

Single cell & single molecule analysis of cancer

Michael Schatz

October 22, 2015

JHU Genomics Symposium



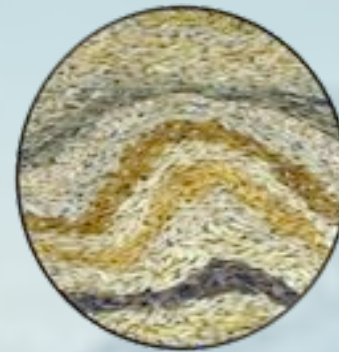
Schatzlab Overview



Human Genetics

Role of mutations in disease

Narzisi *et al.* (2015)
Iossifov *et al.* (2014)



Plant Biology

Genomes & Transcriptomes

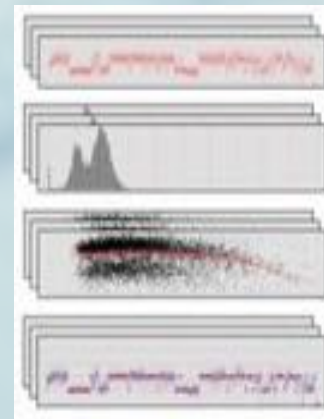
Ming *et al.* (2015)
Schatz *et al.* (2014)



Algorithmics & Systems Research

Ultra-large scale biocomputing

Stevens *et al.* (2015)
Marcus *et al.* (2014)



Single Cell & Single Molecule

CNVs, SVs, & Cell Phylogenetics

Garvin *et al.* (2015)
Goodwin *et al.* (2015)

Outline

1. Single Molecule Sequencing

Long read sequencing of a breast cancer cell line

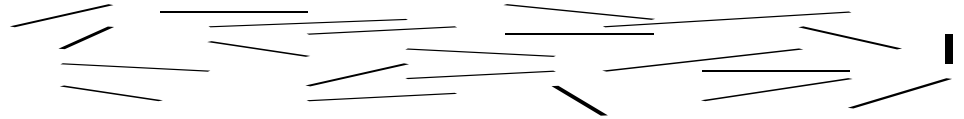
2. Single Cell Copy Number Analysis

Intra-tumor heterogeneity and metastatic progression



Sequence Assembly Problem

1. Shear & Sequence DNA



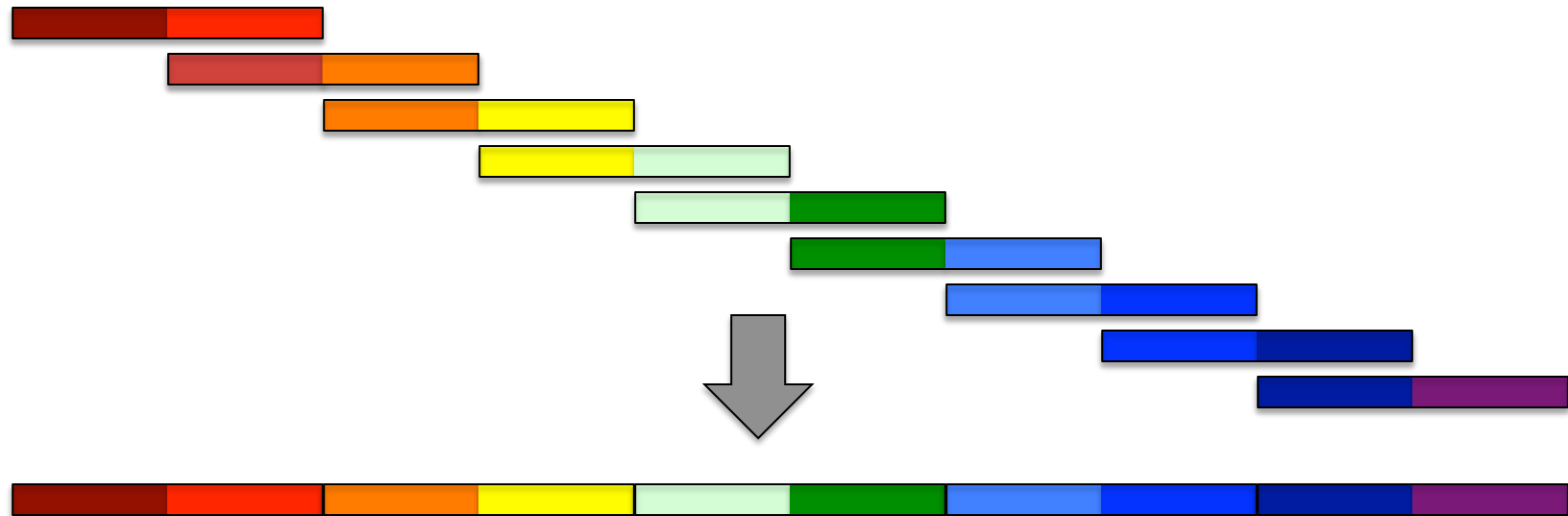
2. Construct assembly graph from overlapping reads

...AGCCTAGGGATGCGCGACACGT

GGATGCGCGACACGT CGCATATCCGGTTTGGTCAACCTCGGACGGAC

CAACCTCGGACGGACCTCAGCGAA...

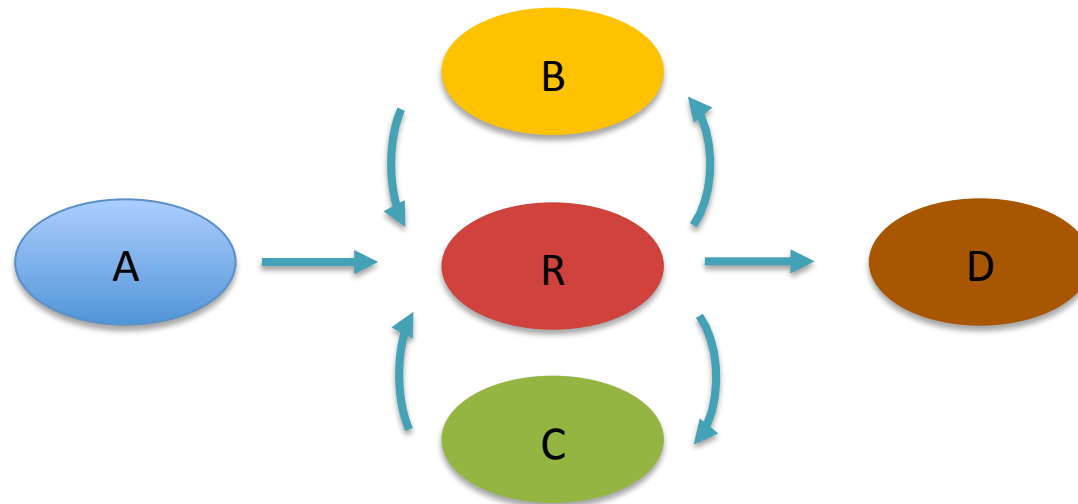
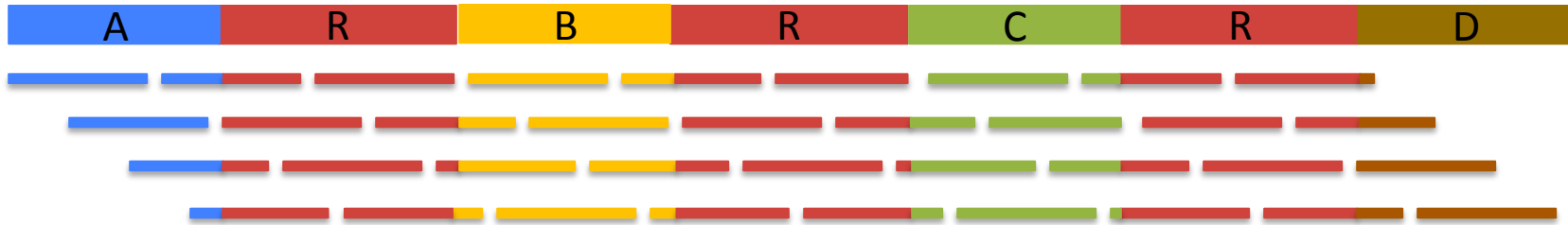
3. Simplify assembly graph



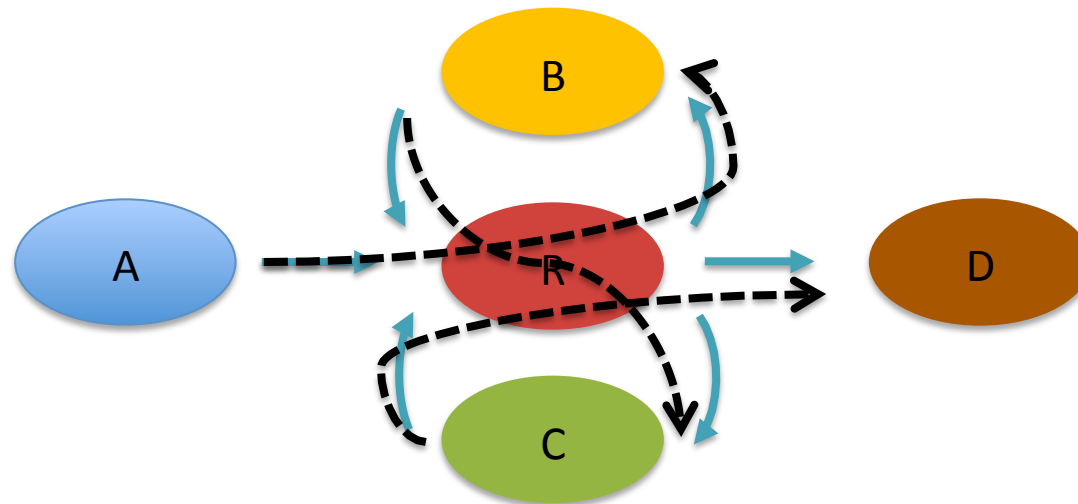
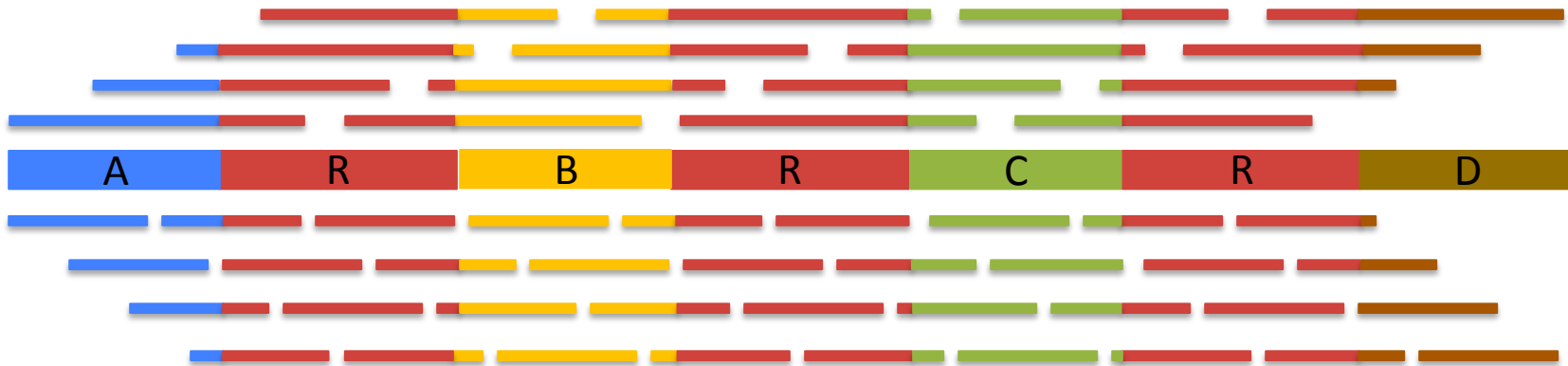
On Algorithmic Complexity of Biomolecular Sequence Assembly Problem

Narzisi, G, Mishra, B, Schatz, MC (2014) *Algorithms for Computational Biology*. Lecture Notes in Computer Science. Vol. 8542

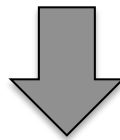
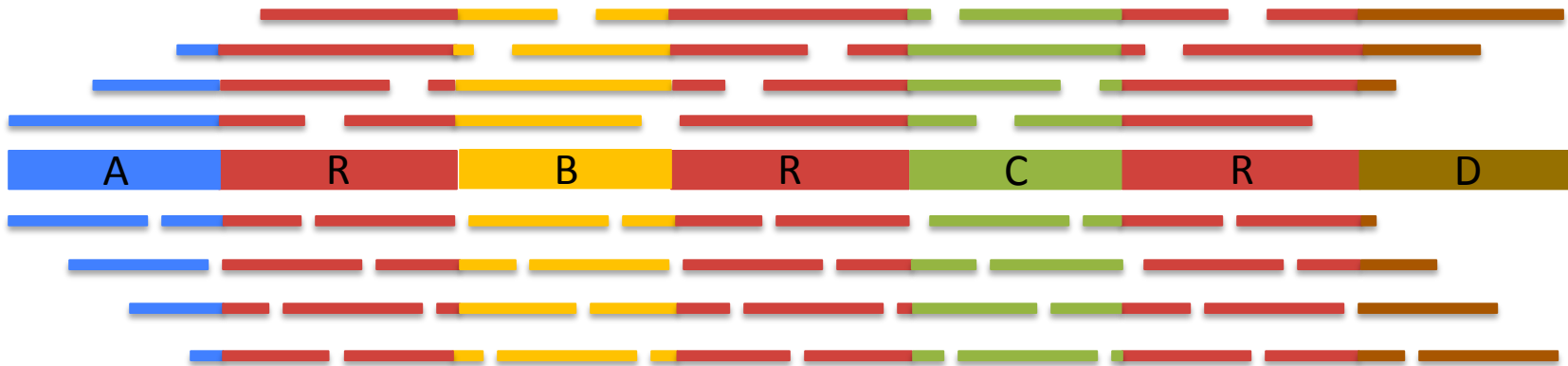
Assembly Complexity



Assembly Complexity



Assembly Complexity



The advantages of SMRT sequencing

Roberts, RJ, Carneiro, MO, Schatz, MC (2013) *Genome Biology*. 14:405

Genomics Arsenal in the Year 2015

Long Read Sequencing: De novo assembly, SV analysis, phasing

Illumina/Moleculo



(Kuleshov et al. 2014)

Pacific Biosciences



(Berlin et al, 2014)

Oxford Nanopore



(Quick et al, 2014)

Long Span Sequencing: Chromosome Scaffolding, SV analysis, phasing

Molecular Barcoding



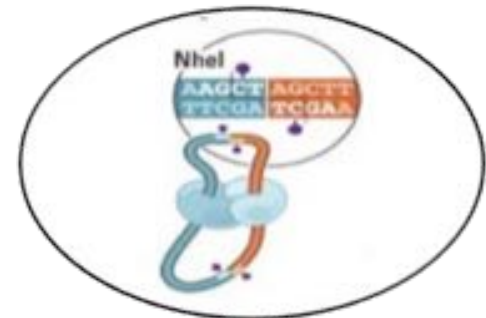
(10Xgenomics.com)

Optical Mapping



(Cao et al, 2014)

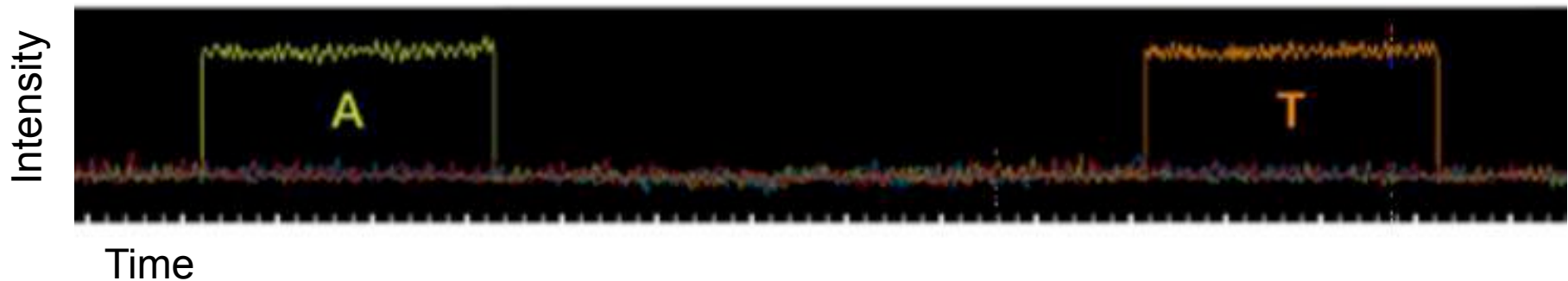
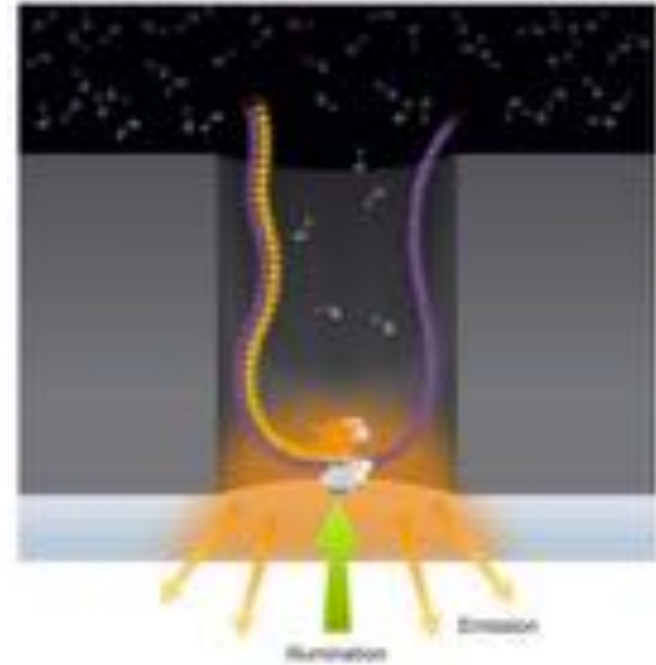
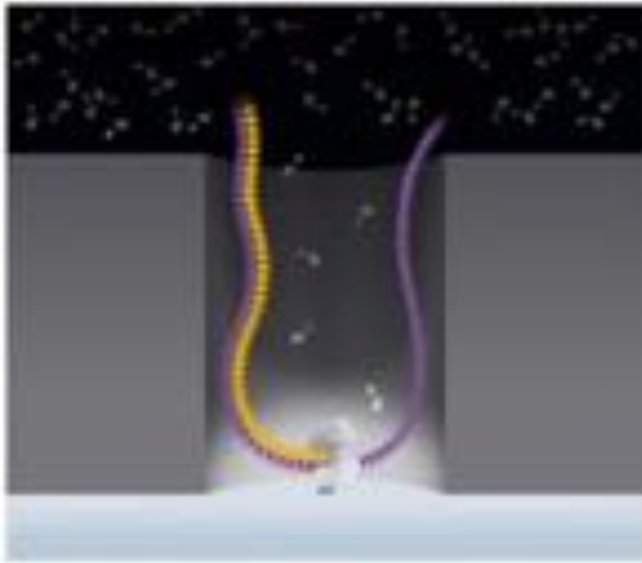
Chromatin Assays



