

Algorithms for single cell and single molecule biology

Michael Schatz

March 27, 2015

Biotech Symposium / Simons Foundation





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

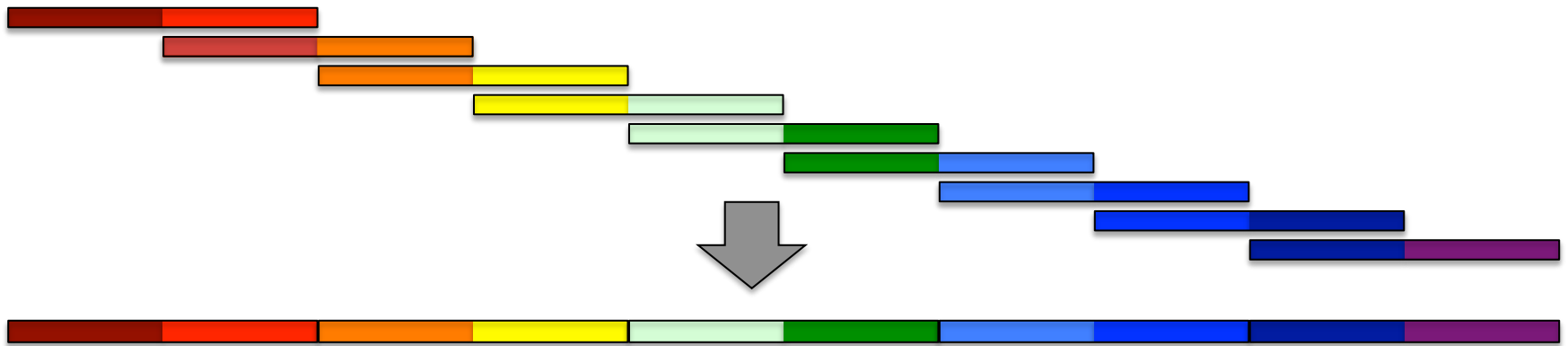
1. Shear & Sequence DNA



2. Construct assembly graph from overlapping reads

...AGCCTAGGGATGCGCGACACGT
GGATGCGCGACACGTCGCATATCCGTTTGGTCAACCTCGGACGGAC
CAACCTCGGACGGACCTCAGCGAA...

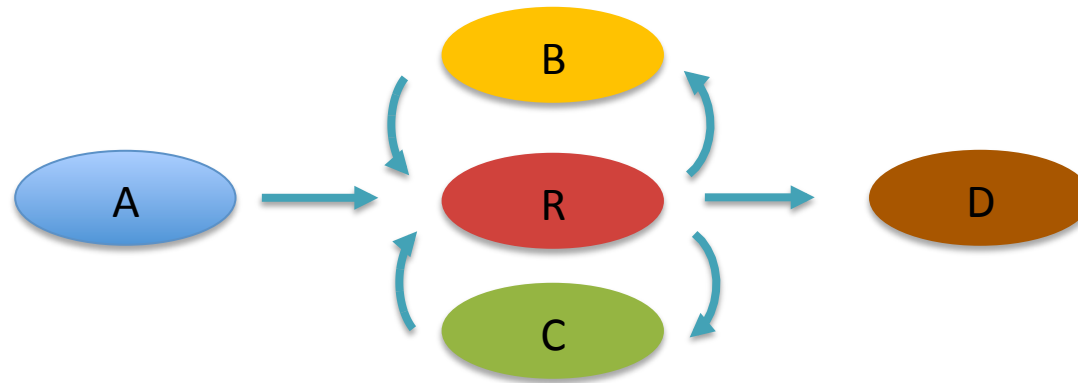
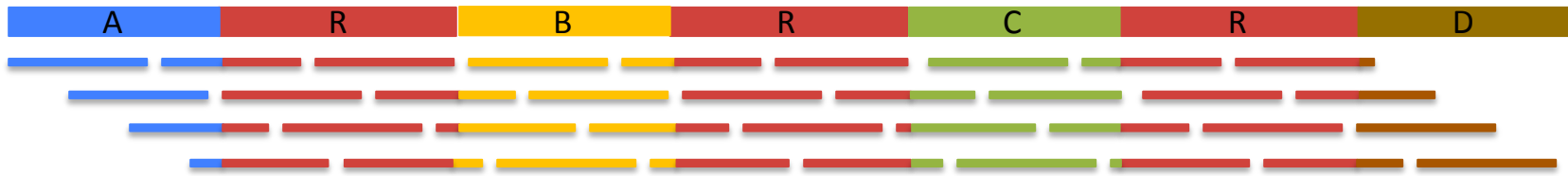
3. Simplify assembly graph



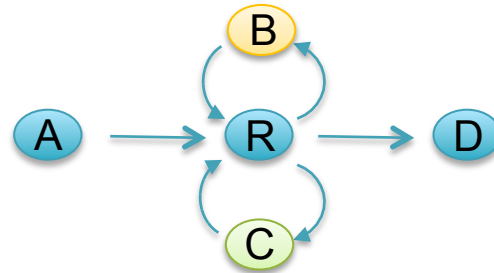
On Algorithmic Complexity of Biomolecular Sequence Assembly Problem

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

Assembly Complexity



Counting Eulerian Tours



AR**B**RCRD
or
ARC**R**BRD

Often an astronomical number of possible assemblies

- Value computed by application of the BEST theorem (Hutchinson, 1975)

$$\mathcal{W}(G, t) = (\det L) \left\{ \prod_{u \in V} (r_u - 1)! \right\} \left\{ \prod_{(u,v) \in E} a_{uv}! \right\}^{-1}$$

$L = n \times n$ matrix with $r_u - a_{uu}$ along the diagonal and $-a_{uv}$ in entry uv

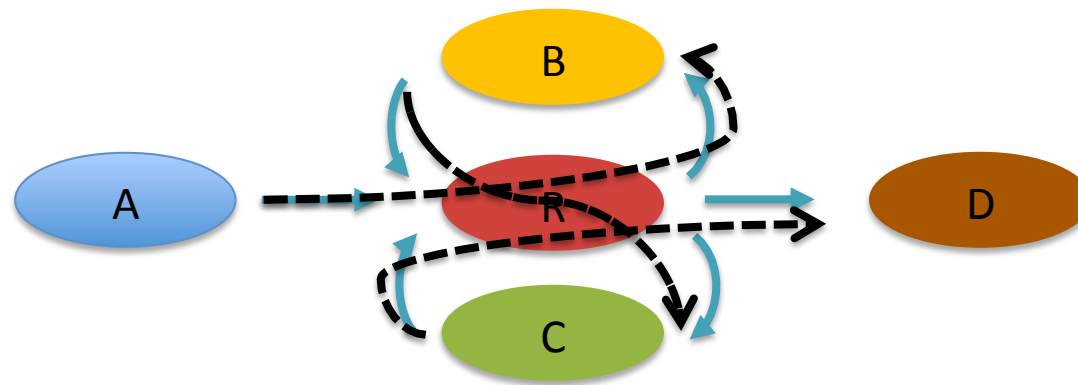
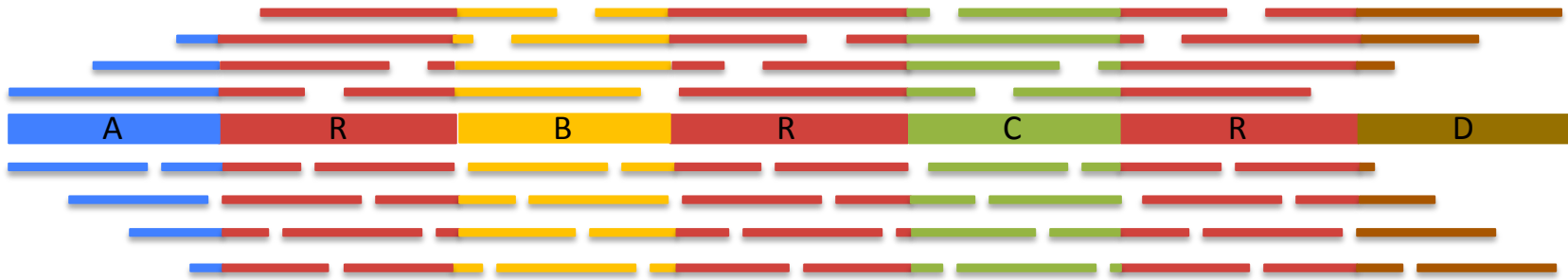
$r_u = d^+(u) + 1$ if $u=t$, or $d^+(u)$ otherwise

a_{uv} = multiplicity of edge from u to v

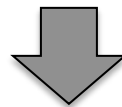
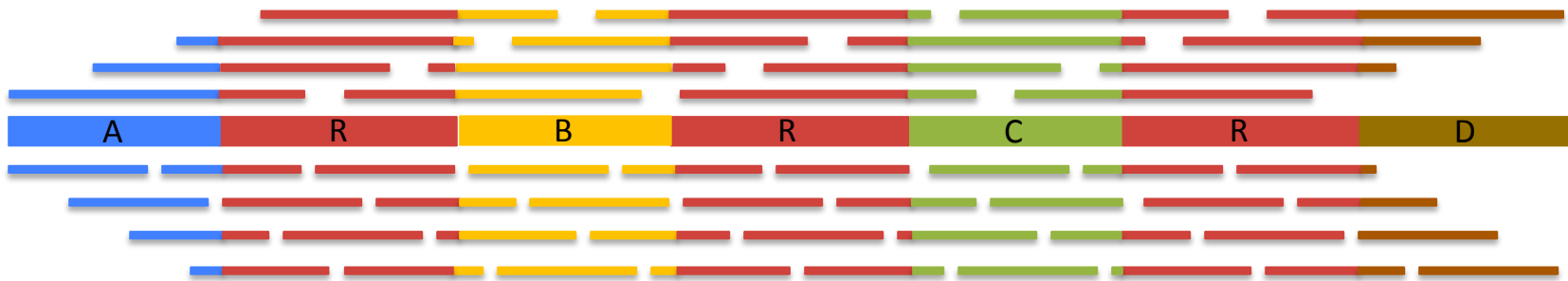
Assembly Complexity of Prokaryotic Genomes using Short Reads.

Kingsford C, Schatz MC, Pop M (2010) *BMC Bioinformatics*. 11:21.

