

# Cancer genomic rearrangements

Marcin Imielinski M.D. Ph.D.

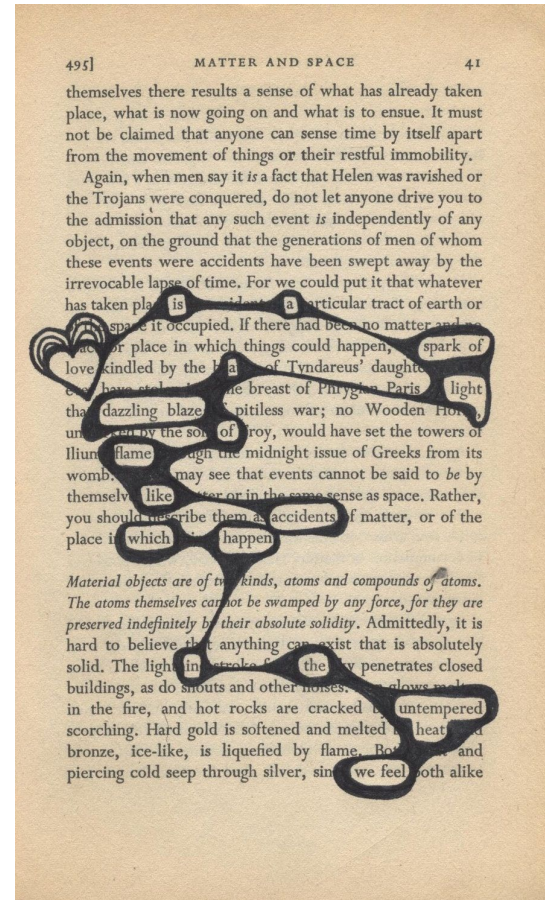
Assistant Professor (Weill Cornell)

Core Member (NYGC)

Clinical and Research Genomics

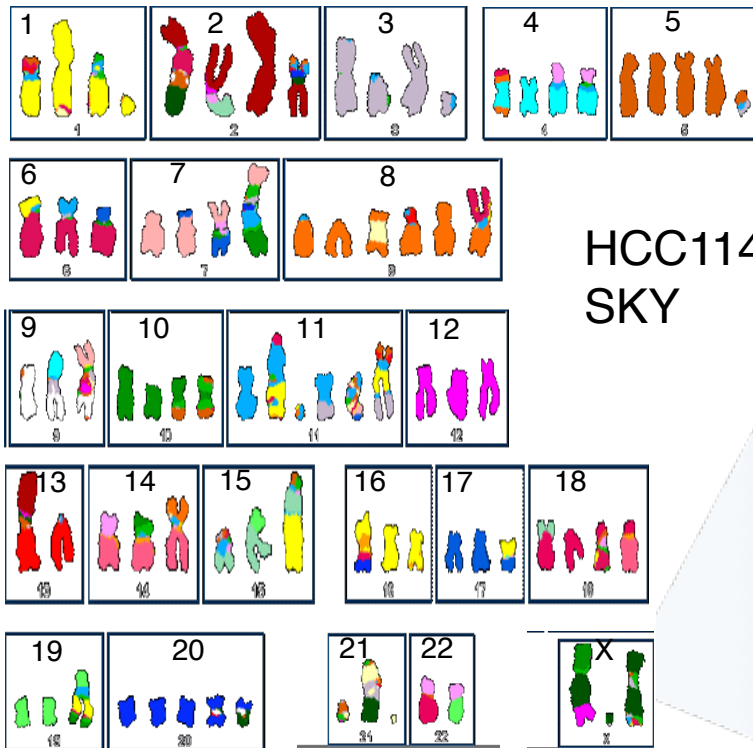
April 4 2018

# If the human genome is a book ...



... cancer is collage poetry

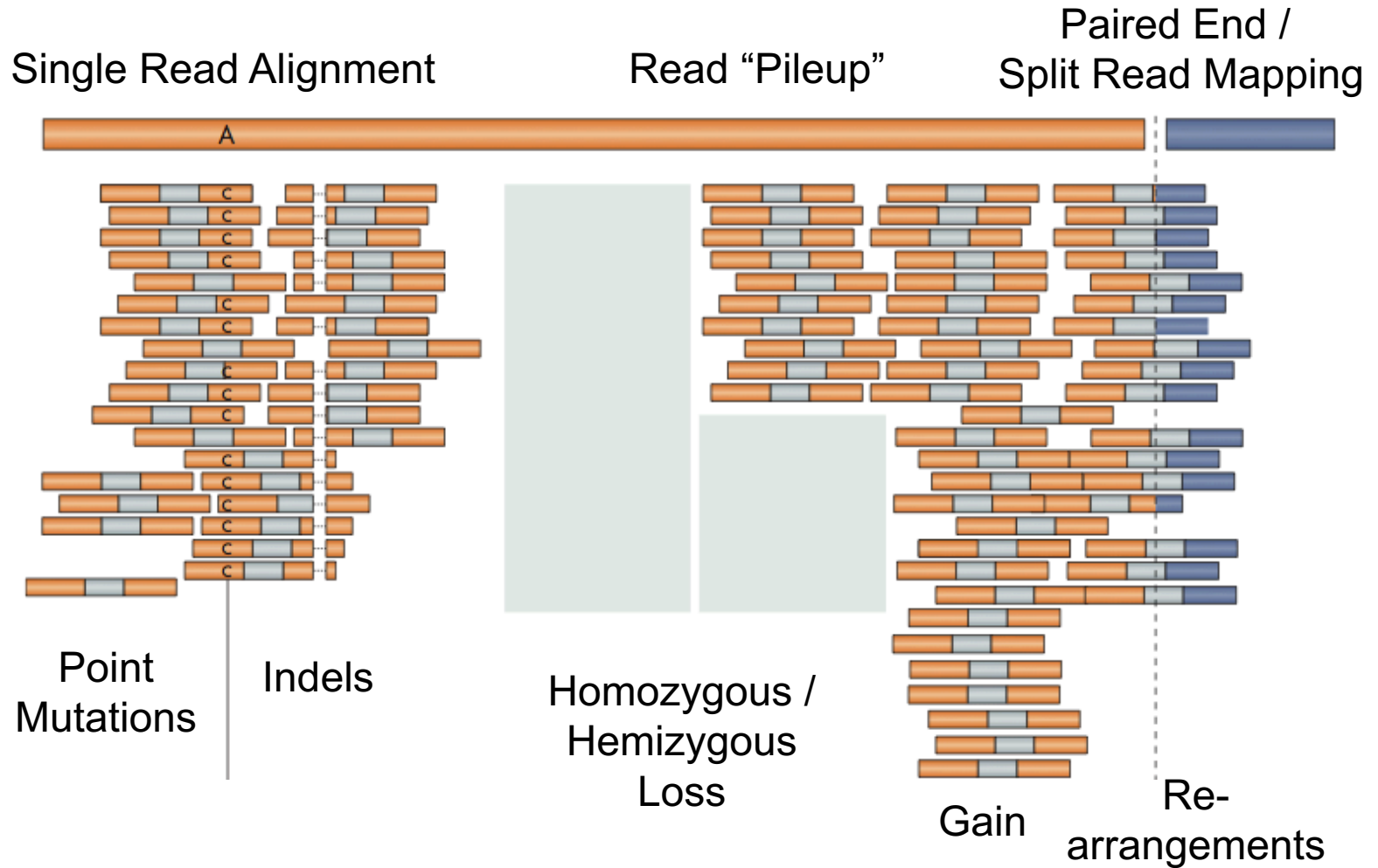
# Chromosomal collages



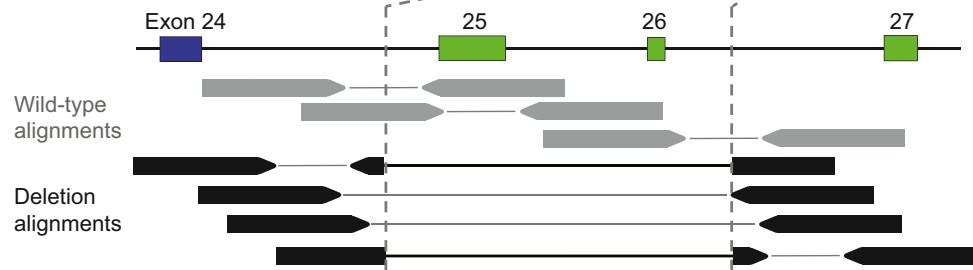
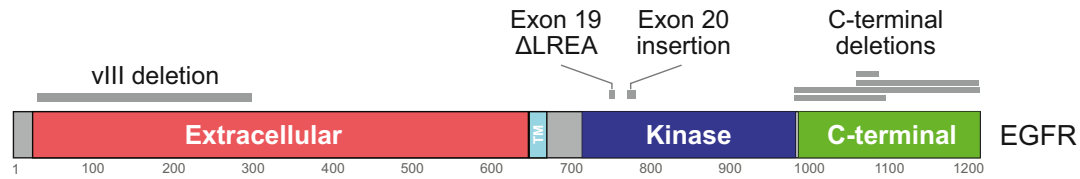
HCC1143  
SKY



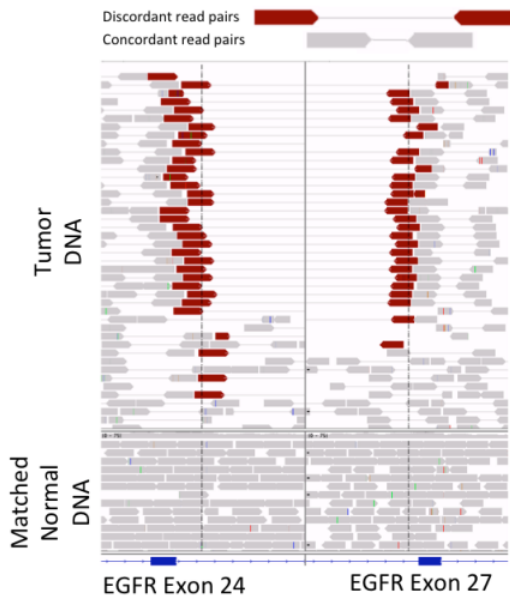
# Chromosomal shreds



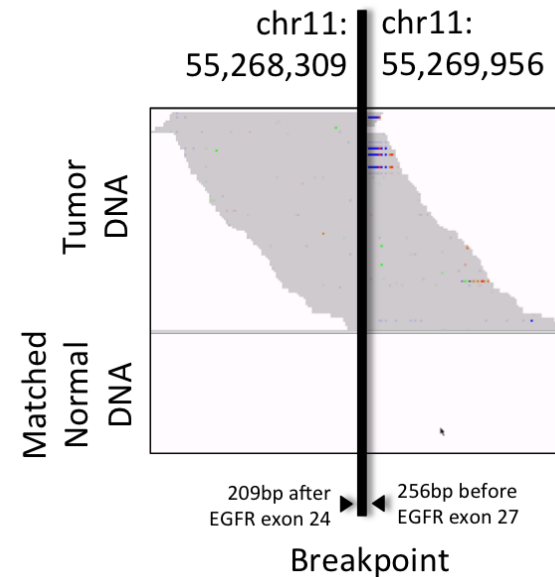
# Junction detection from tiny eads



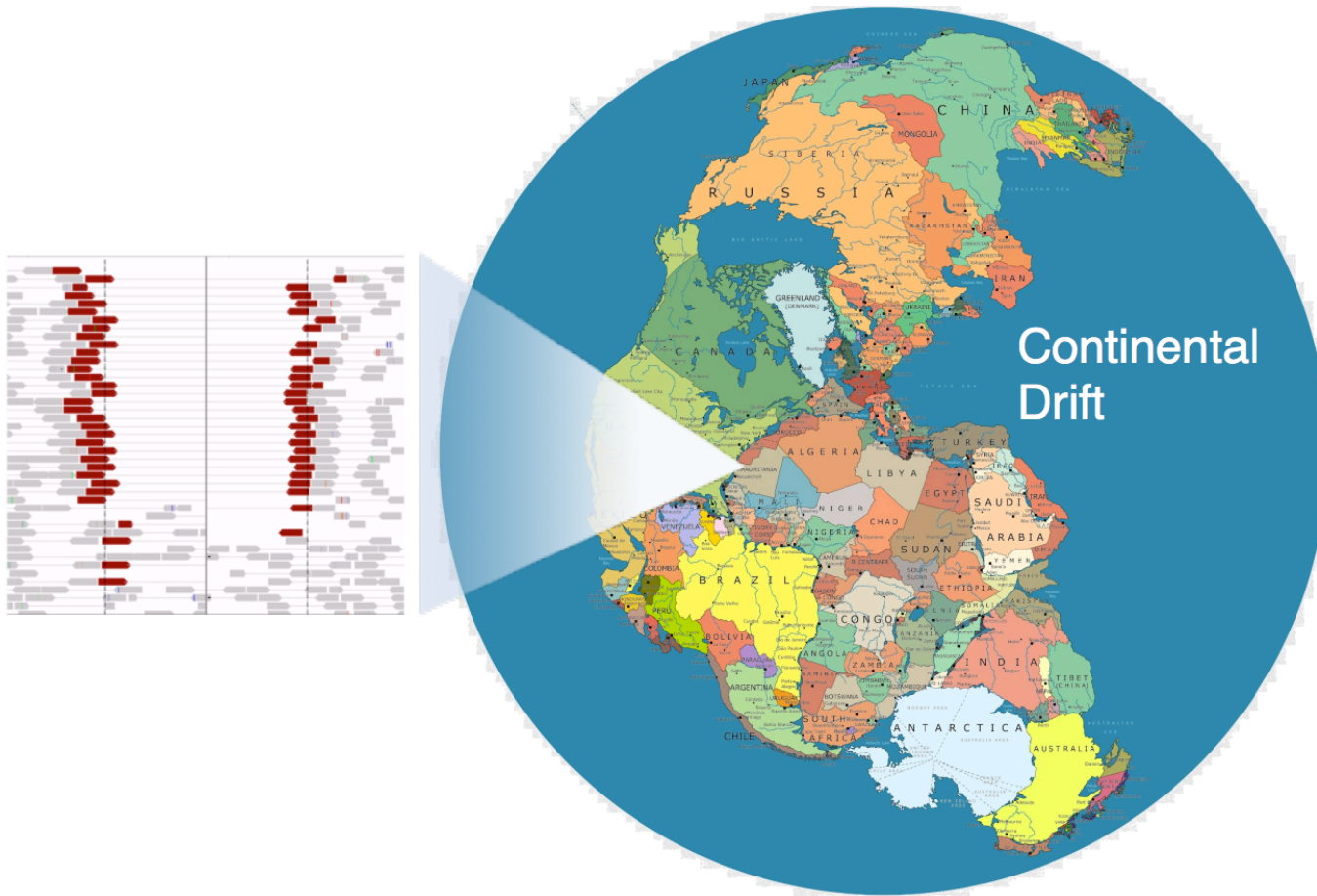
EGFR deletion paired end mapping



EGFR deletion Split read mapping



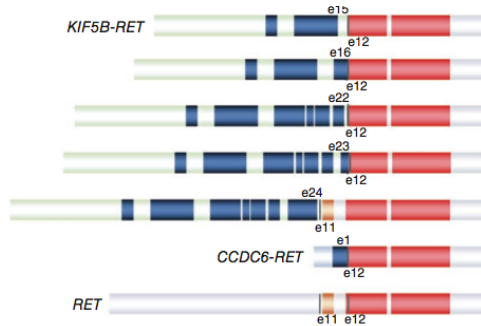
# Chromosomal Pangaea



The Earth  
(500 million years ago)

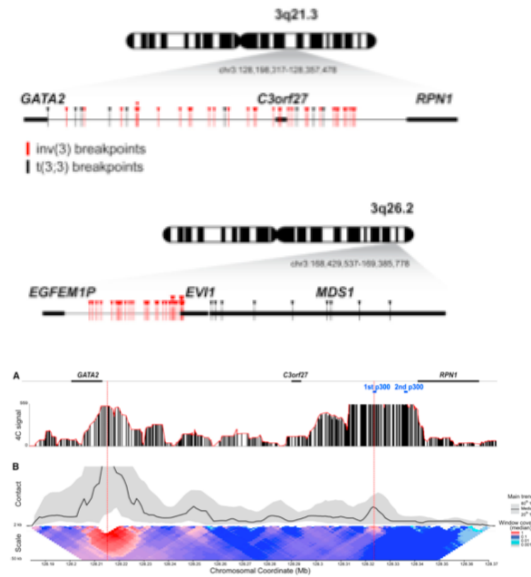
# Studying cancer genome structure: motivation

## Gene Fusions



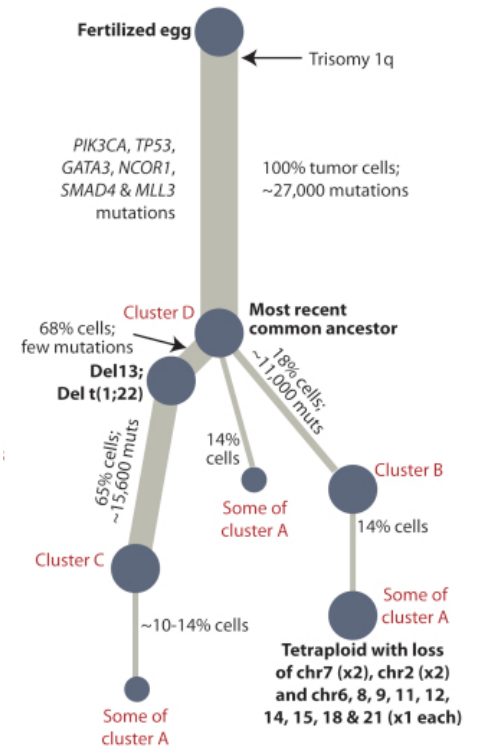
Takeuchi .. Ishikawa  
(Nature Medicine 2012)