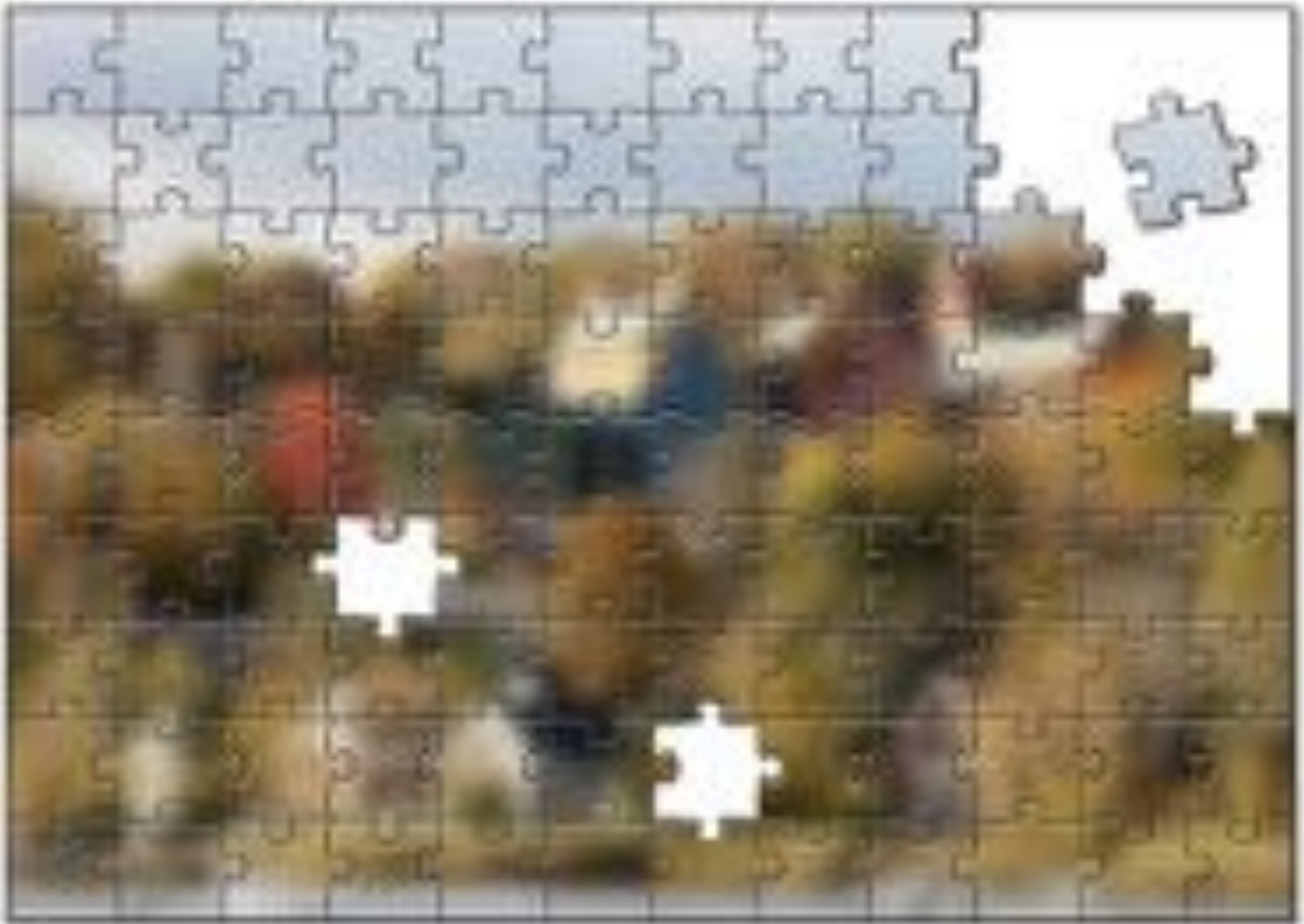
(Putnam et al, 2015)

PacBio SMRT Sequencing

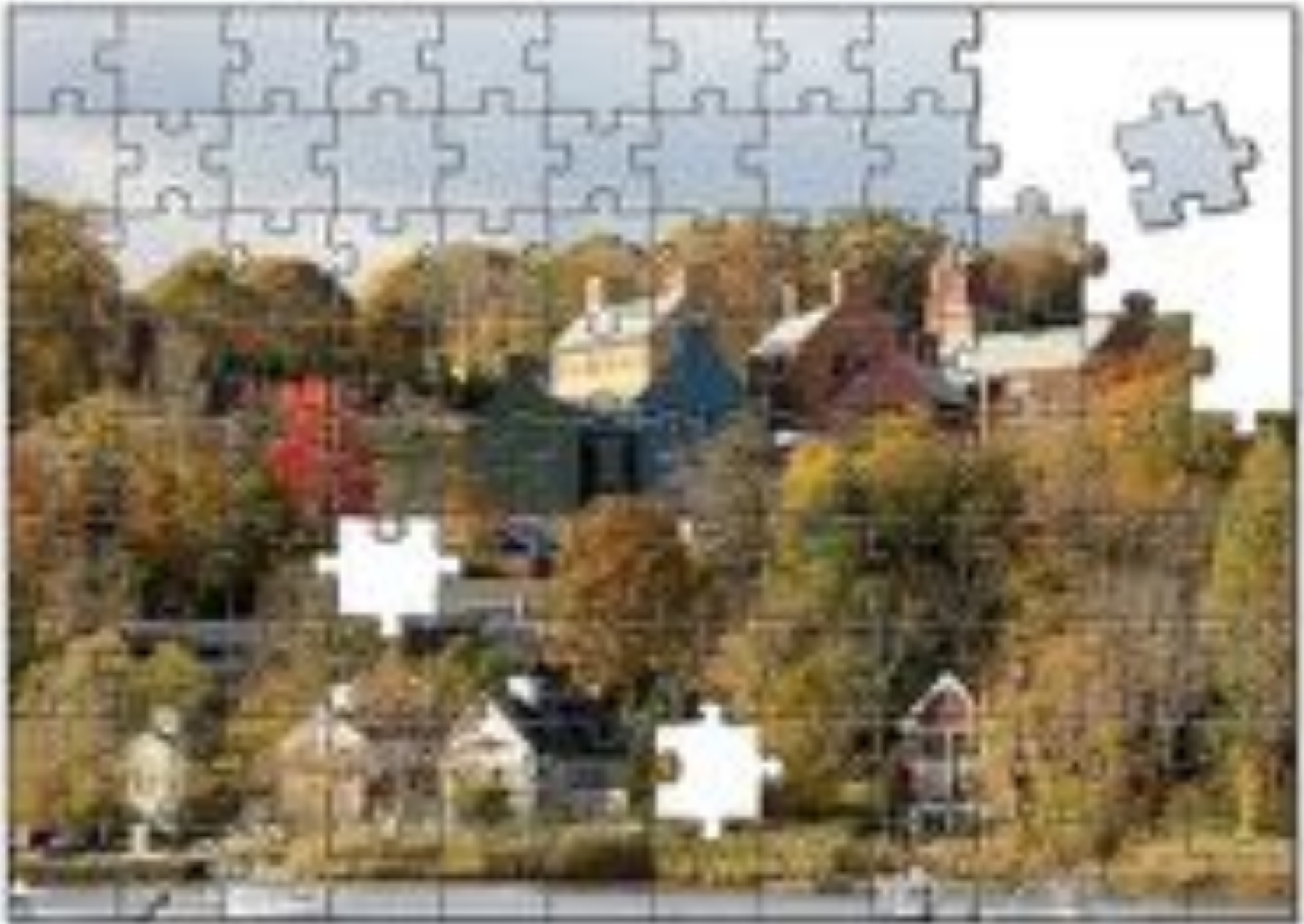
Imaging of fluorescently phospholinked labeled nucleotides as they are incorporated by a polymerase anchored to a Zero-Mode Waveguide (ZMW).



Single Molecule Sequences



“Corrective Lens” for Sequencing



PacBio Assembly Algorithms

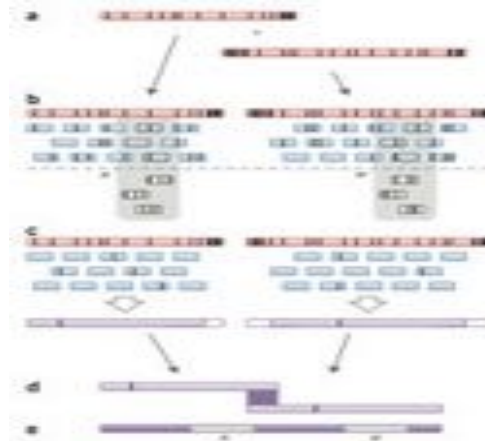
PBJelly



**Gap Filling
and Assembly Upgrade**

English *et al* (2012)
PLOS One. 7(11): e47768

PacBioToCA & ECTools



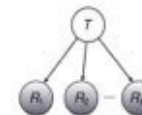
**Hybrid/PB-only Error
Correction**

Koren, Schatz, *et al* (2012)
Nature Biotechnology. 30:693–700

HGAP/MHAP & Quiver



$$\Pr(\mathbf{R} | T) = \prod_k \Pr(R_k | T)$$



Quiver Performance Results Comparison to Reference Genome (<i>M. ruber</i> ; 3.1 MB; SMRT® Cells)		
	Initial Assembly	Quiver Consensus
QV	43.4	54.5
Accuracy	99.99540%	99.99964%
Differences	141	11

**PB-only Correction &
Polishing**

Chin *et al* (2013)
Nature Methods. 10:563–569

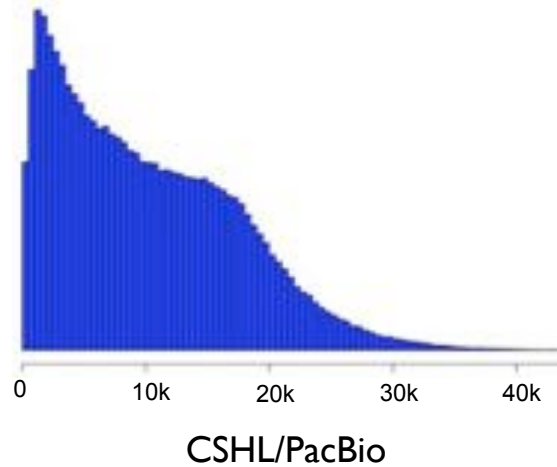
< 5x

PacBio Coverage

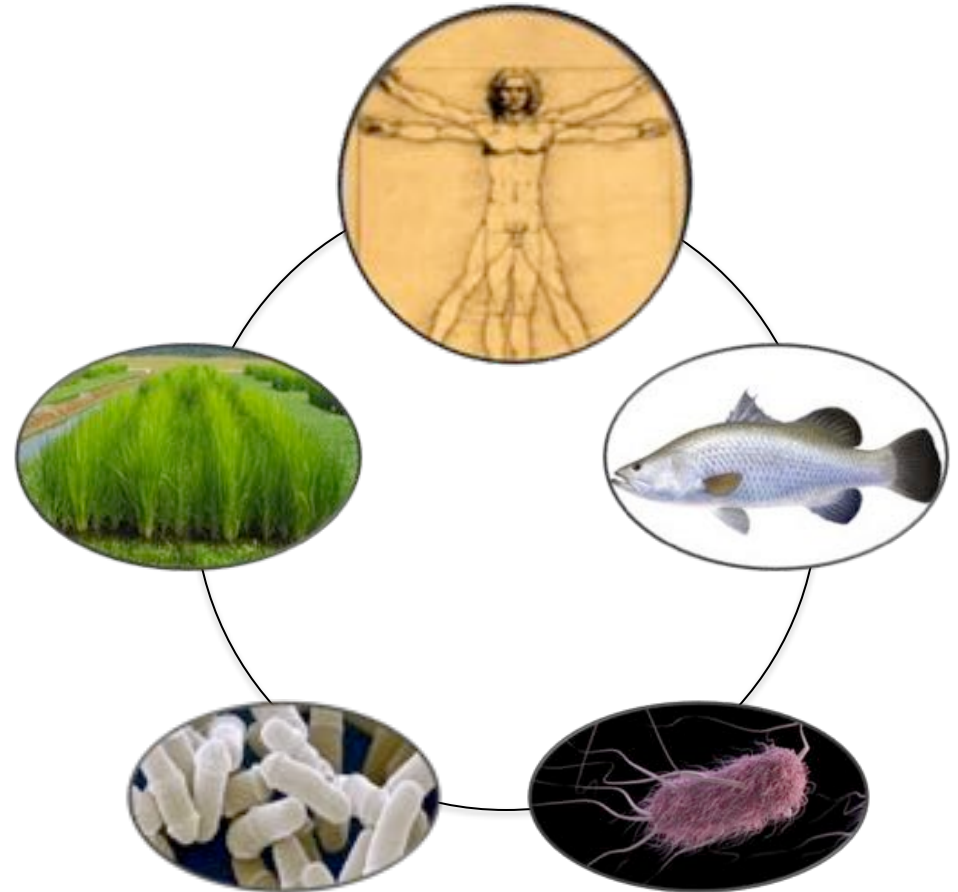
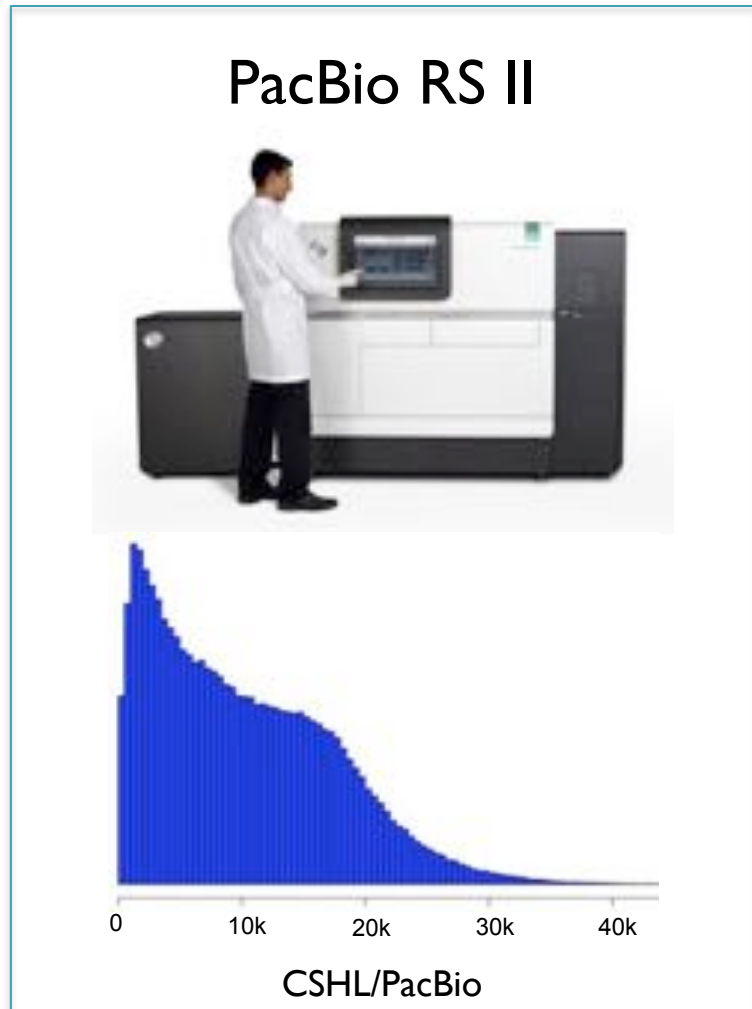
> 50x