Assembly Complexity



Assembly Complexity



The advantages of SMRT sequencing

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

N50 size

Def: 50% of the genome is in contigs as large as the N50 value

Example: 1 Mbp genome 50%



N50 size = 30 kbp

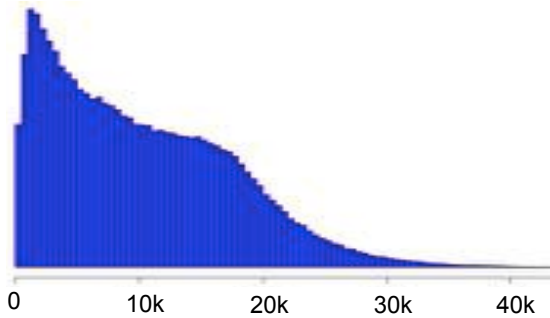
$(300\text{k} + 100\text{k} + 45\text{k} + 45\text{k} + 30\text{k} = 520\text{k} \geq 500\text{kbp})$

A larger N50 is indicative of improvement in every dimension:

- Better resolution of genes and flanking regulatory regions
- Better resolution of transposons and other complex sequences
- Better resolution of chromosome organization

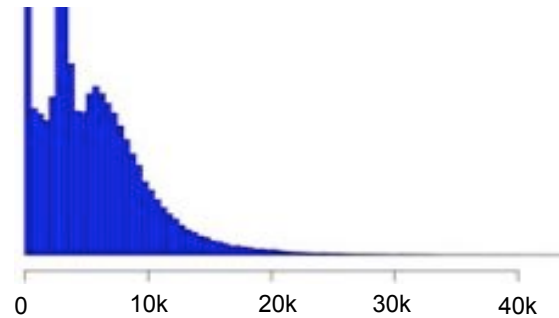
Single Molecule Sequencing

PacBio RS II



CSHL/PacBio

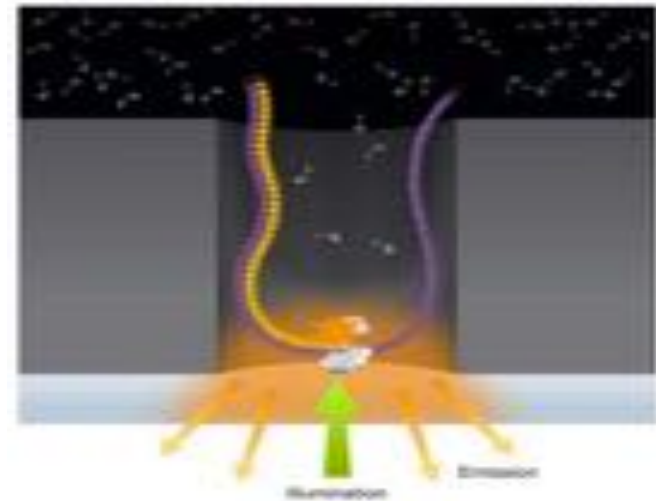
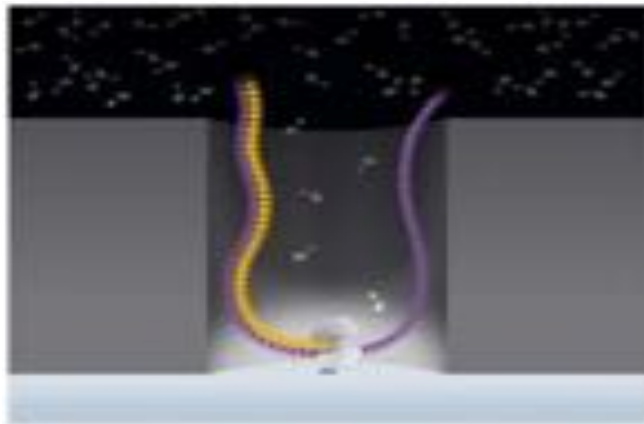
Oxford Nanopore



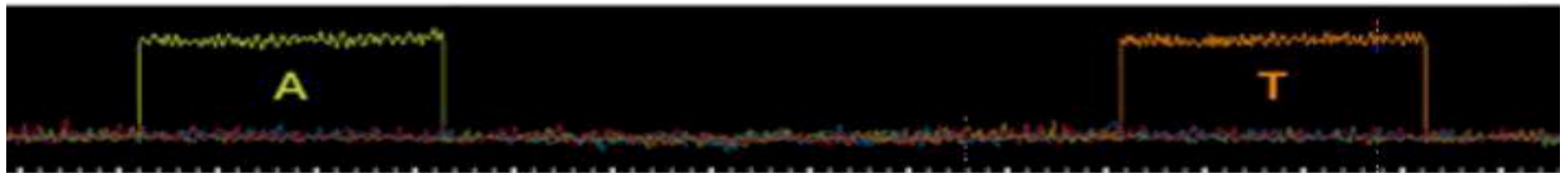
CSHL/ONT

PacBio SMRT Sequencing

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



Intensity



Time

http://www.pacificbiosciences.com/assets/files/pacbio_technology_background.pdf

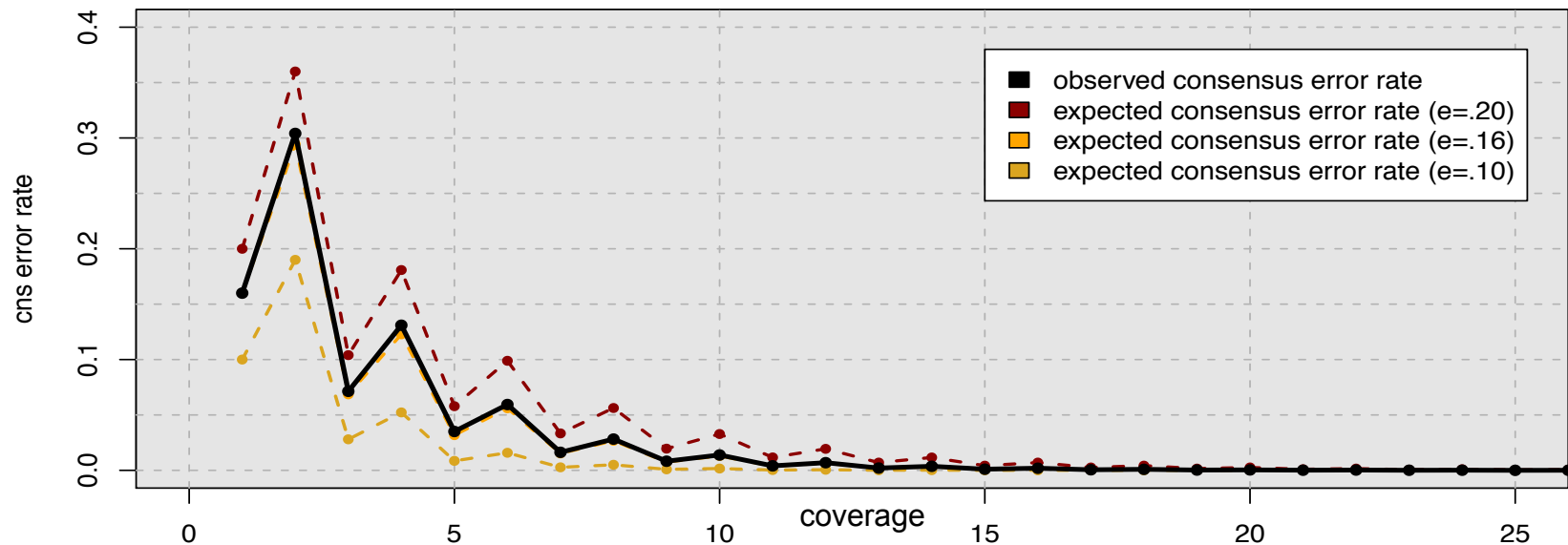
Single Molecule Sequences



“Corrective Lens” for Sequencing



Consensus Accuracy and Coverage



Coverage can overcome random errors

- Dashed: error model; Solid: observed accuracy

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

$$CNS\ Error = \sum_{i=\lceil c/2 \rceil}^c \binom{c}{i} (e)^i (1-e)^{n-i}$$

PacBio Assembly Algorithms

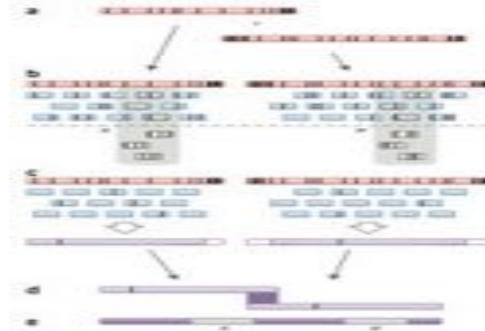
PBJelly



Gap Filling

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

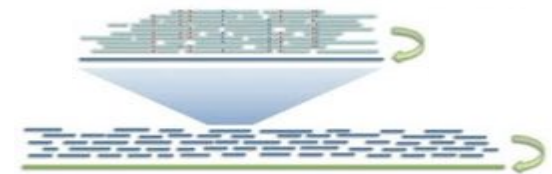
PacBioToCA & ECTools



Hybrid Error Correction

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

HGAP & 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

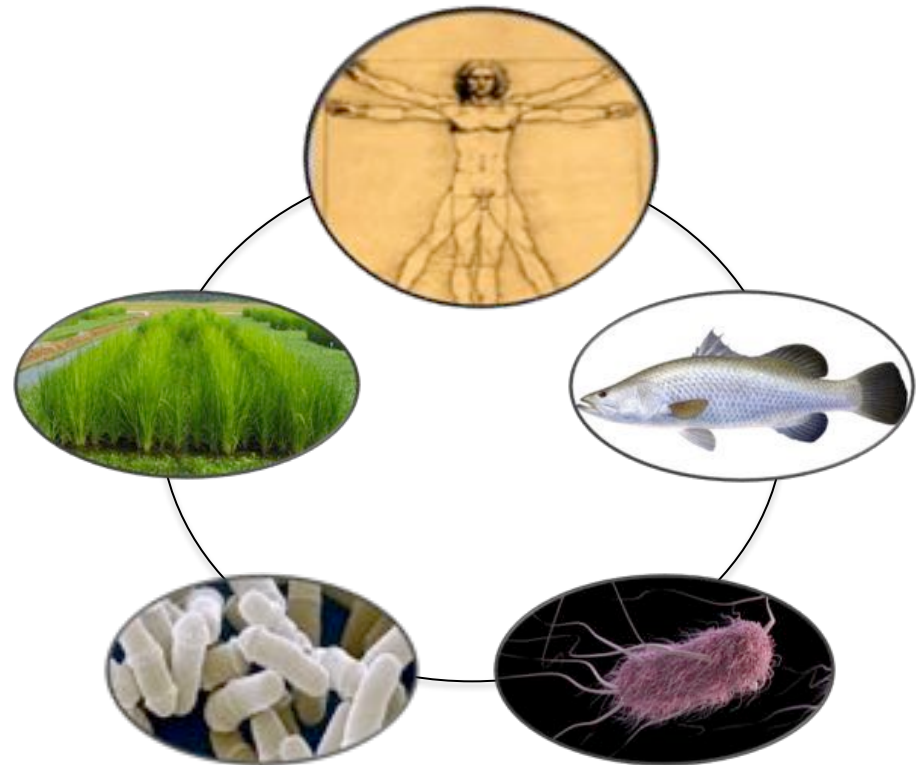
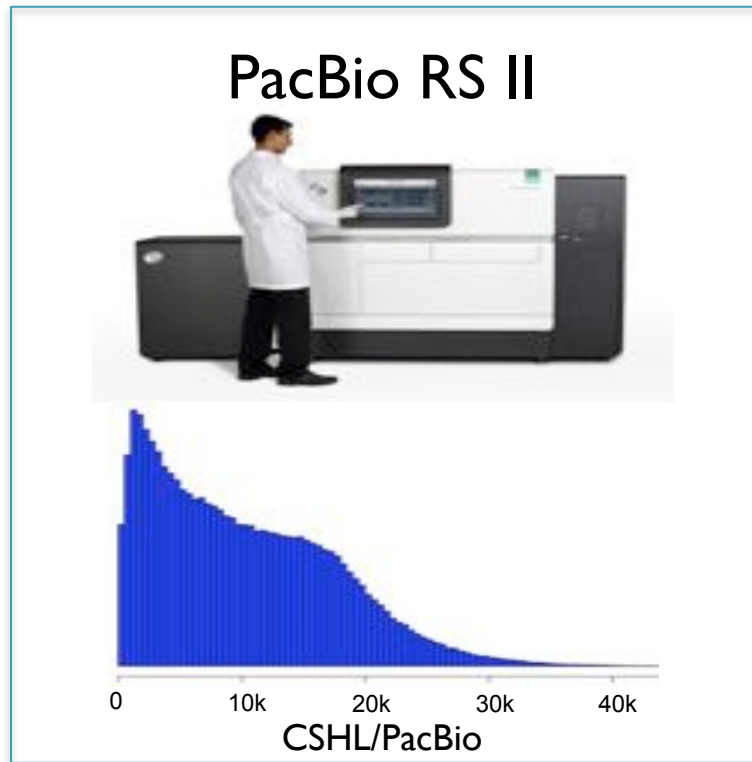
Chin *et al* (2013)
Nature Methods. 10:563–569