## Noncoding Rearrangements



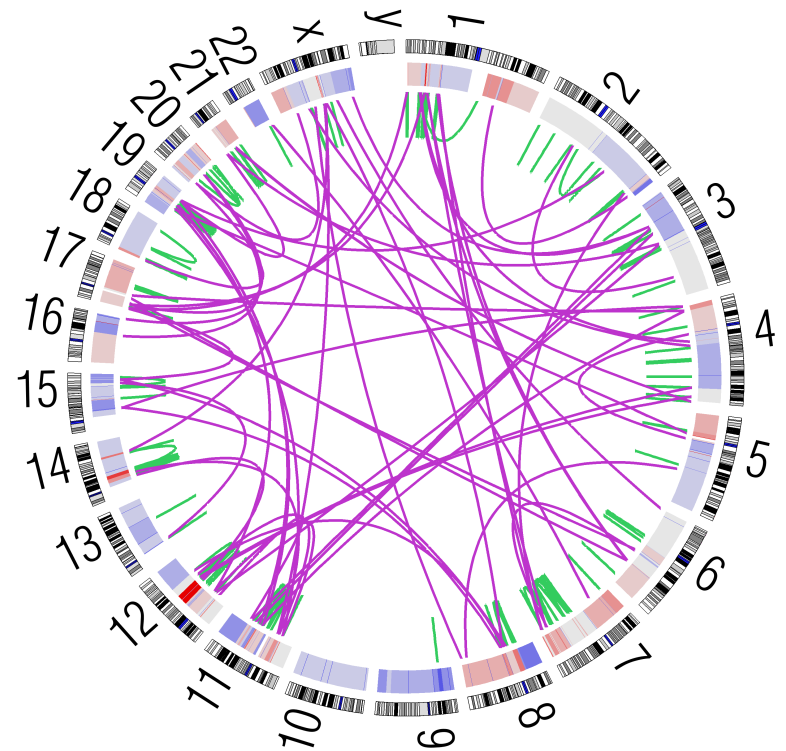
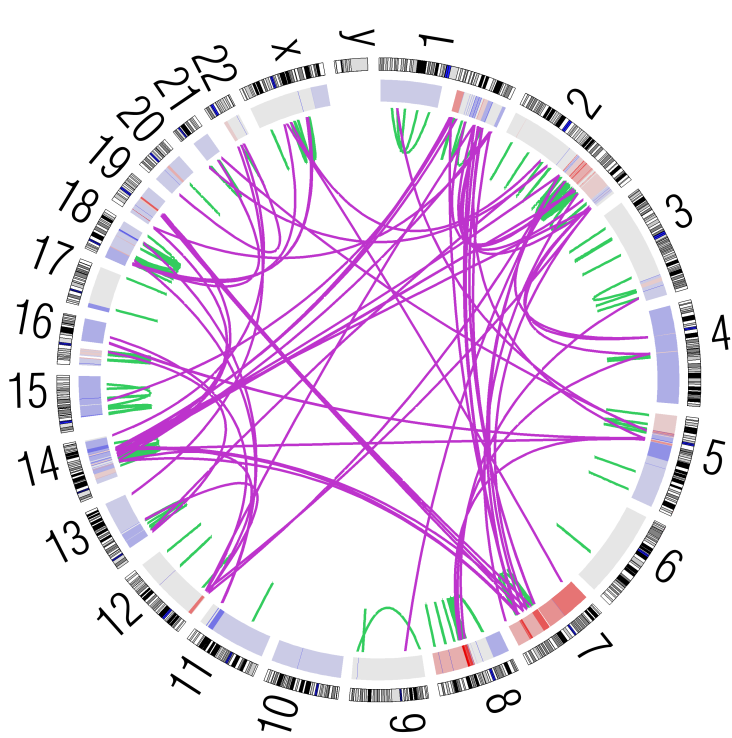
Groschel .. Delwel  
(Cell 2014)

## Mutational processes



Nik-Zainal .. Stratton  
(Cell 2012)

# Which is the smoker?

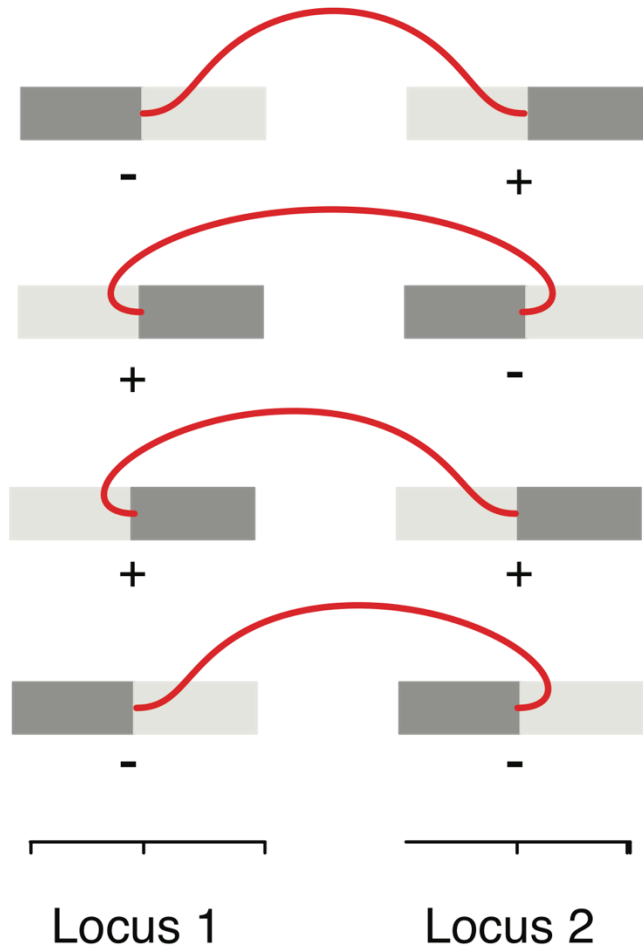




# Circos .. so beautiful

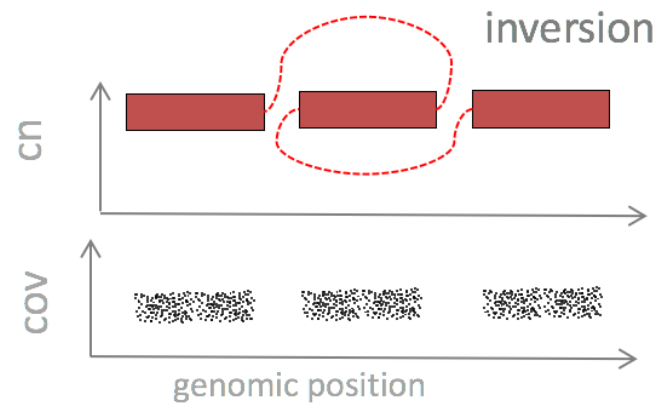
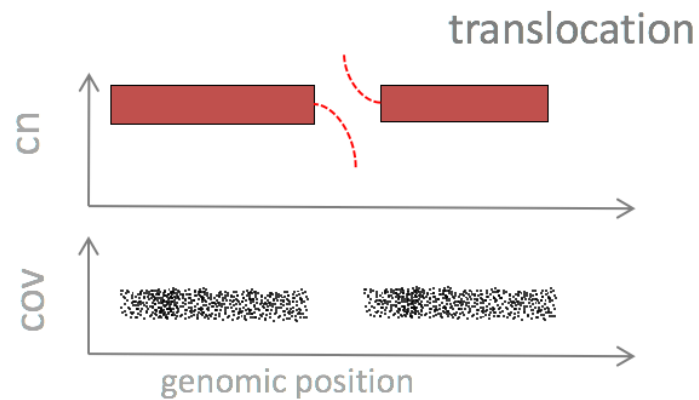
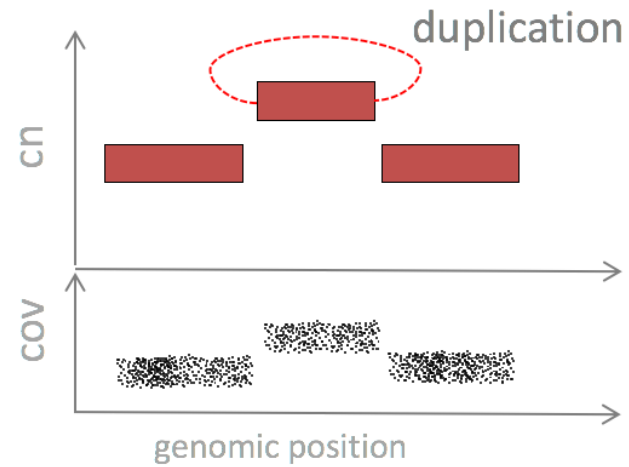
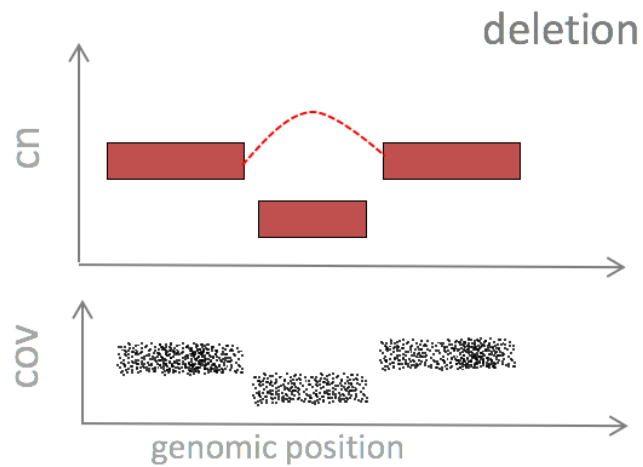