3rd Gen Long Read Sequencing

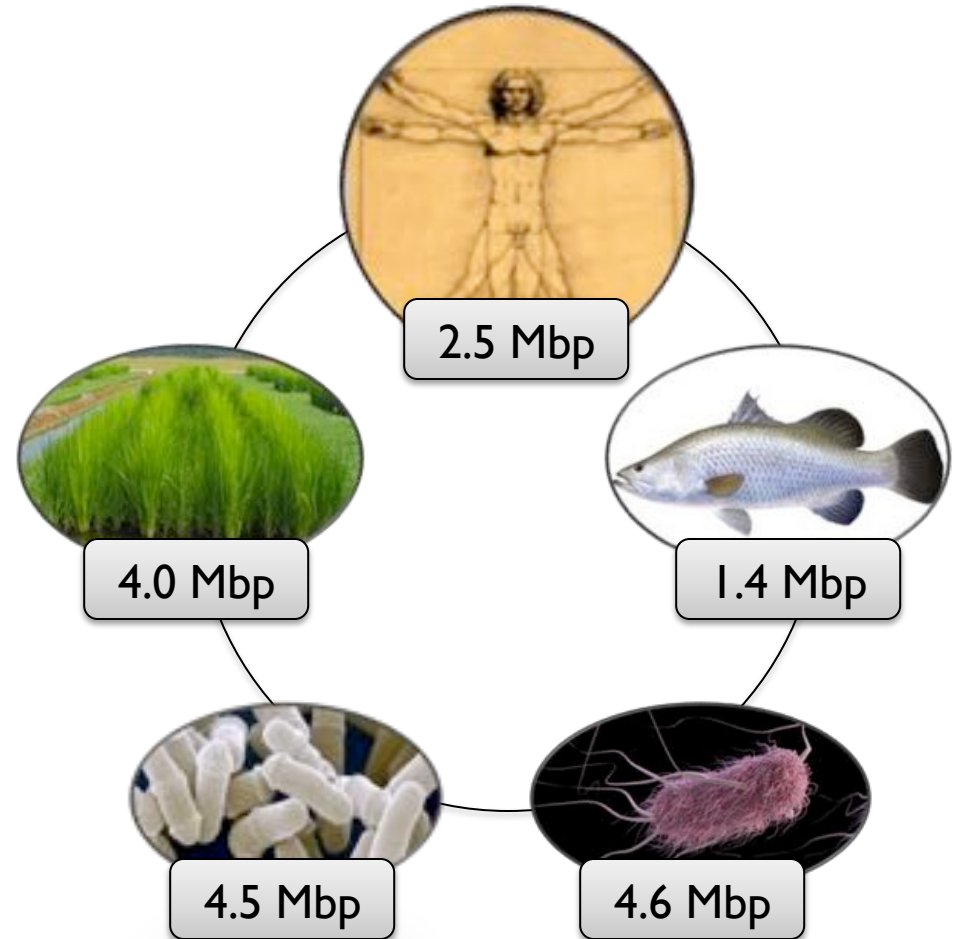
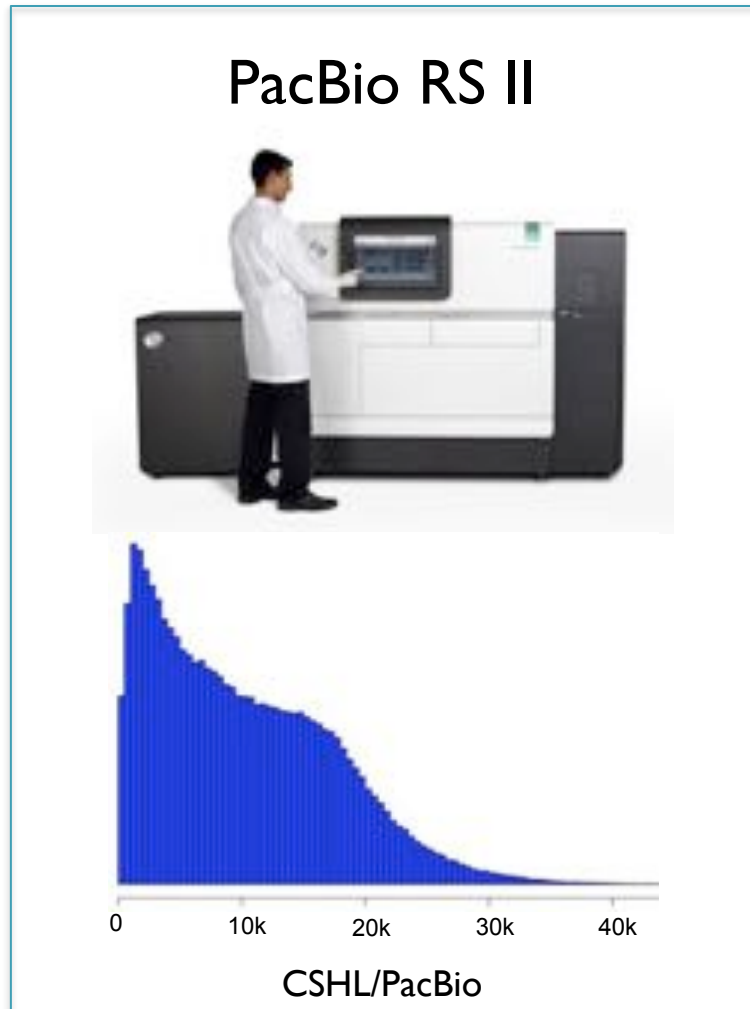
PacBio RS II



3rd Gen Long Read Sequencing

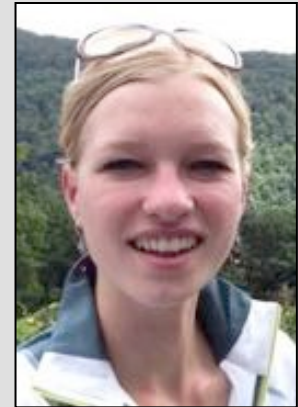


3rd Gen Long Read Sequencing

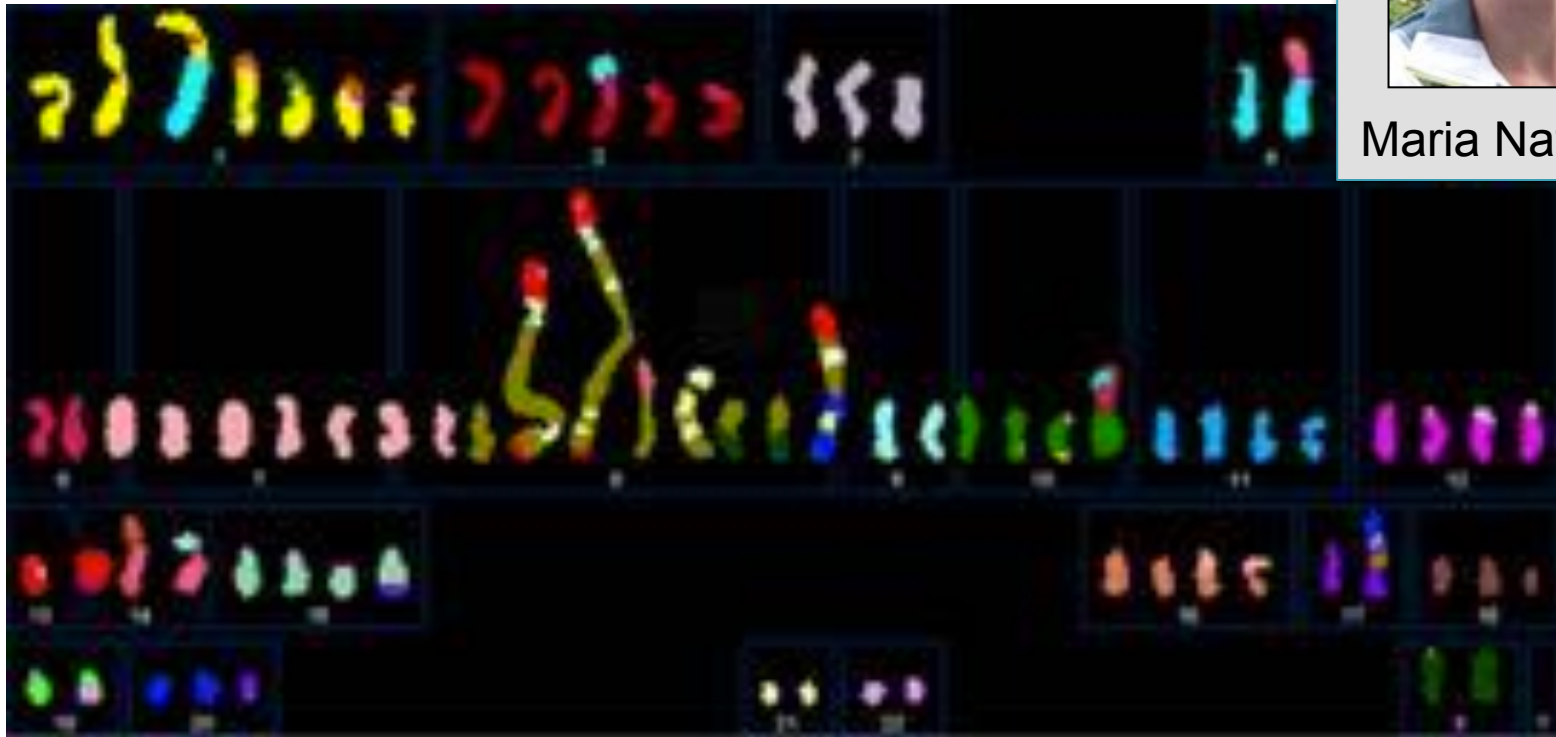


SK-BR-3

Most commonly used Her2-amplified breast cancer



Maria Nattestad

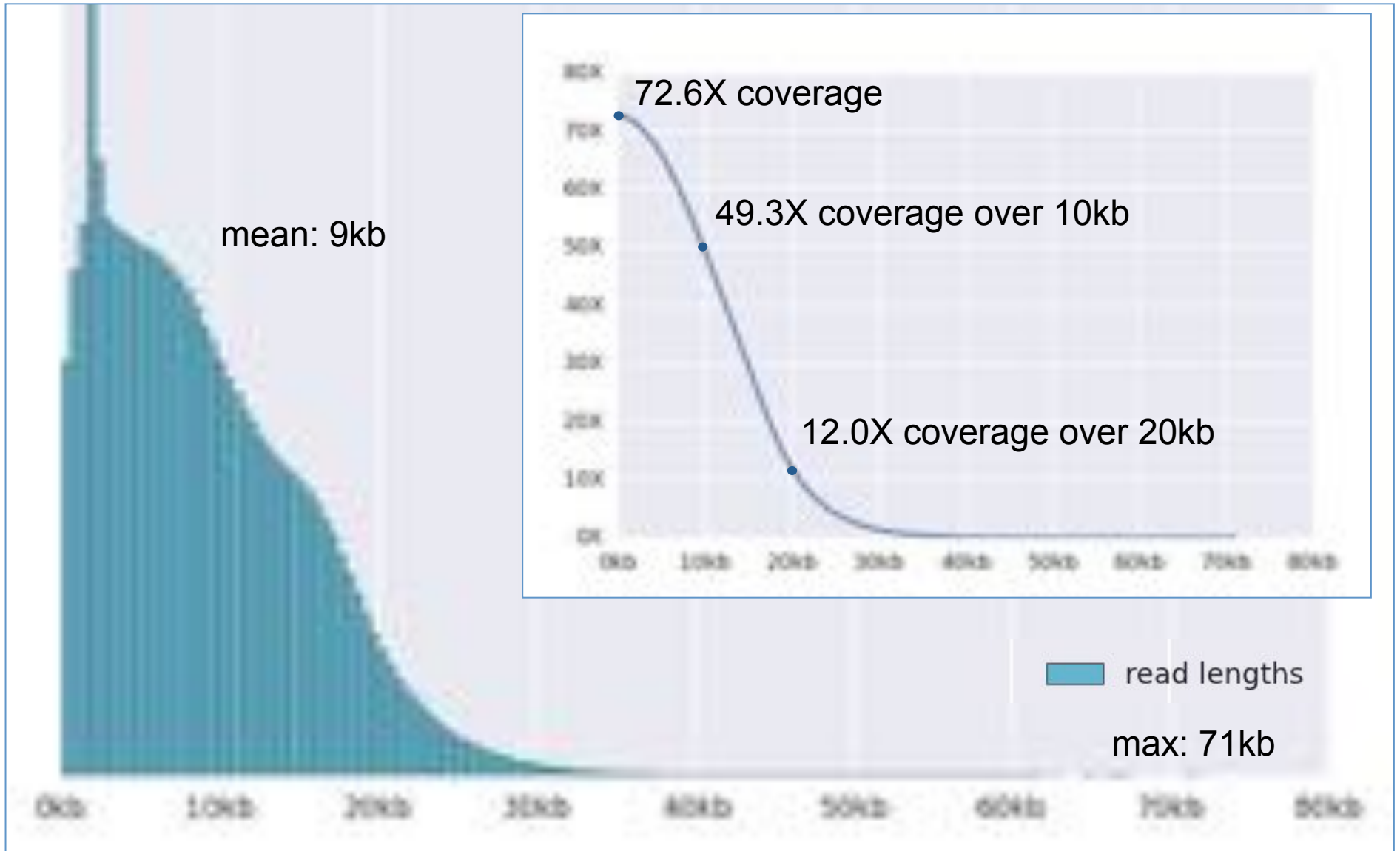


(Davidson et al, 2000)

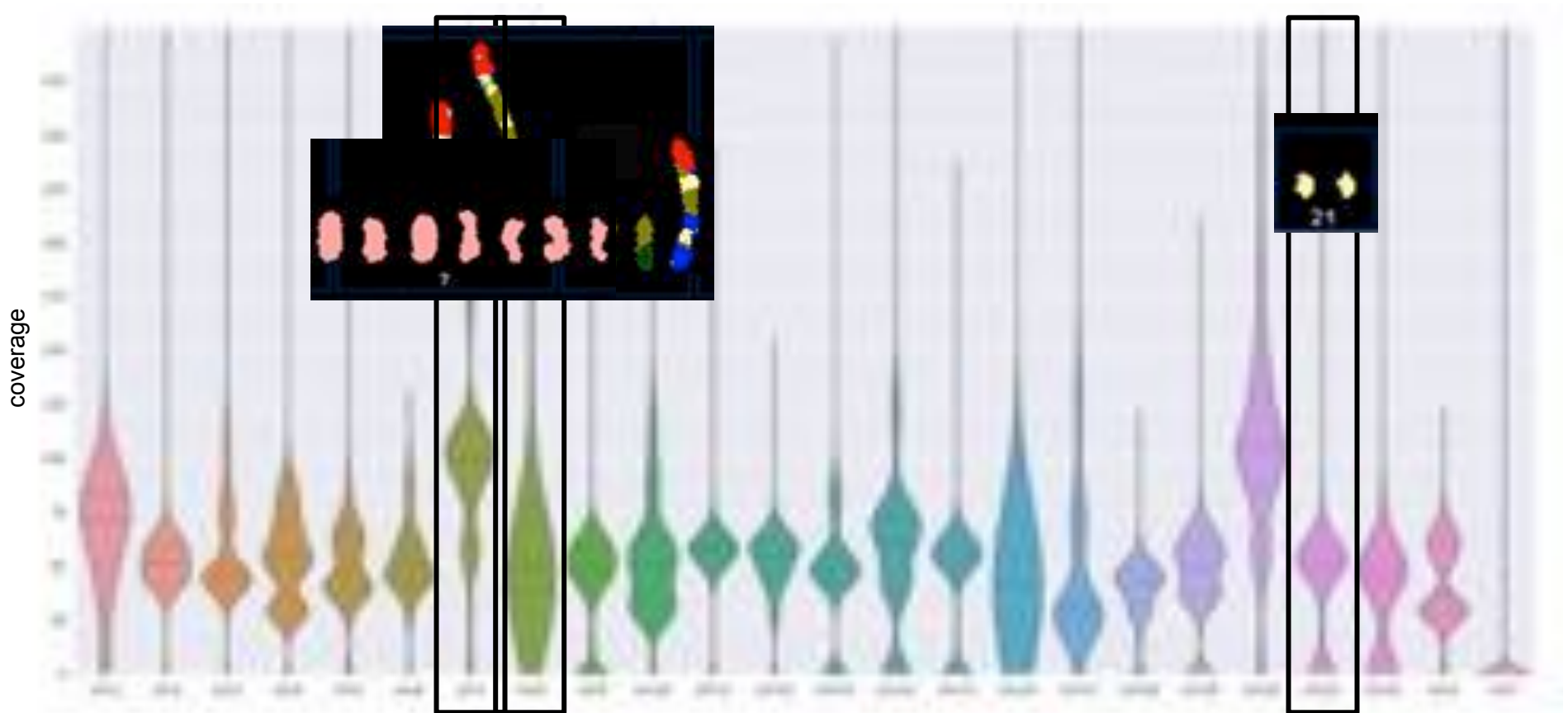
Can we resolve the complex structural variations, especially around Her2?