< 5x

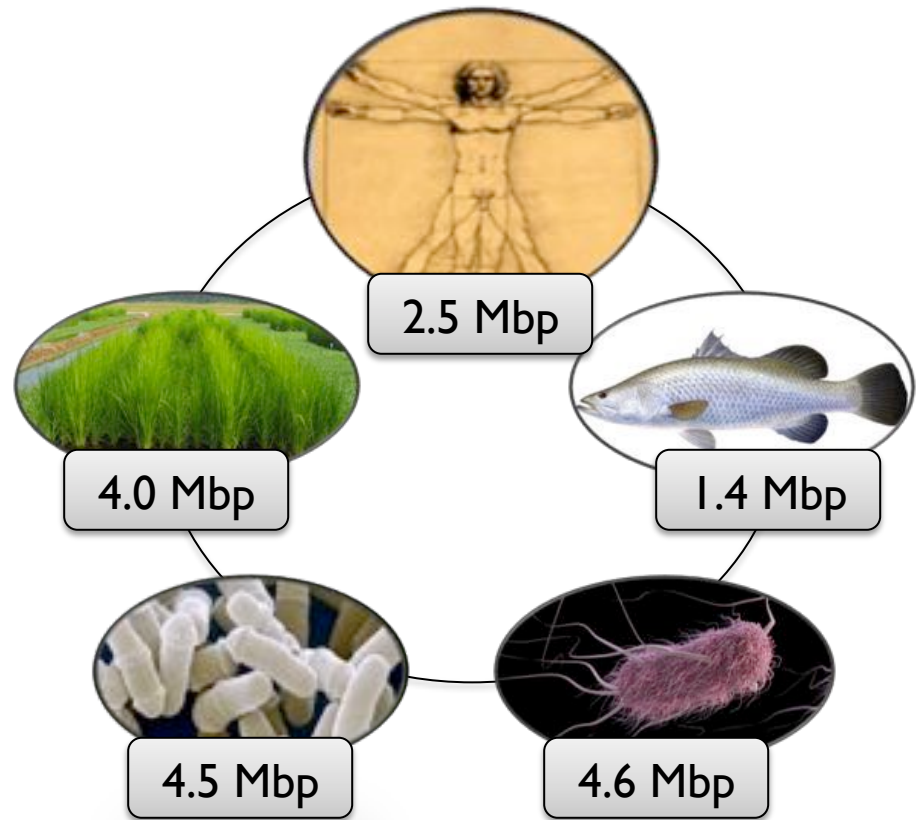
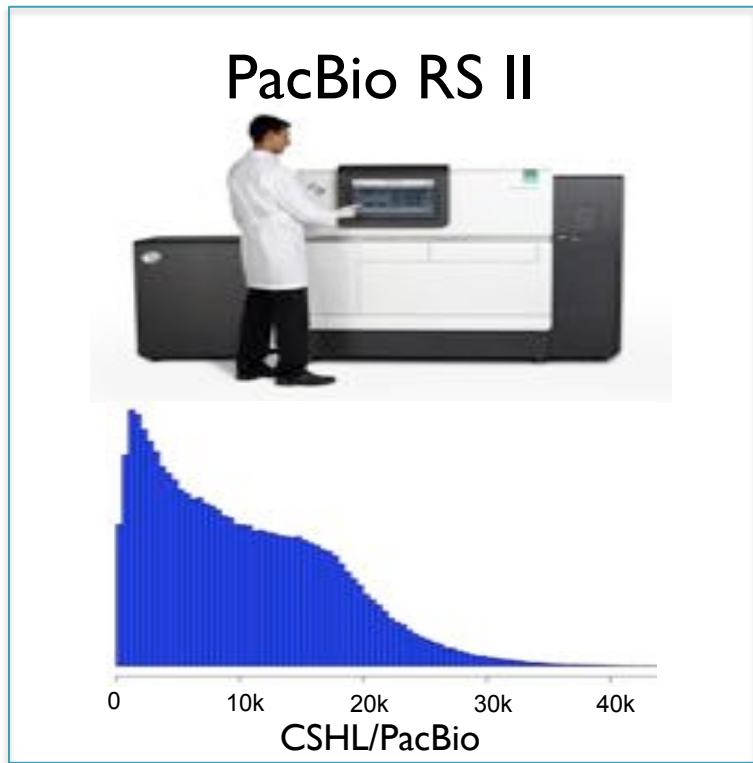
PacBio Coverage

> 50x

PacBio Sequencing



PacBio Sequencing



Her2 amplified breast cancer

Breast cancer

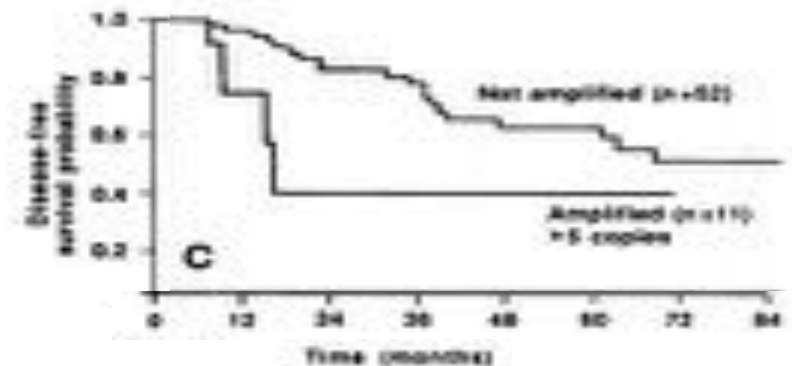
- About 12% of women will develop breast cancer during their lifetimes
- ~230,000 new cases every year (US)
- ~40,000 deaths every year (US)

Statistics from American Cancer Society and Mayo Clinic.

Recurrence and metastasis from Gonzalez-Angulo, et al, 2009.

Her2 amplified breast cancer

- 20% of breast cancers
- 2-3X recurrence risk
- 5X metastasis risk



(Adapted from Slamon et al, 1987)

SK-BR-3

Most commonly used Her2+ breast cancer cell line

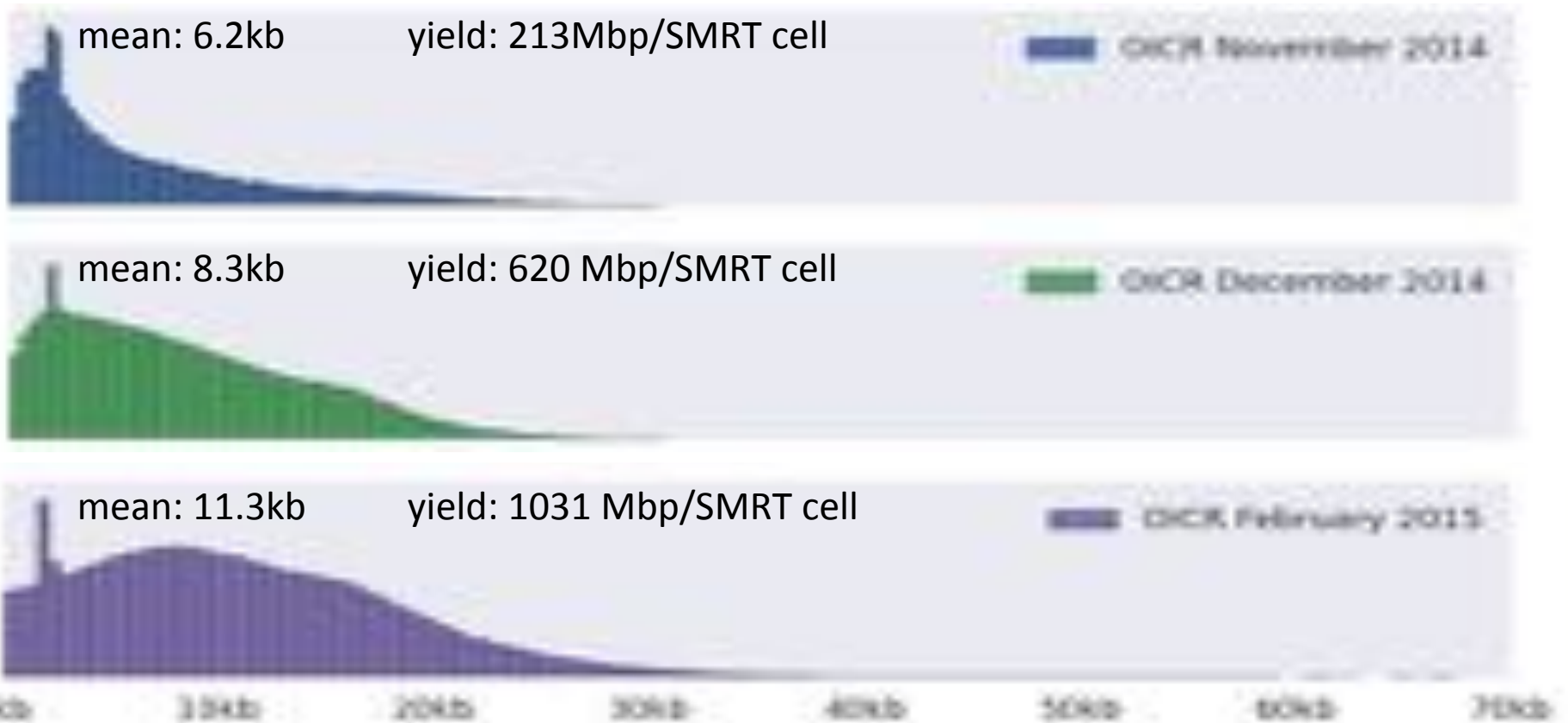


(Davidson et al, 2000)

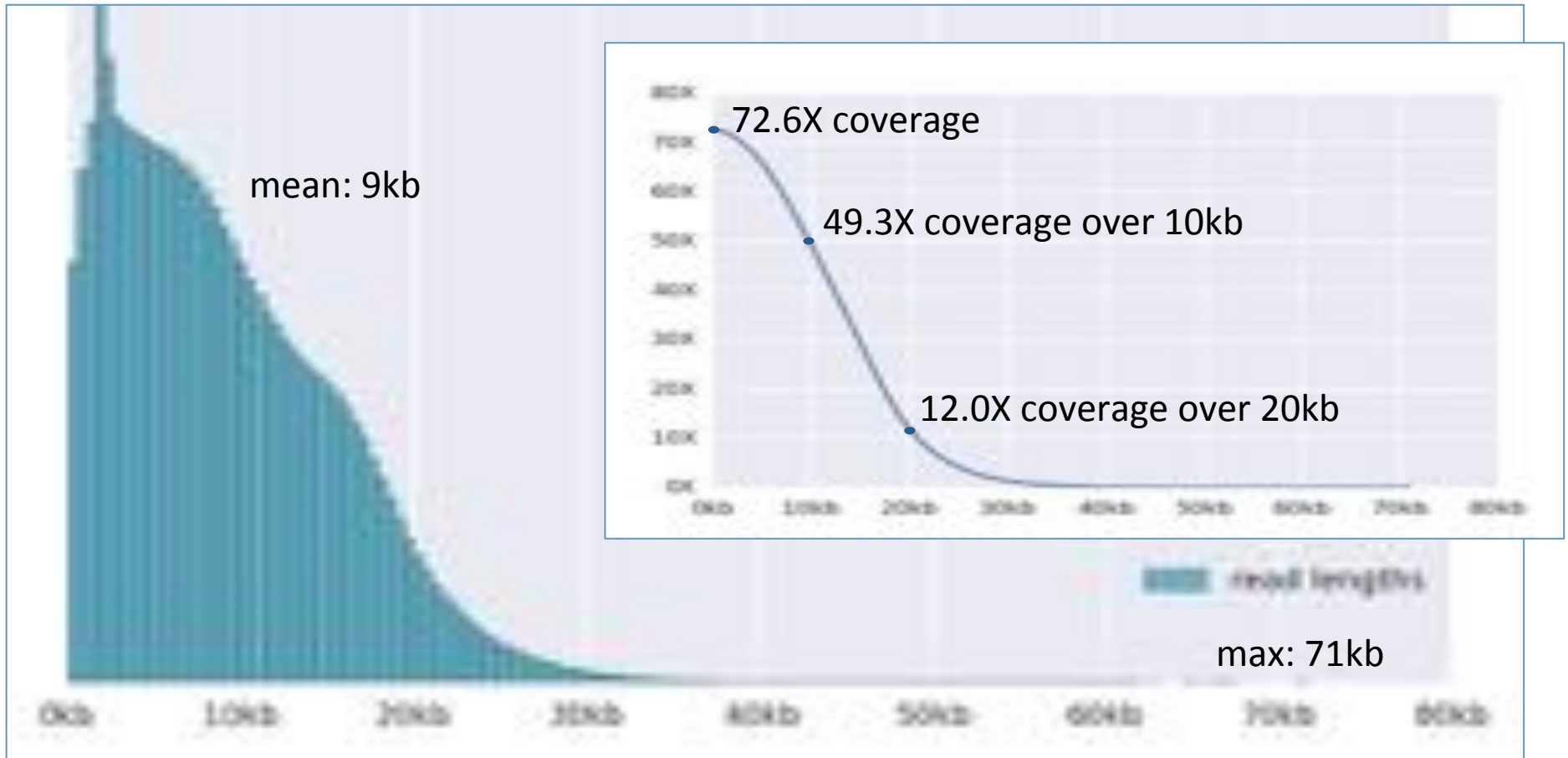
Can we resolve the complex structural variations, especially around Her2?