# Junction: the “atomic unit” of a genomic rearrangement

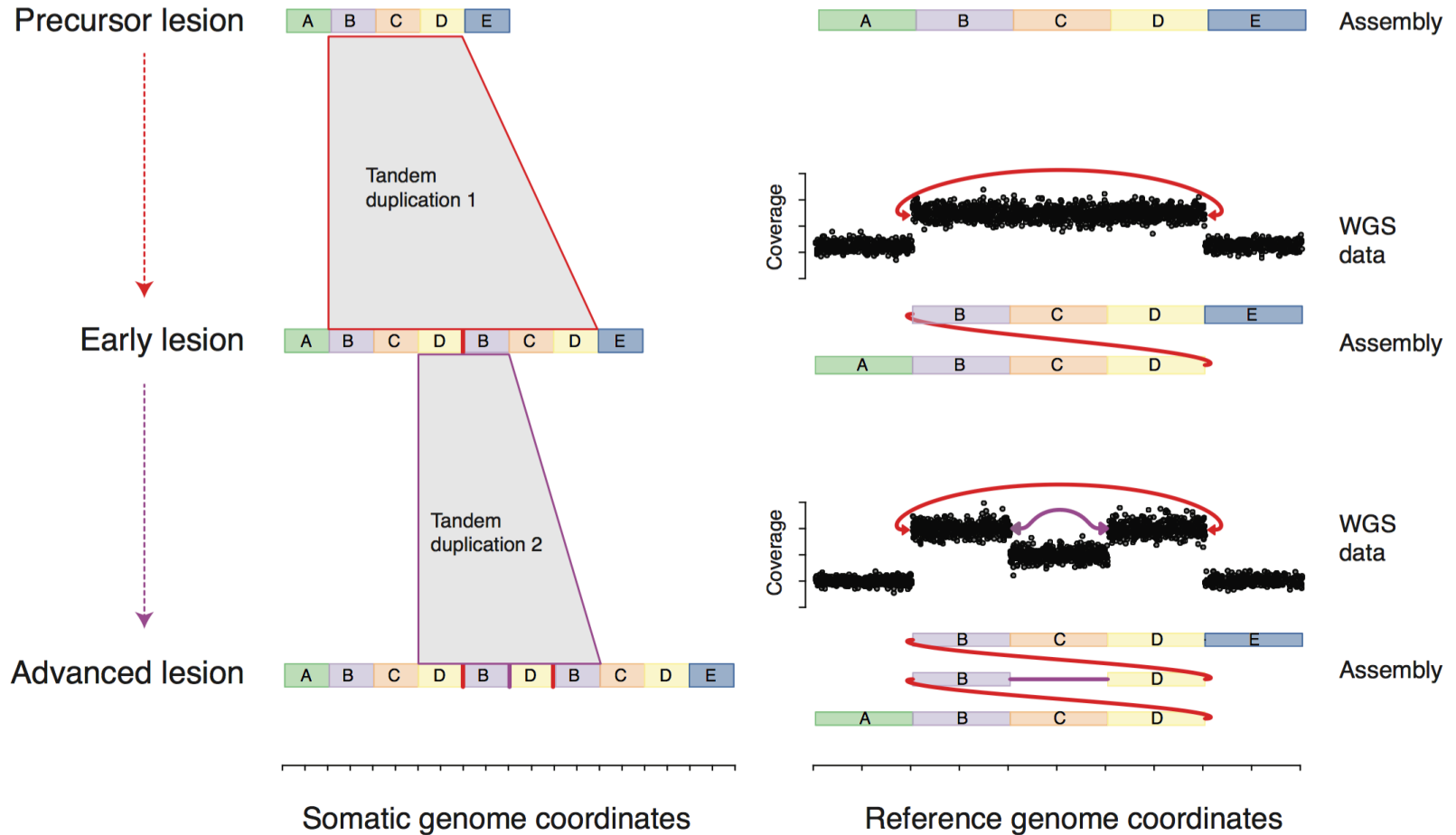


Junction =  
pair of  
locations  
AND  
orientations

# Can it be all so simple?

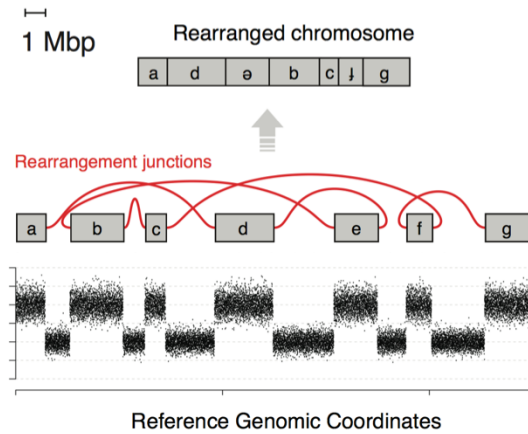


# What is an “event”?

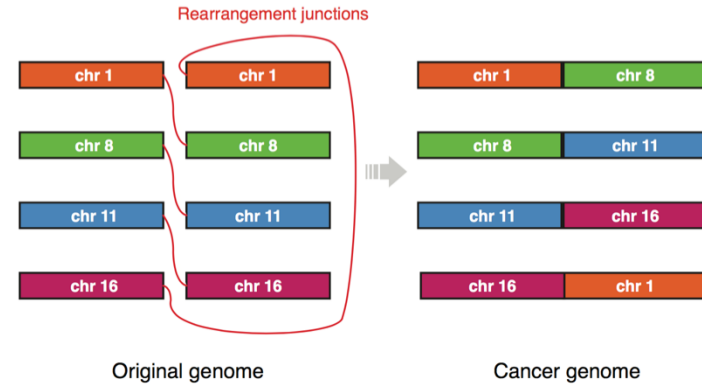


# Complex structural variation in cancer

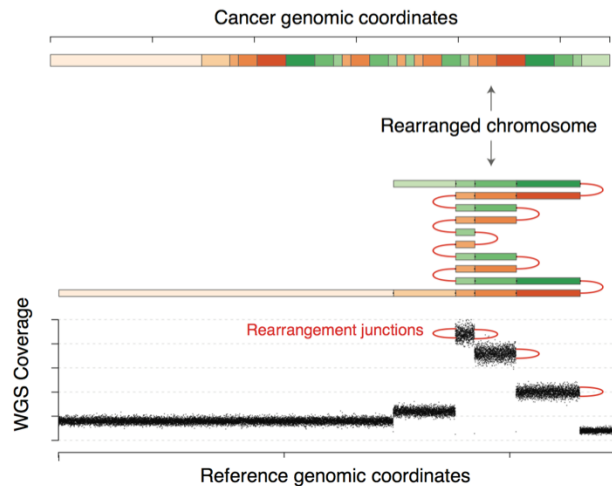
## Chromothripsis



## Chromoplexy

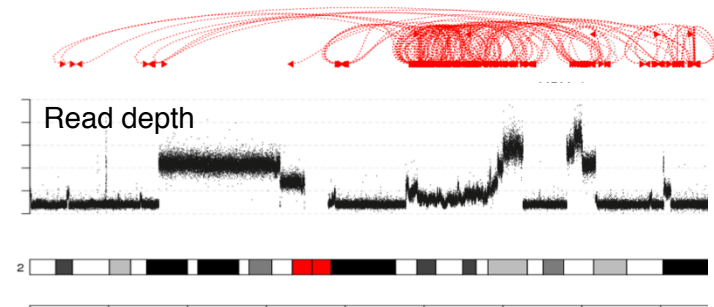


## Breakage-fusion-bridge

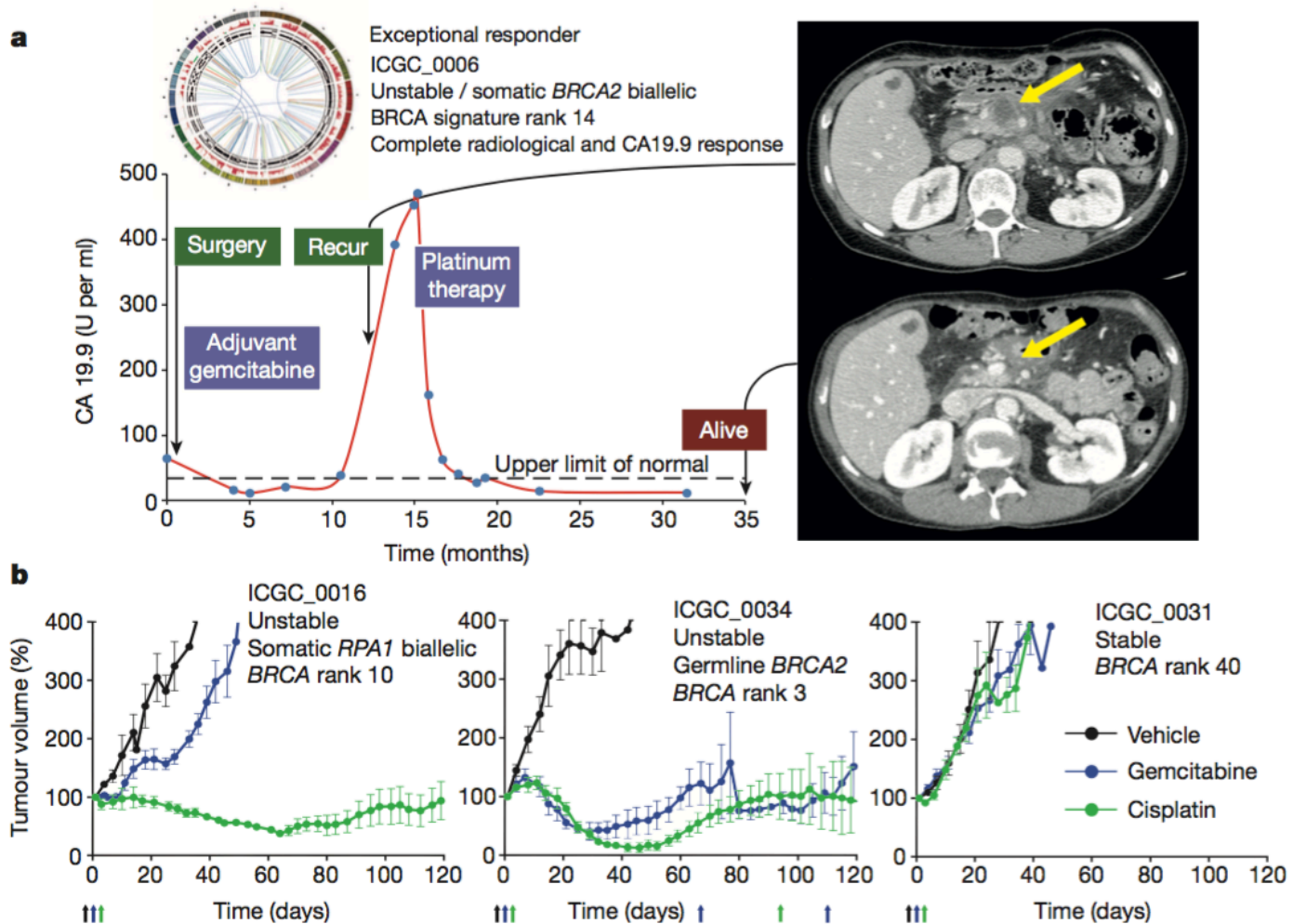


??????

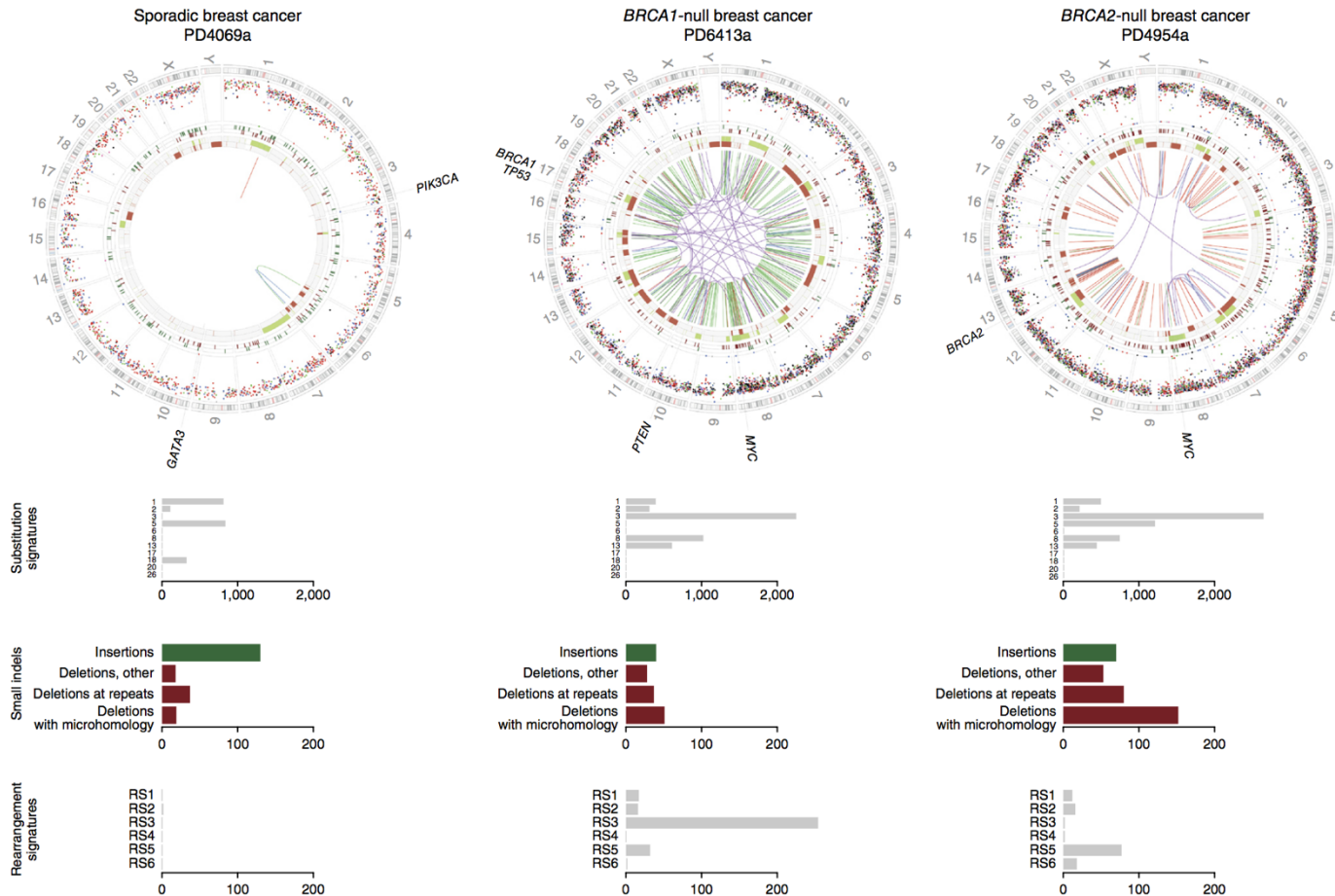
Junctions



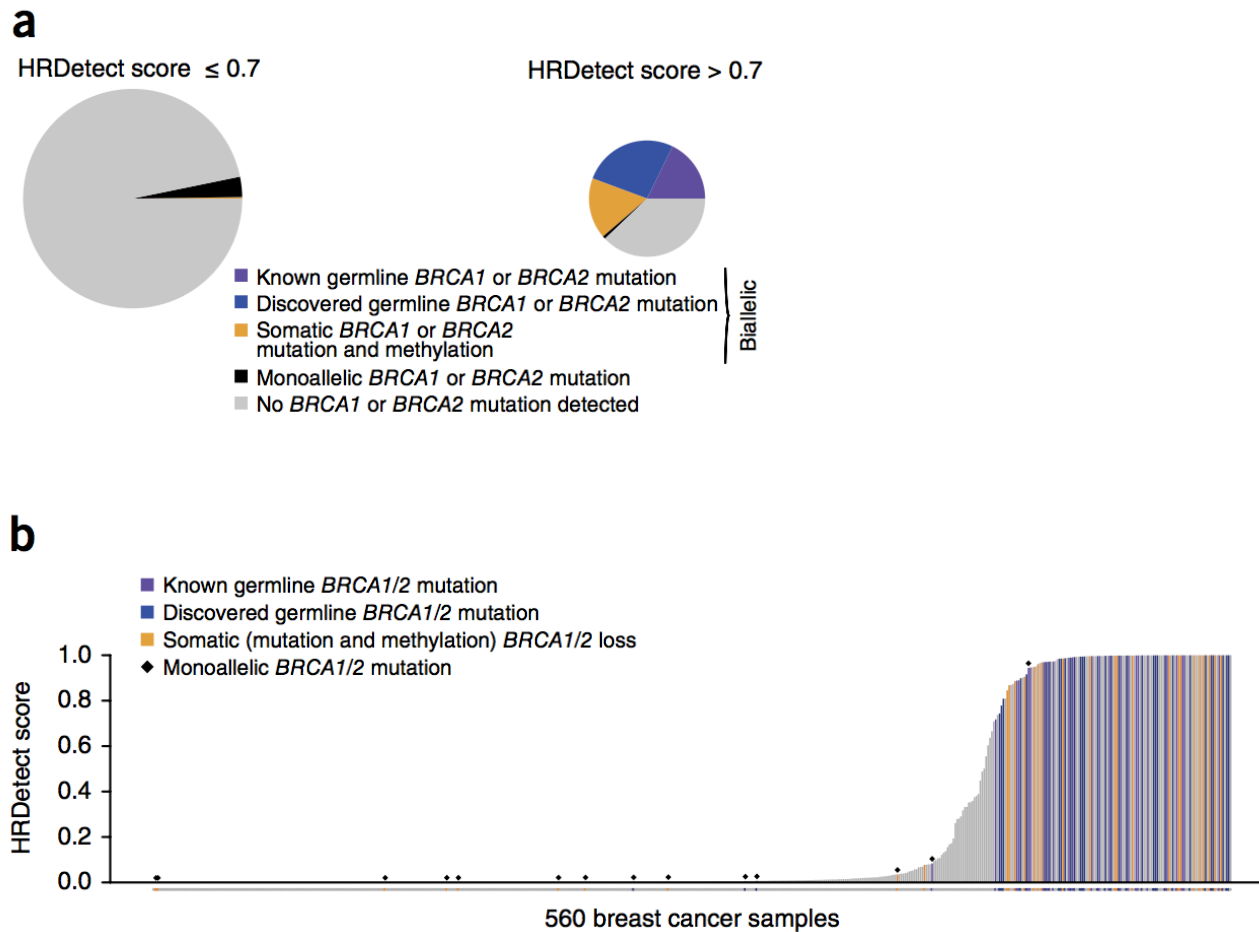
# Clinical consequences of rearrangement signatures: exceptional chemotherapy response



# Clinical consequences of rearrangement signatures: “BRCAness phenotype”



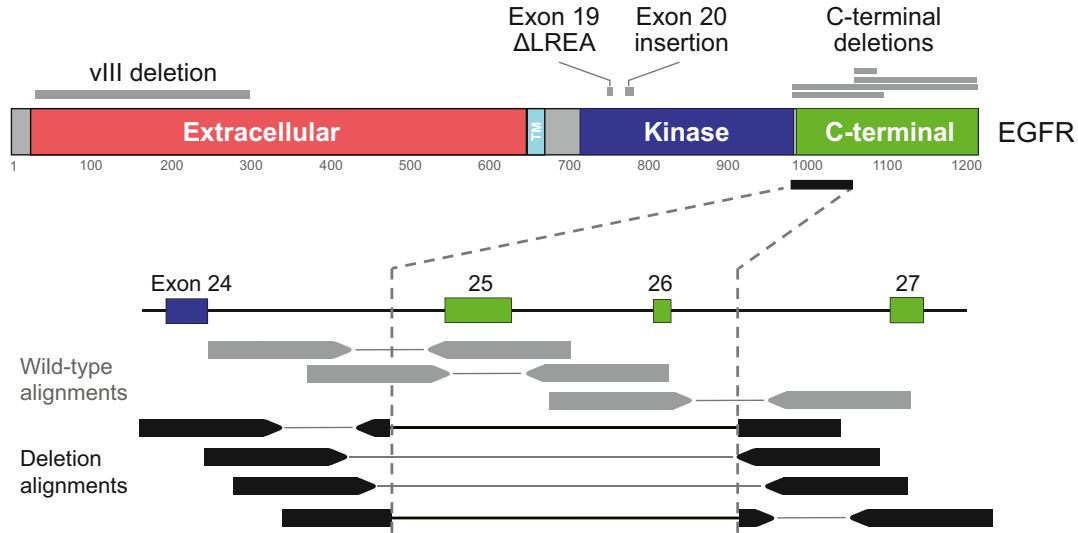
# Clinical consequences of rearrangement signatures: “BRCAness phenotype”