Ongoing collaboration between CSHL and OICR to *de novo* assemble the complete cell line genome with PacBio long reads

PacBio read length distribution



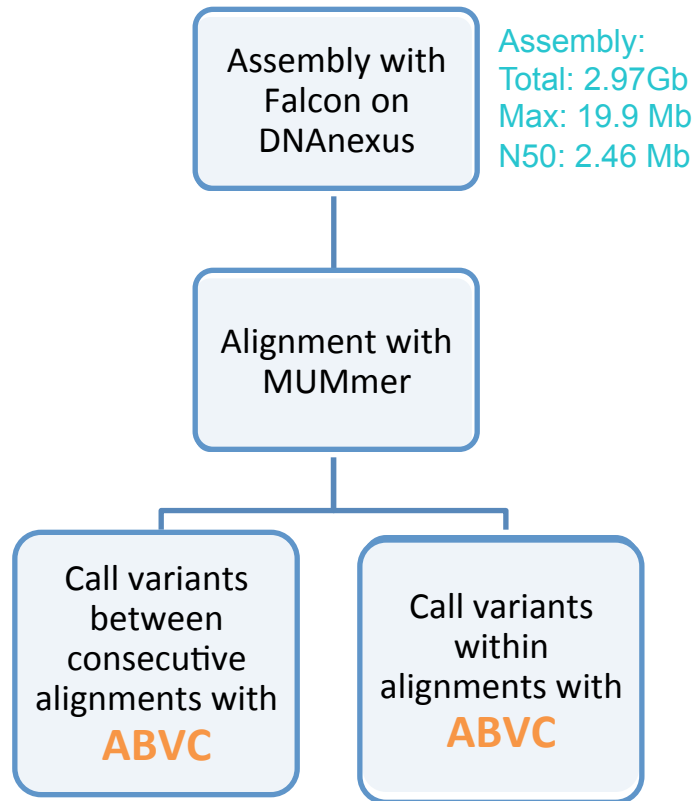
Genome Wide Coverage Analysis



Genome-wide coverage averages around 54X
Coverage per chromosome varies greatly as expected from previous karyotyping results

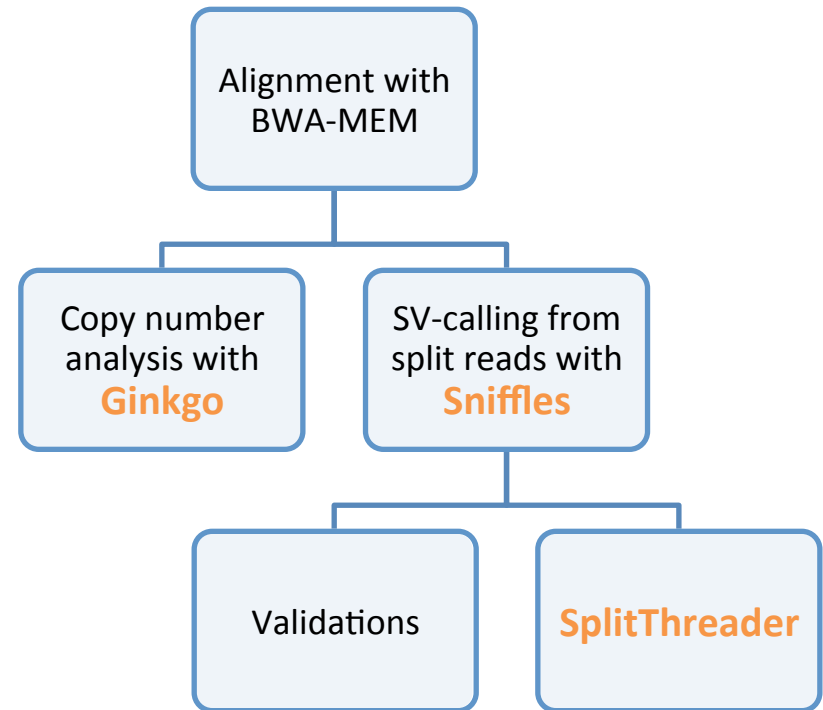
Structural Variation Analysis

Assembly-based



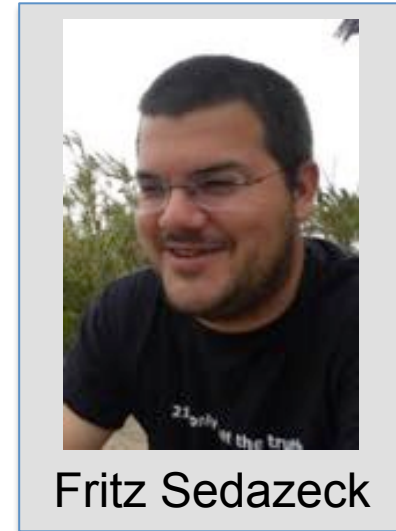
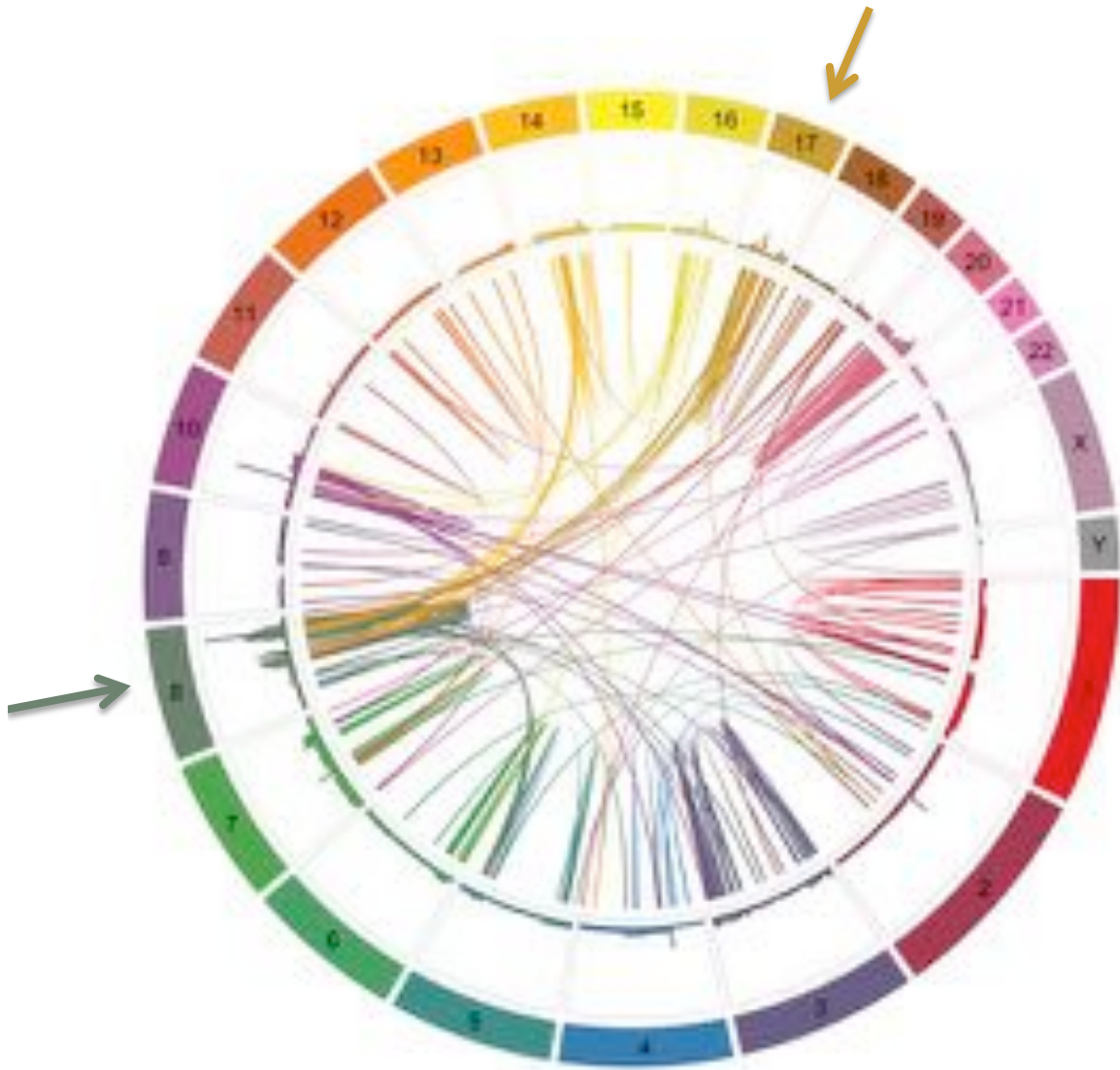
~ 11,000 local variants
50 bp < size < 10 kbp

Split-Read based



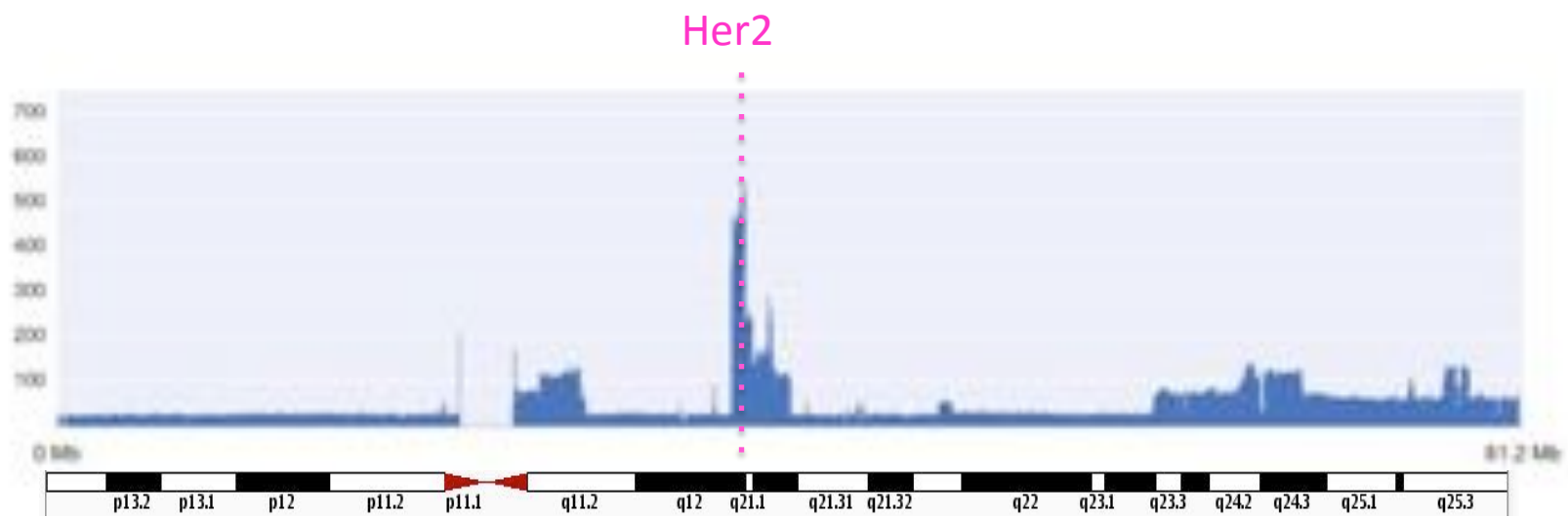
350 long-range variants
(>10kb distance)

Long Range Variations in SK-BR-3



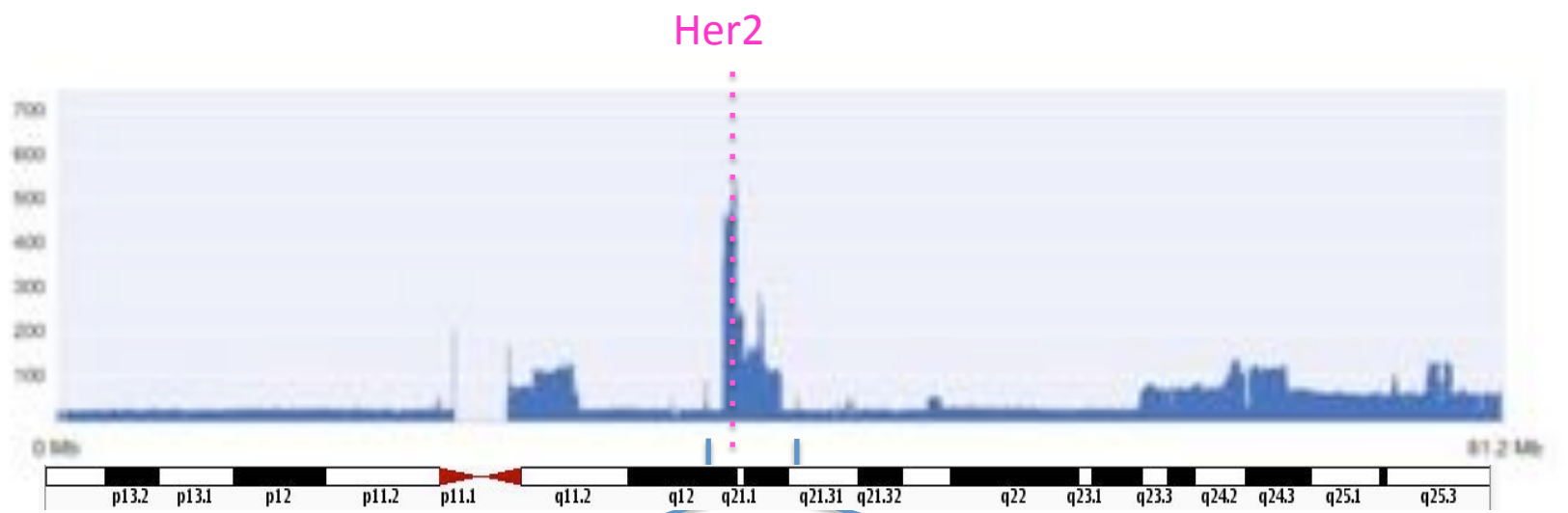
Analysis by Sniffles

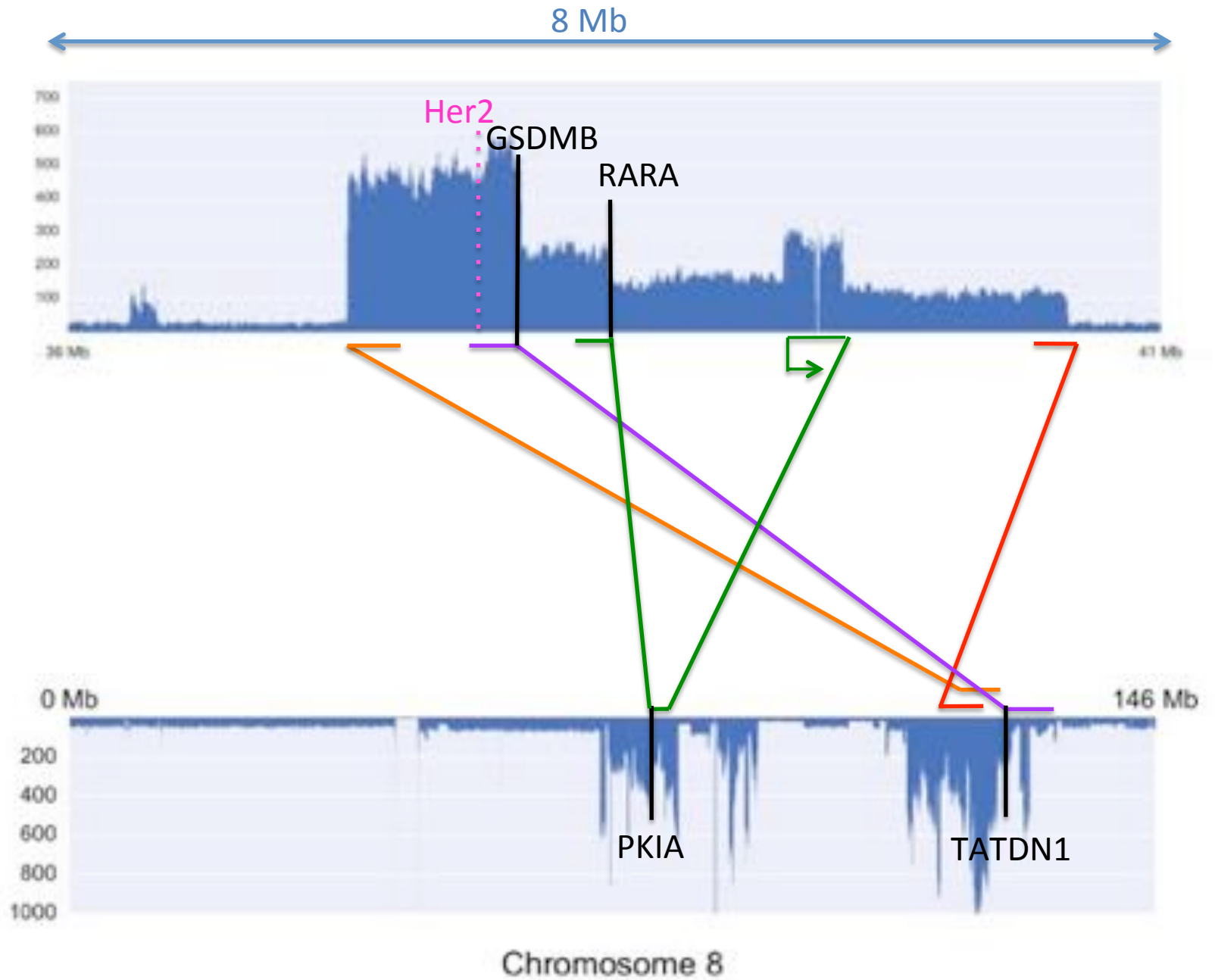
- 350 variants \geq 10kbp
- Requires 10 split reads broken within a 200 bp interval on both sides of the translocation



← Chr 17: 83 Mb →

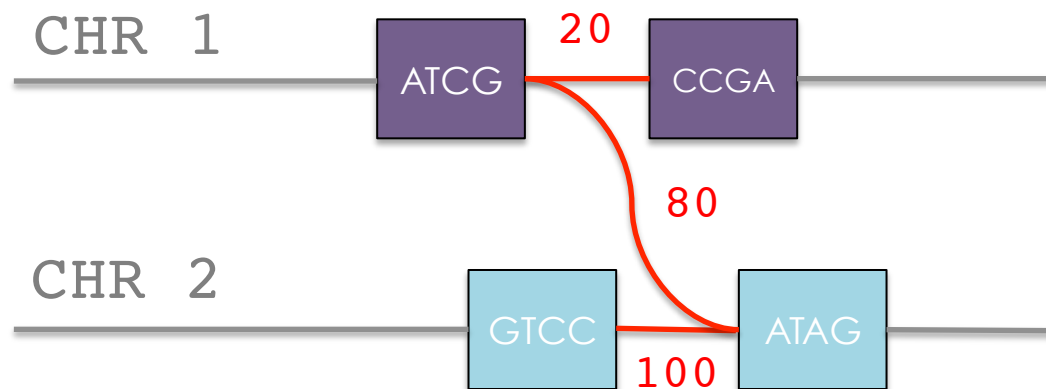
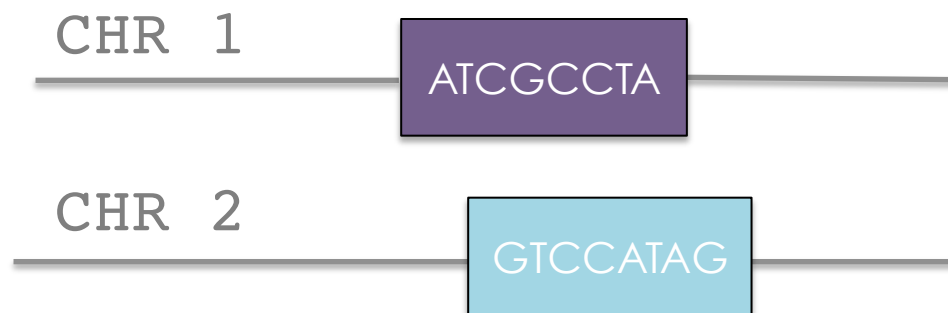
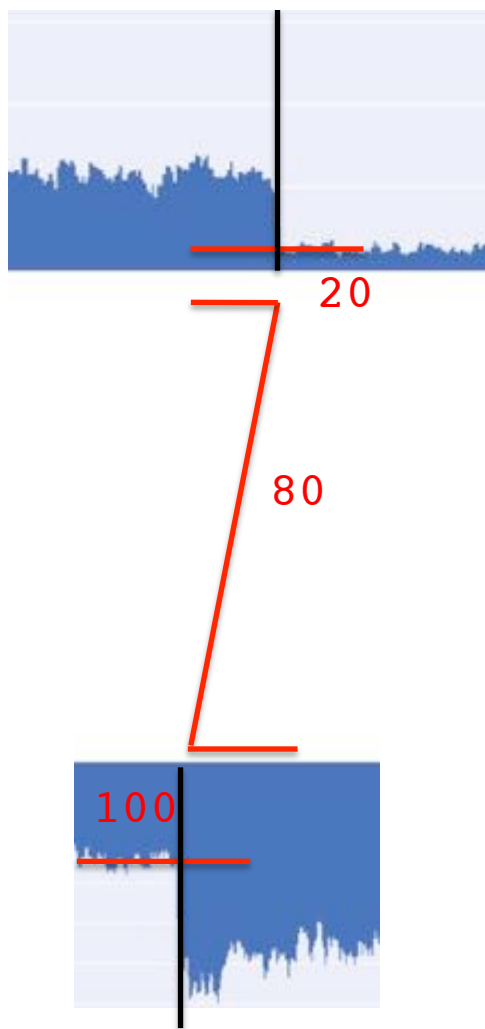
← 8 Mb →