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

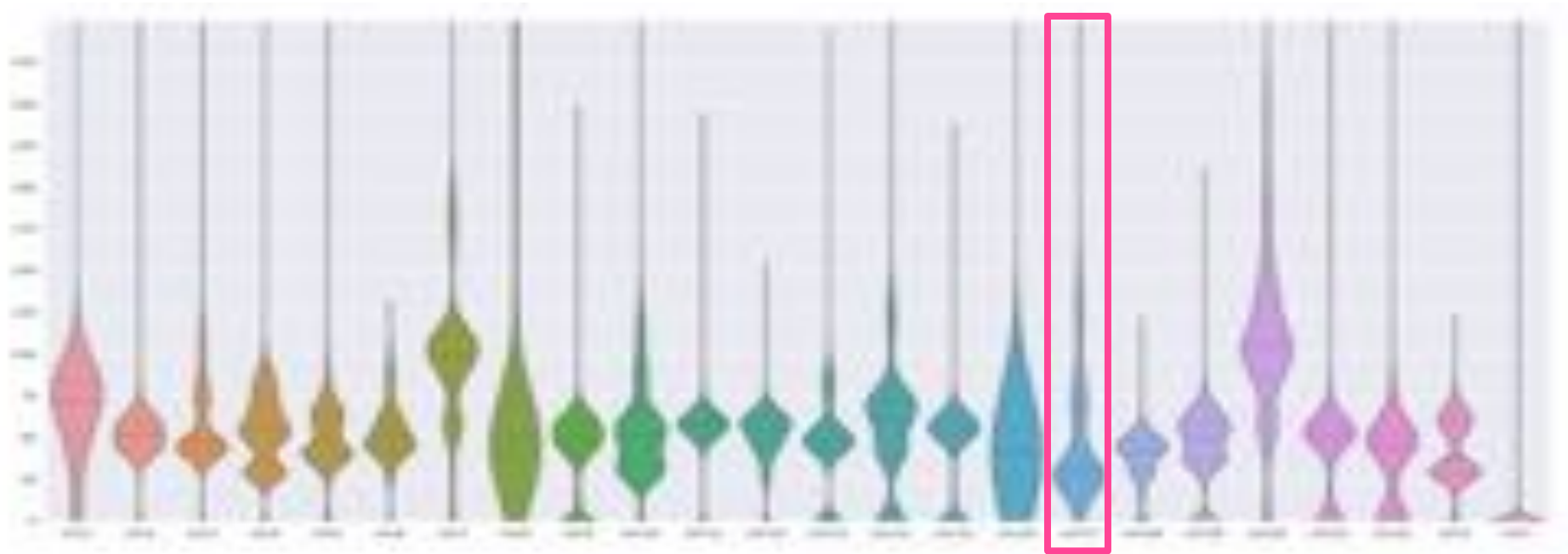
Improving SMRTcell Performance



PacBio read length distribution

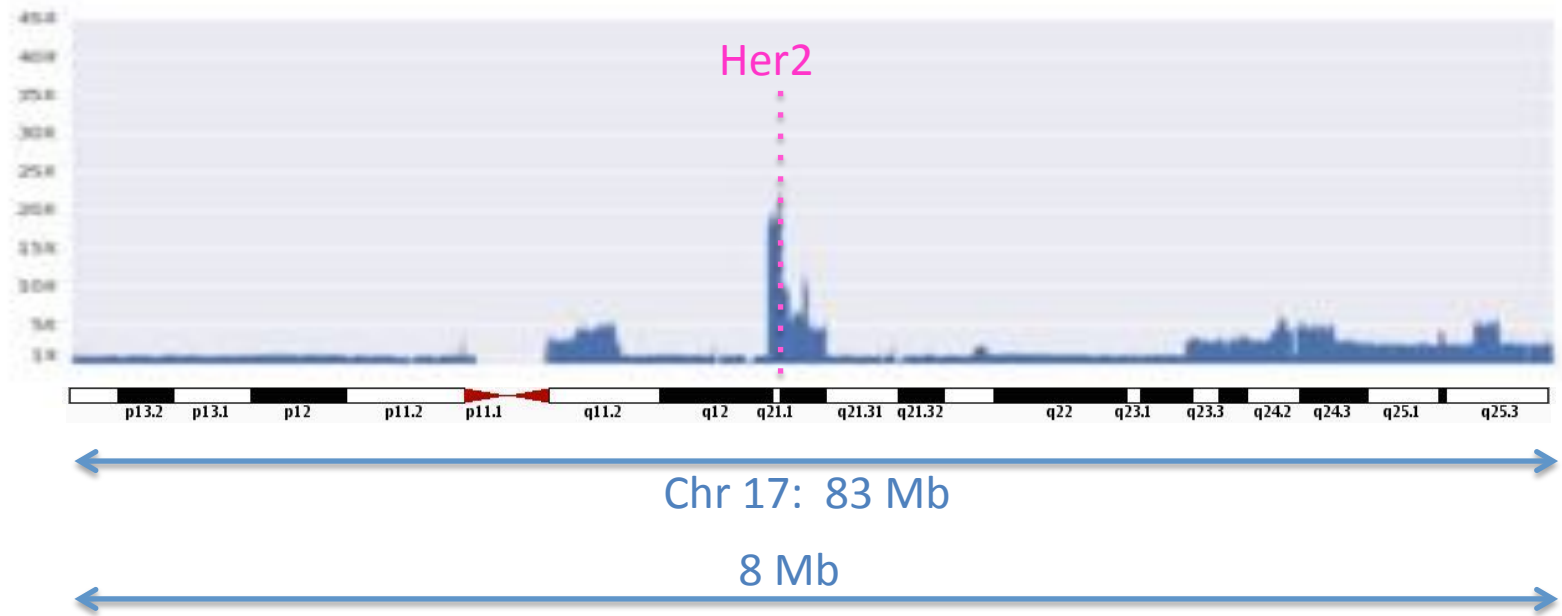


Genome-wide alignment coverage

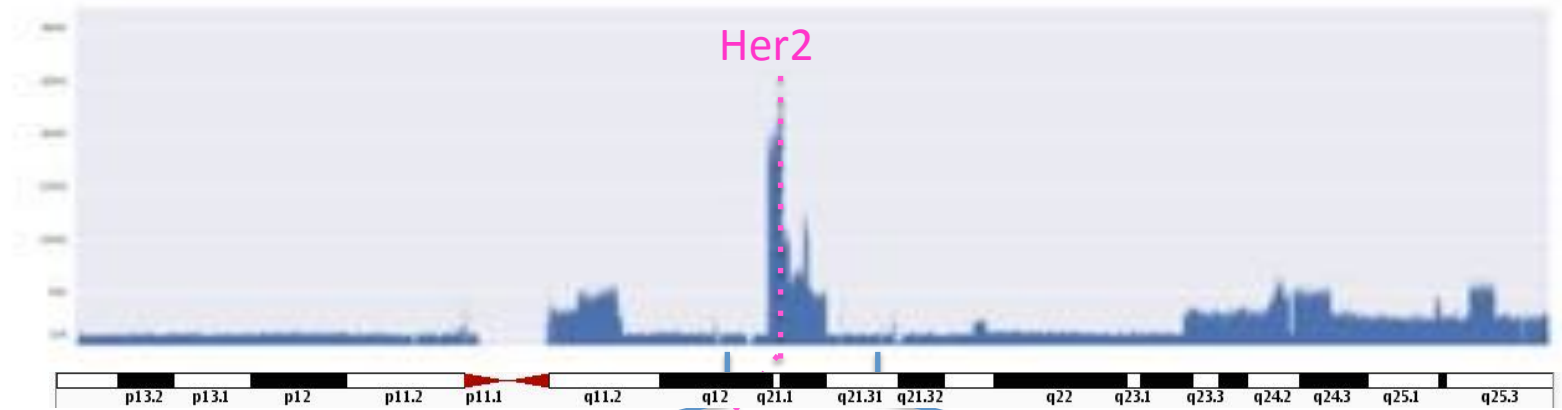


Genome-wide coverage averages around 54X
Coverage per chromosome greatly varies

PacBio

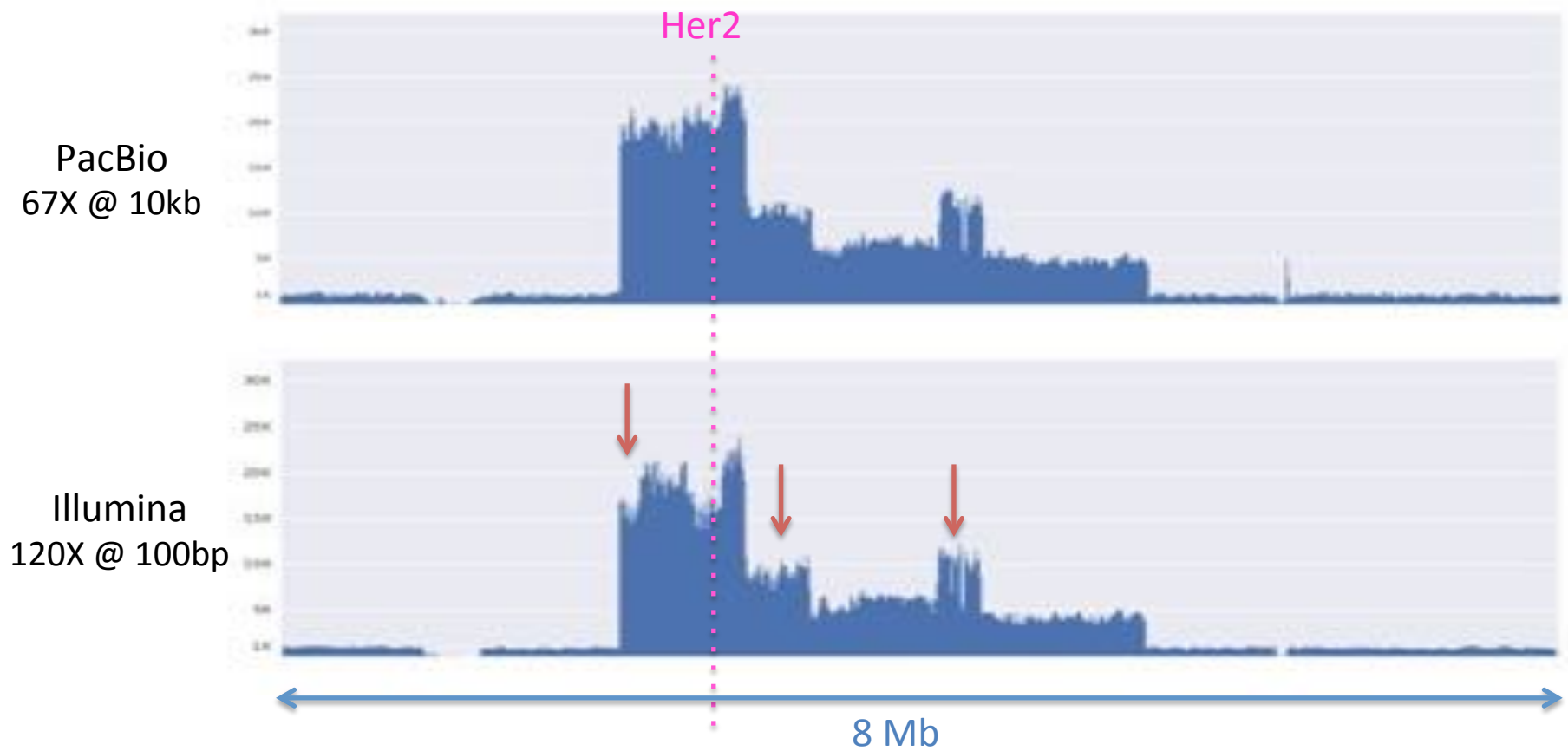


PacBio



PacBio





PacBio and Illumina coverage values are highly correlated but Illumina shows greater variance because of poorly mapping reads