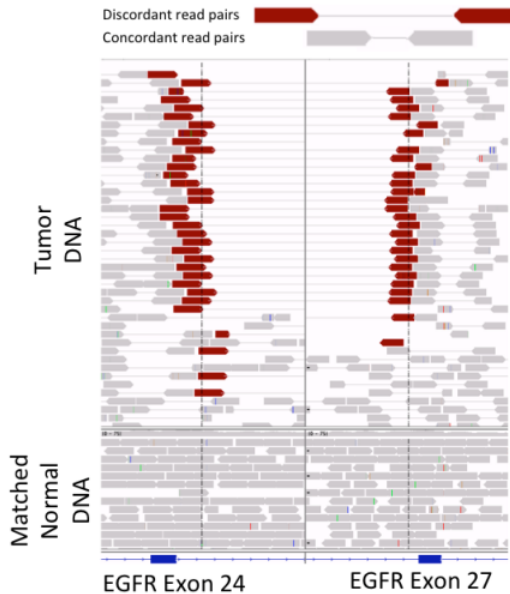


# Standard WGS

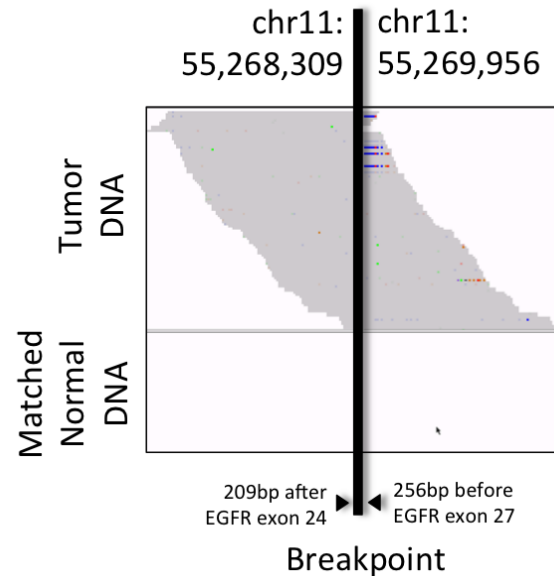
## Paired-end rearrangement mapping



EGFR deletion paired end mapping

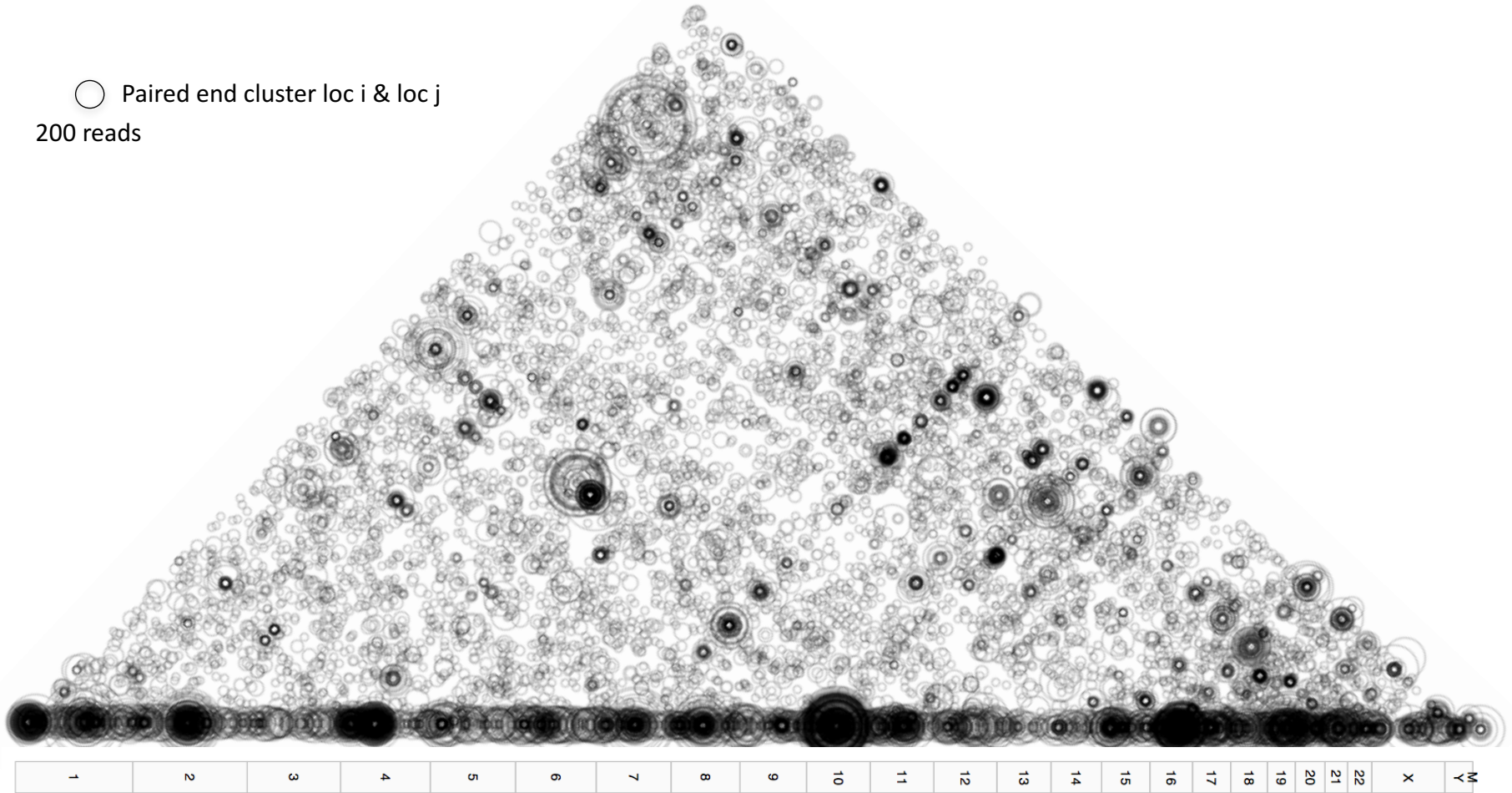


EGFR deletion Split read mapping



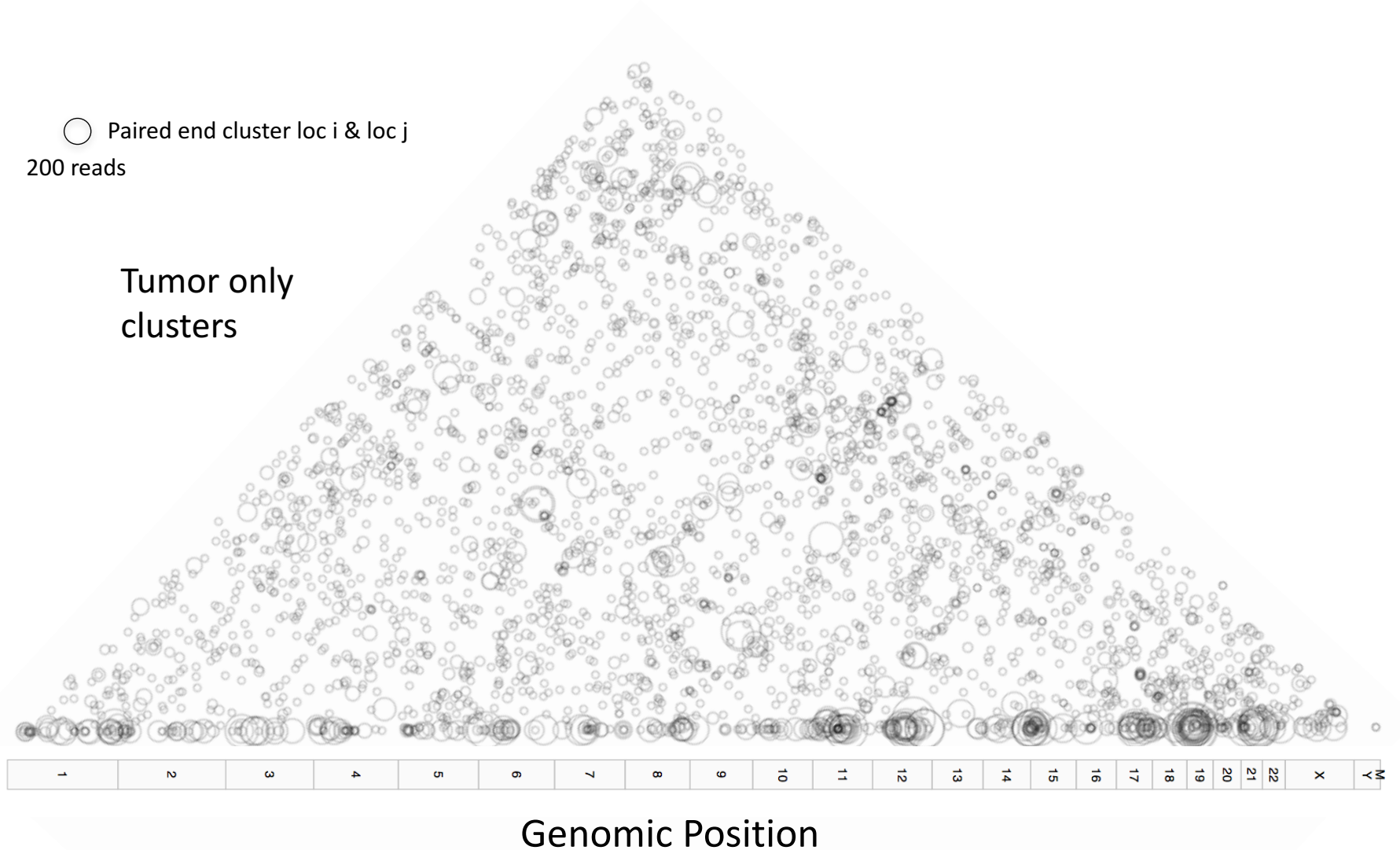
# Challenges: signal vs noise in rearrangement analysis

○ Paired end cluster loc i & loc j  
200 reads



Genomic Position

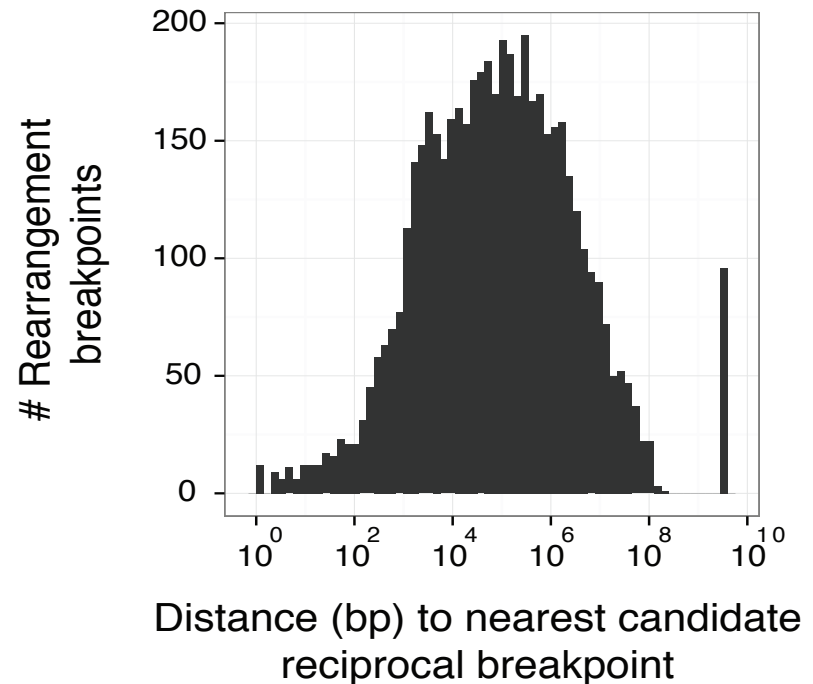
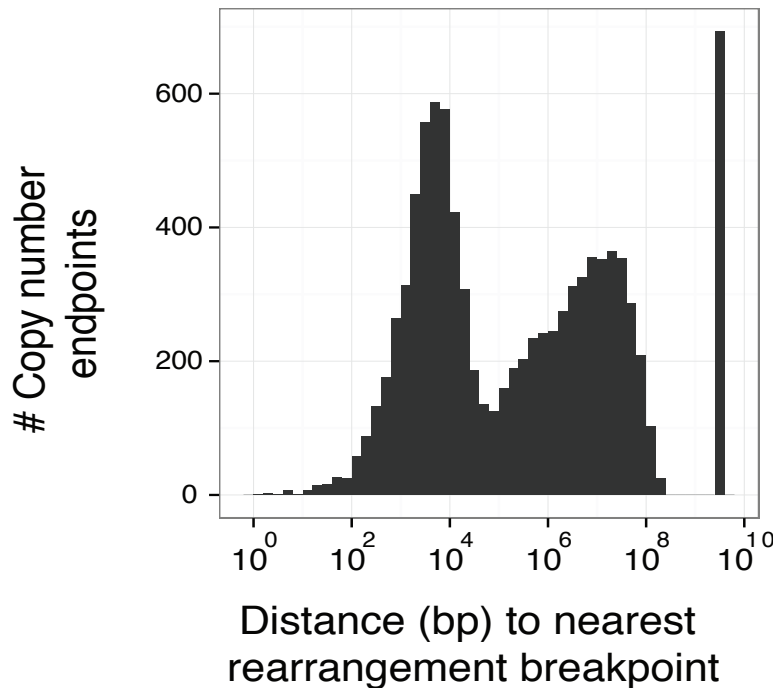
# Challenges: signal vs noise in rearrangement analysis



# ALERT:

Copy number and rearrangement data ... don't agree!

Data from 80 lung cancer whole genome T/N pairs

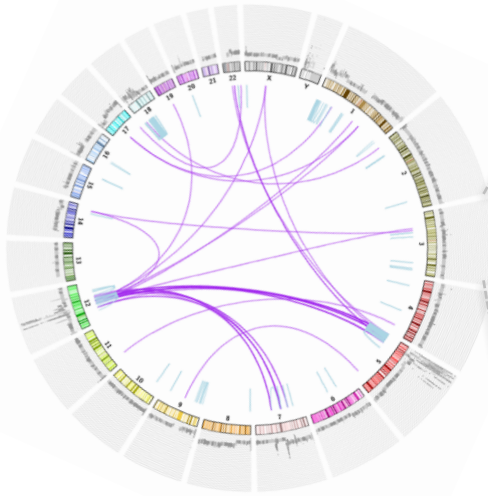


- Separately inferred → Inconsistent
- Over-segmentation of WGS data
- Unmapped and false positive junctions

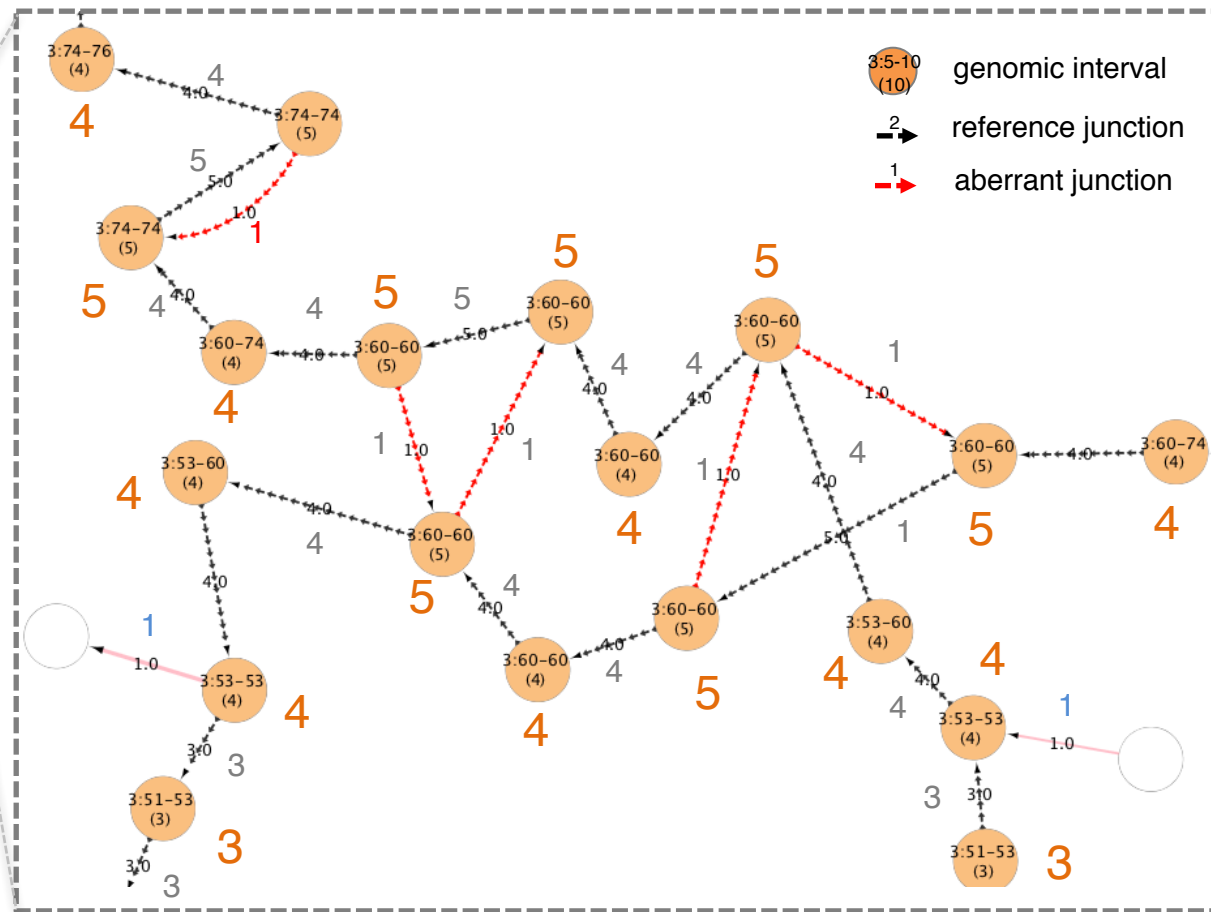
# JaBbA:

## From junctions to balanced assembly graphs

Whole genome sequence



*minimize*  
 $\epsilon^T \Sigma^{-1} \epsilon + \lambda_e (\|\delta\|_2 + \|\eta\|_2)$   
*subject to*  
 $v_i + \gamma = \mu_i \beta + \epsilon_i$   
 $\frac{v^T L}{\|L\|_1} + \gamma = \mu_0 \beta$   
 $(B_{ev})^T e + \delta - v = 0$   
 $B_{ve} e + \eta - v = 0$   
 $\tau_{min} \leq v^T L / \|L\|_1 \leq \tau_{max}$   
 $v, e, \gamma, \beta, \delta, \eta \geq 0$



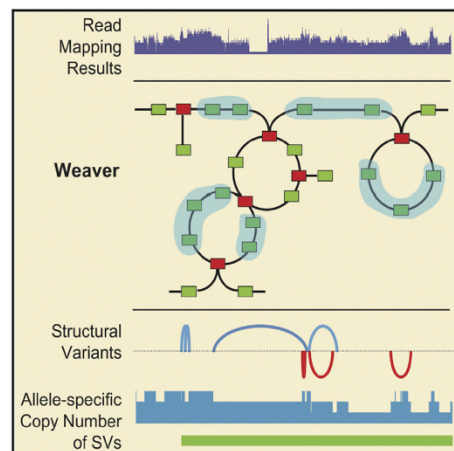
JaBbA graph

# Other “graph callers”

## Cell Systems

### Allele-Specific Quantification of Structural Variations in Cancer Genomes

#### Graphical Abstract



#### Authors

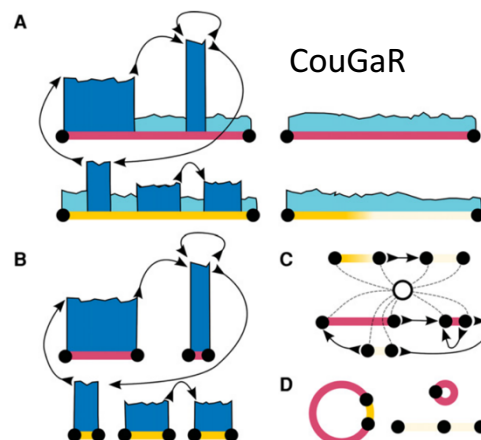
Yang Li, Shiguo Zhou,  
David C. Schwartz, Jian Ma

Correspondence  
jianma@cs.cmu.edu

#### In Brief

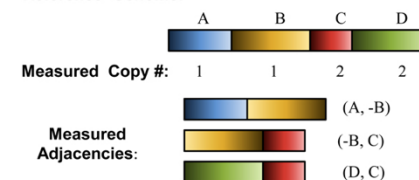
A new algorithm that quantifies allele-specific structural variations can greatly improve the analysis of complex genomic alterations in cancer.

#### Article



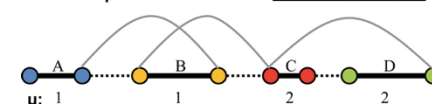
**Figure 1.** Overview of CouGaR algorithm. Tumor and normal samples are processed through a five-step algorithm. (A) We identify regions that are potentially amplified (dark blue) across two different chromosomes (red and yellow lines) in the tumor samples (left two contigs) compared to normal samples (right two contigs). We compute depth of coverage (DOC) information and cluster discordant read pairs to represent novel (with respect to hg19) adjacencies in the genome. (B) We identify continuous regions of amplification in the tumor genome using an HMM and DOC information from both tumor and normal samples. (C) We add a single super-source/-sink node, and using a min-cost circulation algorithm, we solve for the copy count of each region in the tumor genome. (D) Finally, a minimal set of circular and linear contigs that explain the coverage is found by formulating an integer programming problem that puts a penalty term on the number of unique contigs used.