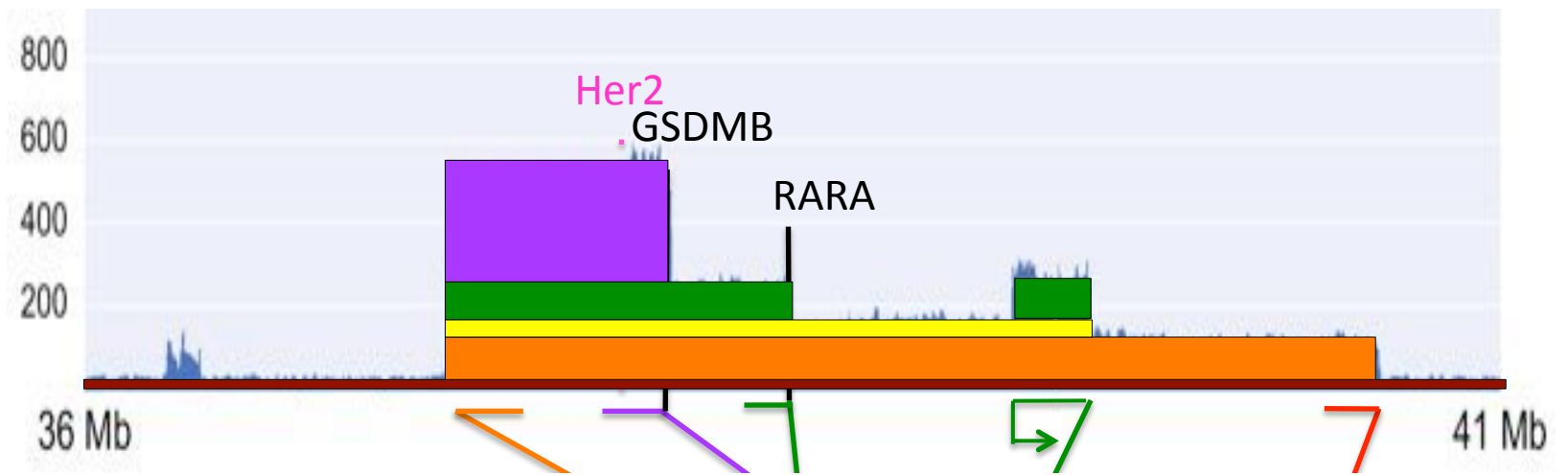




SplitThreader

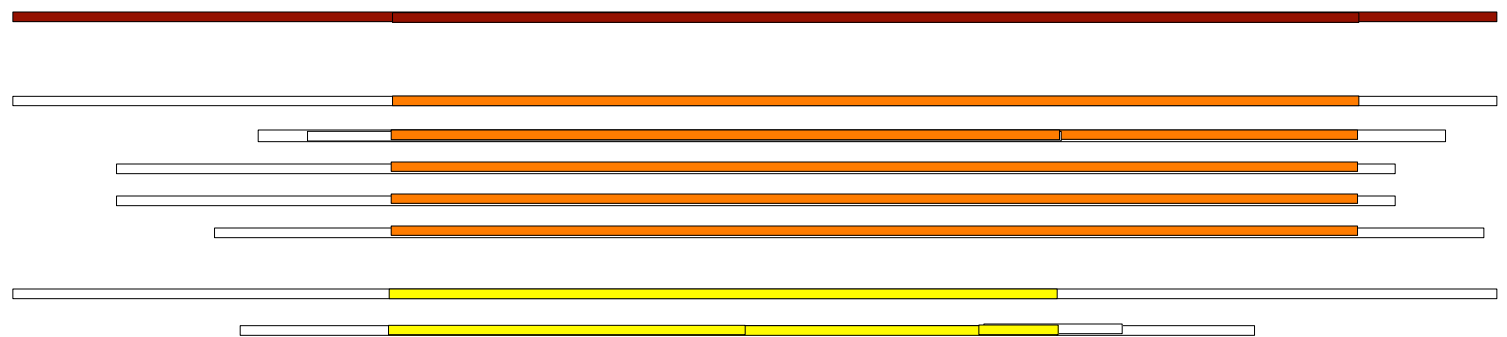
Graphical threading to retrace complex history of rearrangements in cancer genomes





Chr 17

Chr 8



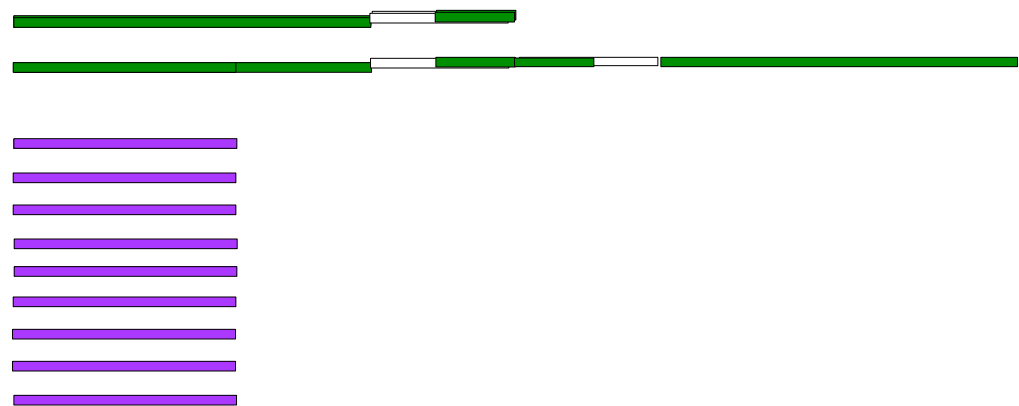
1. Healthy chromosome 17

2. Translocation into chromosome 8

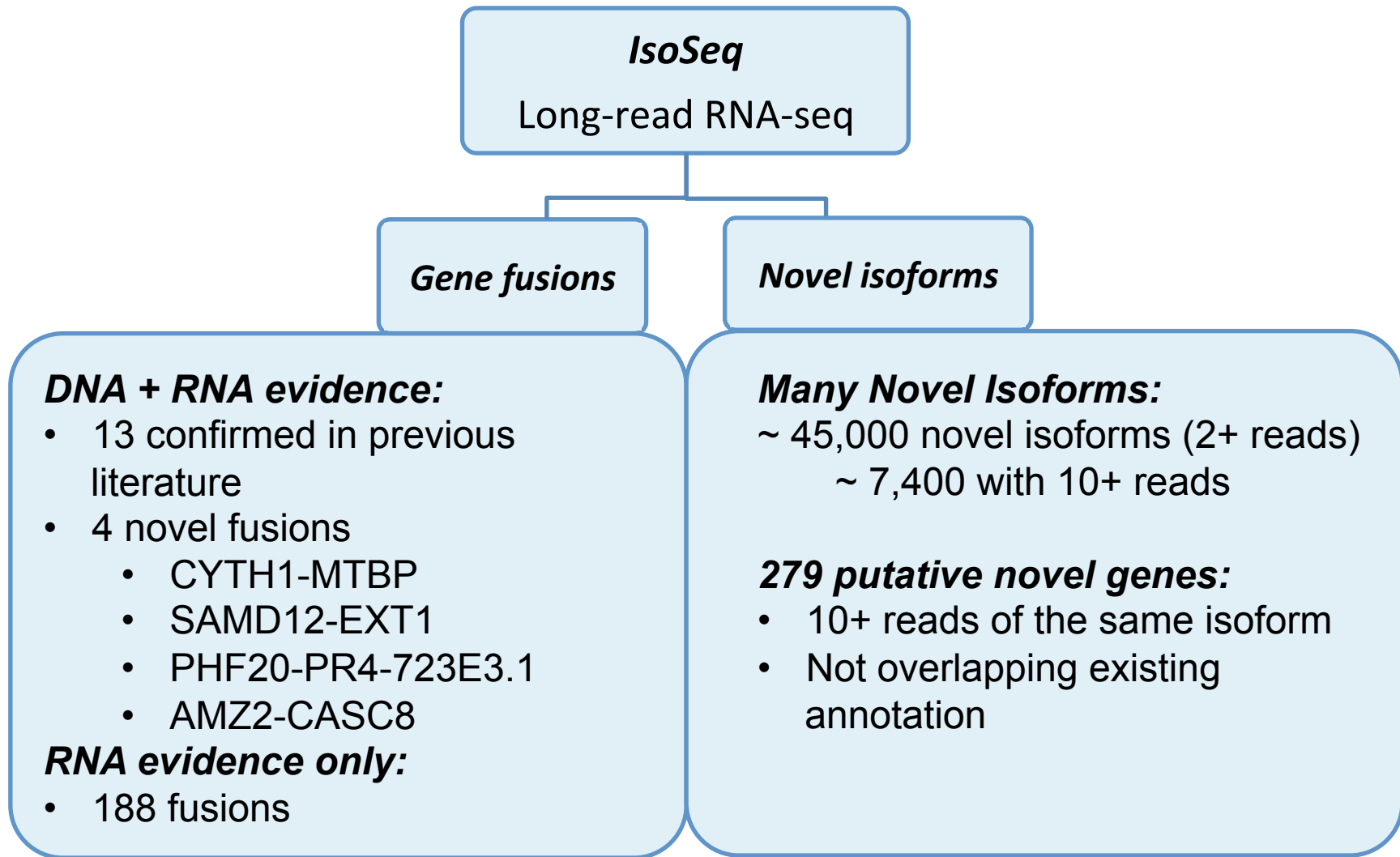
3. Translocation within chromosome 8

4. Complex variant and inverted duplication within chromosome 8

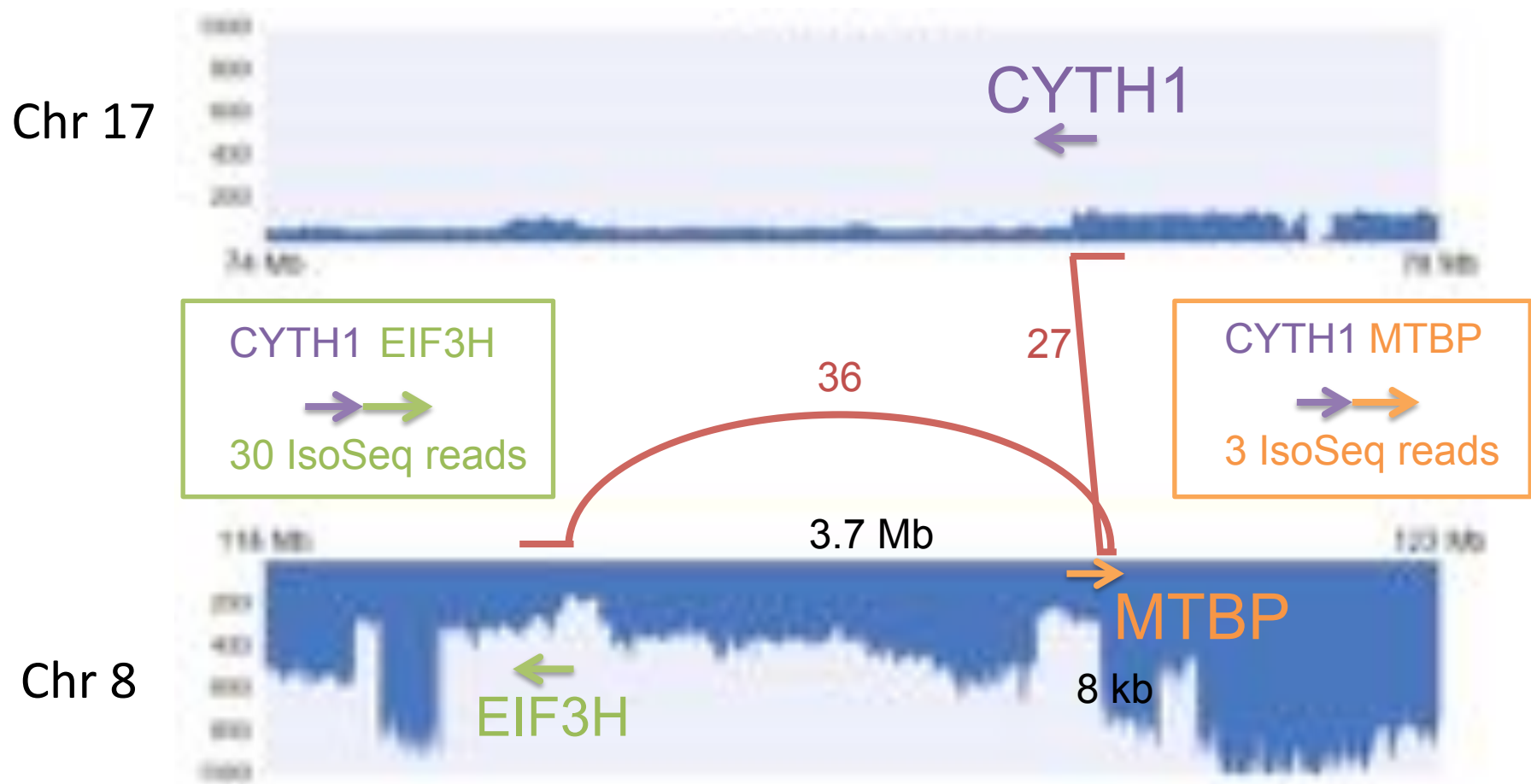
5. Translocation within chromosome 8



Transcriptome analysis with IsoSeq



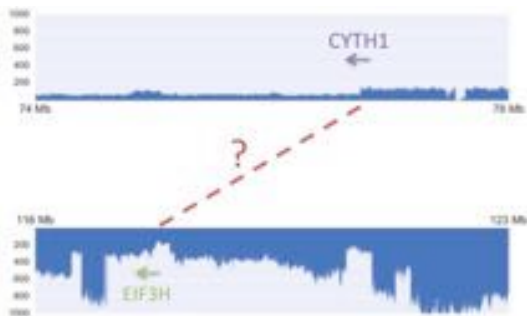
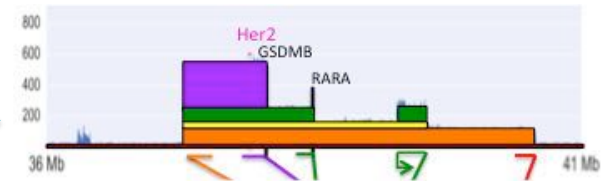
CYTH1-EIF3H gene fusion



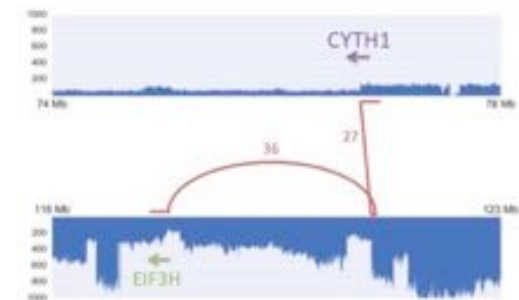
The genome informs the transcriptome



Explain amplifications



Trace gene fusions

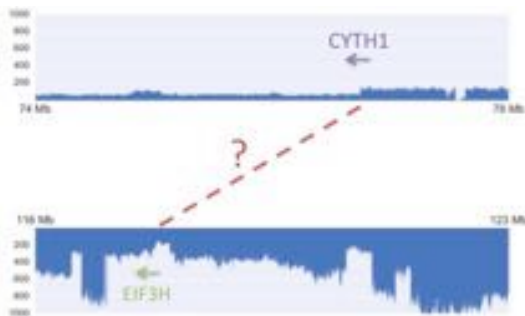
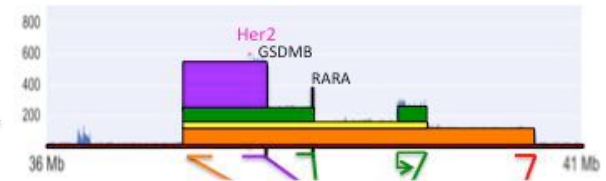


Data and additional results: <http://schatzlab.cshl.edu/data/skbr3/>

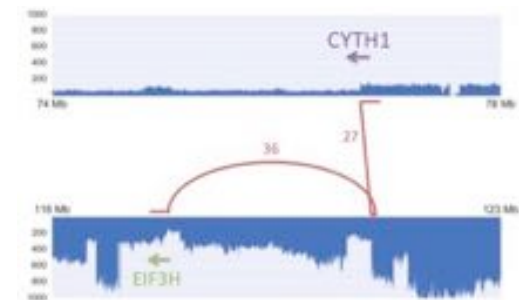
The genome informs the transcriptome ... and informs the prognosis



Explain amplifications



Trace gene fusions



Data and additional results: <http://schatzlab.cshl.edu/data/skbr3/>

PacBio Roadmap



PacBio RS II

\$750k instrument cost
1895 lbs

~\$75k / human @ 50x



SMRTcell

150k Zero Mode Waveguides
~10kb average read length
~1 GB / SMRTcell
~\$500 / SMRTcell

PacBio Roadmap



PacBio Sequel

\$350k instrument cost
841 lbs

~\$15k / human @ 50x



SMRTcell v2

1M Zero Mode Waveguides
~15kb average read length
~10 GB / SMRTcell
~\$1000 / SMRTcell