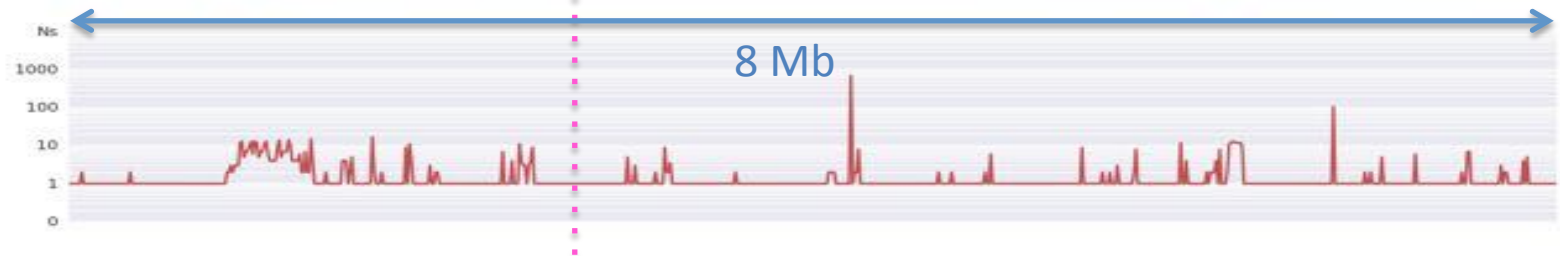
PacBio
67X @ 10kb



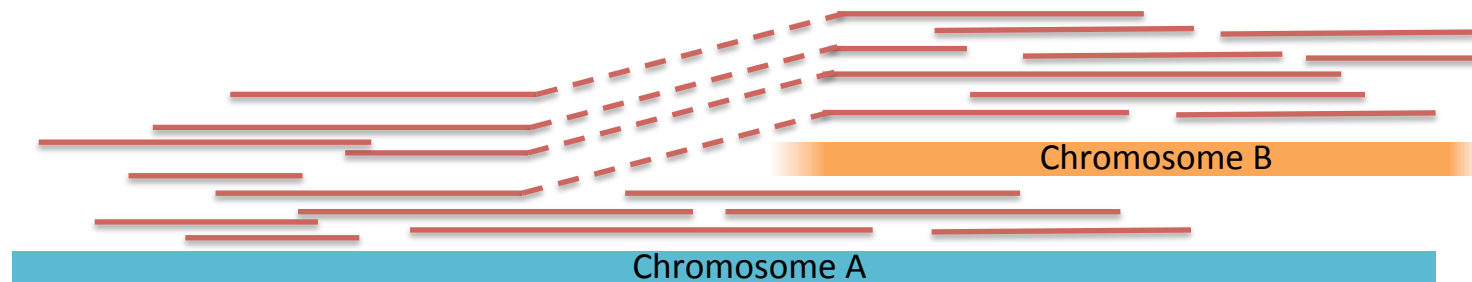
Illumina
120X @ 100bp



Repeats
21-mers



Structural variant discovery with long reads



1. Alignment-based split read analysis: Efficient capture of most events

BWA-MEM + Lumpy

2. Local assembly of regions of interest: In-depth analysis with *base-pair precision*

Localized HGAP + Celera Assembler + MUMmer

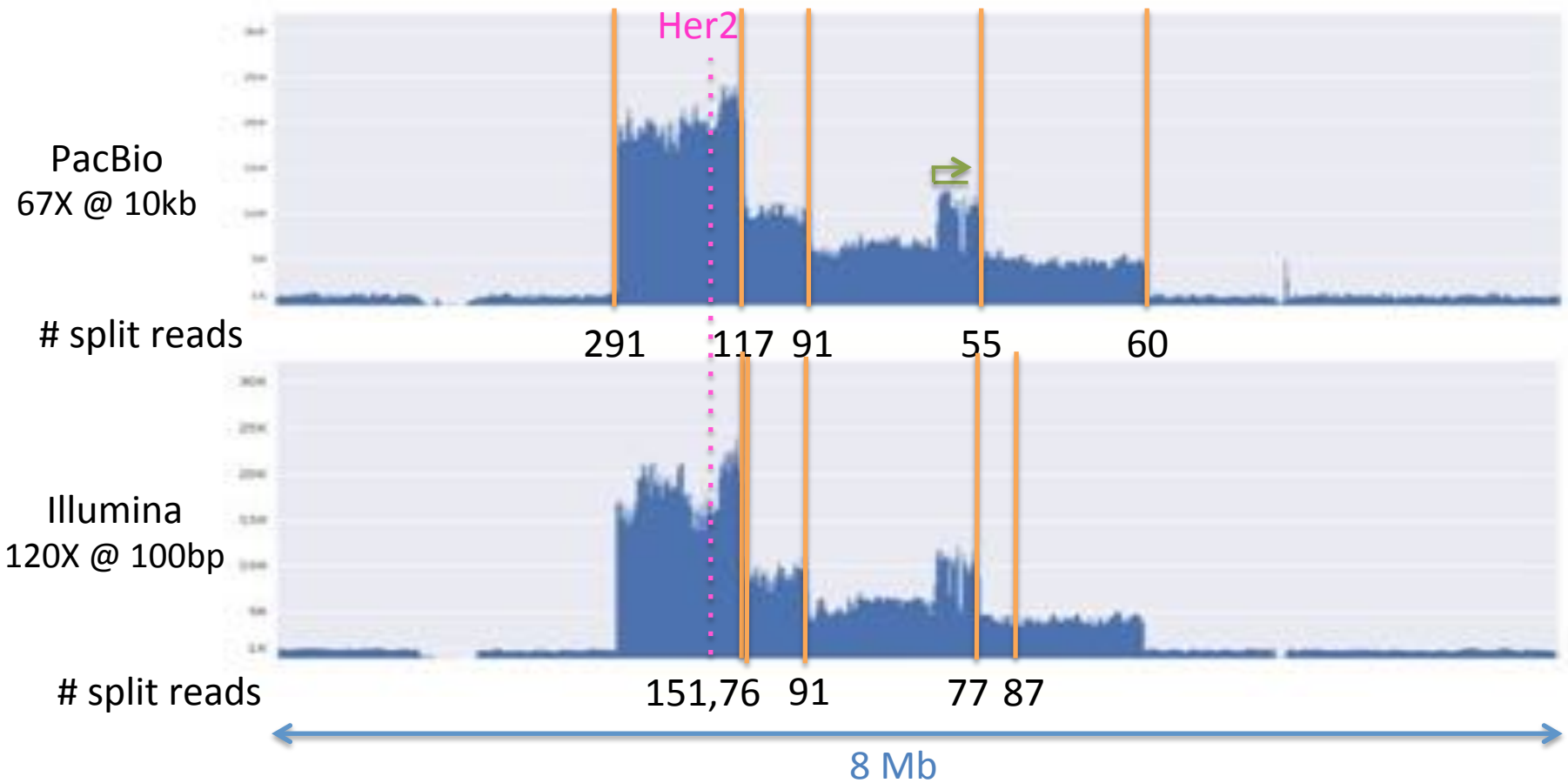
3. Whole genome assembly: In-depth analysis including *novel sequences*