#### Reference Genome:



#### PREGO

#### Interval-Adjacency Graph:



#### Reconstructed Cancer Genome:



Li .. Ma

Cell Systems 2016

Dzamba .. Brudno

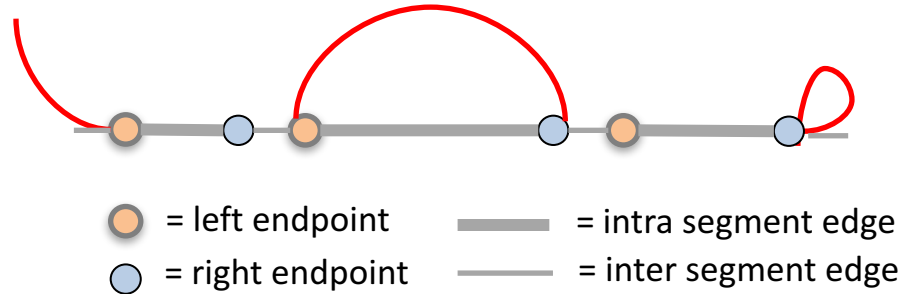
Genome Research 2016

Osper .. Raphael

BMC Bioinf 2012

# Graph representation of whole genomes

## Bidirected Graph



- Nodes represent left and right **sides** of intervals
- Undirected edges of two flavors (intra and inter segment)
- **Paths must be**

/-

# Stranded adjacency matrix A

$$A = \begin{array}{cc} & \begin{array}{c} + \\ - \end{array} \\ \begin{array}{c} + \\ - \end{array} & \begin{array}{|c|c|} \hline A_1 & A_2 \\ \hline A_3 & A_1^T \\ \hline \end{array} \end{array}$$

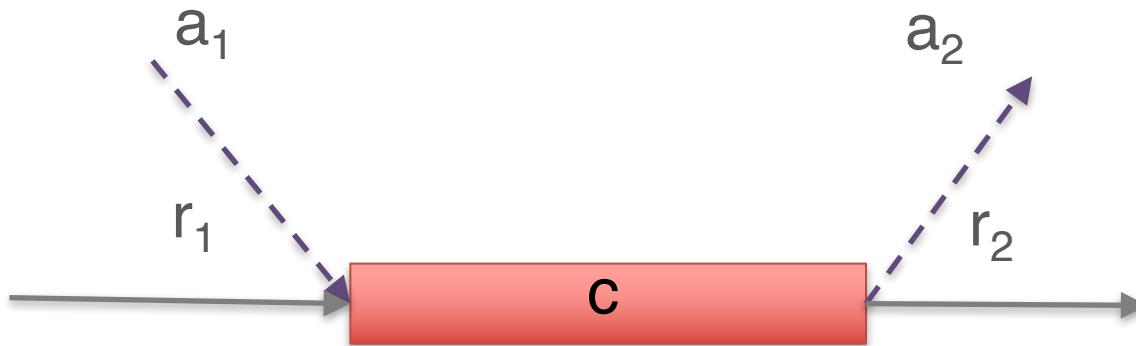
$a_{ij} = a_{\bar{j}\bar{i}}$   
= number of copies of junctions joining intervals  $i$  and  $j$

$$A_2 = A_2^T \quad A_3 = A_3^T$$



# JaBbA (Junction Balance Analysis):

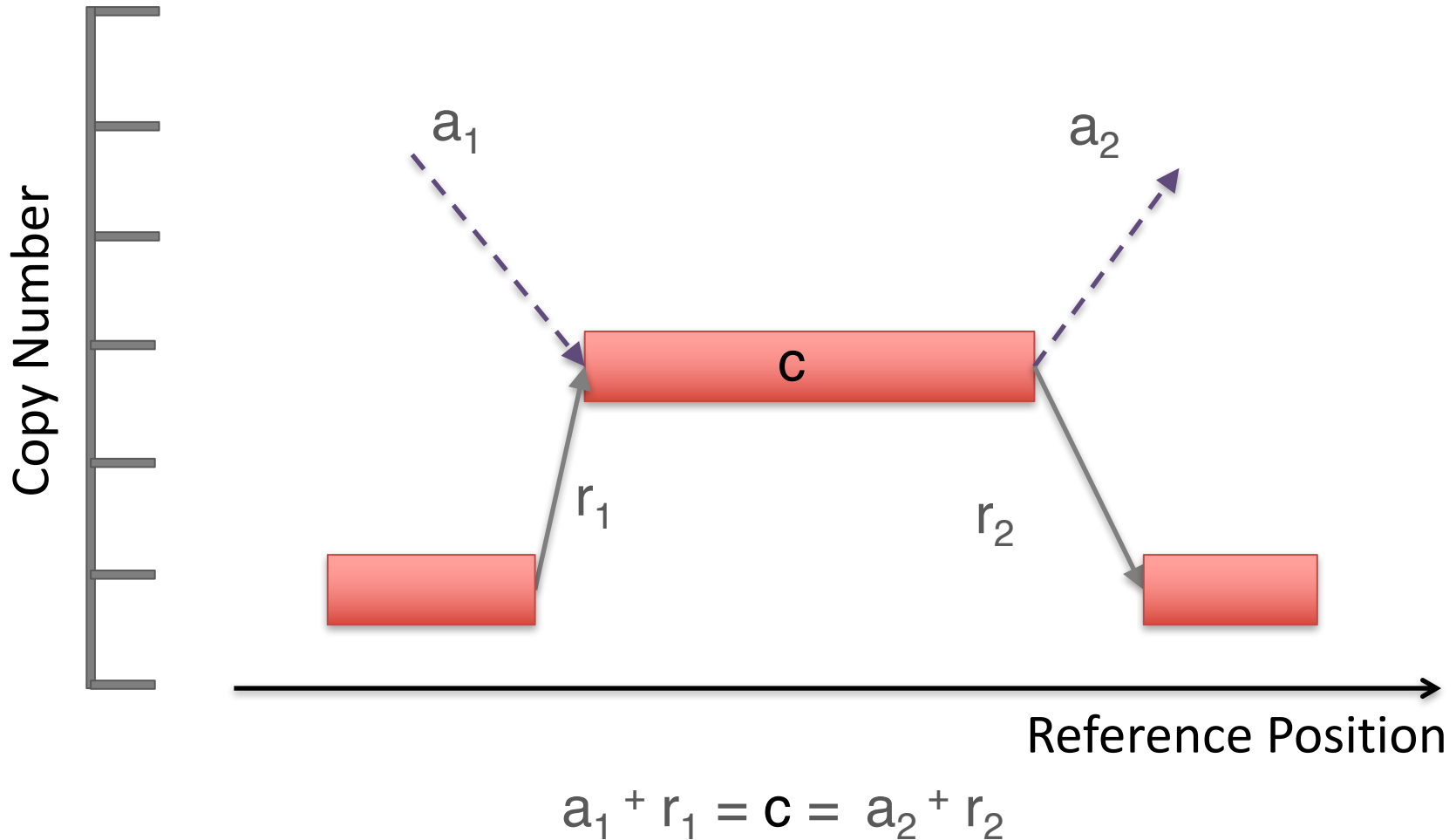
Integrating rearrangements and copy state



$$a_1 + r_1 = c = a_2 + r_2$$

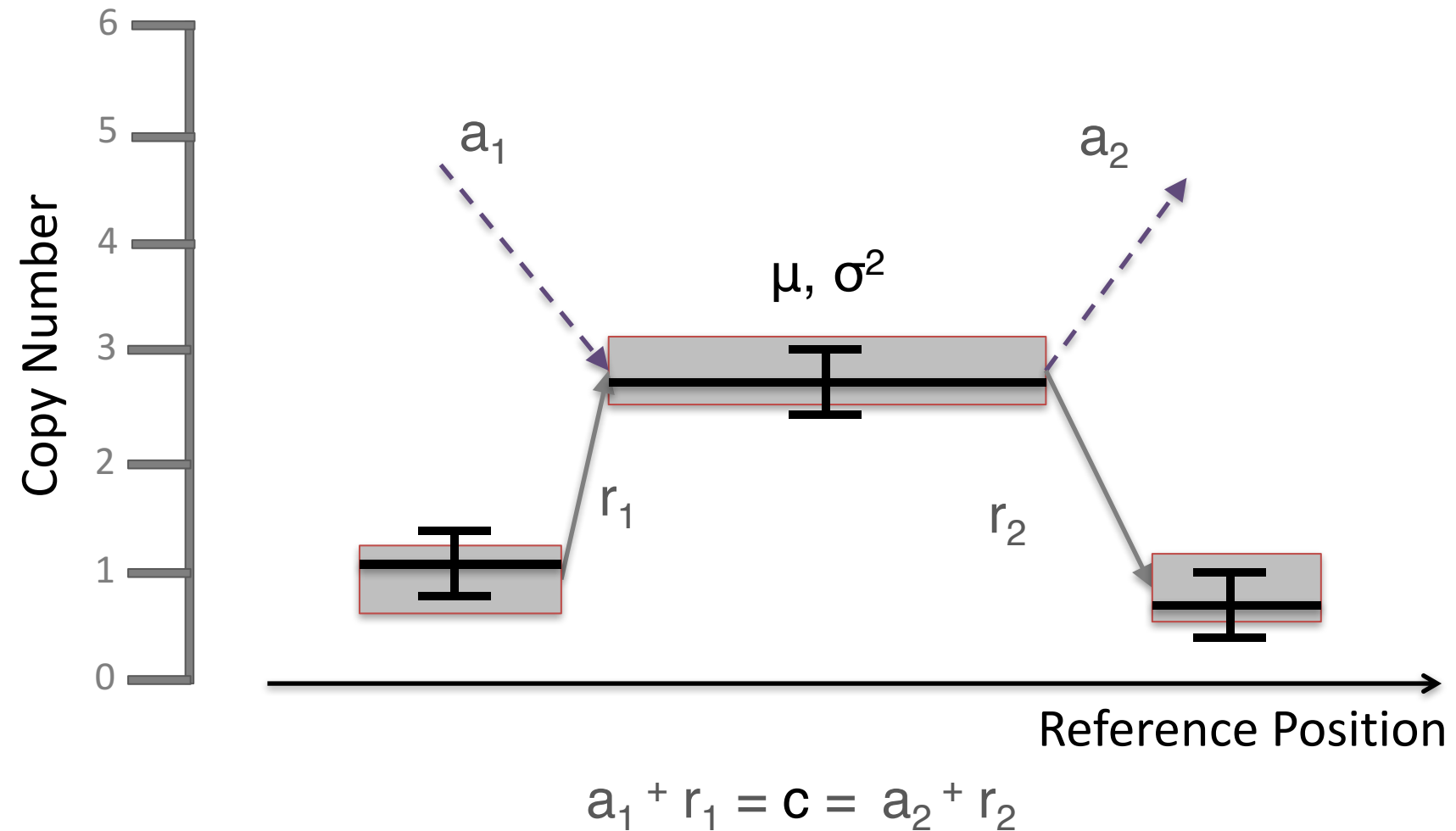
# JaBbA (Junction Balance Analysis):

Integrating rearrangements and copy state



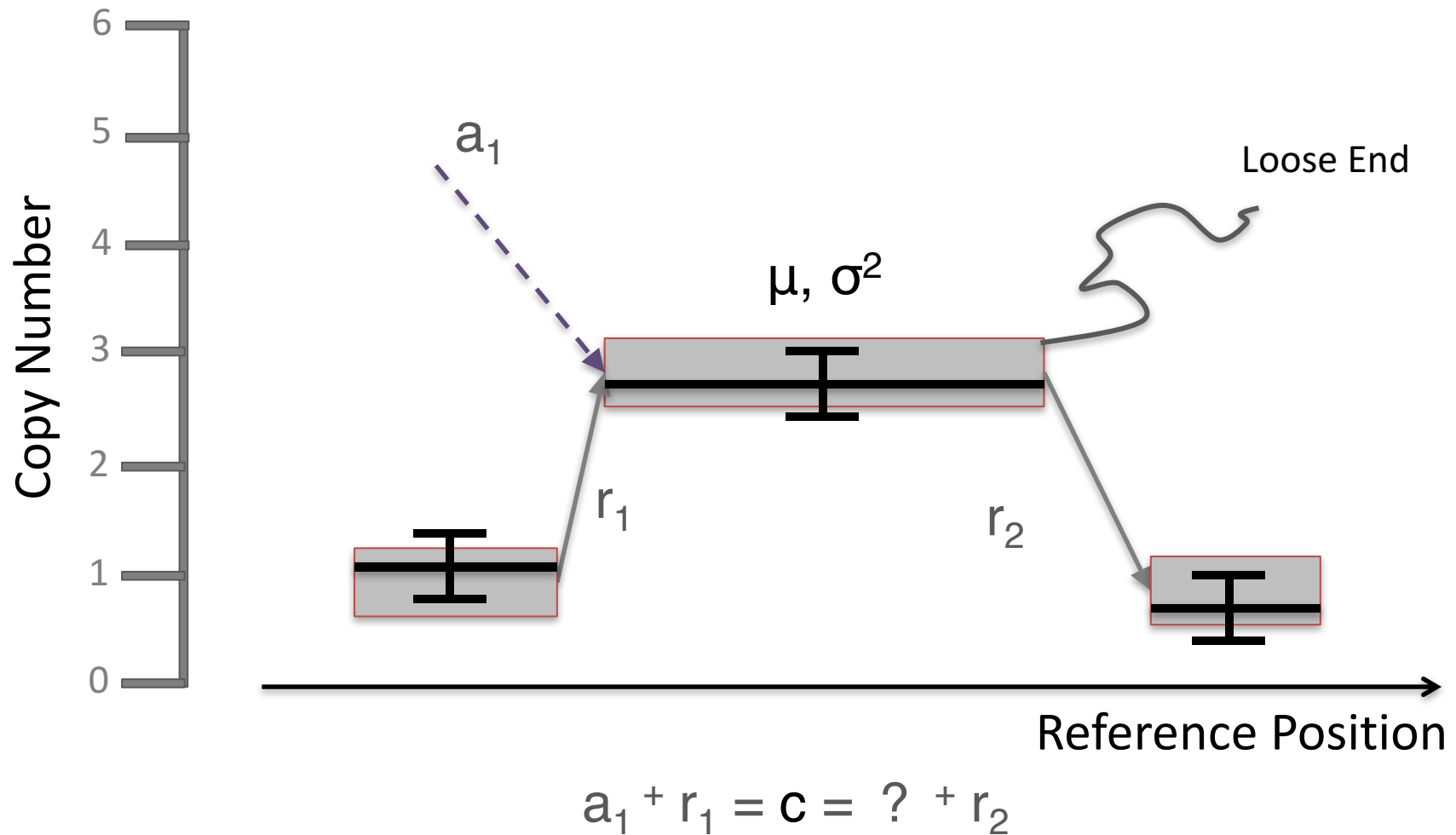
# JaBbA (Junction Balance Analysis):

Challenge: Noisy coverage data



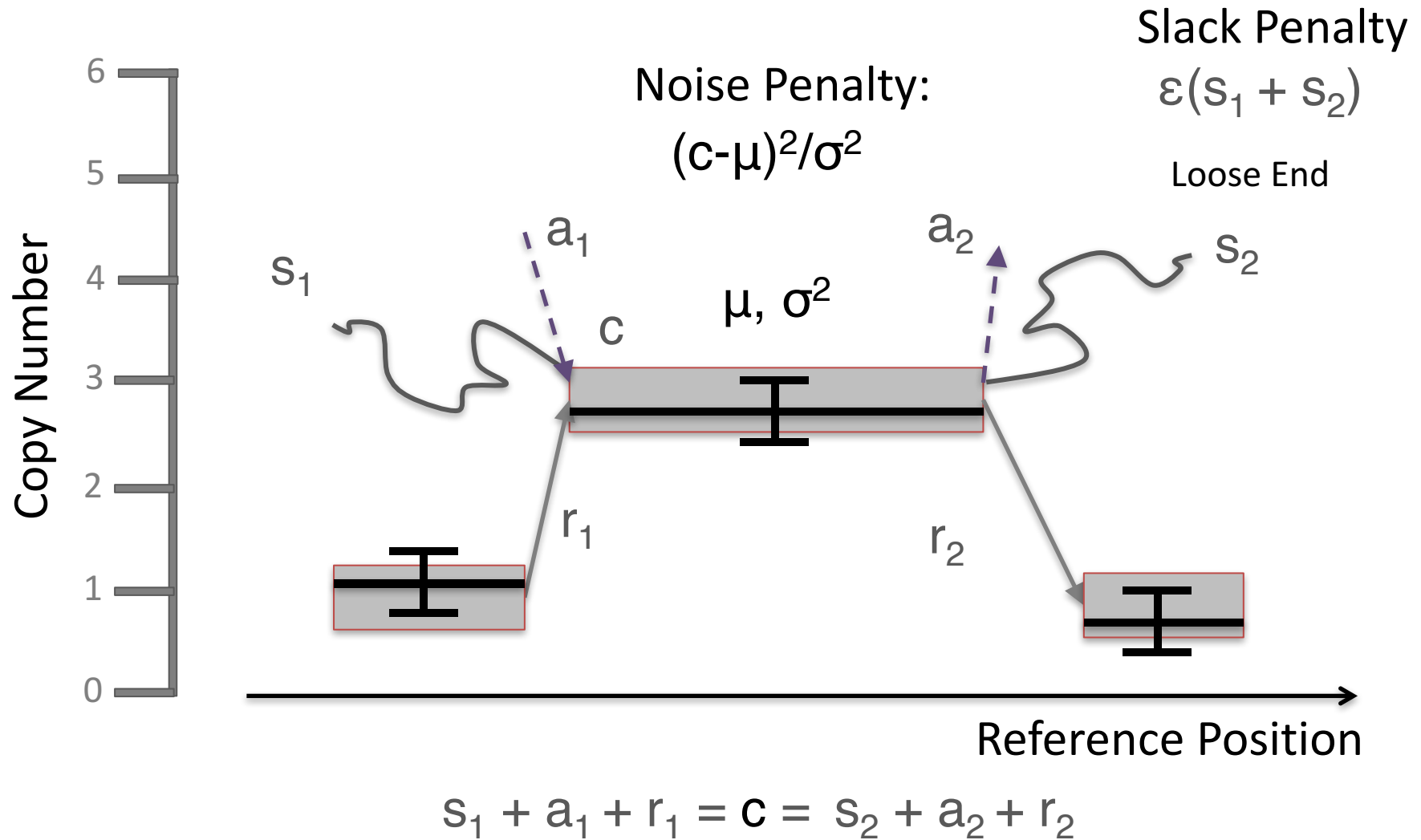
# JaBbA (Junction Balance Analysis):

