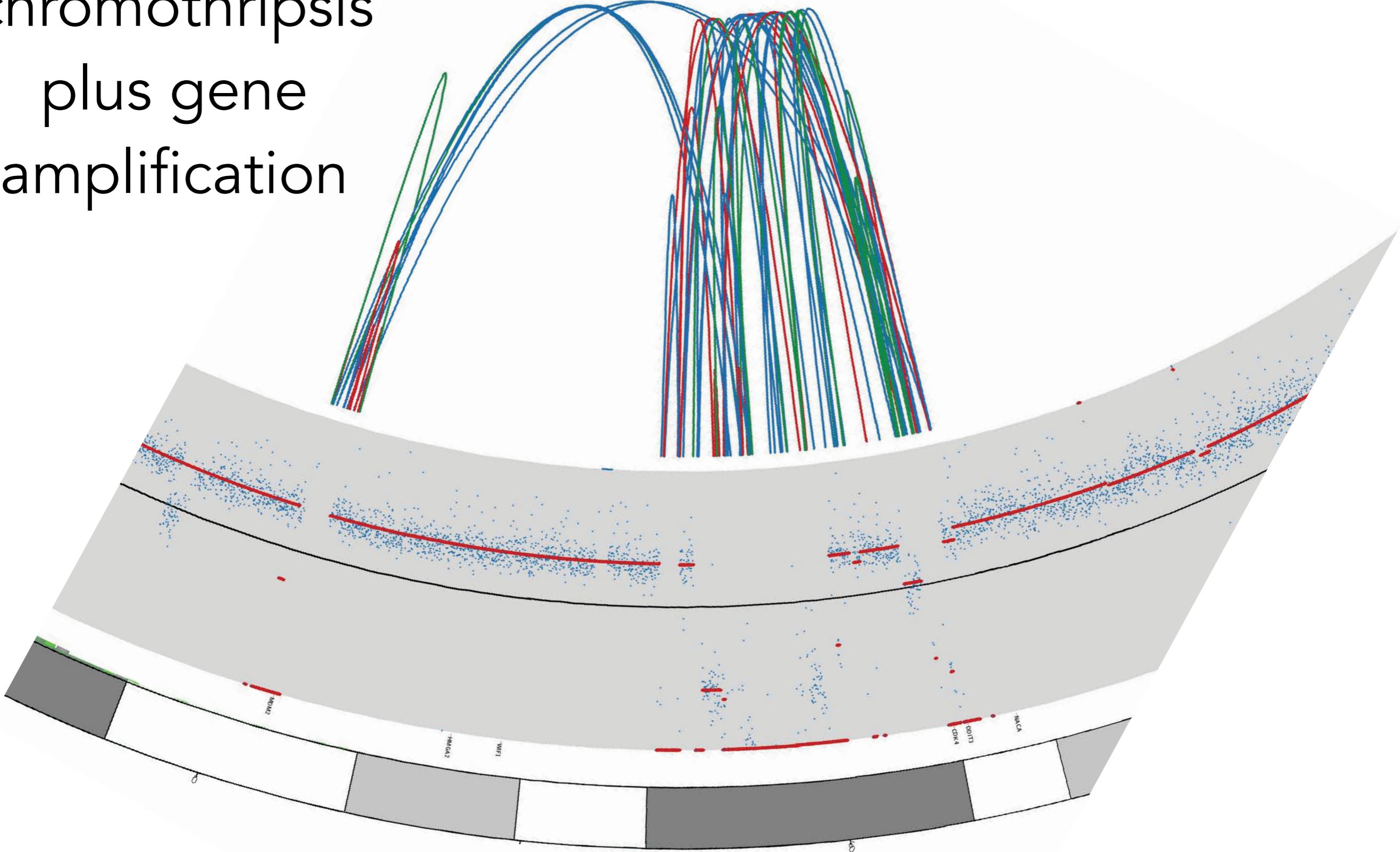
The classic multifocal chromothripsis



Represents 253/296 breaks in this tumor; this and 4 other events account for 271/296 breaks