DNAnexus-enabled version of Falcon

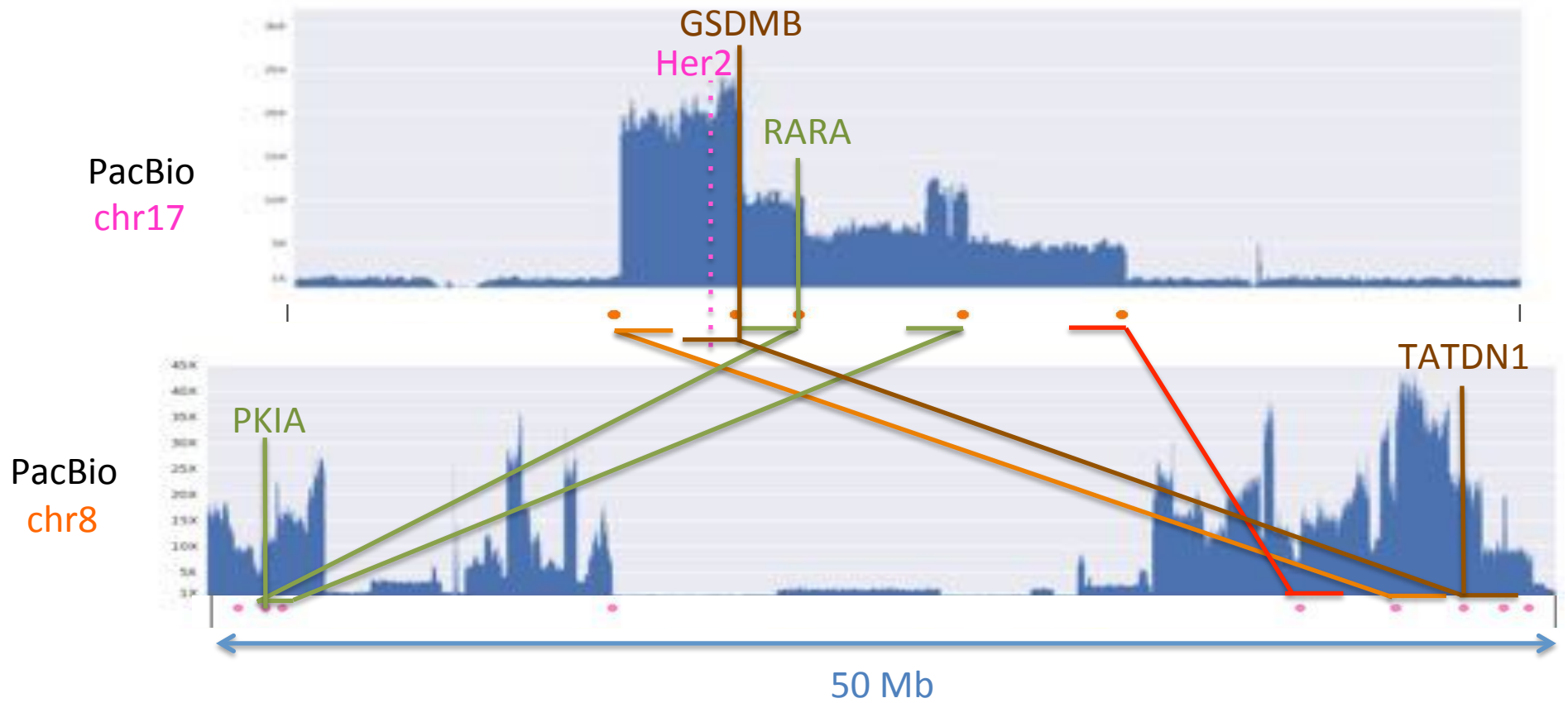
Total Assembly: 2.64Gbp

Contig N50: 2.56 Mbp

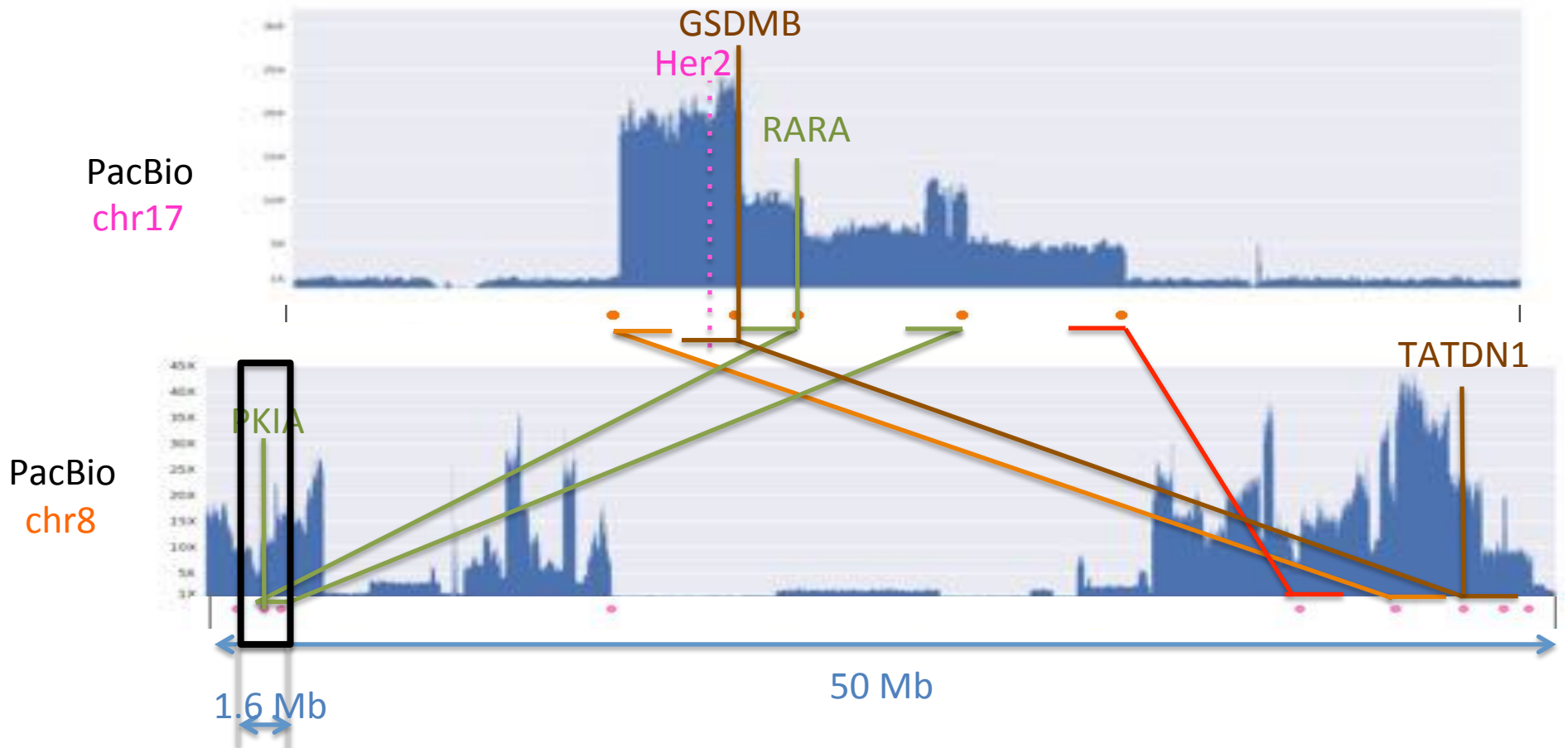
Max Contig: 23.5Mbp



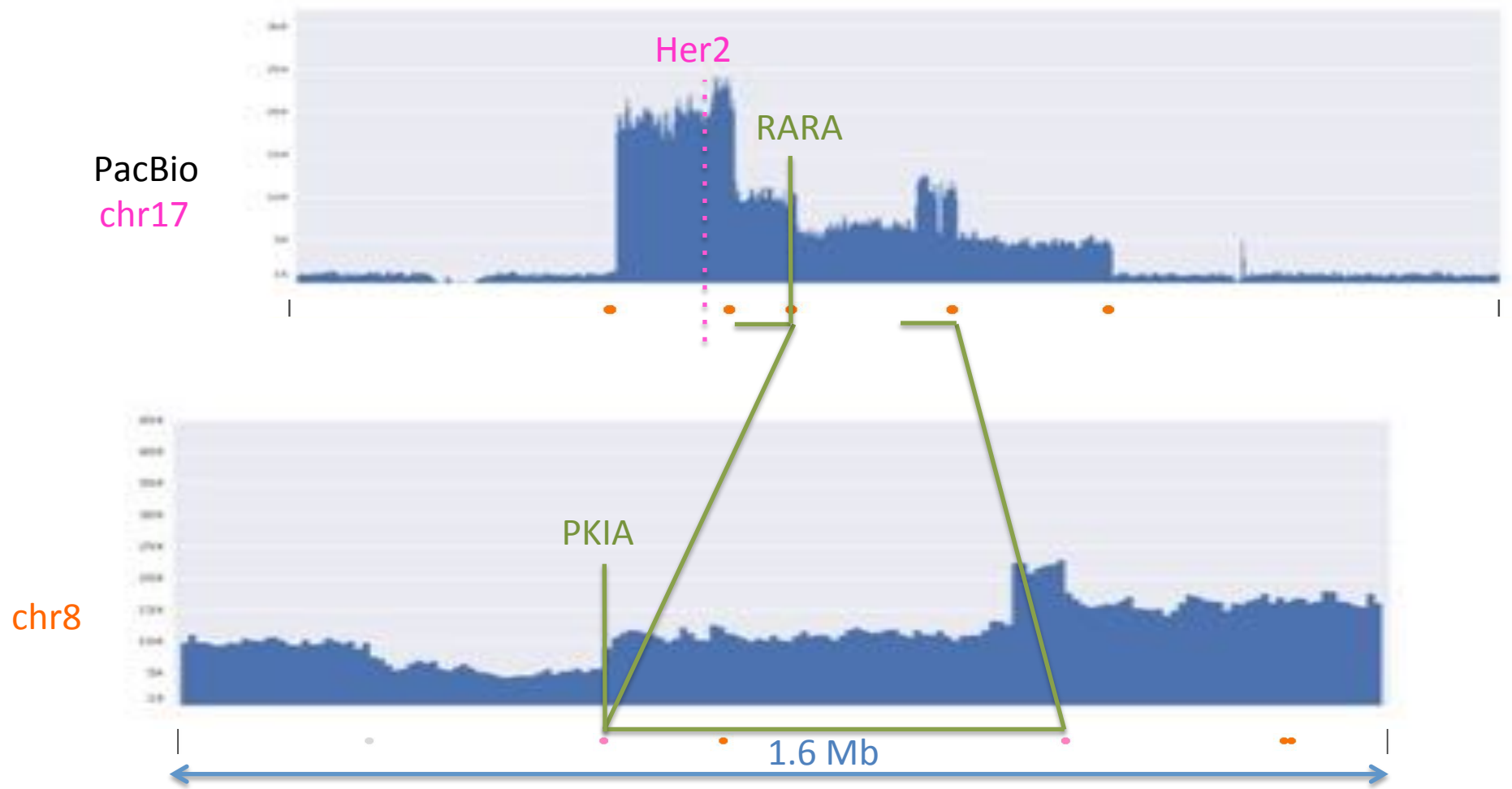
Green arrow indicates an inverted duplication.
False positive and missing Illumina calls due to mis-mapped reads (especially low complexity).



Confirmed both known gene fusions in this region

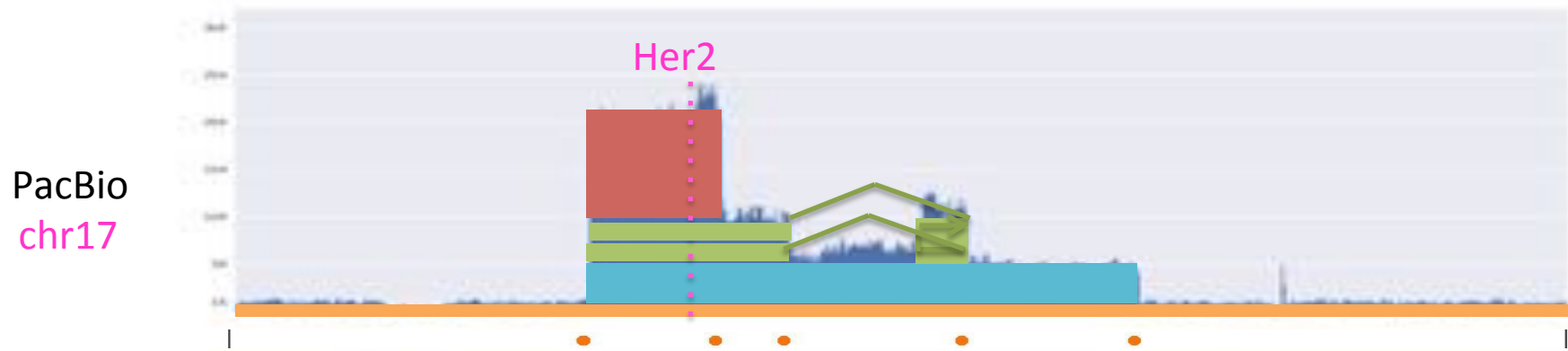


Confirmed both known gene fusions in this region



Joint coverage and breakpoint analysis to discover underlying events

Cancer lesion Reconstruction



By comparing the proportion of reads that are spanning or split at breakpoints we can begin to infer the history of the genetic lesions.