Challenge: Missing rearrangements

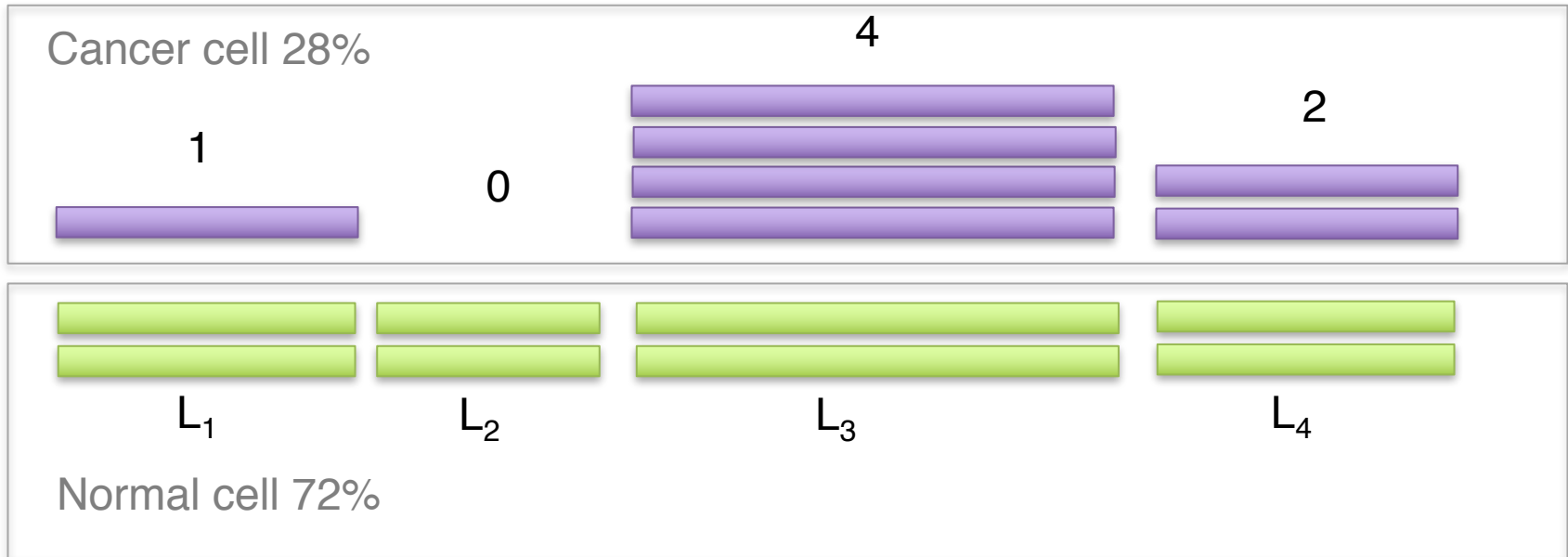


# JaBbA (Junction Balance Analysis):

Statistical model



# Transforming analog fragment density to digital copy number



Density of fragments aligning to interval 3 is

$$\frac{(4 \times 0.28 + 2 \times 0.72) \times L_3}{0.28 \times (L_1 + 4L_3 + L_4) + 0.72 \times 2 \times (L_1 + L_2 + L_3 + L_4)}$$

# Transforming analog fragment density to digital copy number

$$\frac{\mu_i}{\mu^T L} = \frac{\alpha v_i + 2(1 - \alpha)}{\alpha v^T L + 2(1 - \alpha) \|L\|_1}$$

$\mu \in \mathbb{R}^n$  Vector of fragment densities across n intervals (data)

$L \in \mathbb{Z}^n$  Vector of n interval widths (data)

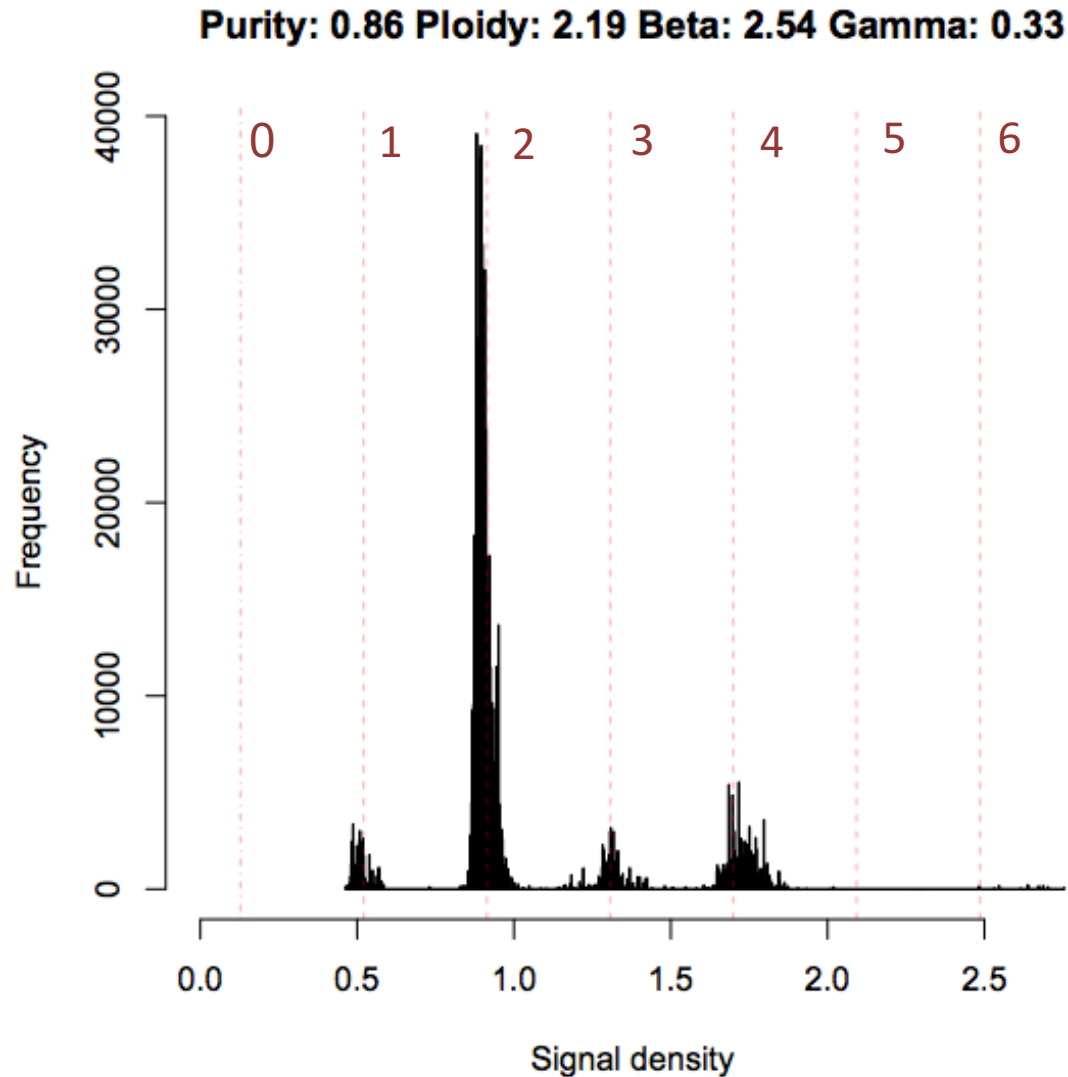
$v \in \mathbb{Z}^n$  Vector of n interval copy numbers (inferred)

$\alpha \in [0,1]$  Tumor cell fraction (purity, inferred)

Let  $\gamma = \frac{2(1 - \alpha)}{\alpha}$        $\beta = \frac{\alpha v^T L + 2(1 - \alpha) \|L\|_1}{\alpha \mu^T L}$

$\rightarrow v_i + \gamma = \beta \mu_i$

# Transforming analog fragment density to digital copy number



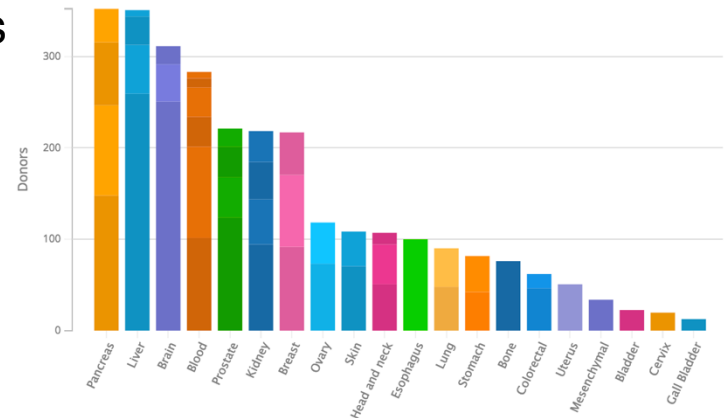
Example: lung  
adenocarcinoma





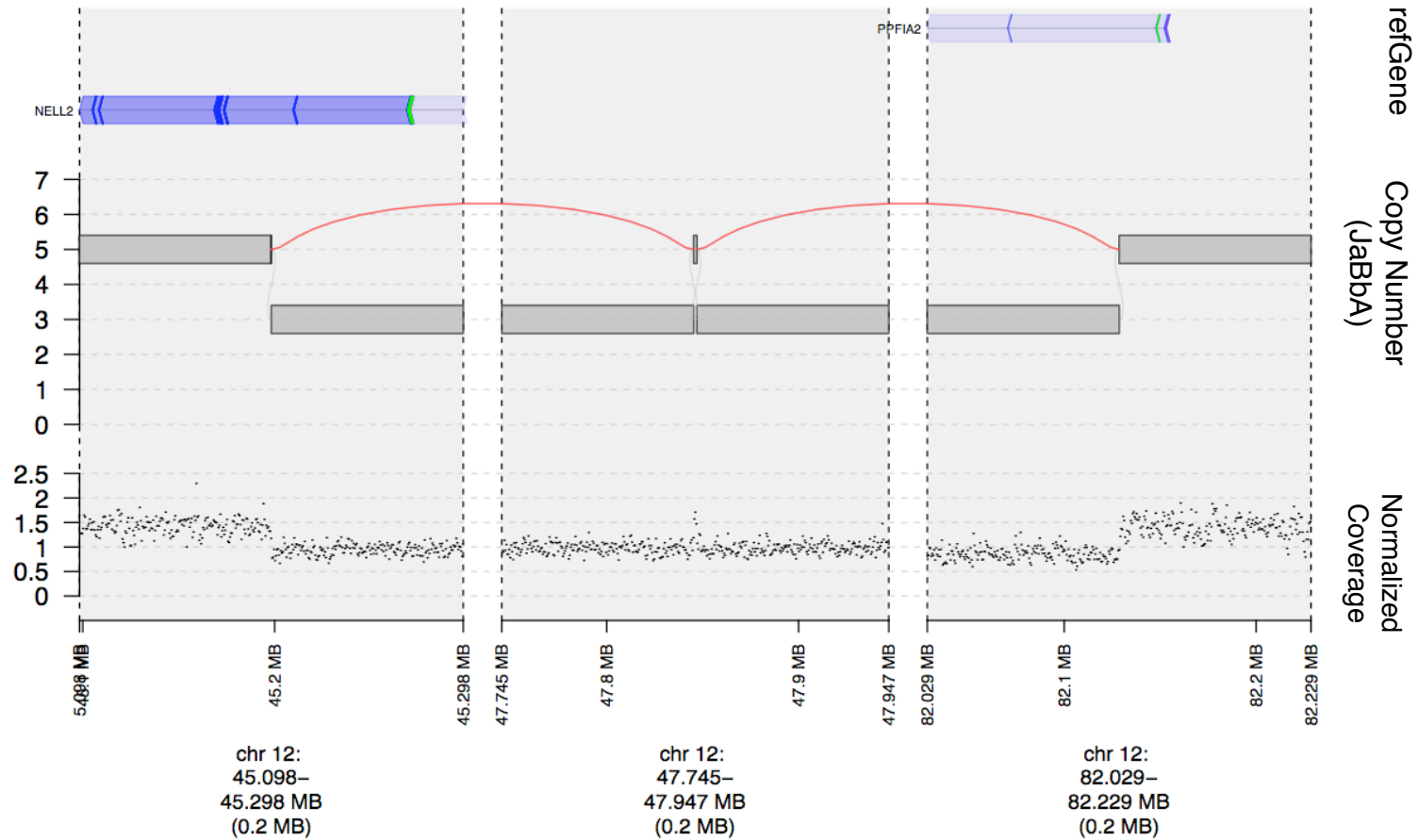
- 2834 tumor and matched normal whole genome sequences across 30 cancer types and 48 projects
- 1.5 Petabytes of raw data + downstream analytic pipelines
- 13 analysis working groups, including PCAWG-6 (structural variation dataset)

Donor Distribution by Primary Site  
48 projects and 20 primary sites



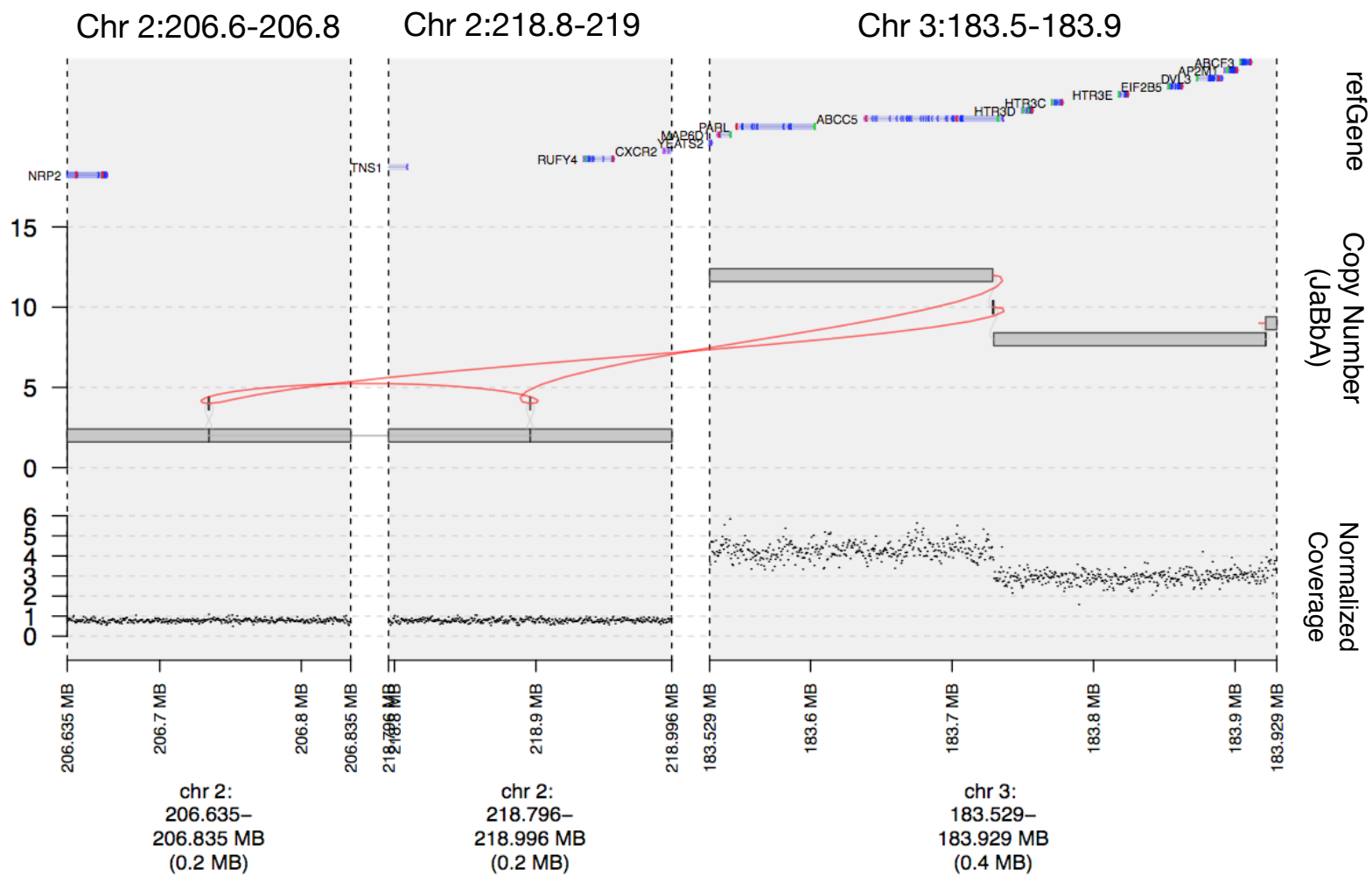
- 330 WGS cell lines across 16 cancer types
- collaboration with Mahmoud Ghandi and Jesse Boehm at Broad Institute