GBM_6_T ; chainID=1 ; numBreaks=253

chromothripsis plus gene amplification



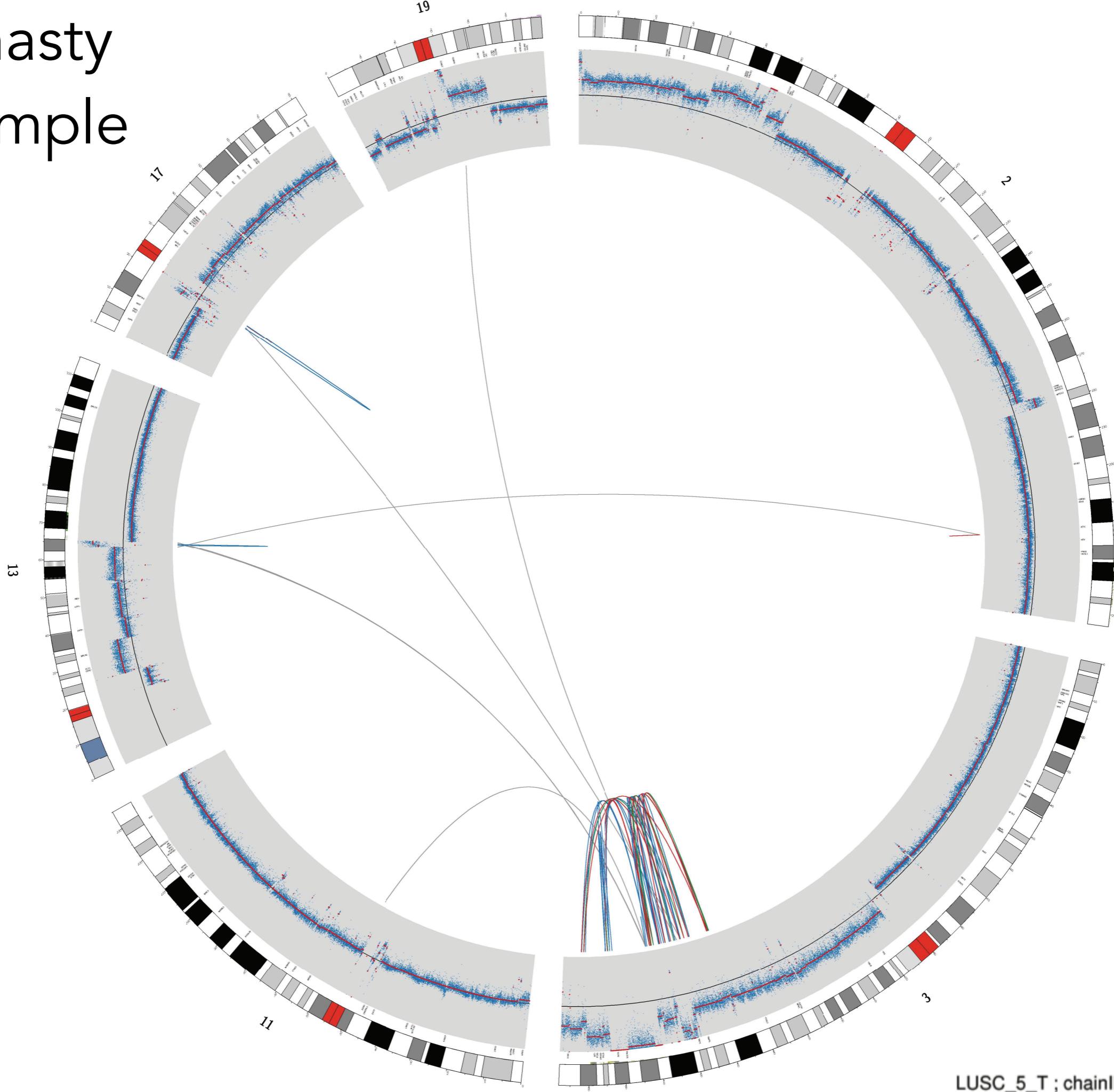
12

2mb

amplifies mdm2, cdk4, DDIT3

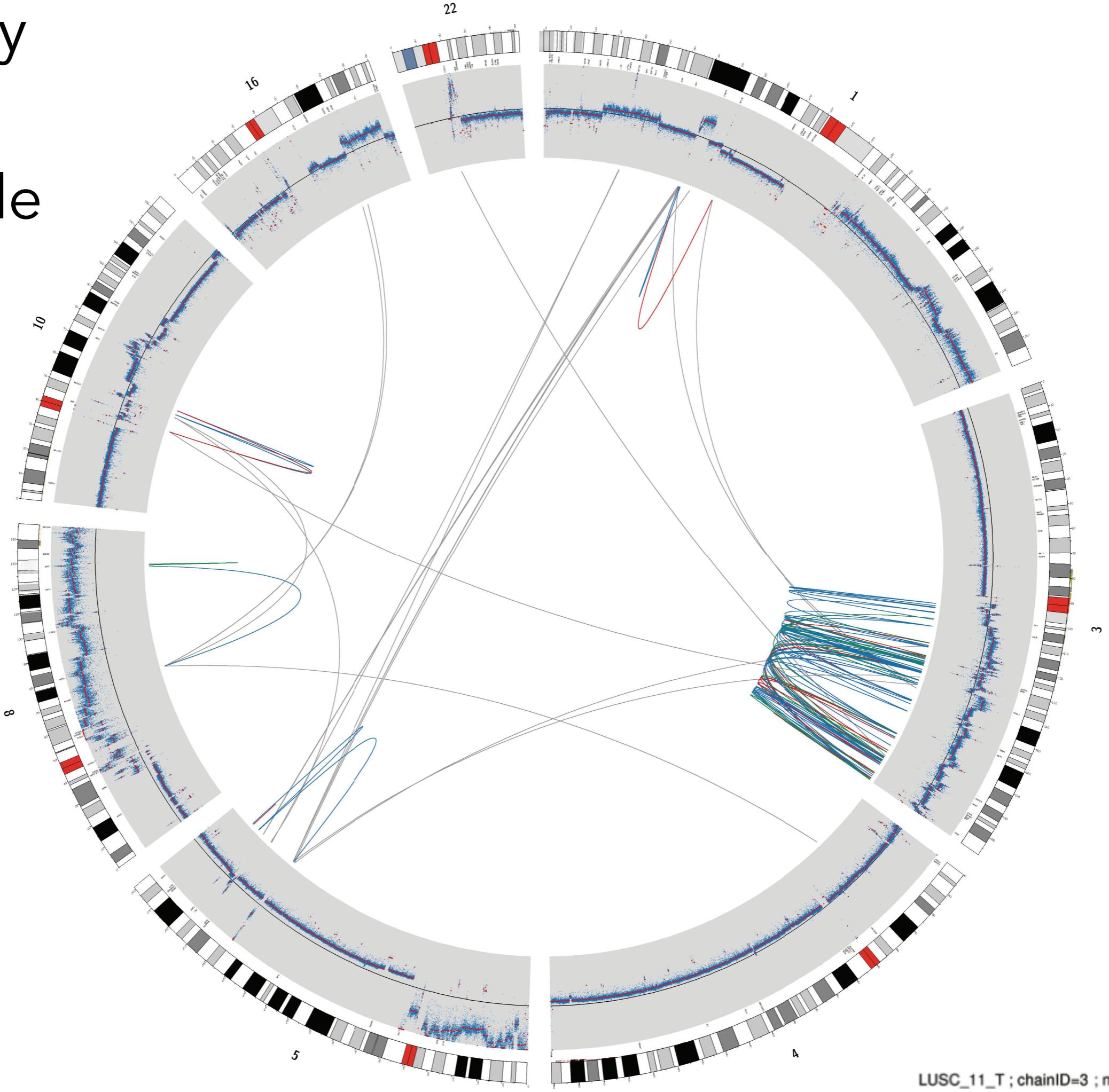
GBM_19_T ; chainID=2 ; numBreaks=56

A nasty example



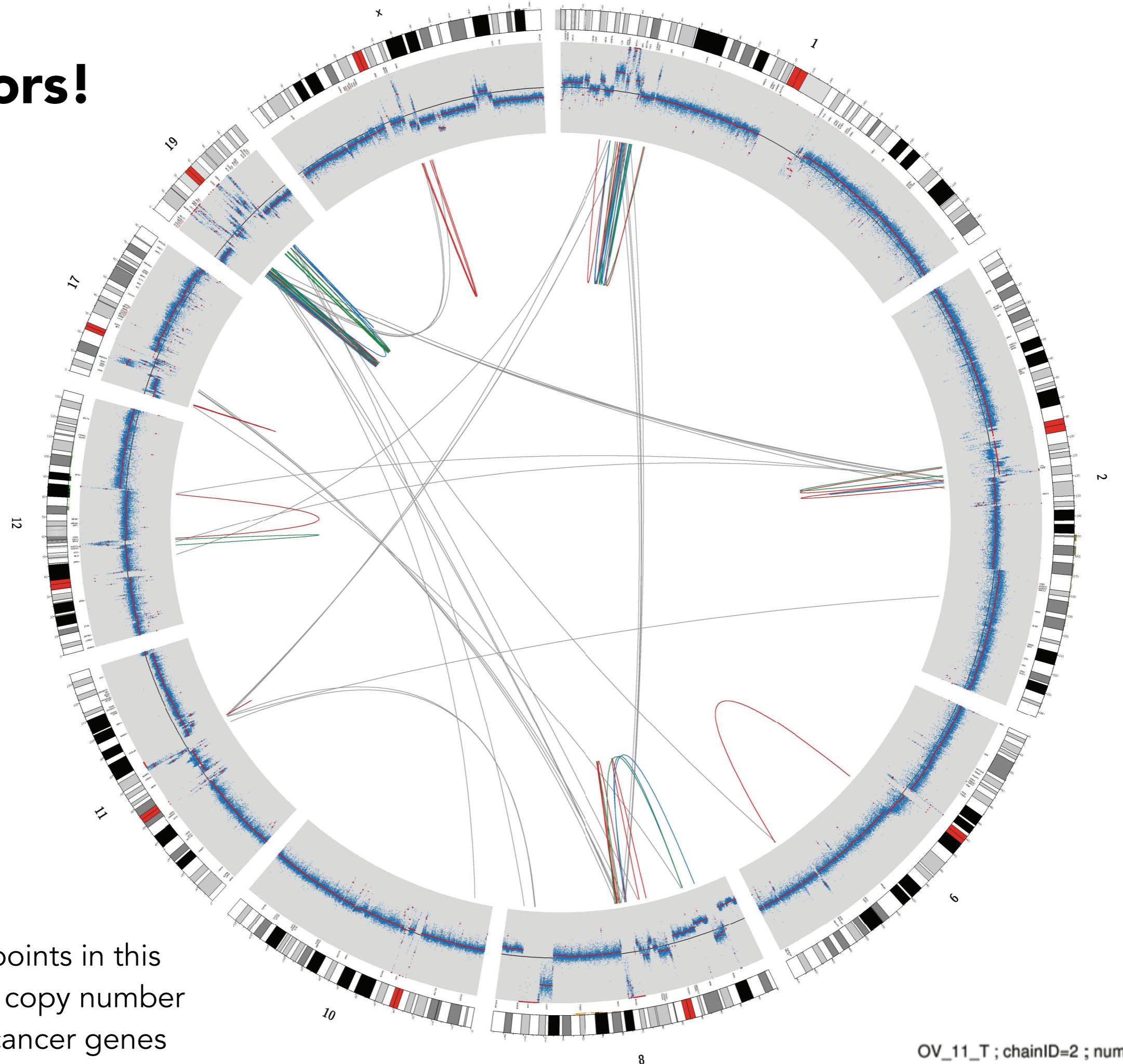
LUSC_5_T ; chainID=5 ; numBreaks=71

A really nasty example



LUSC_11_T ; chainID=3 ; numBreaks=102

Horrors!