1. Healthy diploid genome
2. Original translocation into chromosome 8
3. Duplication, inversion, and inverted duplication within chromosome 8
4. Final duplication from within chromosome 8

SKBR3 Oncogene Analysis

Known missense mutation in p53: **R175H**

Arg

Reference ATCTGAGCAGCGCTCATGGTGGGGGCAGCGCCTCACAACCTCCGTCATGTGCTGTGACTGCTT
Illumina ATCTGAGCAGCGCTCATGGTGGGGGCAG**T**GCCTCACAACCTCCGTCATGTGCTGTGACTGCTT
PacBio ATCTGAGCAGCGCTCATGGTGGGGGCAG**T**GCCTCACAACCTCCGTCATGTGCTGTGACTGCTT

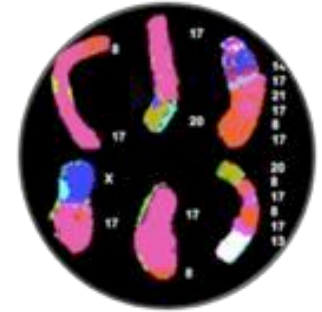
His

Oncogene amplifications	
ErbB2 (Her2)	≈20X
MYC	≈27X
MET	≈8X

**Genetic Lesion
History Analysis
Underway**

Known Gene fusions		Confirmed by PacBio reads?
TATDN1	GSDMB	Yes
RARA	PKIA	Yes
ANKHD1	PCDH1	Yes
CCDC85C	SETD3	Yes
SUMF1	LRRFIP2	Yes
WDR67 (TBC1D31)	ZNF704	Yes
DHX35	ITCH	Yes
NFS1	PREX1	Yes *read-through transcription
CYTH1	EIF3H	Yes *nested inside 2 translocations

SK-BR-3 Her2+ Breast Cancer Reference Genome



Released all data pre-publication to accelerate breast cancer research:

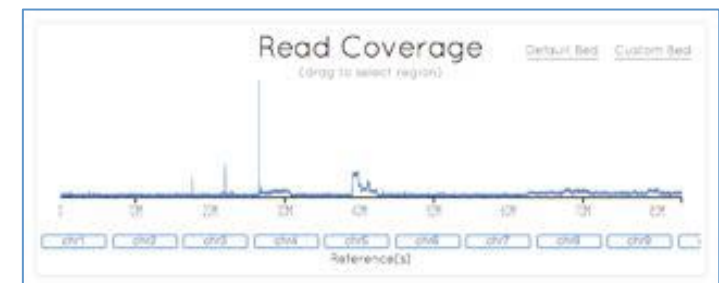
<http://schatzlab.cshl.edu/data/skbr3/>

Available *today* under the Toronto Agreement:

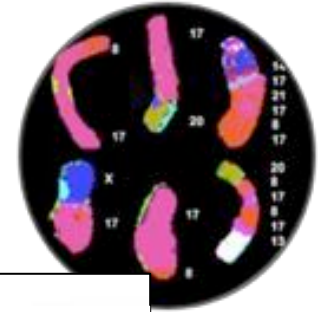
- Fastq & BAM files of aligned reads
- Interactive Coverage Analysis with BAM.IOBIO
- Whole genome assembly

Available soon

- Whole genome methylation analysis
- Full-length cDNA Transcriptome analysis
- Comparison to single cell analysis of >100 individual cells



SK-BR-3 Her2+ Breast Cancer Reference Genome



Released

Available

- Fast
- Interactive
- Whole

Available

- Whole
- Full-length cDNA transcriptome analysis
- Comparison to single cell analysis of >100 individual cells

bioRxiv
beta
THE PREPRINT SERVER FOR BIOLOGY

New Results

Error correction and assembly complexity of single molecule sequencing reads.

Hayan Lee, James Gurtowski, Shinjae Yoo, Shoshana Marcus, W. Richard McCombie, Michael Schatz