Oxford Nanopore



MinION

\$2k / instrument
1 GB / day
~\$300k / human @ 50x



PromethION

\$75k / instrument
>>100GB / day
??? / human @ 50x

Oxford Nanopore sequencing, hybrid error correction, and de novo assembly of a eukaryotic genome

Goodwin, S, Gurtowski, J, Ethe-Sayers, S, Deshpande, P, Schatz MC, McCombie, WR (2015) Genome Research doi: 10.1101/gr.191395.115

Our Destiny



Outline

1. Single Molecule Sequencing

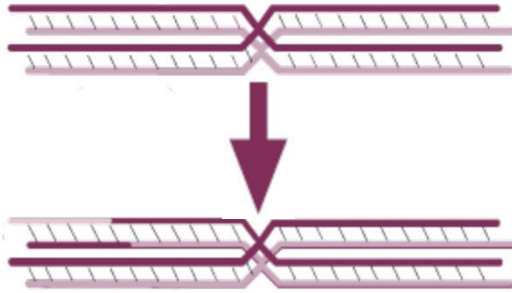
Long read sequencing of a breast cancer cell line

2. Single Cell Copy Number Analysis

Intra-tumor heterogeneity and metastatic progression



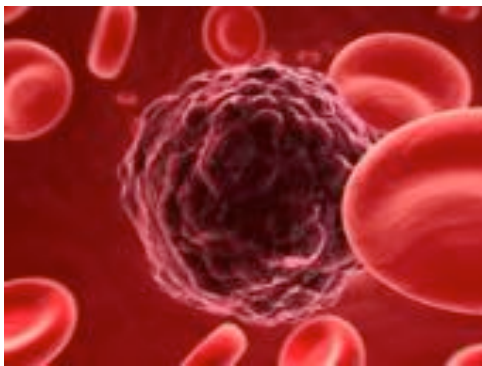
Single Cell Sequencing



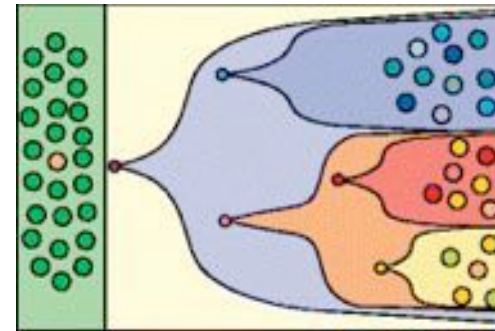
Recombination /
Crossover in germ cells



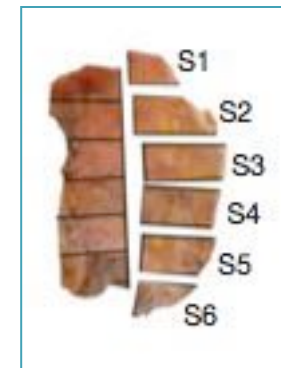
Neuronal mosaicism



Circulating tumor cells

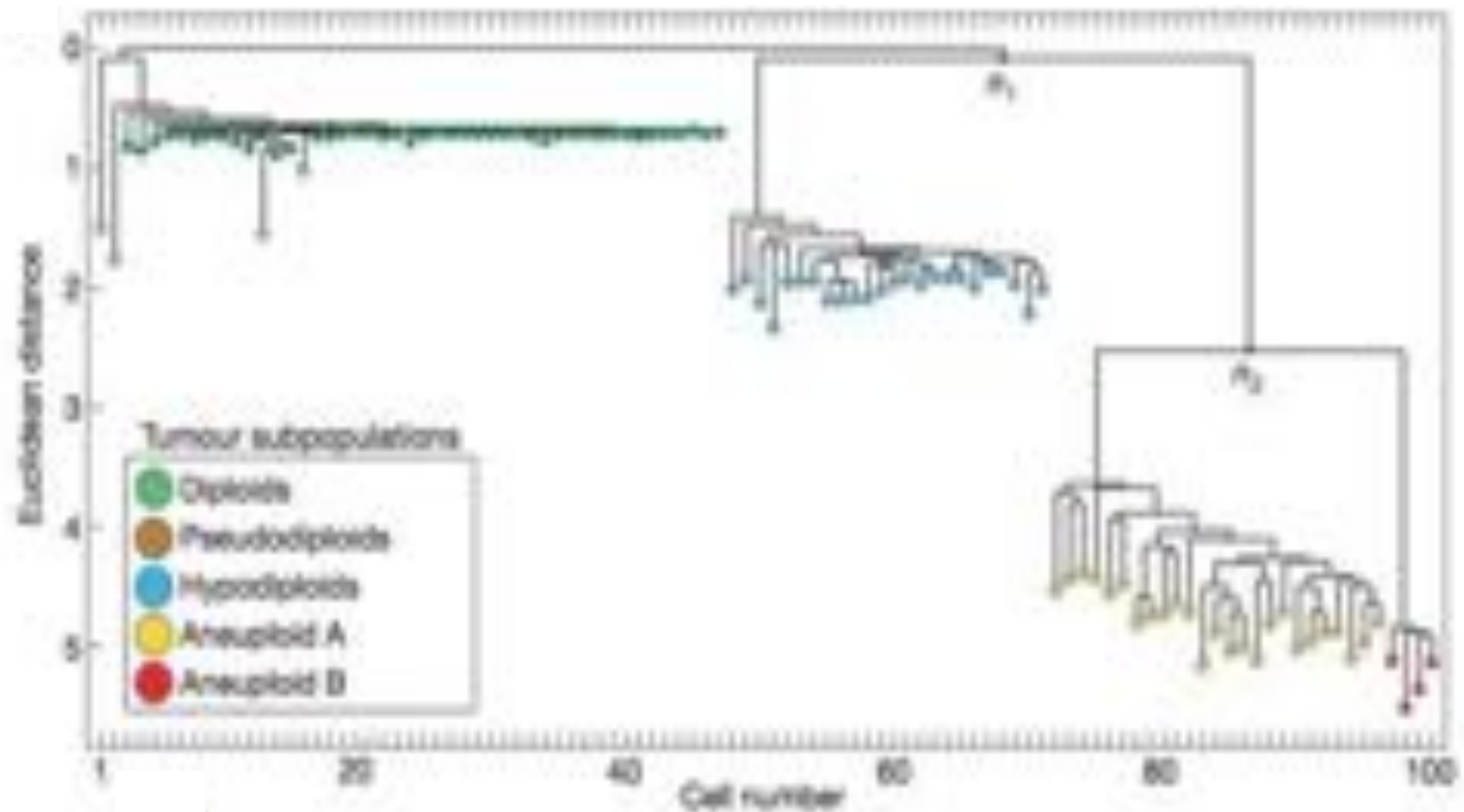


Clonal Evolution
in tumors

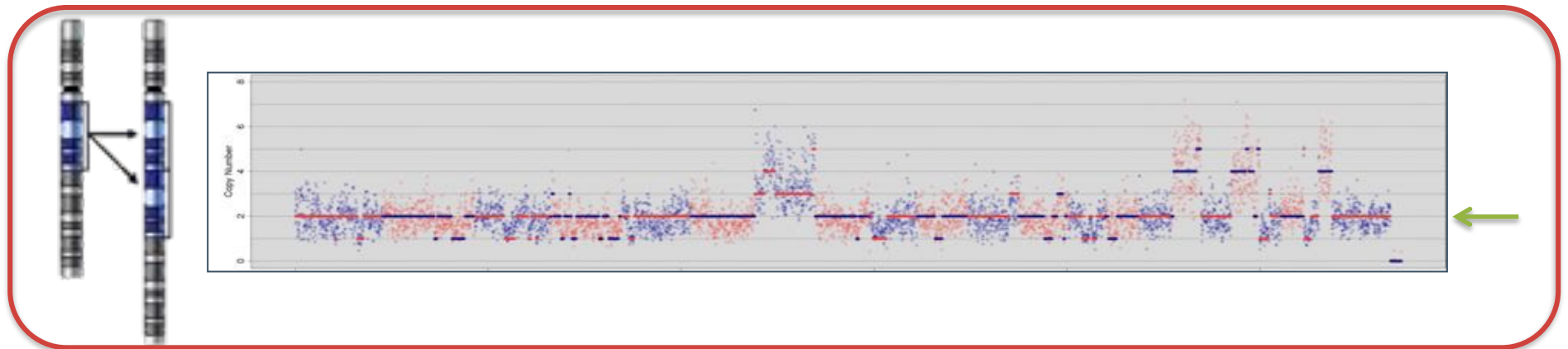
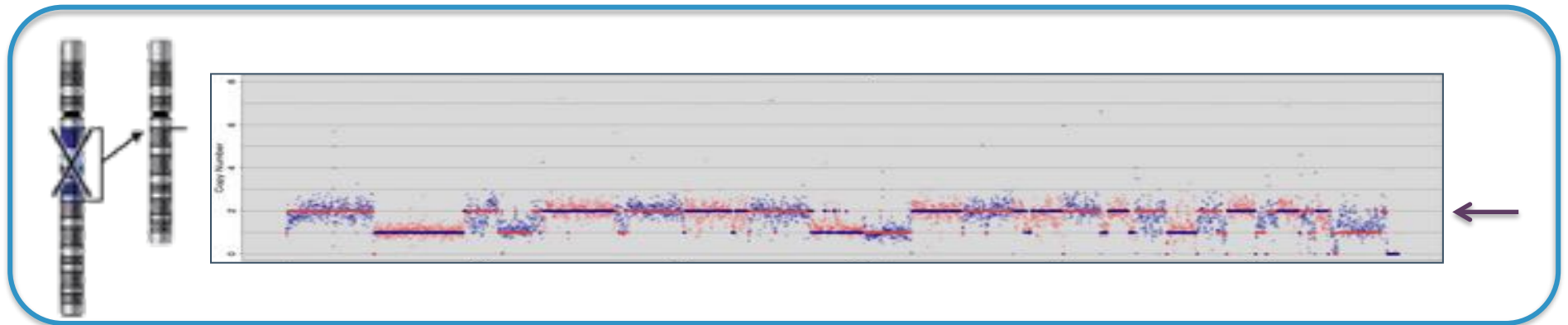
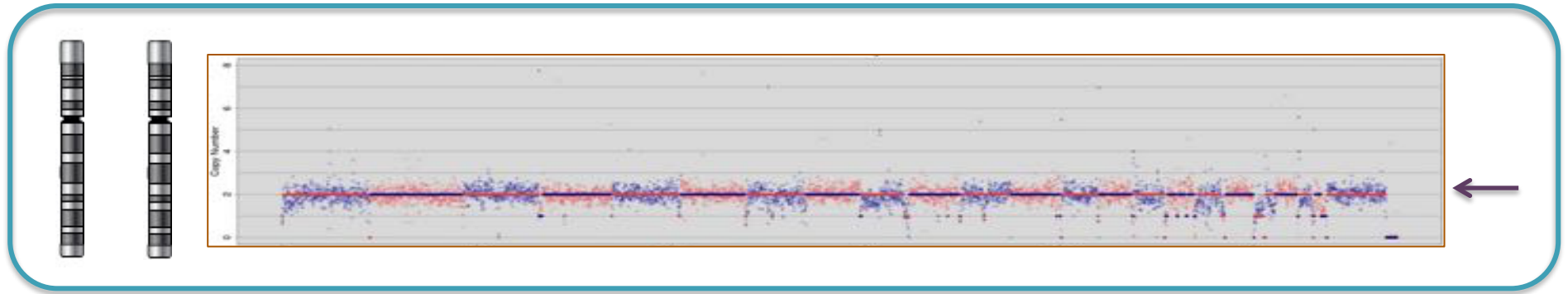


Tumour evolution inferred by single-cell sequencing

Nicholas Navin^{1,2}, Jude Kendall¹, Jennifer Troge³, Peter Andrews¹, Linda Rodgers¹, Jeanne McIndoo¹, Kerry Cook¹, Asya Stepansky¹, Dan Levy¹, Diane Esposito¹, Lakshmi Muthuswamy³, Alex Krasnitz¹, W. Richard McCombie¹, James Hicks¹ & Michael Wigler¹



Copy-number Profiles

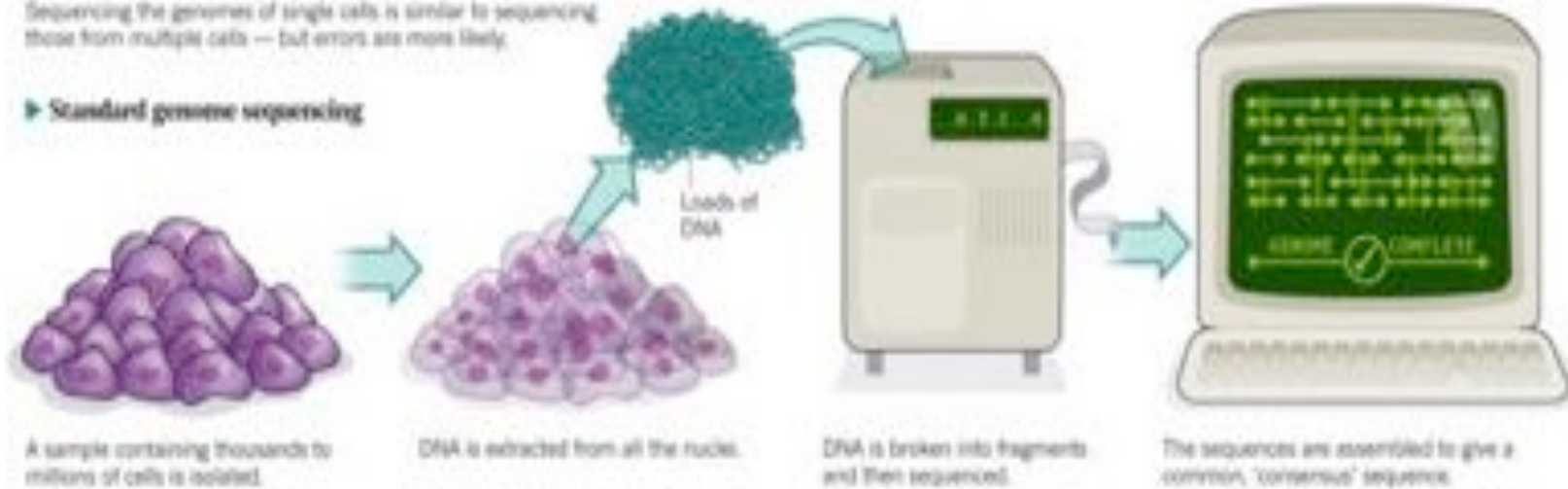


Whole Genome Amplification

ONE GENOME FROM MANY

Sequencing the genomes of single cells is similar to sequencing those from multiple cells — but errors are more likely.

► Standard genome sequencing



Whole Genome Amplification

ONE GENOME FROM MANY

Sequencing the genomes of single cells is easier to sequencing those from multiple cells ... but errors are more likely

In Standard genome sequencing



A sample containing thousands to millions of cells is pooled



DNA is extracted from all the cells



DNA is broken into fragments and then sequenced



The sequences are assembled to give a common 'consensus' sequence

In Single-cell sequencing



A single cell is difficult to isolate, but it can be done mechanically or with an advanced cell sorter



The DNA is extracted and amplified during which a few errors may occur

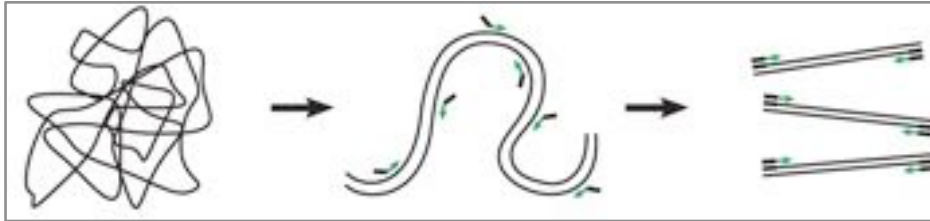


Amplified DNA is sequenced



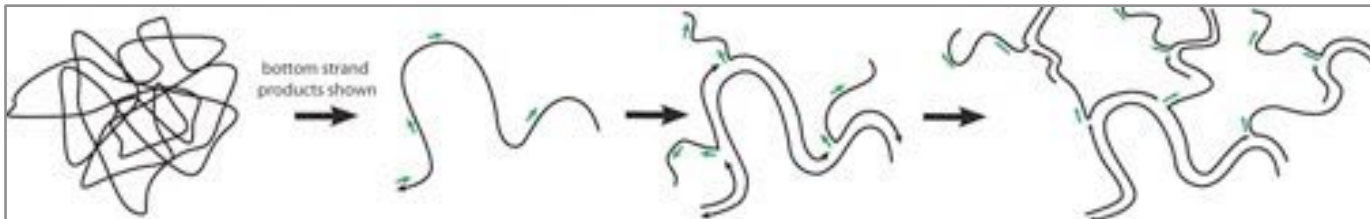
Direct individual reads are more likely to contain errors (SNPs) than the assembled consensus gene

Whole Genome Amplification Techniques



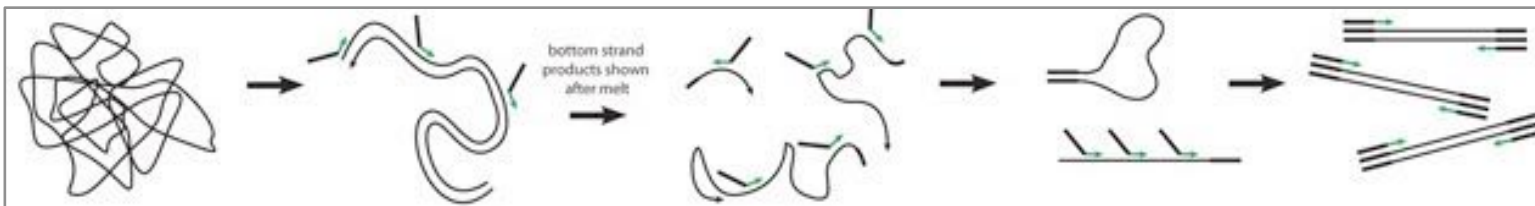
DOP-PCR: Degenerate Oligonucleotide Primed PCR