93/132 breakpoints in this sample; alters copy number of numerous cancer genes

OV_11_T ; chainID=2 ; numBreaks=93

Summary

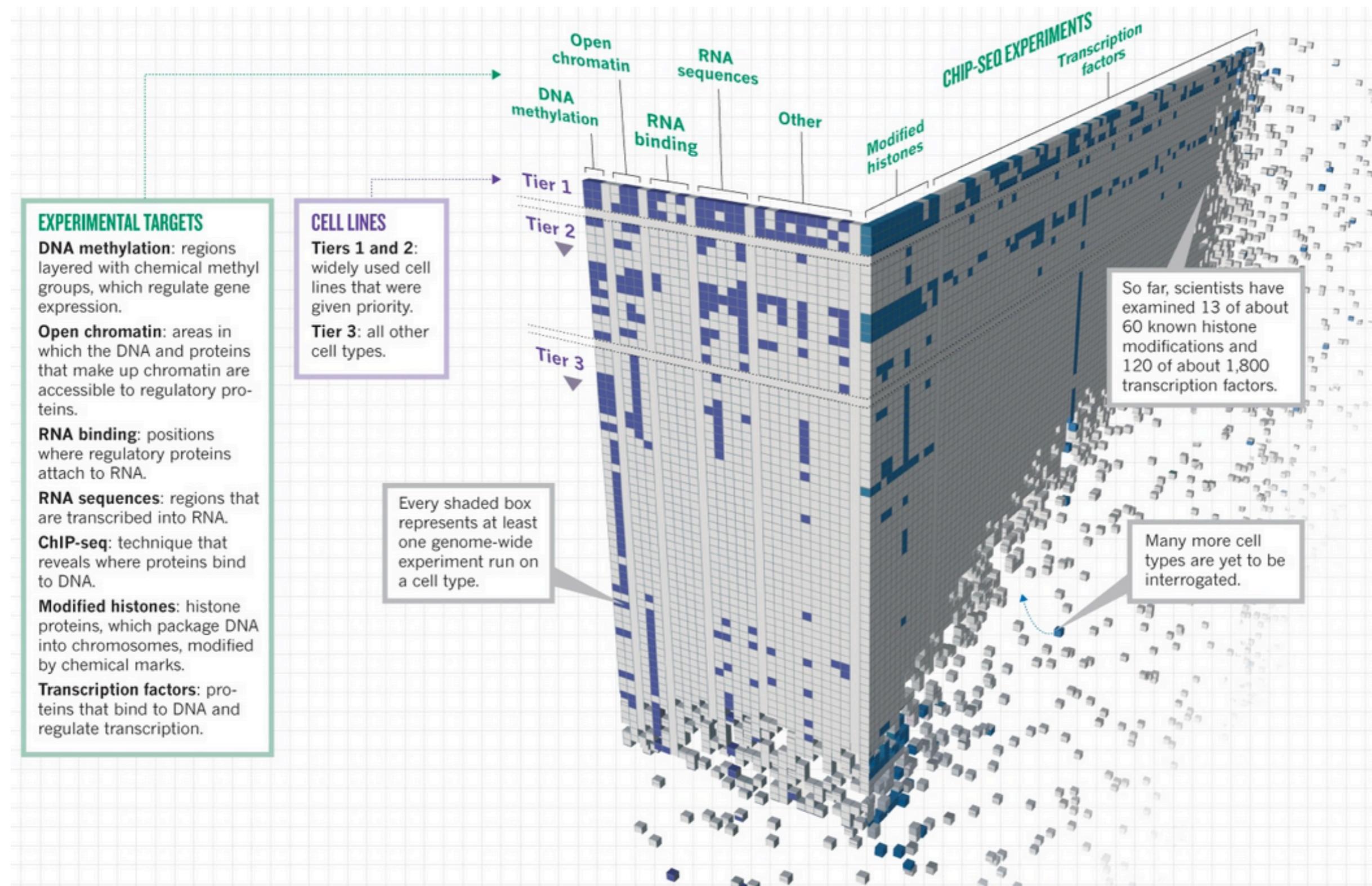
- Pooling data from multiple samples greatly improves SV discovery. Especially for identifying variants that are private to a single sample (e.g. tumors)
- Complex SVs are quite common in tumors.
- Many appear to be chromothripsis.
- 70% of glioblastomas have very complex rearrangements

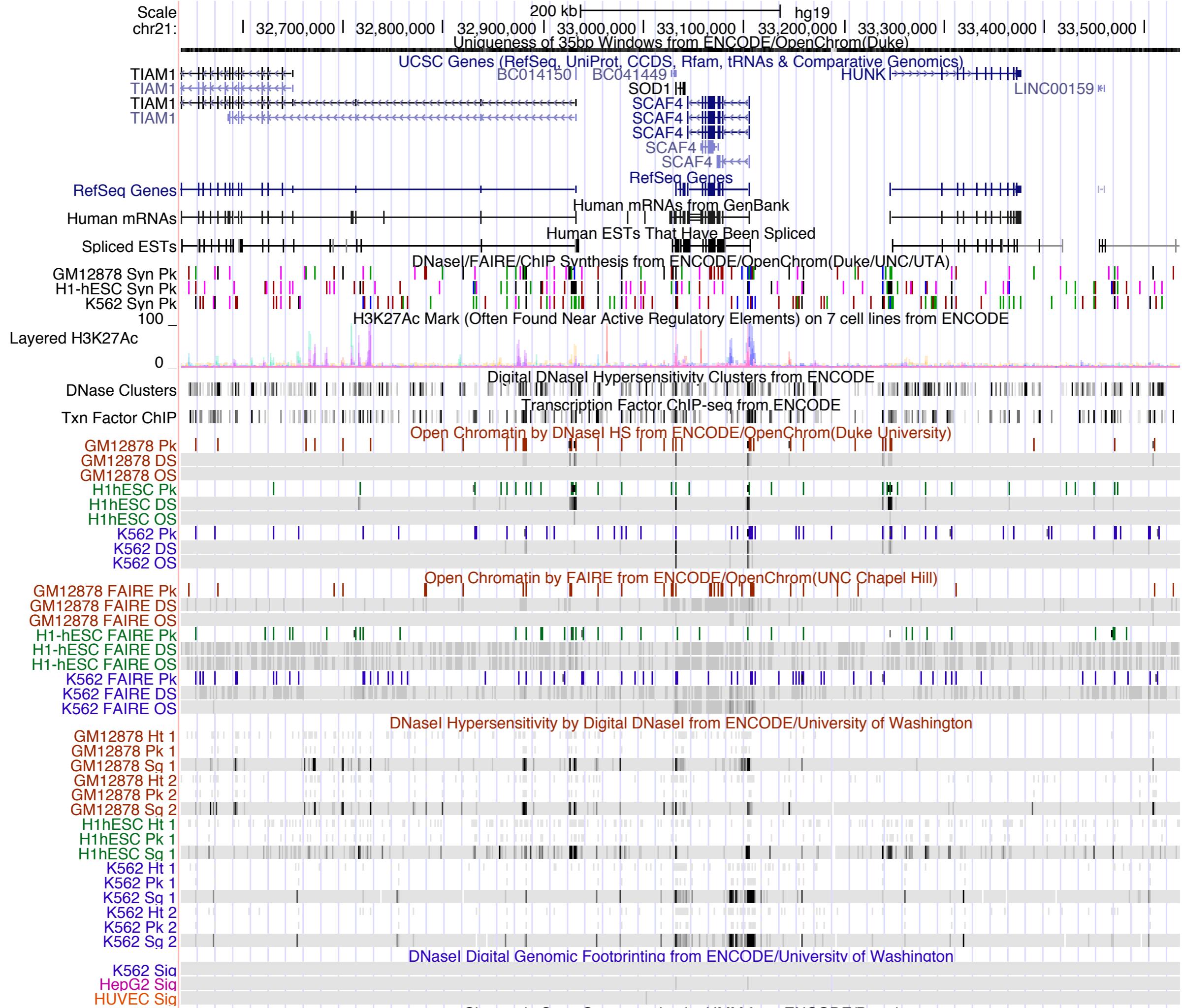
Questions we are tackling now:

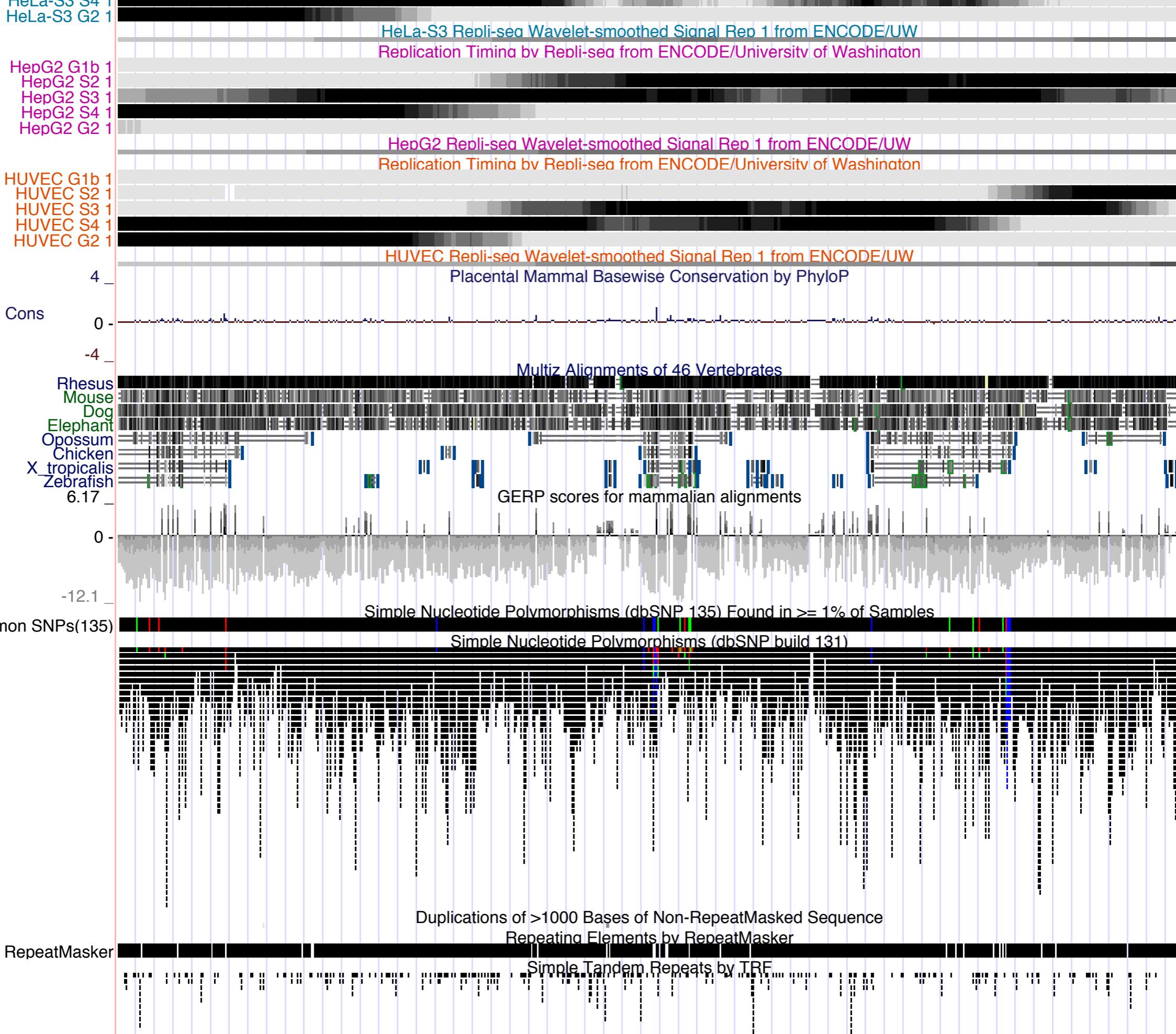
1. How do complex rearrangements arise?
 - single event or stepwise accumulation?
 - look at copy number states.
2. Are complex rearrangements early events in tumorigenesis, or the tardy consequence of genomic instability?
 - look at allele frequencies and heterogeneity.

Now for something different.

How do we make sense of complex datasets to gain insights into genome biology?





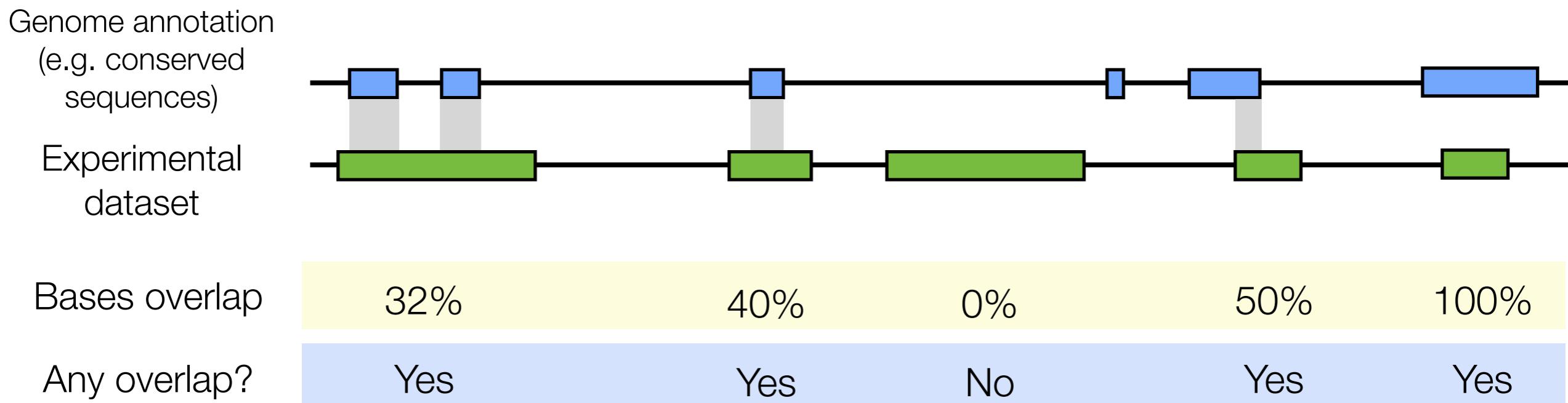


This is a hard yet important problem.

- Understand the function (or lack) of every base pair 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**

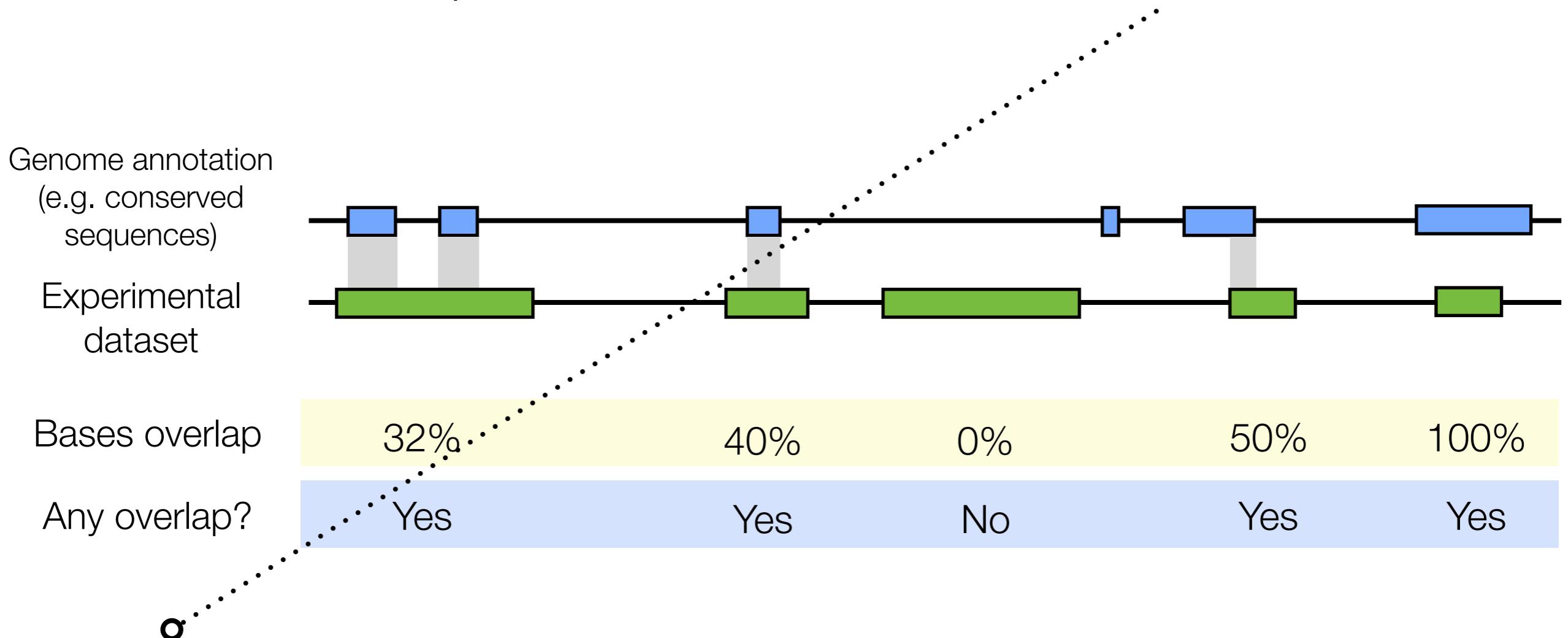
How do we detect genomic co-association?