# Long range coupling of copy changes



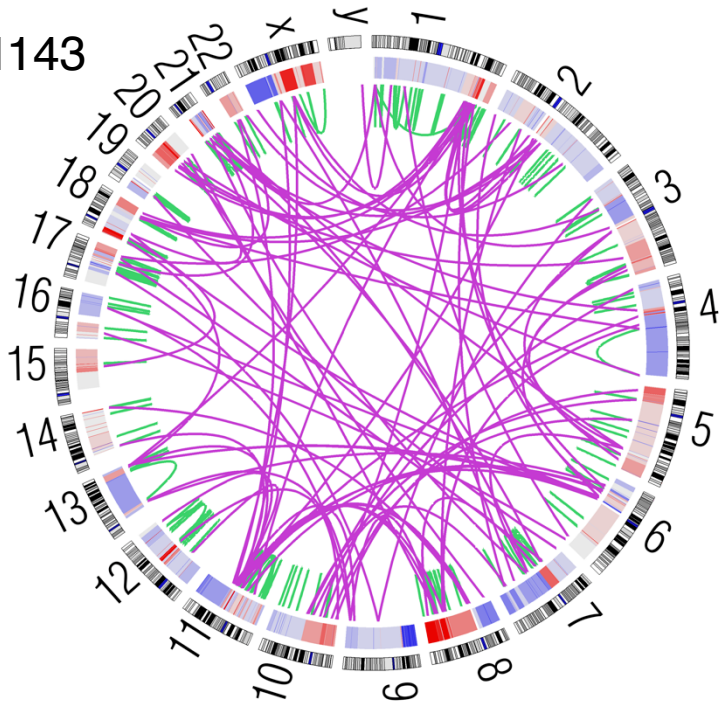
Long distance coupling of copy changes through rearrangement junctions

# “Kidnapped” loci (Lung adenocarcinoma)

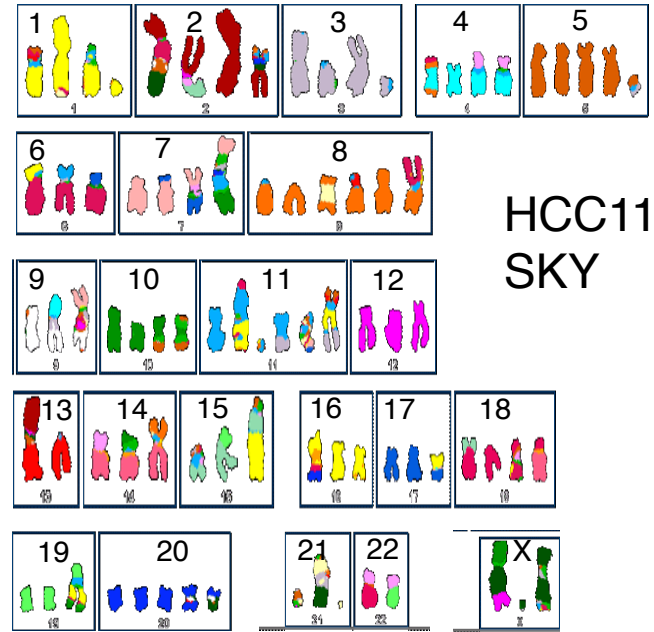


# WGS vs. cytogenetics

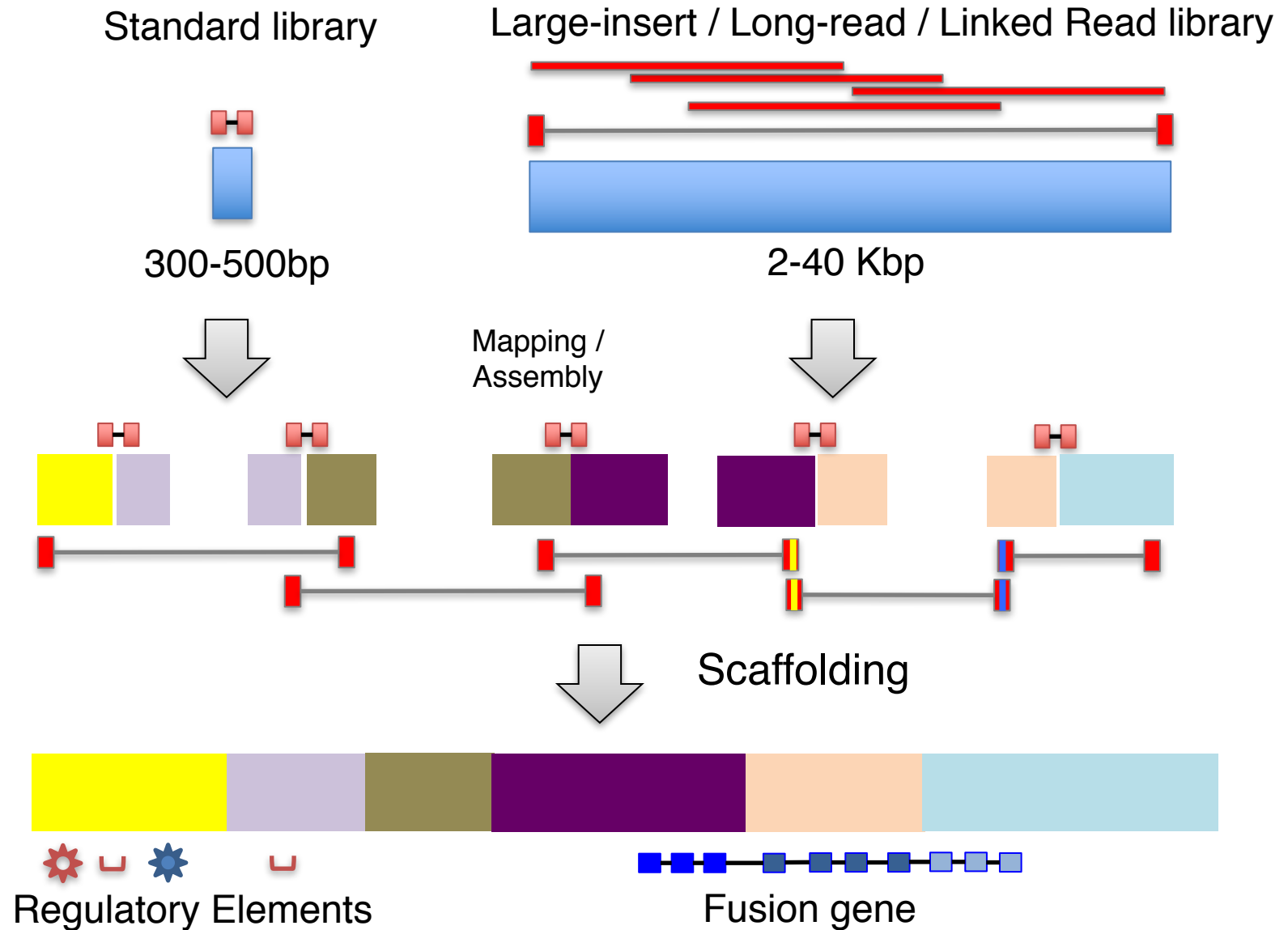
HCC1143  
WGS



!=



# “Shattered” cancer data

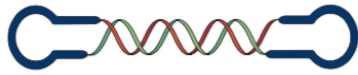


# Long-read sequencing

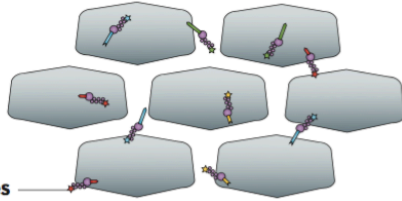
12-20 Kbp reads, \$85-400/Gb\*\*

## Aa Pacific Biosciences

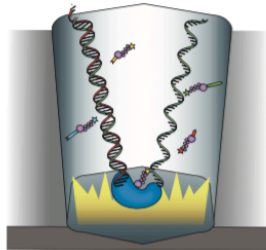
**SMRTbell template**  
Two hairpin adapters allow continuous circular sequencing



**ZMW wells**  
Sites where sequencing takes place



**Labelled nucleotides**  
All four dNTPs are labelled and available for incorporation



**Modified polymerase**  
As a nucleotide is incorporated by the polymerase, a camera records the emitted light

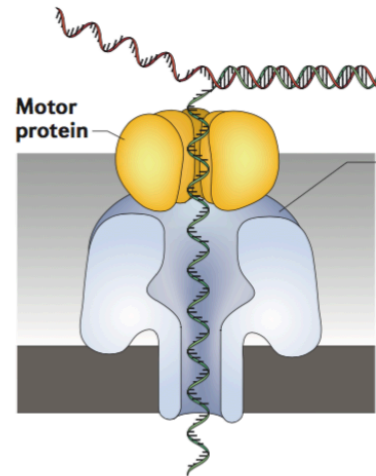
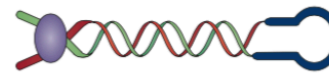
**PacBio output**  
A camera records the changing colours from all ZMWs; each colour change corresponds to one base



200-900 Kbp, \$100-180/Gb\*\*

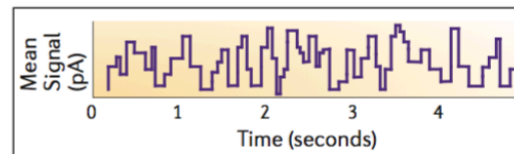
## Ab Oxford Nanopore Technologies

**Leader-Hairpin template**  
The leader sequence interacts with the pore and a motor protein to direct DNA, a hairpin allows for bidirectional sequencing



**Alpha-hemolysin**  
A large biological pore capable of sensing DNA

**Current**  
Passes through the pore and is modulated as DNA passes through



**ONT output (squiggles)**  
Each current shift as DNA translocates through the pore corresponds to a particular k-mer

\*\*<https://blog.genohub.com/2017/06/16/pacbio-vs-oxford-nanopore-sequencing/>

# Linked-read whole genome sequencing (10x genomics)

100 Kbp “synthetic long reads” \$7/Gb

## Emulsion PCR

Arbitrarily long DNA is mixed with beads loaded with barcoded primers, enzyme and dNTPs



## GEMs

Each micelle has 1 barcode out of 750,000



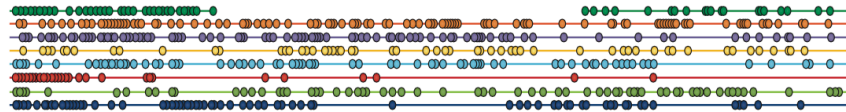
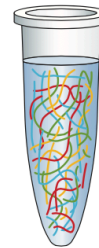
## Amplification

Long fragments are amplified such that the product is a barcoded fragment ~350 bp



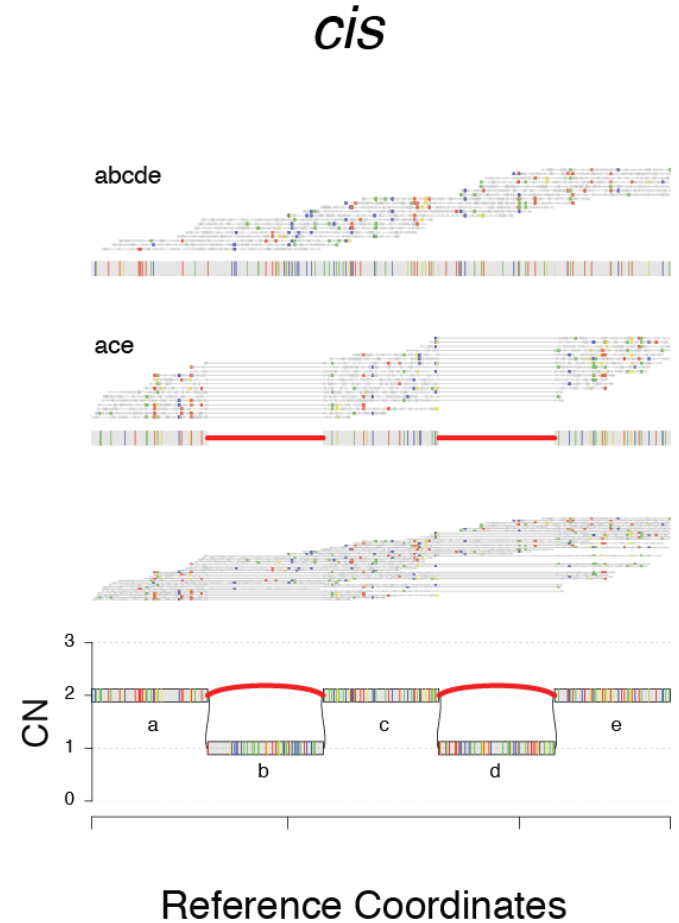
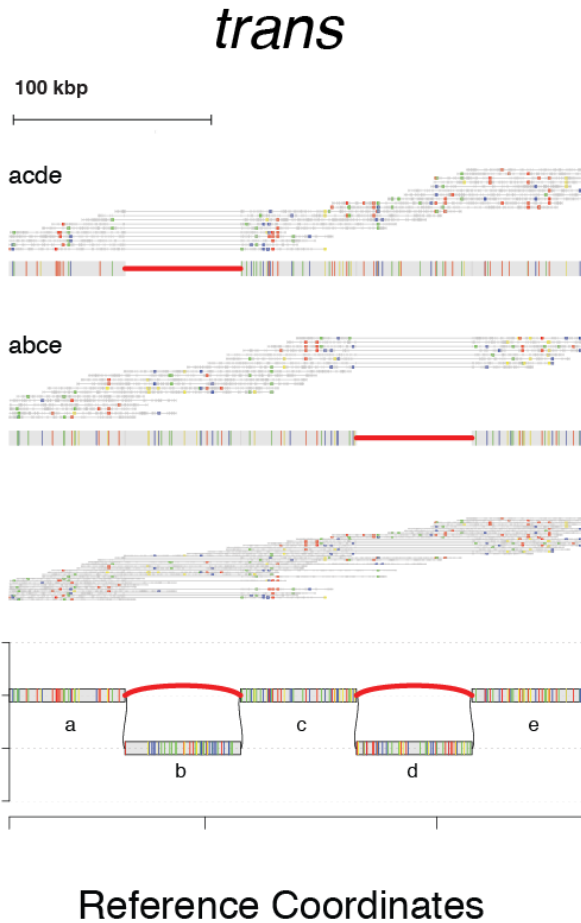
## Pooling

The emulsion is broken and DNA is pooled, then it undergoes a standard library preparation



# Phasing rearrangements with 10X

SNV      Junction  
A T C G      —

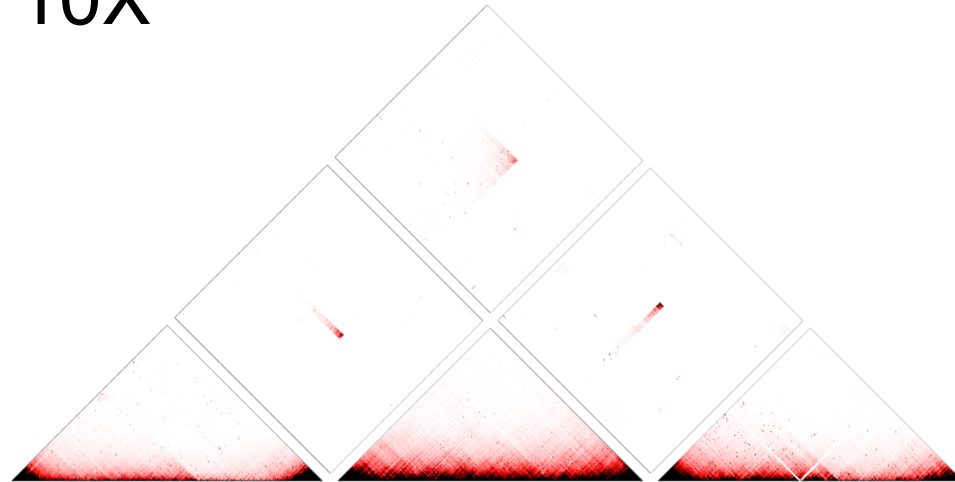
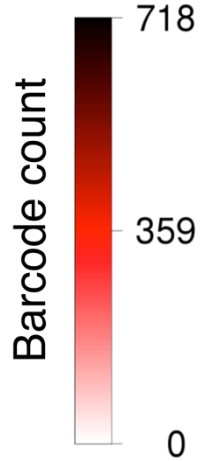