Telenius et al. (1992) Genomics



MDA: Multiple Displacement Amplification

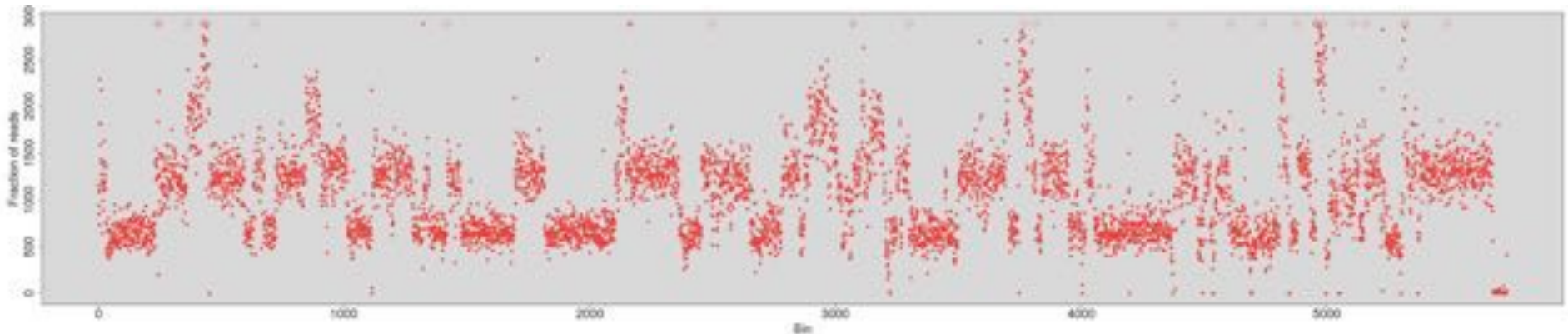
Dean et al. (2002) PNAS



MALBAC: Multiple Annealing and Looping Based Amplification Cycles

Zong et al. (2012) Science

Data are noisy

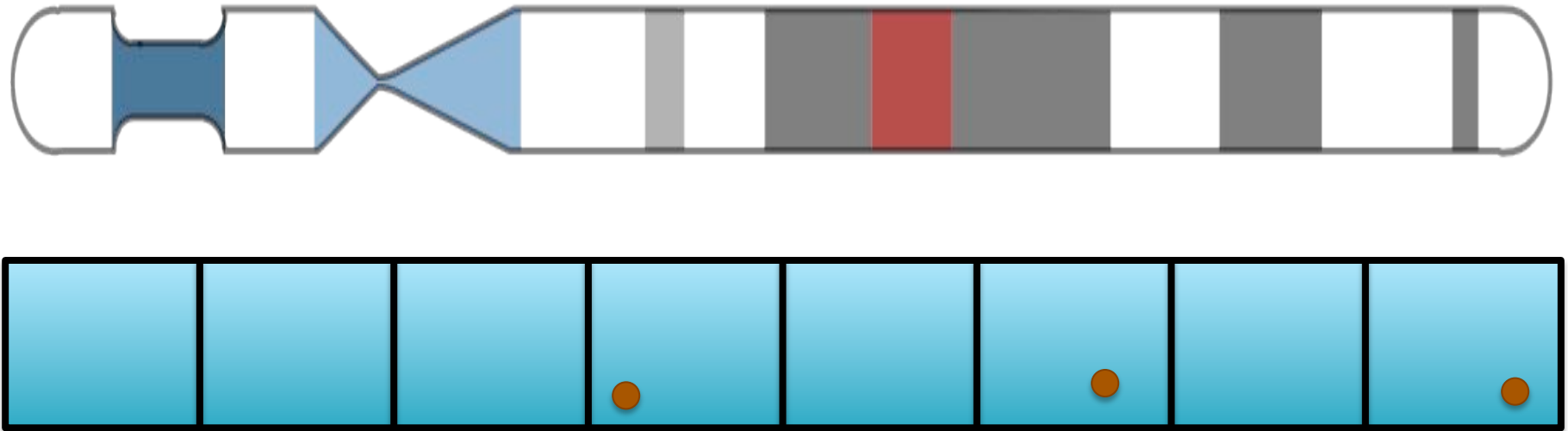


Potential for biases at every step

- WGA: Non-uniform amplification
- Library Preparation: Low complexity, read duplications, barcoding
- Sequencing: GC artifacts, short reads
- Computation: mappability, GC correction, segmentation, tree building

Coverage is too sparse and noisy for SNP analysis,
-> requires special processing

I) Binning

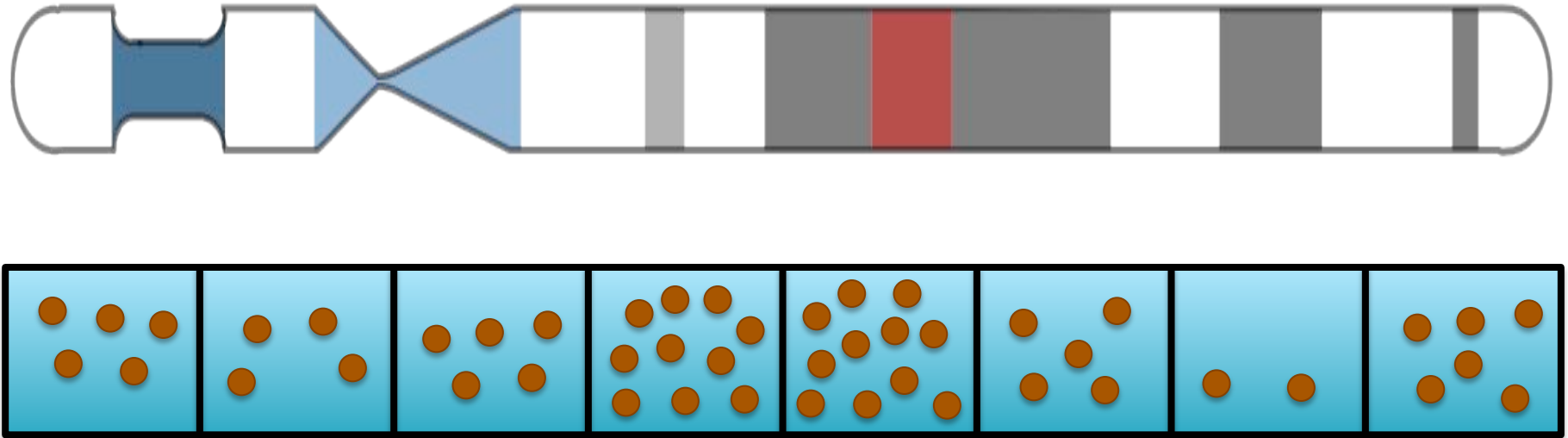


Single Cell CNV analysis

- Divide the genome into “bins” with $\sim 50 - 100$ reads / bin
- Map the reads and count reads per bin

Use uniquely mappable bases to establish bins

I) Binning

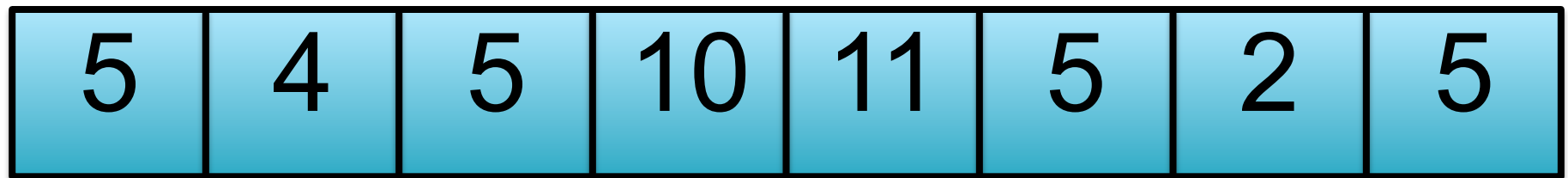


Single Cell CNV analysis

- Divide the genome into “bins” with $\sim 50 - 100$ reads / bin
- Map the reads and count reads per bin

Use uniquely mappable bases to establish bins

I) Binning

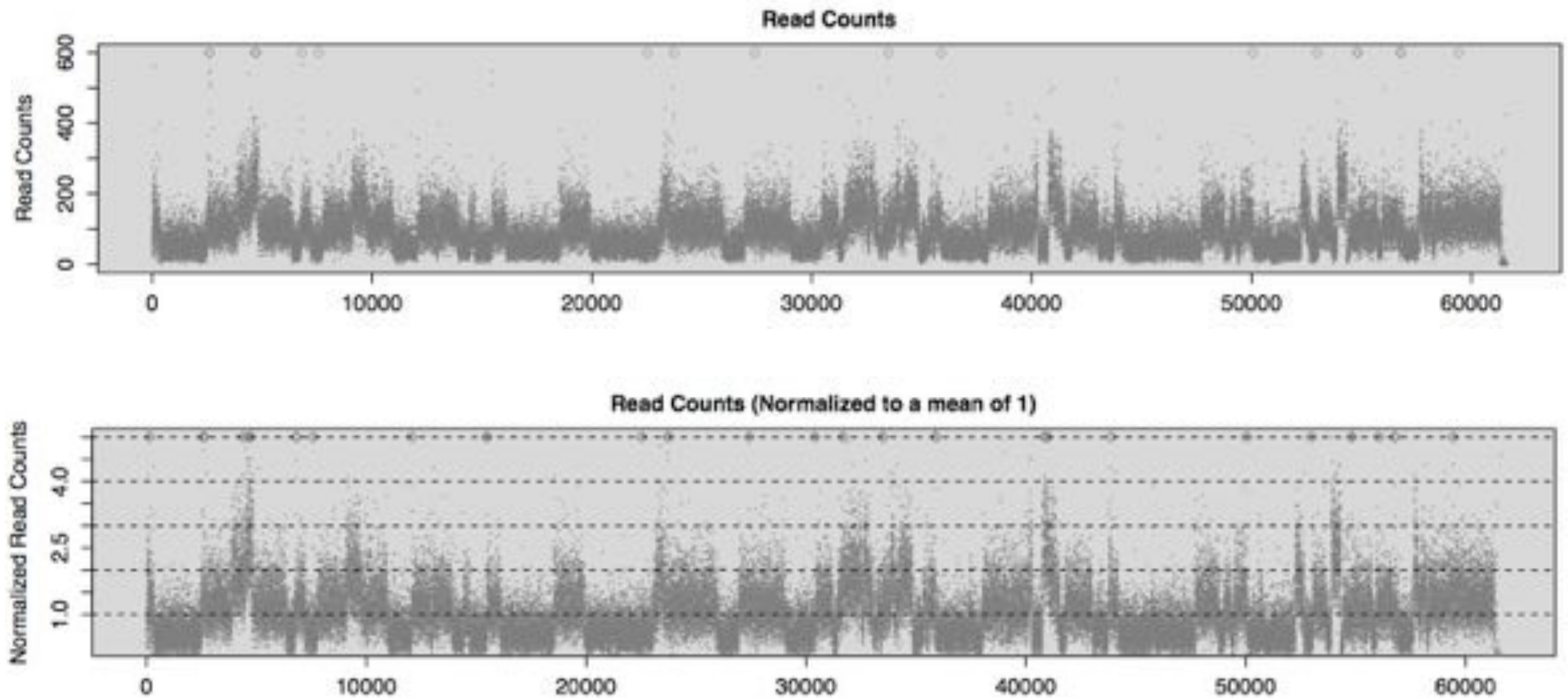


Single Cell CNV analysis

- Divide the genome into “bins” with ~50 – 100 reads / bin
- Map the reads and count reads per bin

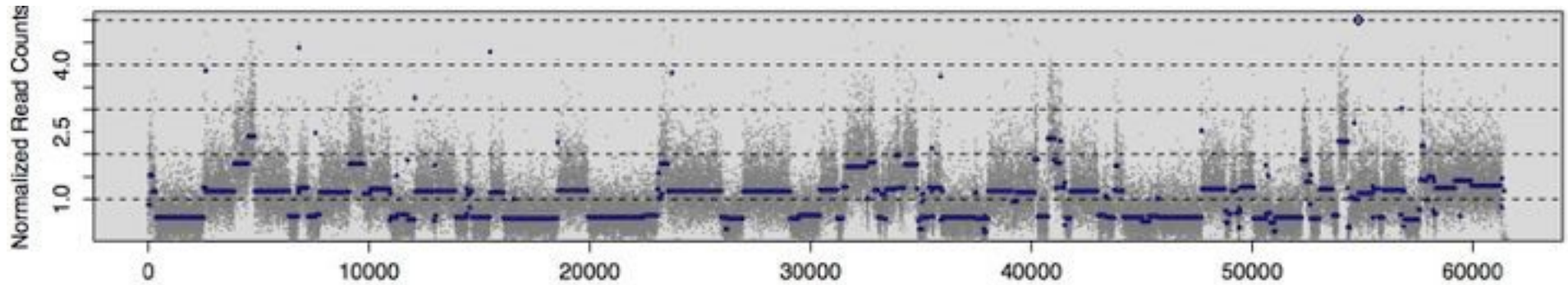
Use uniquely mappable bases to establish bins

2) Normalization

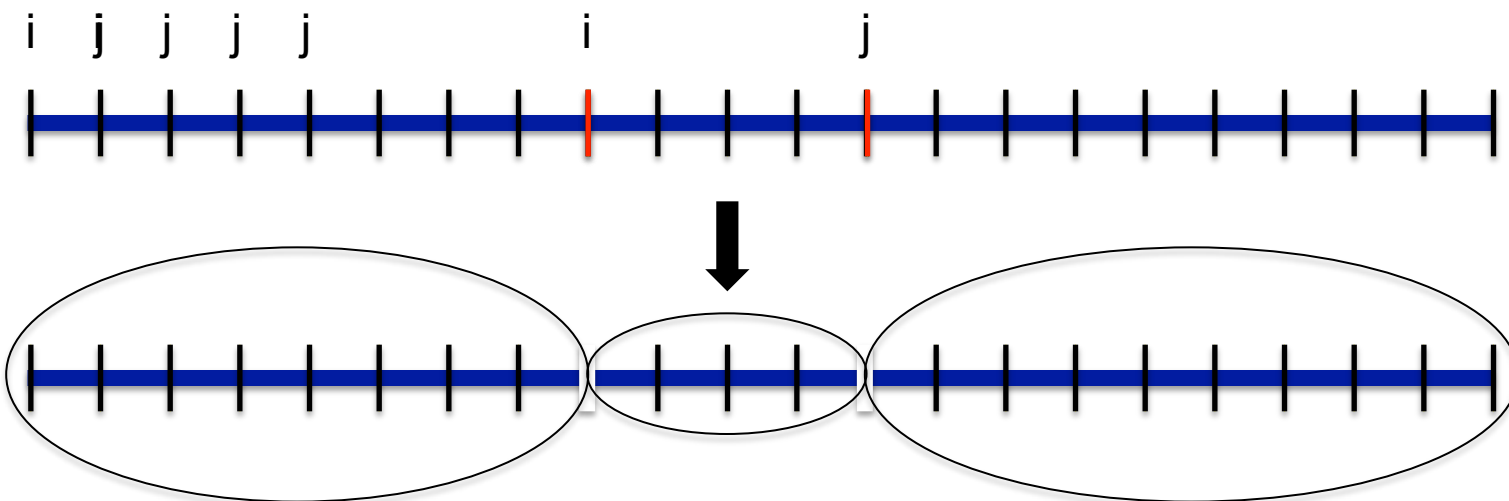


Also correct for mappability, GC content, amplification biases

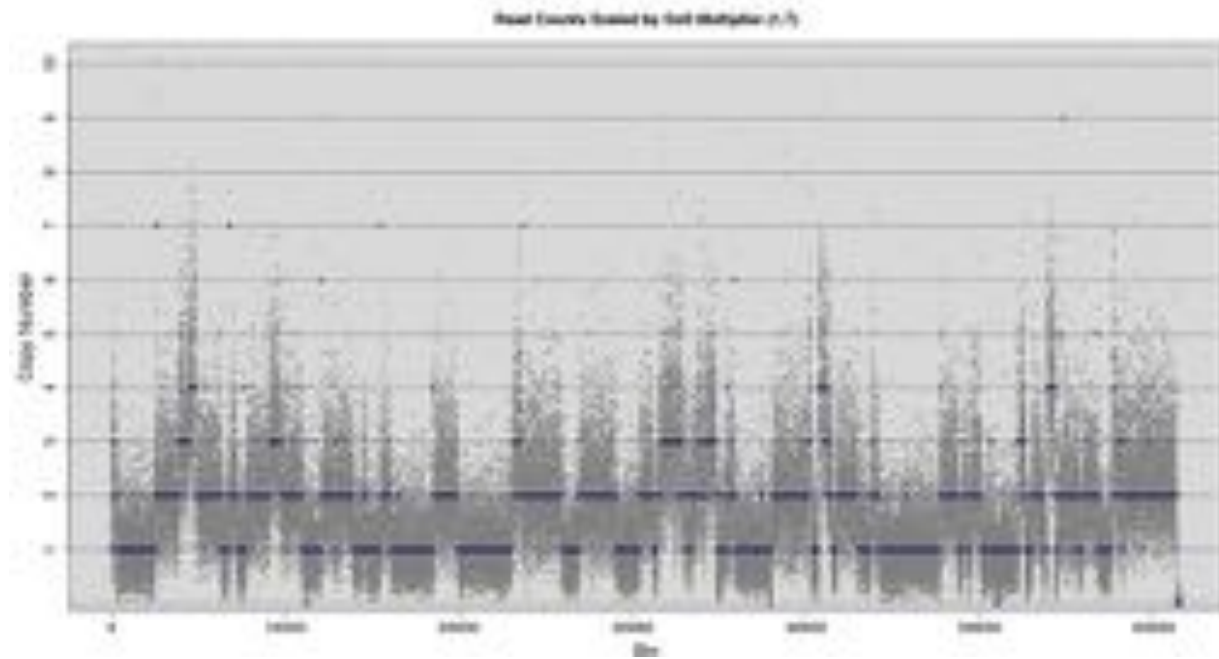
3) Segmentation



Circular Binary Segmentation (CBS)