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



How do we detect genomic co-association?

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

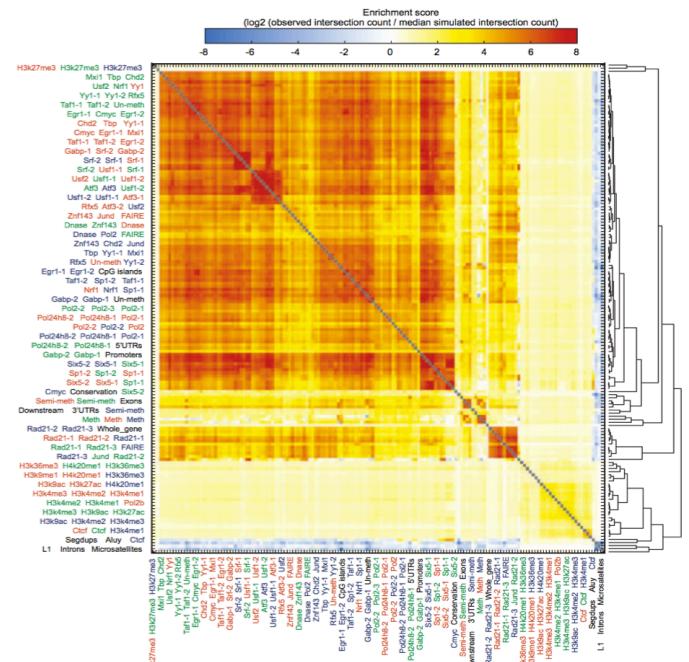


How do we develop a proper null expectation in order to reduce the dimensionality of the data to an informative statistic?

How do we develop a proper null expectation?

1. Monte-Carlo simulation

**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



How do we develop a proper null expectation?

1. Monte-Carlo simulation

2. Block bootstrap sampling

Peter Bickel)

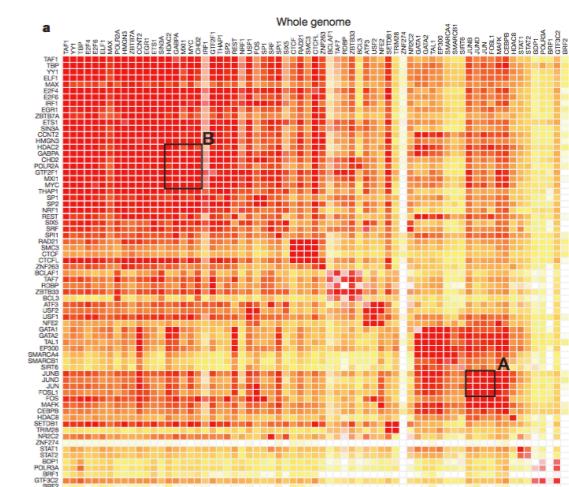
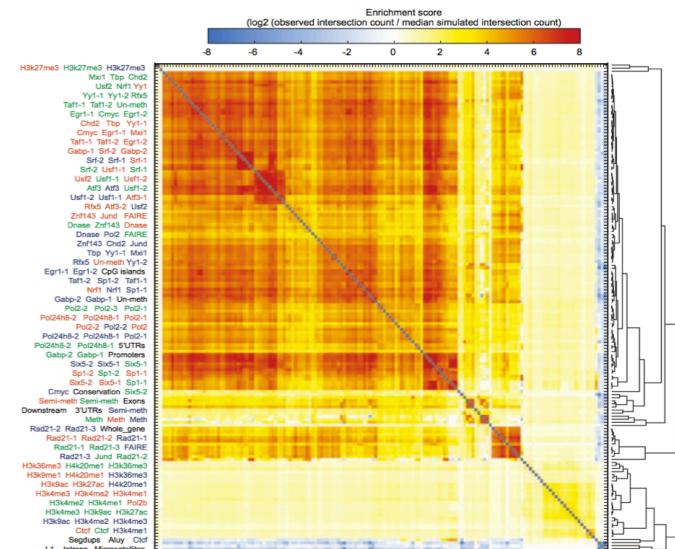
Binary Interval Search (BITS): A Scalable Algorithm for Counting Interval Intersections

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, USA

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



ARTICLE

[doi:10.1038/nature11247](https://doi.org/10.1038/nature11247)