# Phasing genomic kidnappings with 10X

10X Chromium library barcode overlap

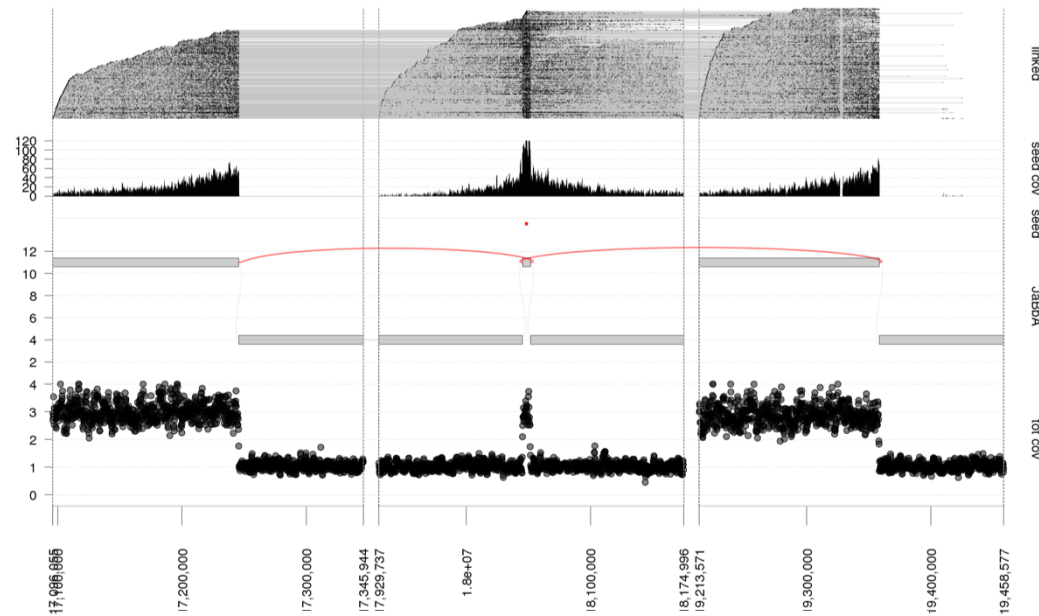


10X Linked reads

Linked read coverage

Assembly graph

Coverage



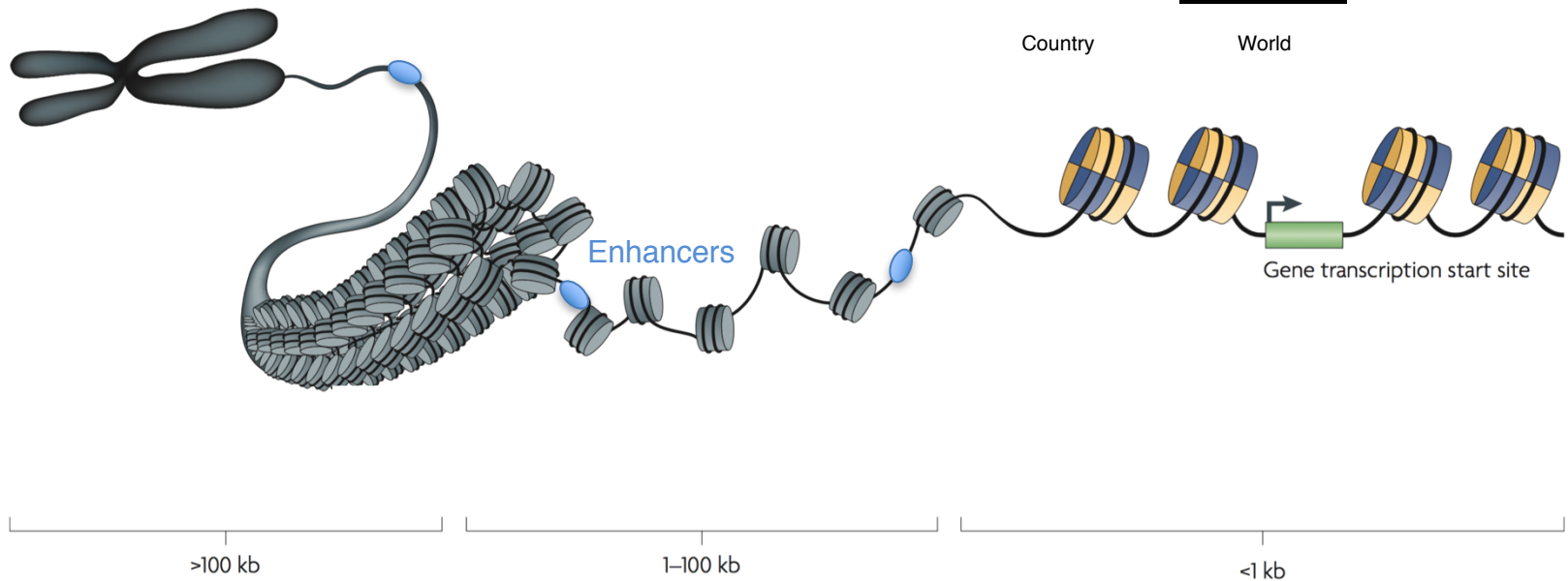
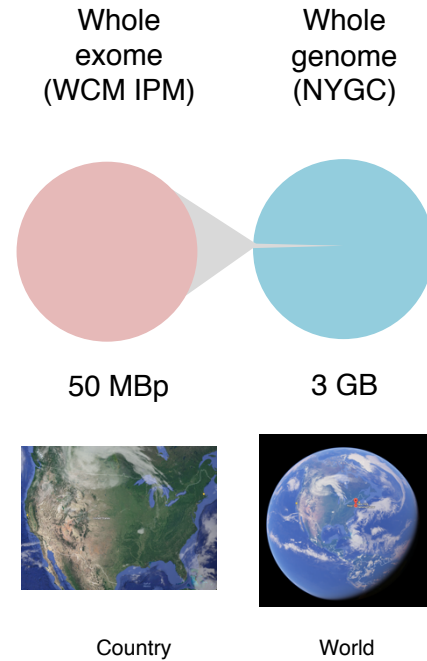
Dis-contiguous  
Reference loci

chr 21:  
17.03-  
17.45 MB  
(0.42 MB)

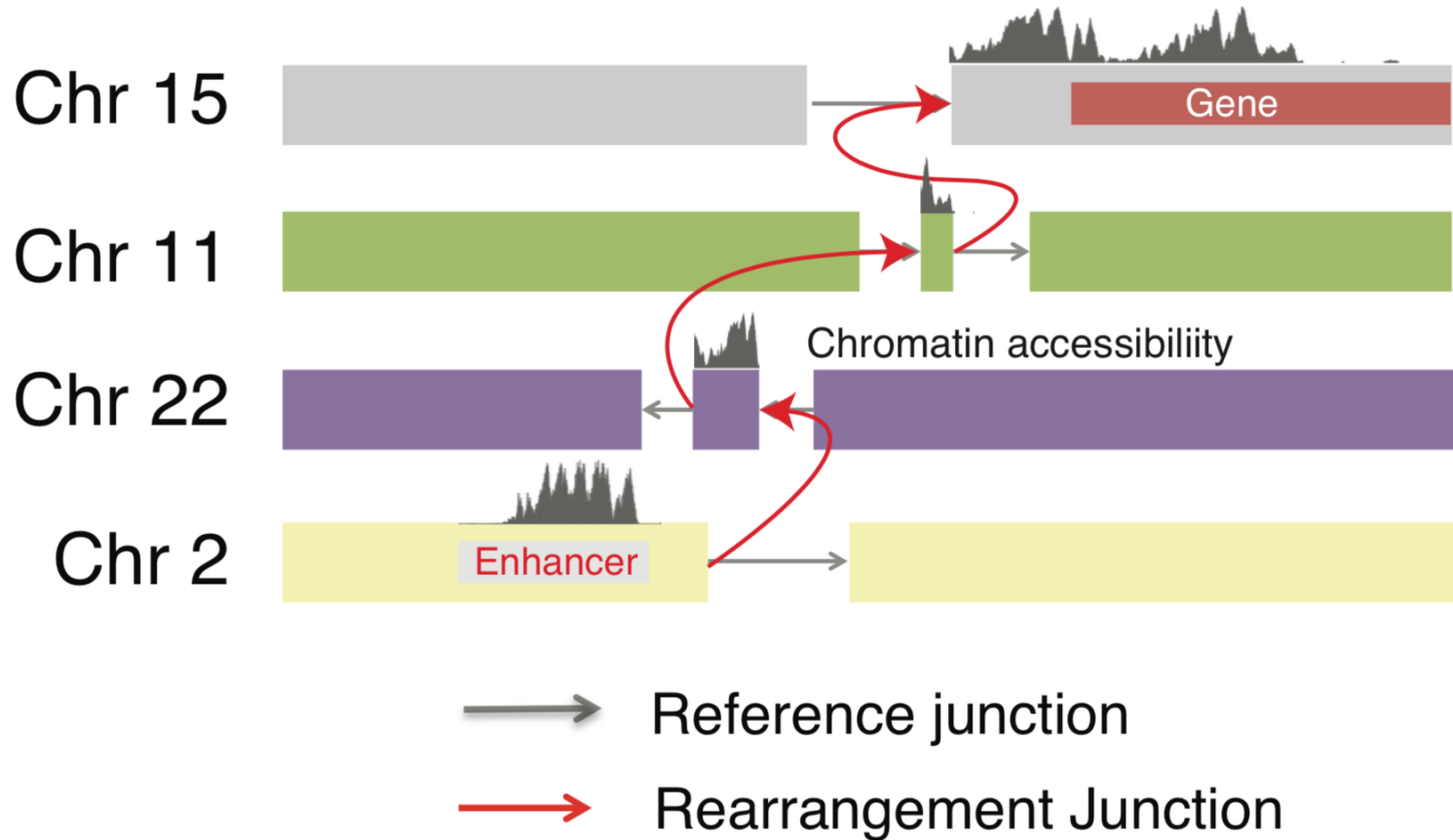
chr 21:  
17.845-  
18.252 MB  
(0.41 MB)

chr 21:  
19.14-  
19.56 MB  
(0.42 MB)

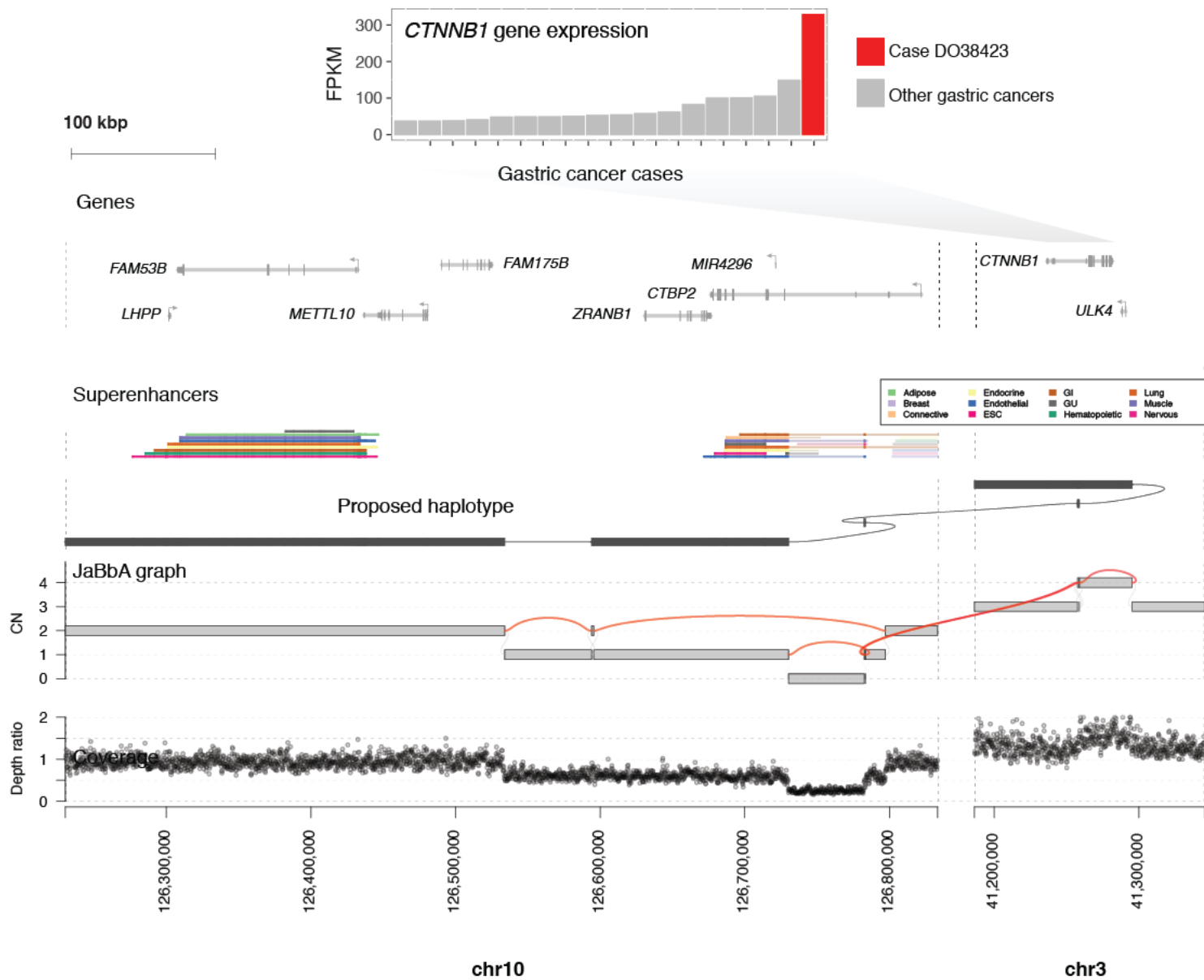
# Diving into the ocean of regulatory DNA



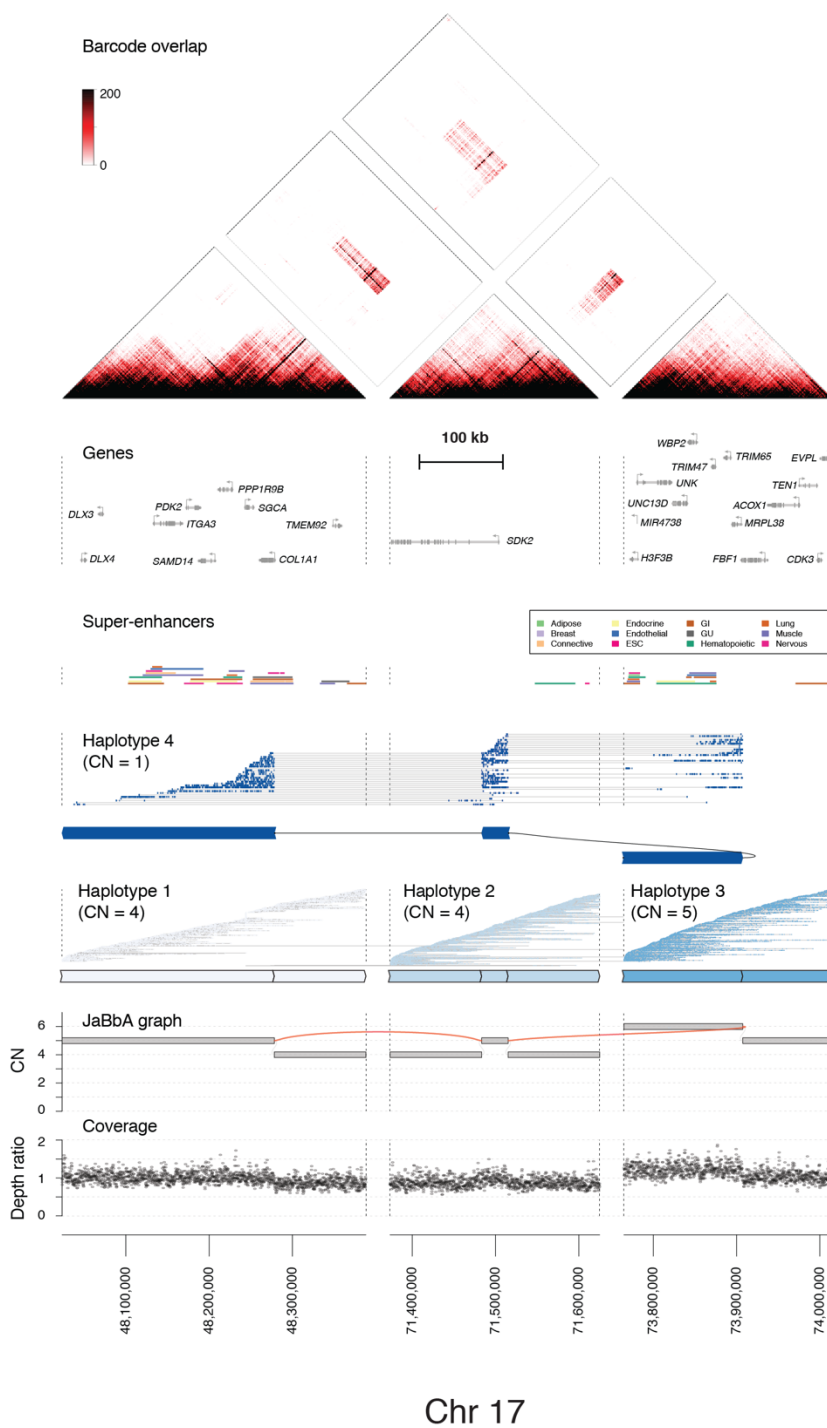
# Enhancer perturbations



# Enhancer hijacking in PCAWG



# 10X phasing of enhancer hijacking in IPM bladder cancer case



[m ski@m skilab.org](mailto:m ski@m skilab.org)  
<http://github.com/m skilab>



Evan Biederstedt (research specialist)

Julie Behr (Tri-I CBM PhD student)

Aditya Deshpande (Tri-I CBM PhD student)

Zoran Gajic (undergrad)

Kofi Gyan (Tri-I CBM PhD student)

Kevin Hadi (WCM PBSB PhD student)

Khagay Nagdimov (intern)

Joel Rosiene (medical student)

Huasong Tian (research scientist)

Netha Ulahannan (postdoctoral associate)

Trent Walradt (medical student)

Charalampos Xanthopoulos (soft eng)

Xiaotong Yao (Tri-I CBM PhD student)