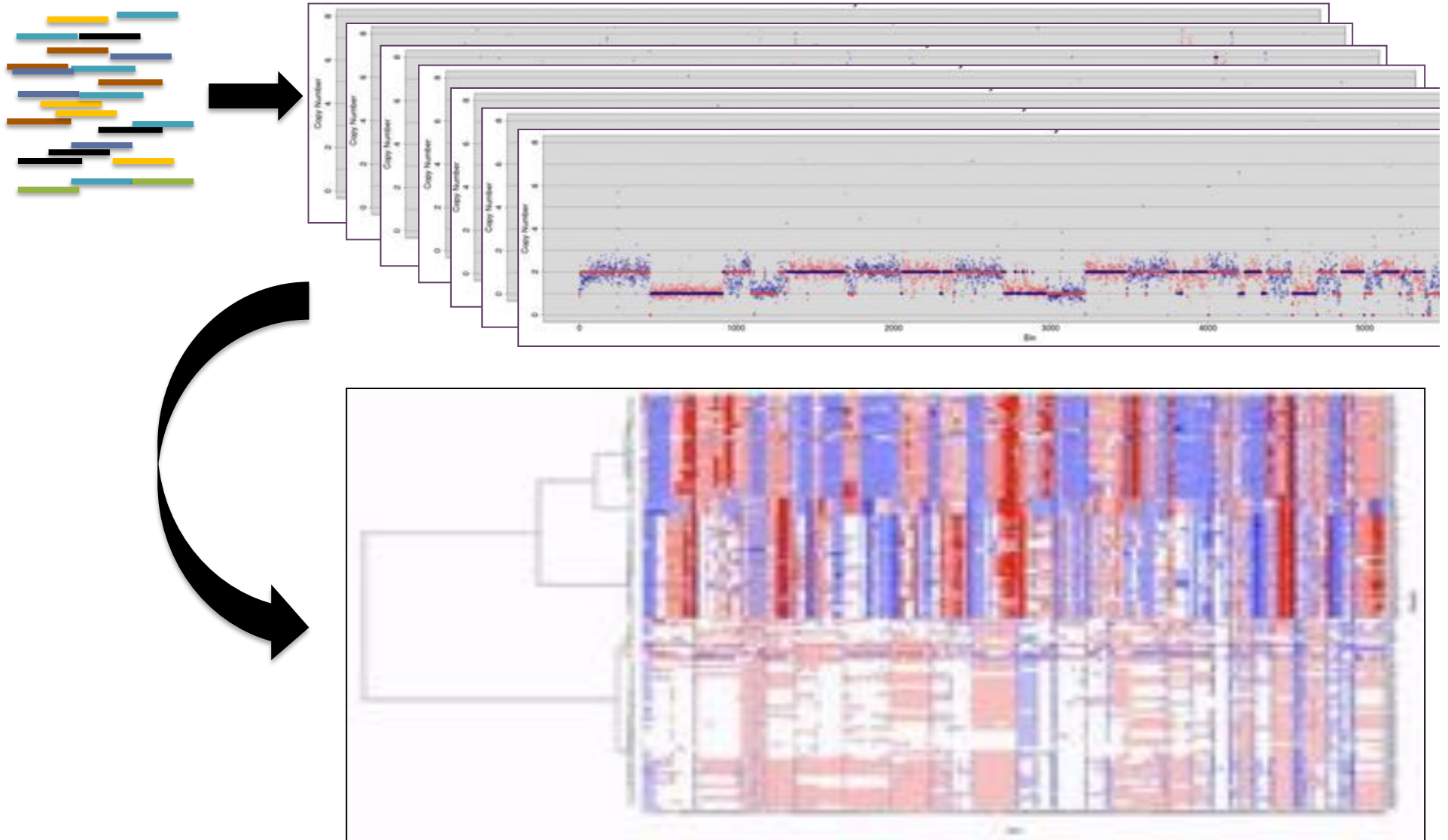


4) Estimating Copy Number



$$CN = \operatorname{argmin} \left\{ \sum_{i,j} (\hat{Y}_{i,j} - Y_{i,j})^2 \right\}$$

5) Cells to Populations



Ginkgo

<http://qb.cshl.edu/ginkgo>



Interactive Single Cell CNV analysis & clustering

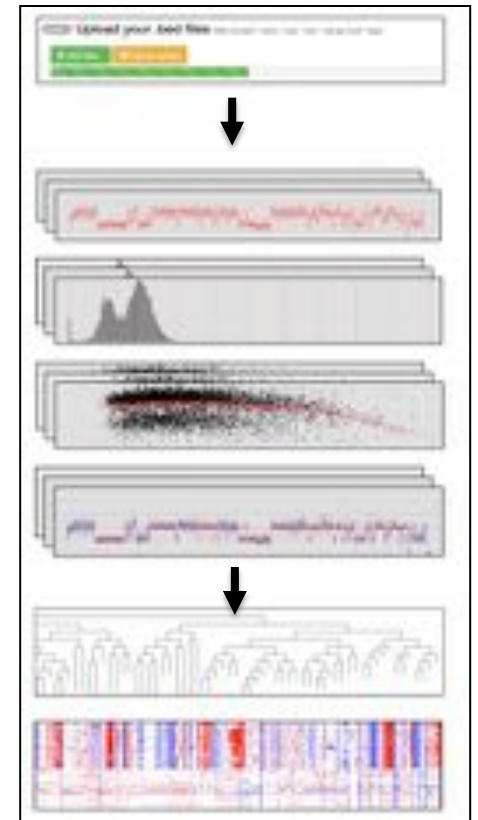
- Easy-to-use, web interface, parameterized for binning, segmentation, clustering, etc
- Per cell through project-wide analysis in any species

Compare MDA, DOP-PCR, and MALBAC

- DOP-PCR shows superior resolution and consistency

Available for collaboration

- Analyzing CNVs with respect to different clinical outcomes
- Extending clustering methods, prototyping scRNA

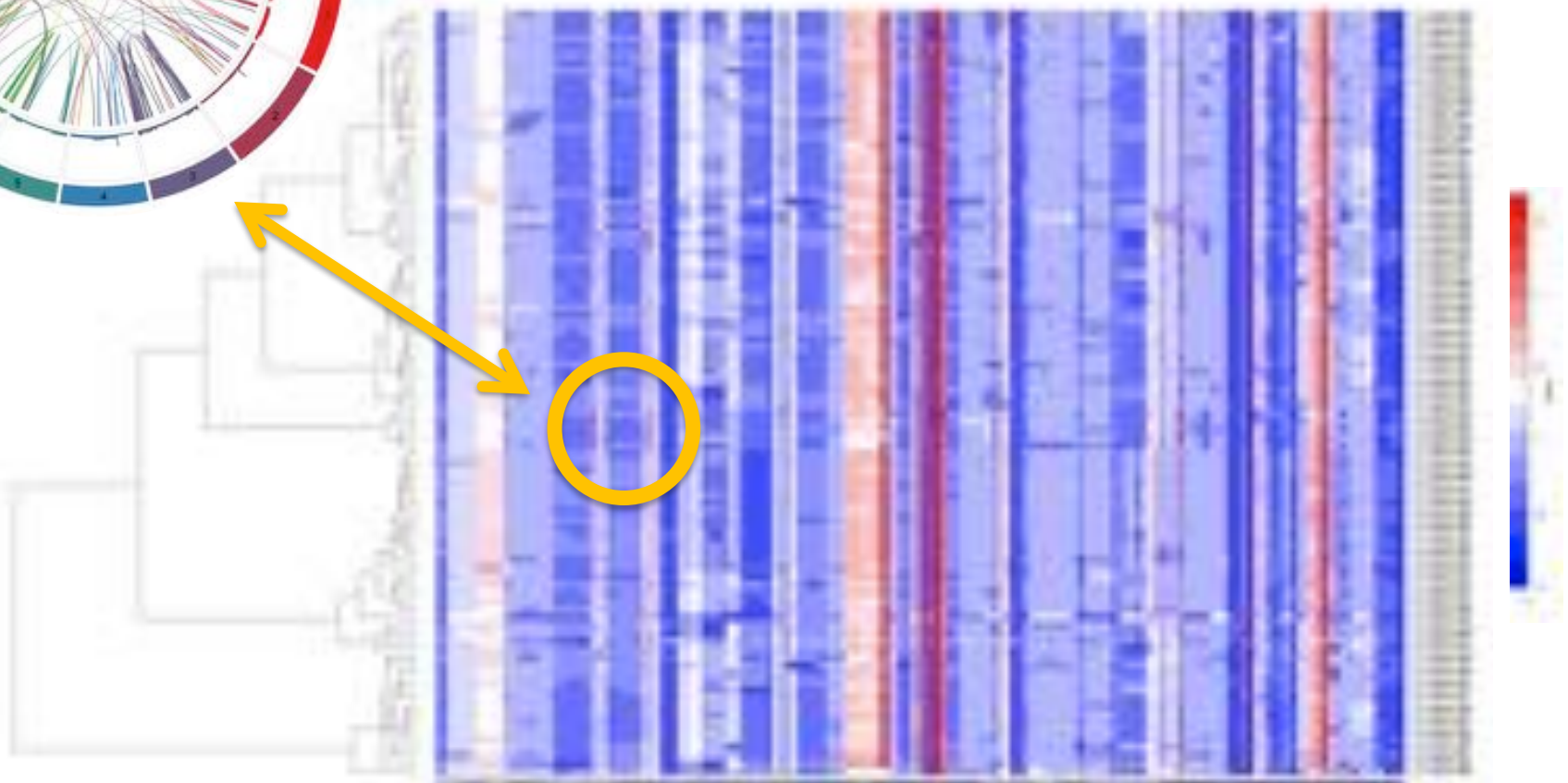
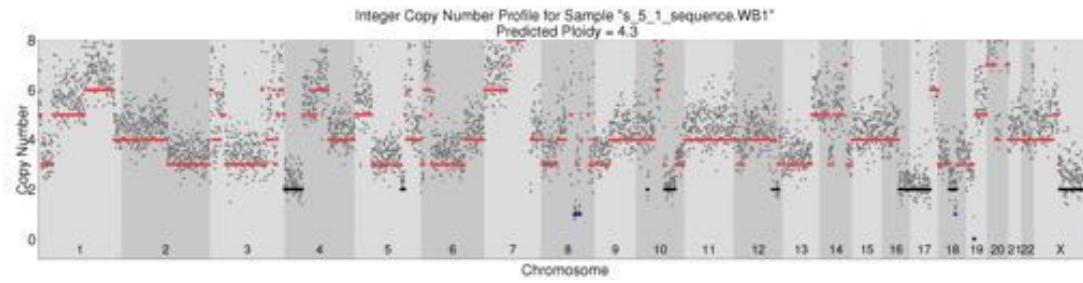


Interactive analysis and assessment of single-cell copy-number variations.

Garvin T, Aboukhalil R, Kendall J, Baslan T, Atwal GS, Hicks J, Wigler M, Schatz MC (2015)

Nature Methods doi:10.1038/nmeth.3578

CNVs in 100 SK-BR-3 Cells



Understanding Genome Structure & Function

Single Molecule Sequencing

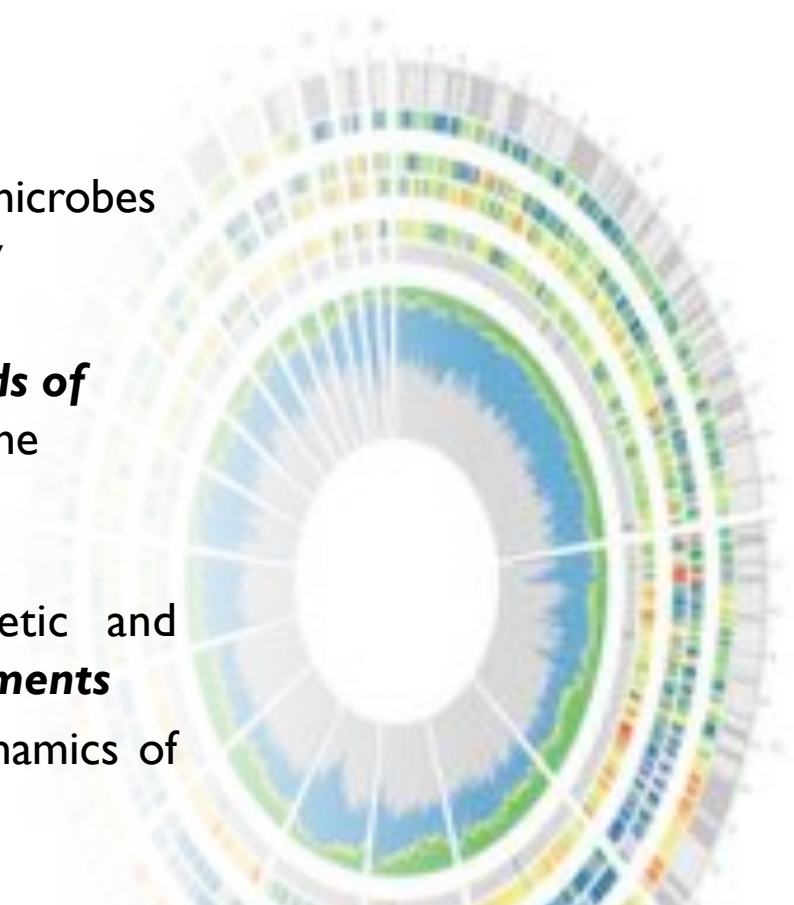
- Now have the ability to **perfectly assemble** microbes and many small eukaryotes, **reference quality** assemblies of larger eukaryotes
- Using this technology to find **10s of thousands of novel structural variations** per human genome

Single Cell Sequencing

- Exciting technologies to probe the genetic and molecular **composition of complex environments**
- We have only begun to explore the rich dynamics of genomes, transcriptomes, and epigenomics

These advances give us incredible power to study how genomes mutate and evolve

With several new biotechnologies in hand, we are now largely limited only by our quantitative power to make comparisons and find patterns



Acknowledgements

Schatz Lab

Rahul Amin
Han Fang
Tyler Gavin
James Gurtowski
Hayan Lee
Zak Lemmon
Giuseppe Narzisi
Maria Nattestad
Aspyn Palatnick
Srividya
Ramakrishnan
Fritz Sedlazeck
Rachel Sherman
Greg Vurture
Alejandro Wences

CSHL

Hannon Lab
Gingeras Lab
Jackson Lab
Hicks Lab
Iossifov Lab
Levy Lab
Lippman Lab
Lyon Lab
Martienssen Lab
McCombie Lab
Tuveson Lab
Ware Lab
Wigler Lab

SBU

Skiena Lab
Patro Lab

Cornell

Susan McCouch
Lyza Maron
Mark Wright

OICR

John McPherson
Karen Ng
Timothy Beck
Yogi Sundaravadanam

NYU

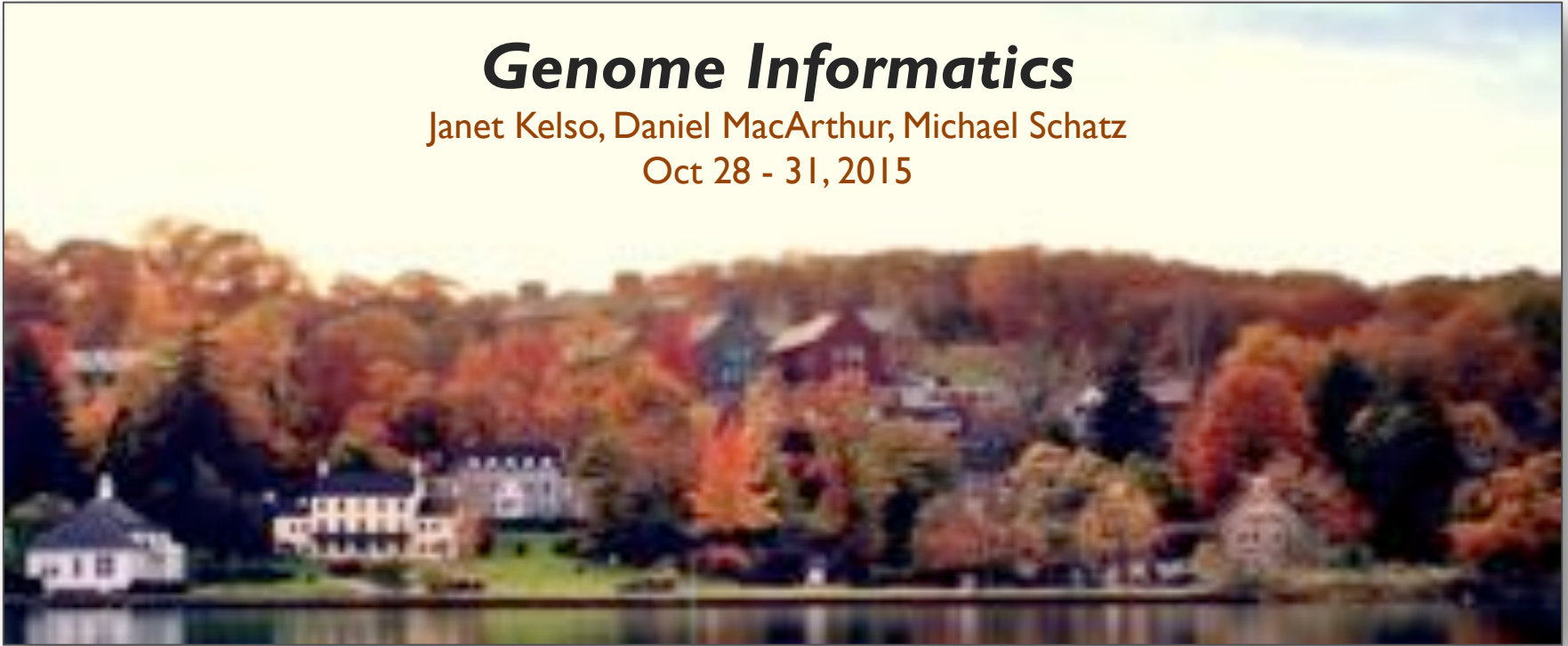
Jane Carlton
Elodie Ghedin



Genome Informatics

Janet Kelso, Daniel MacArthur, Michael Schatz

Oct 28 - 31, 2015



Thank you

<http://schatzlab.cshl.edu>

@mike_schatz