An integrated encyclopedia of DNA elements in the human genome

The ENCODE Project Consortium*

NON PARAMETRIC METHODS FOR GENOMIC INFERENCE

BY PETER J. BICKEL*,†, NATHAN BOLEY*,†, JAMES B. BROWN*,†,
HAIYAN HUANG*,†, NANCY R. ZHANG*,‡

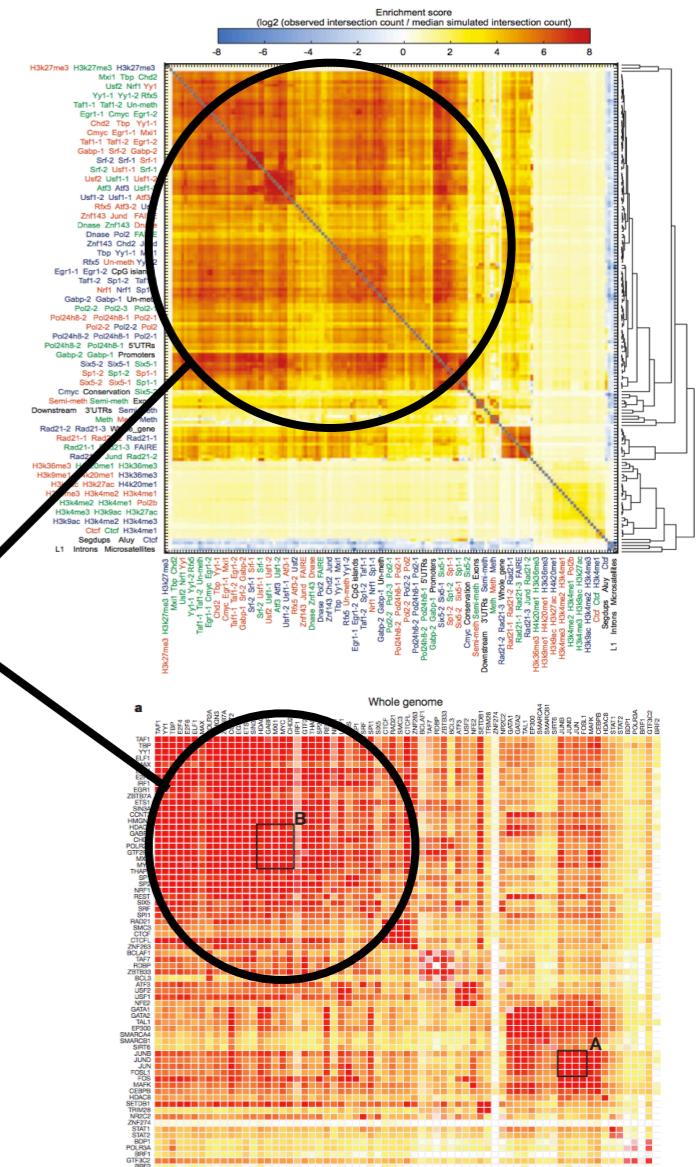
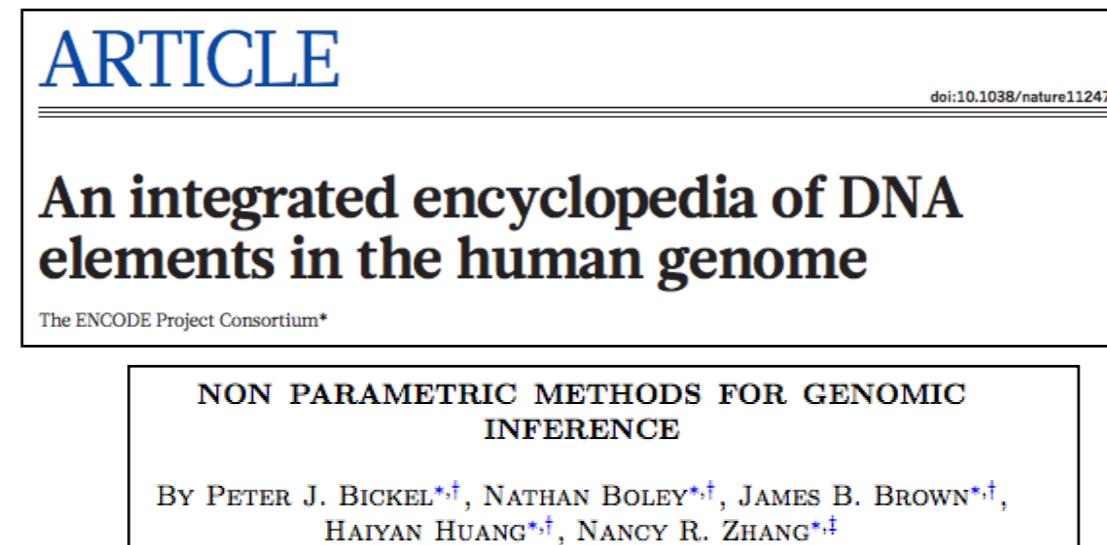
How do we develop a proper null expectation?

1. Monte-Carlo simulation

**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

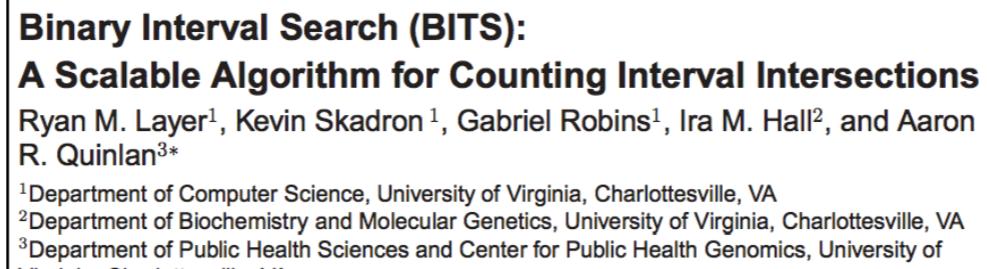
So called TF “HOT regions”
Gerstein et al; modENCODE; 2010

2. Block bootstrap sampling (Peter Bickel)



How do we develop a proper null expectation?

1. Monte-Carlo simulation



2. Block bootstrap sampling (Peter Bickel)

So called TF “HOT regions”
Gerstein et al; modENCODE; 2010



NON PARAMETRIC METHODS FOR GENOMIC INFERENCE

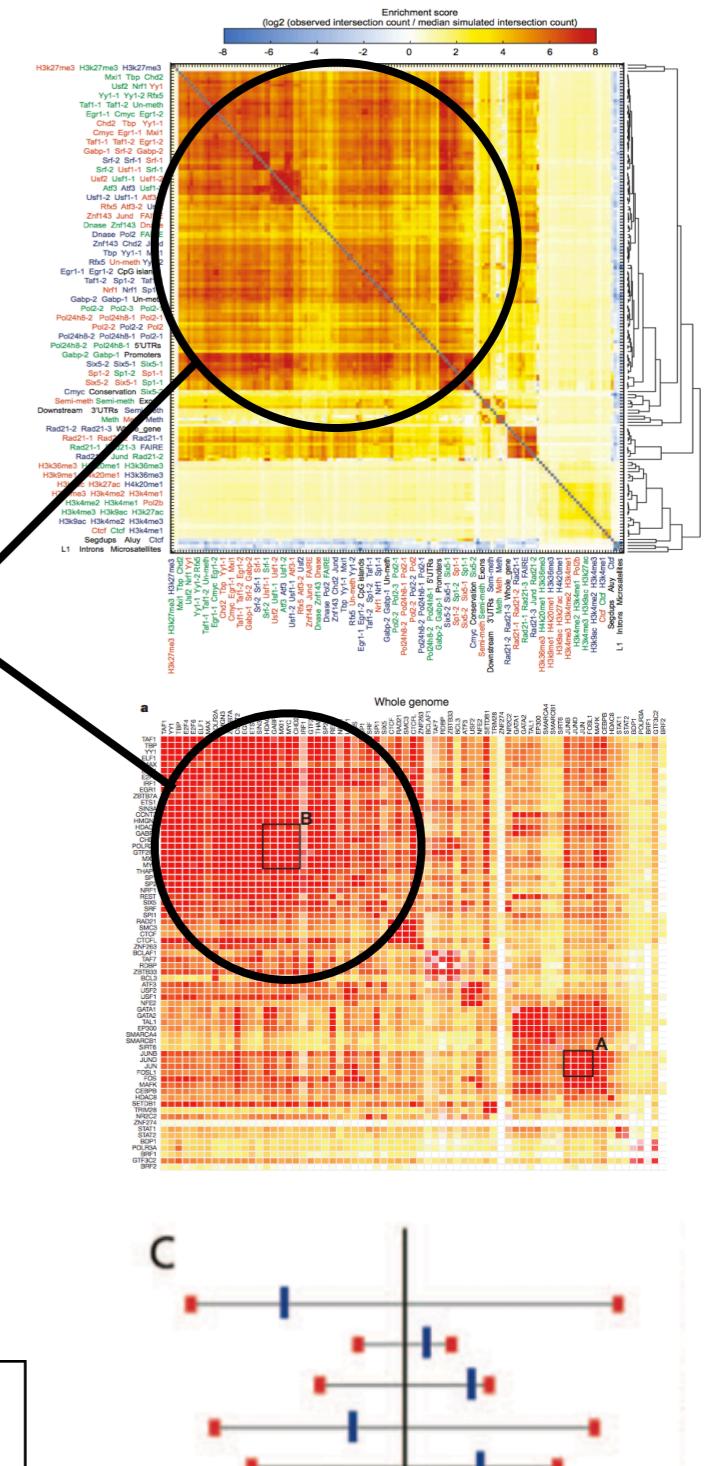
BY PETER J. BICKEL*,†, NATHAN BOLEY*,†, JAMES B. BROWN*,†,
HAIYAN HUANG*,†, NANCY R. ZHANG*,‡

3. Spatial methods (Sarah Wheelan)

Exploring Massive, Genome Scale Datasets with the GenometriCorr Package

Alexander Favorov^{1,2,3*}, Loris Mularoni^{1,3,4a}, Leslie M. Cope¹, Yulia Medvedeva^{2,3,4b},
Andrey A. Mironov^{4,5}, Vsevolod J. Makeev^{2,3}, Sarah J. Wheelan^{1*}

¹ Department of Oncology, Division of Biostatistics and Bioinformatics, Johns Hopkins University School of Medicine, Baltimore, Maryland, United States of America,
²Vavilov Institute of General Genetics, Russian Academy of Sciences, Moscow, Russia, ³ Research Institute of Genetics and Selection of Industrial Microorganisms, Moscow, Russia, ⁴ Department of Bioengineering and Bioinformatics, Moscow State University, Moscow, Russia, ⁵ Institute for Information Transmission Problems, Russian Academy of Sciences, Moscow, Russia

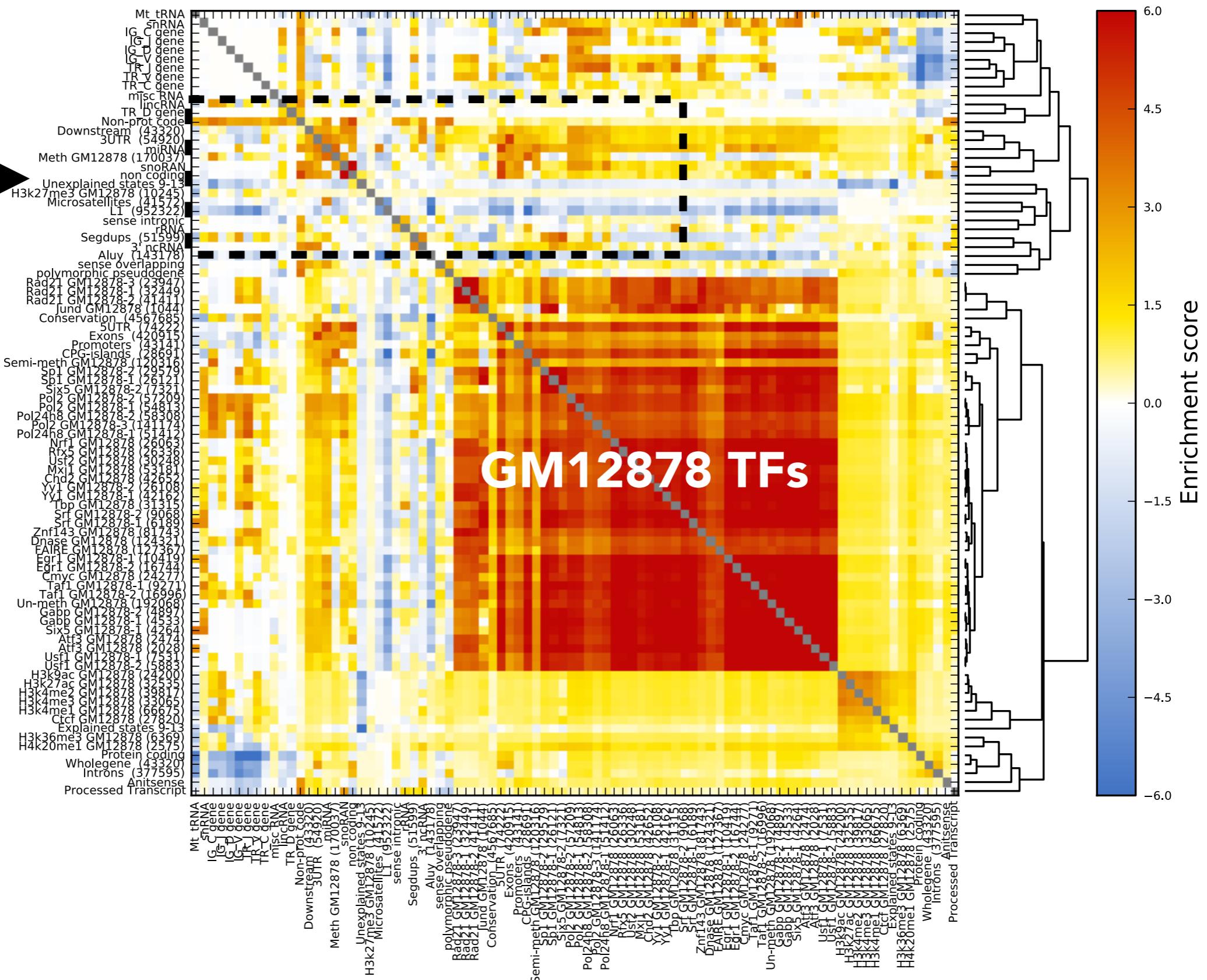


query is randomly distributed with respect to the reference

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

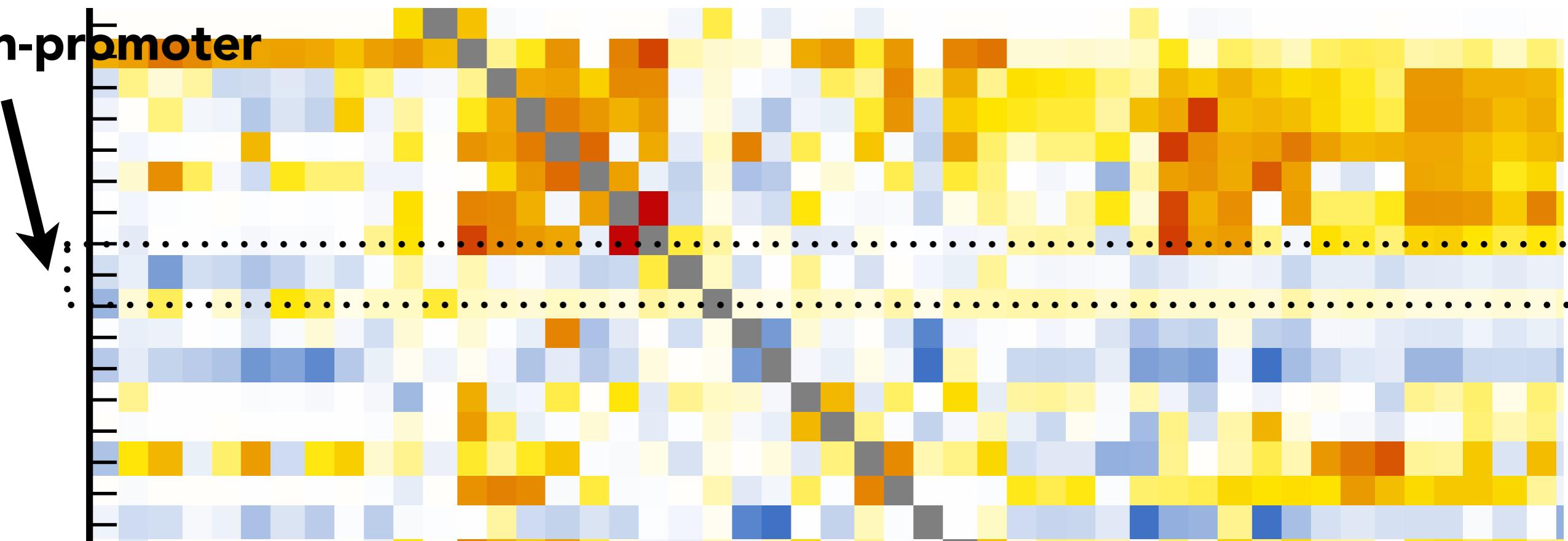
What is the role of conserved, non-coding (CNC) elements?

CNC
elements



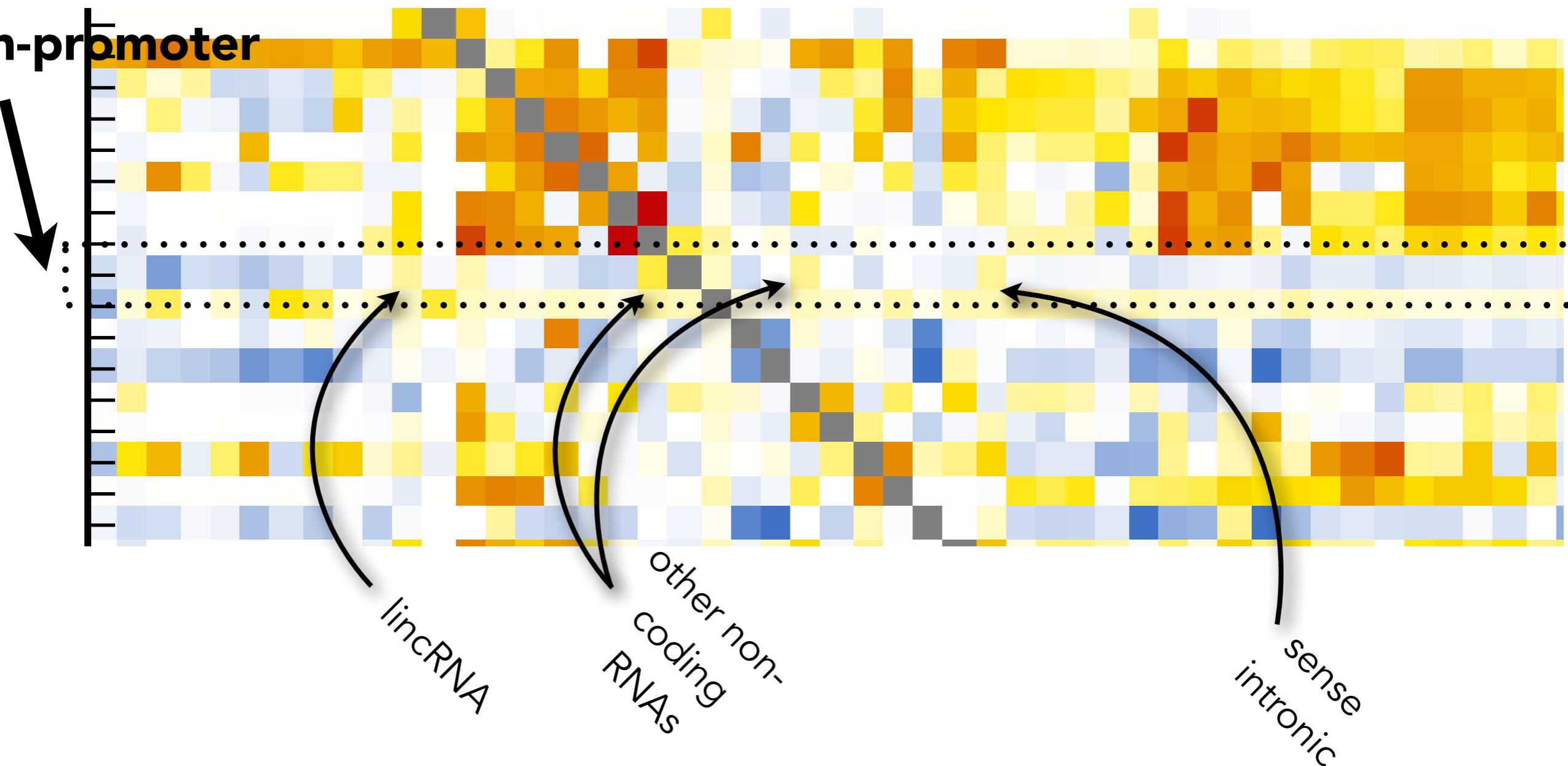
Making sense of conserved, yet non-coding, and non-“functional” elements

**Conserved,
non-coding,
non-enhancer,
non-promoter**



Making sense of conserved, yet non-coding, and non-“functional” elements

**Conserved,
non-coding,
non-enhancer,
non-promoter**



Genome Query Language (GQL)

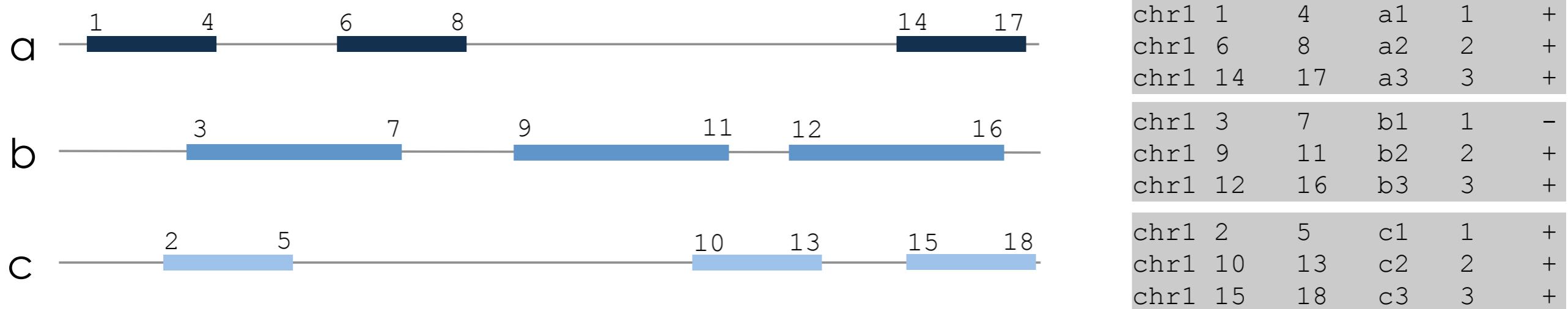
(coming in 2013)

- Writing a new, simple, & expressive language for exploring complex datasets. Based on SQL.
- Set theory, statistics, visualization, data recruitment
- Abstract the language from the engine underneath
- Powerful, fast, and easy to use.
- Built upon bedtools and bedtools2

A GQL example

(find all intervals where at least 2 datasets overlap)

(e.g., find “hot regions” for TF ChIP-seq peaks)



```
> result = SELECT a,b,c WHERE COUNT(>1);  
> PRINT result  
> PLOT result  
> MEASURE result WITH(MONTE_CARLO) AND a,b,c
```

Quinlan lab



Uma Paila, Ph.D.

Postdoctoral Research Associate

udp3f @ virginia.edu

Research Projects and Interests: Investigation of the genetic basis of extreme sensitivity to ionizing radiation; development of new analytical tools for exploring genetic variation identified through next-generation sequencing projects.



John Kubinski
Undergraduate
(Biology)



Neil Kindlon, M.S.

Staff Scientist and Software Engineer

nek3d @ virginia.edu

Research Projects and Interests: Software development for genomic analysis. Structural variation discovery and interpretation using DNA sequencing technologies.

Gift Sinthong
Undergraduate
(Comp. Science)



Ryan Layer

Graduate student

rl6sf @ virginia.edu

Research Interests: Scalable algorithm development for high-throughput genomic analysis; genome data mining and analysis; structural variation discovery and interpretation.

Acknowledgements

Ira Hall

Ankit Malhotra	<i>Univ. of Virginia</i>
Michael Lindberg	<i>Univ. of Virginia</i>
Royden Clark	<i>Univ. of Virginia</i>
Svetlana Sokolova	<i>Univ. of Virginia</i>
Mitchell Leibowitz	<i>Univ. of Virginia</i>

Pat Concannon	<i>Univ. of Virginia</i>
Steve Rich	<i>Univ. of Virginia</i>
Suna Onengut-Gumuscu	<i>Univ. of Virginia</i>
Chris Moskaluk	<i>Univ. of Virginia</i>
Shu-Man Fu	<i>Univ. of Virginia</i>
Gabor Marth	<i>Boston College</i>
Jim Robinson	<i>Broad Institute</i>
Nick Navin	<i>MD Anderson</i>
Kristin Baldwin	<i>Scripps</i>

Nik Krumm	<i>Univ. of Washington</i>
Evan Eichler	<i>Univ. of Washington</i>
Debbie Nickerson	<i>Univ. of Washington</i>
Chris Carlson	<i>Fred Hutchinson CRC</i>
Mark Rieder	<i>Univ. of Washington</i>
Josh Smith	<i>Univ. of Washington</i>
Peter Sudmant	<i>Univ. of Washington</i>

Funding

NHGRI: R01 HG006693-01

NIEHS: R21 ES020521-01

UVA Fund for Excellence in Science and Tech.

UVA Cancer Center Pilot Program