doi: <http://dx.doi.org/10.1101/006395>





Outline

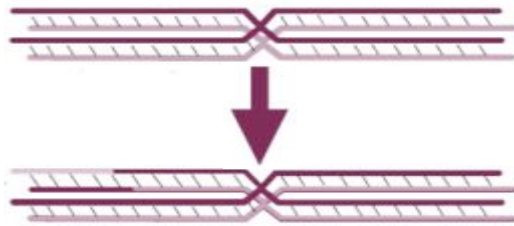
1. Single Molecule Sequencing

Long read sequencing of a breast cancer cell line

2. Single Cell Copy Number Analysis

Intra-tumor heterogeneity and metastatic progression

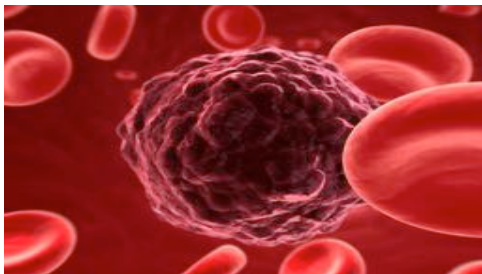
Single Cell Sequencing



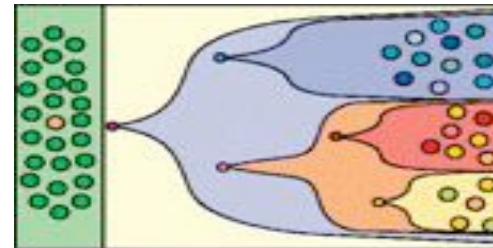
Recombination /
Crossover in germ cells



Neuronal mosaicism



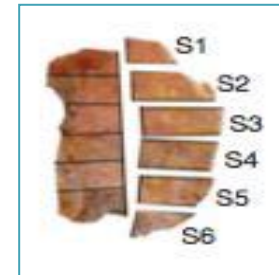
Circulating tumor cells



Clonal Evolution
in tumors

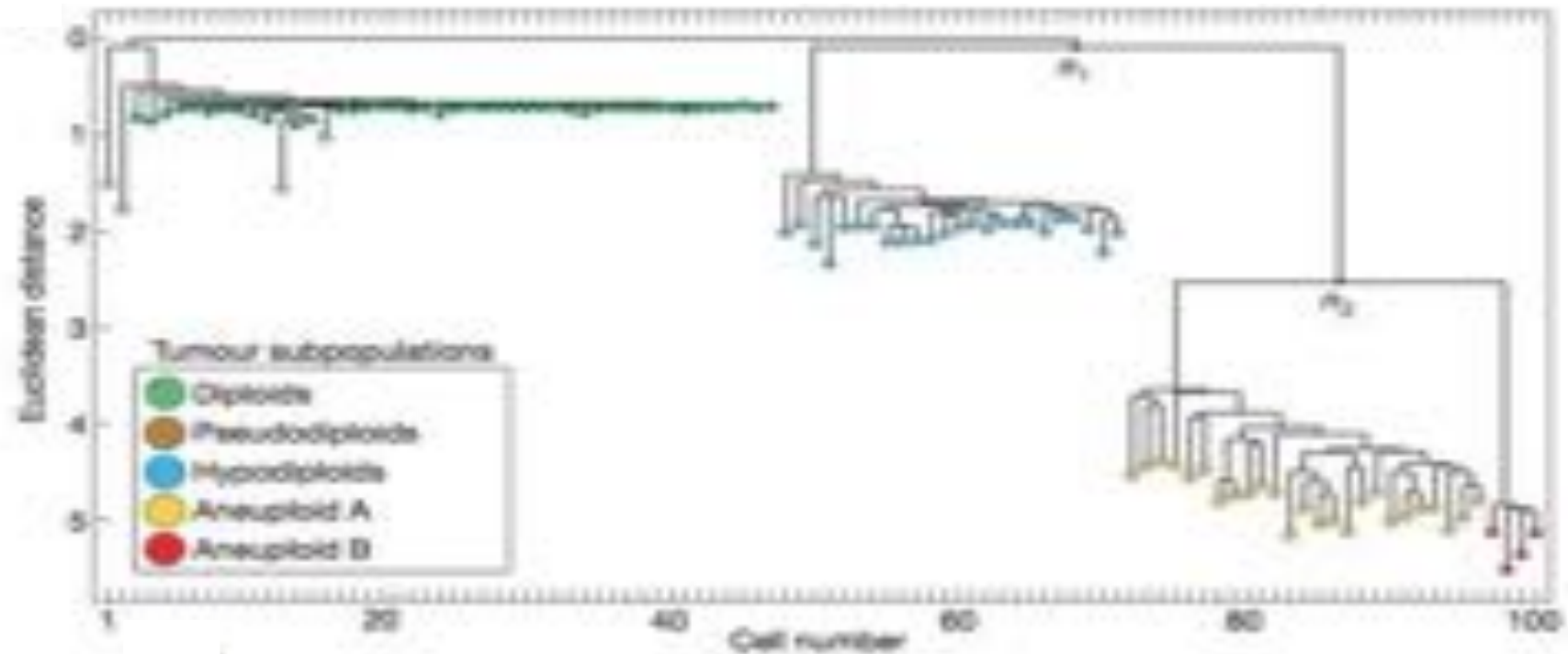
LETTER

doi:10.1038/nature09607

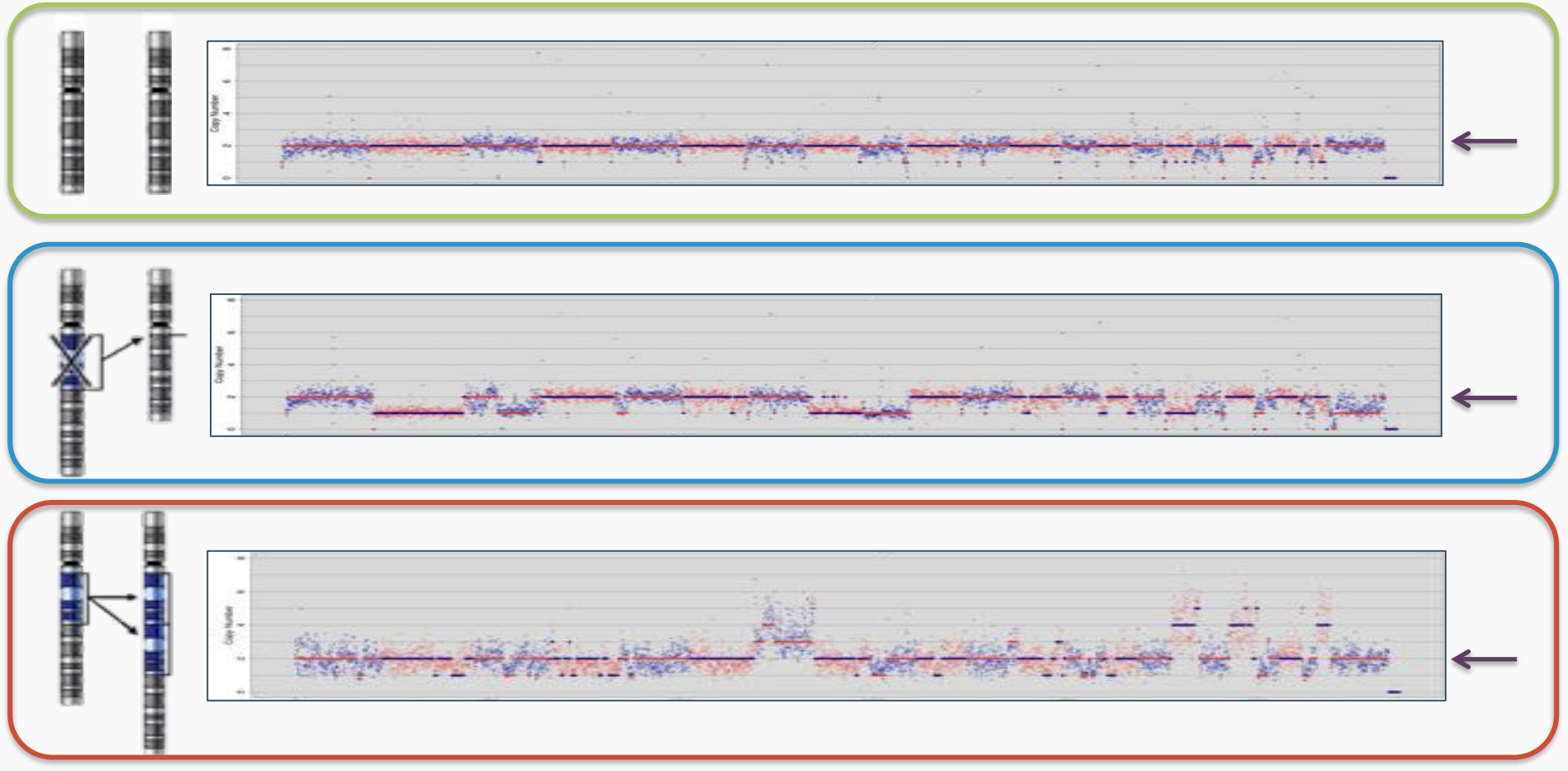


Tumour evolution inferred by single-cell sequencing

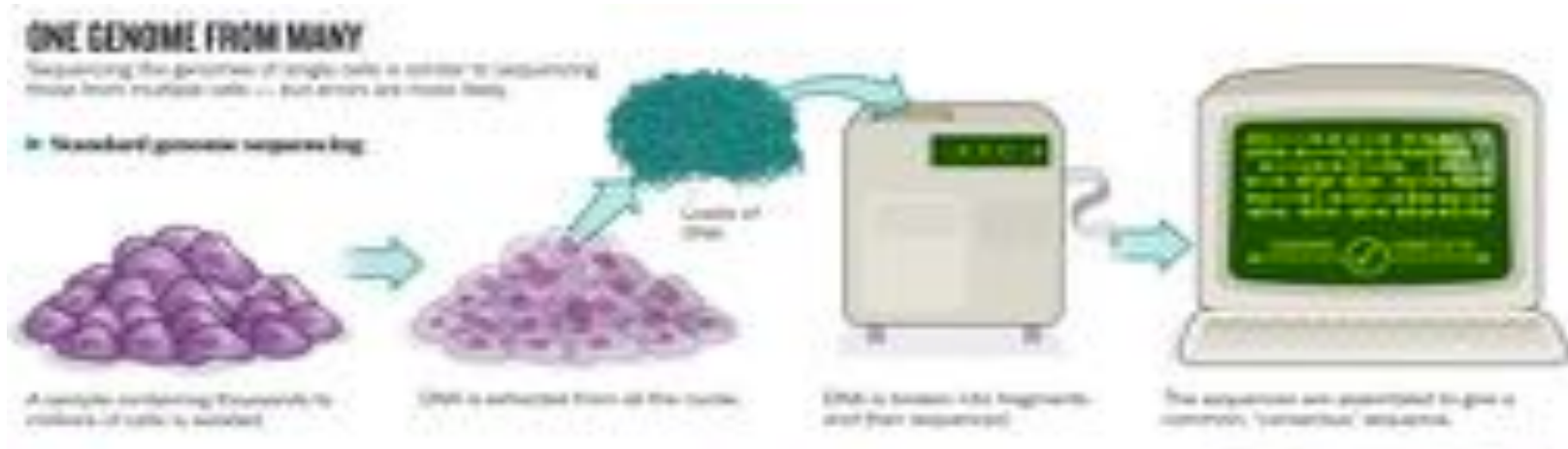
Nicholas Navin^{1,2}, Jude Kendall¹, Jennifer Troge¹, Peter Andrews¹, Linda Rodgers¹, Jeanne McIndoo¹, Kerry Cook¹, Asya Stepanisky¹, Dan Levy³, Diane Esposito³, Lakshmi Muthuswamy³, Alex Krasnitz², W. Richard McComble¹, James Hicks¹ & Michael Wigler¹



Copy-number Profiles



Whole Genome Amplification



Whole Genome Amplification

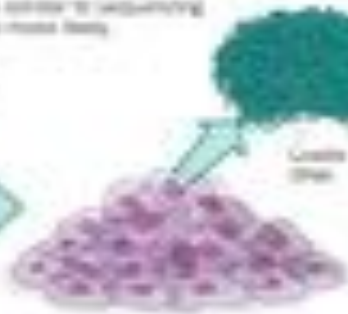
ONE GENOME FROM MANY

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

Standard genome sequencing



A sample containing thousands to millions of cells to be sequenced.



DNA is extracted from all the cells.



DNA is broken into fragments and then sequenced.

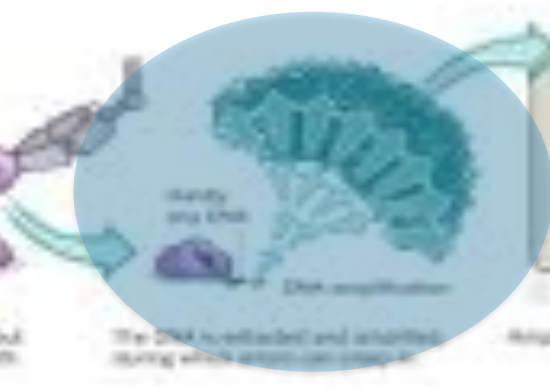


The sequences are presented to give a common, 'consensus' sequence.

Single-cell sequencing



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



The DNA is extracted and amplified during which errors can creep in.



Amplified DNA is sequenced.



Single molecules from the single cells are sequenced separately (difficult), but that sequence can find gaps.

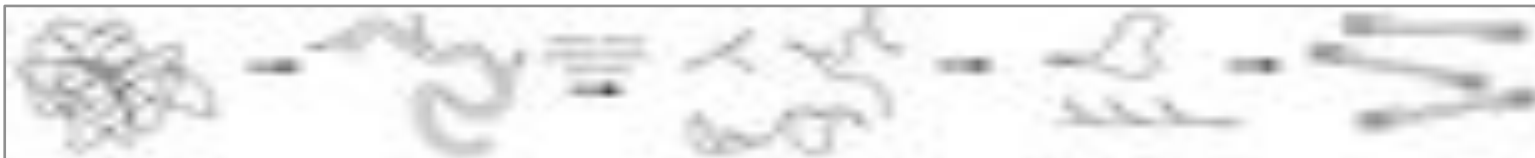
Whole Genome Amplification Techniques



DOP-PCR (Degenerate Oligonucleotide Primed PCR)



MDA (Multiple Displacement Amplification)

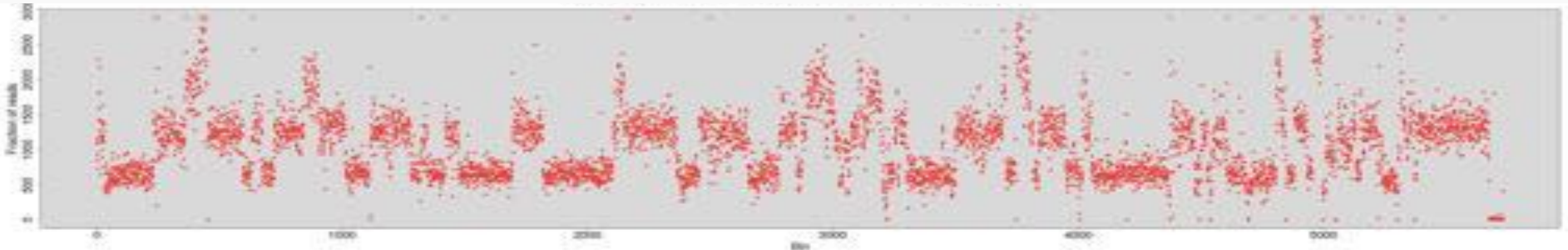


MALBAC (Multiple Annealing and Looping Based Amplification Cycles)

Interactive Analysis and Quality Assessment of Single Cell Copy Number Variations

Garvin, T., Aboukhalil, R. *et al.* (2014) *Under review*

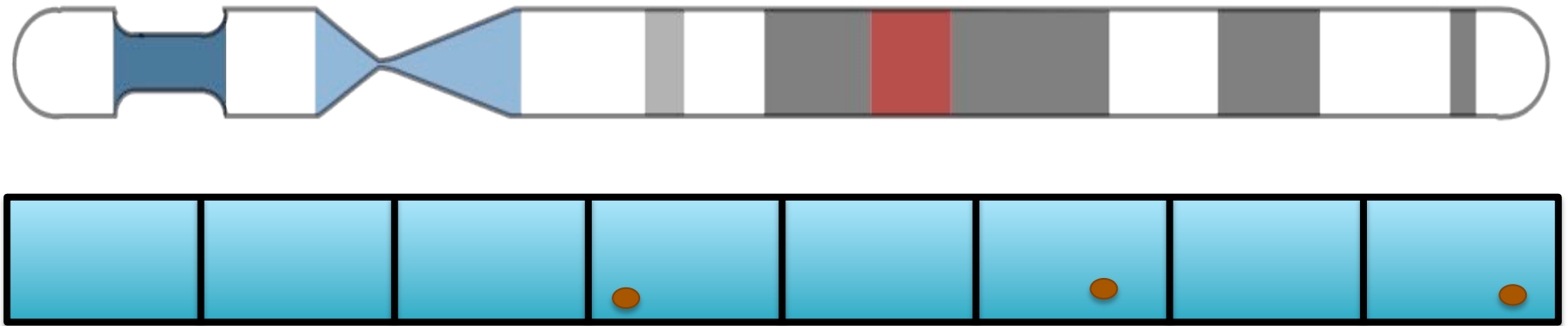
Data are noisy



- Potential for biases at every step
 - WGA: Non-uniform amplification
 - Library Preparation: Low complexity, read duplications, barcoding
 - Sequencing: GC artifacts, short reads
 - Computational analysis: mappability, GC correction, segmentation, tree building

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

I) Binning

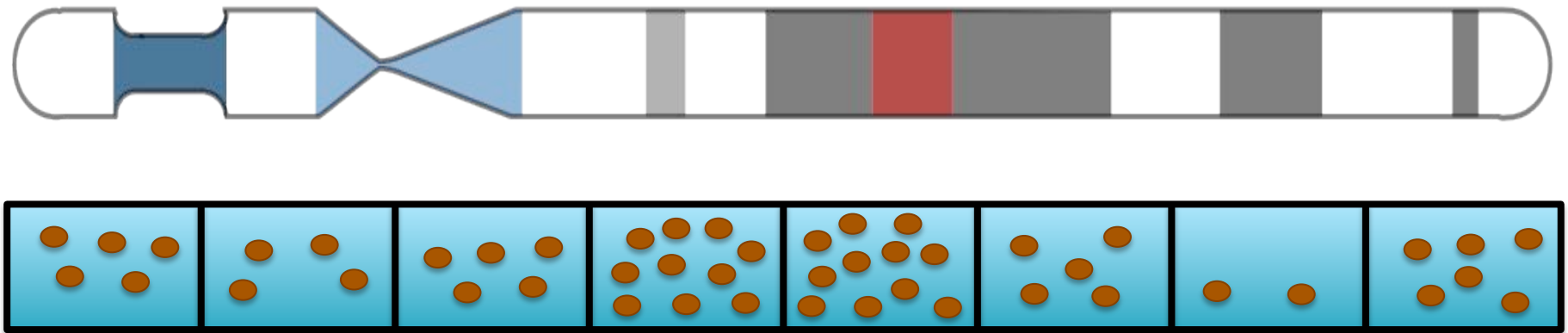


Single Cell CNV analysis

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

Use uniquely mappable bases to establish bins

I) Binning

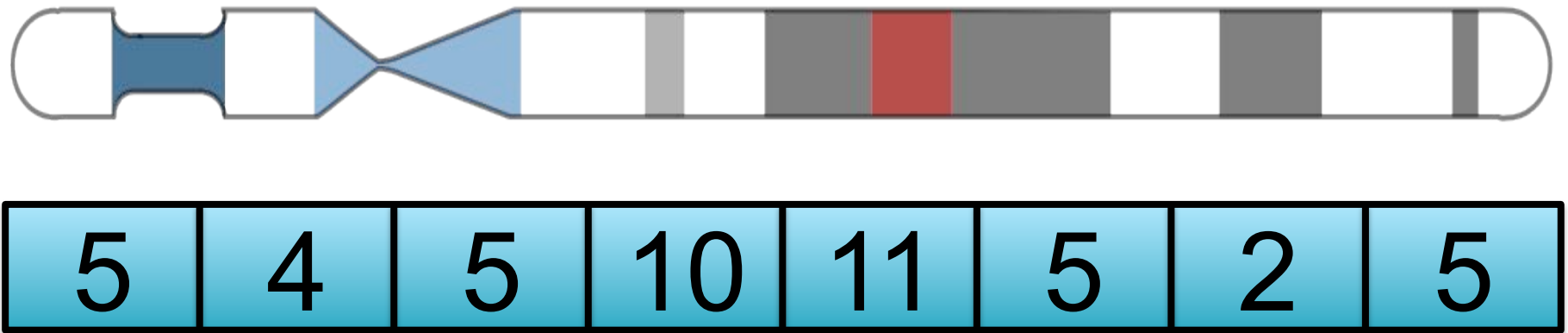


Single Cell CNV analysis

- 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

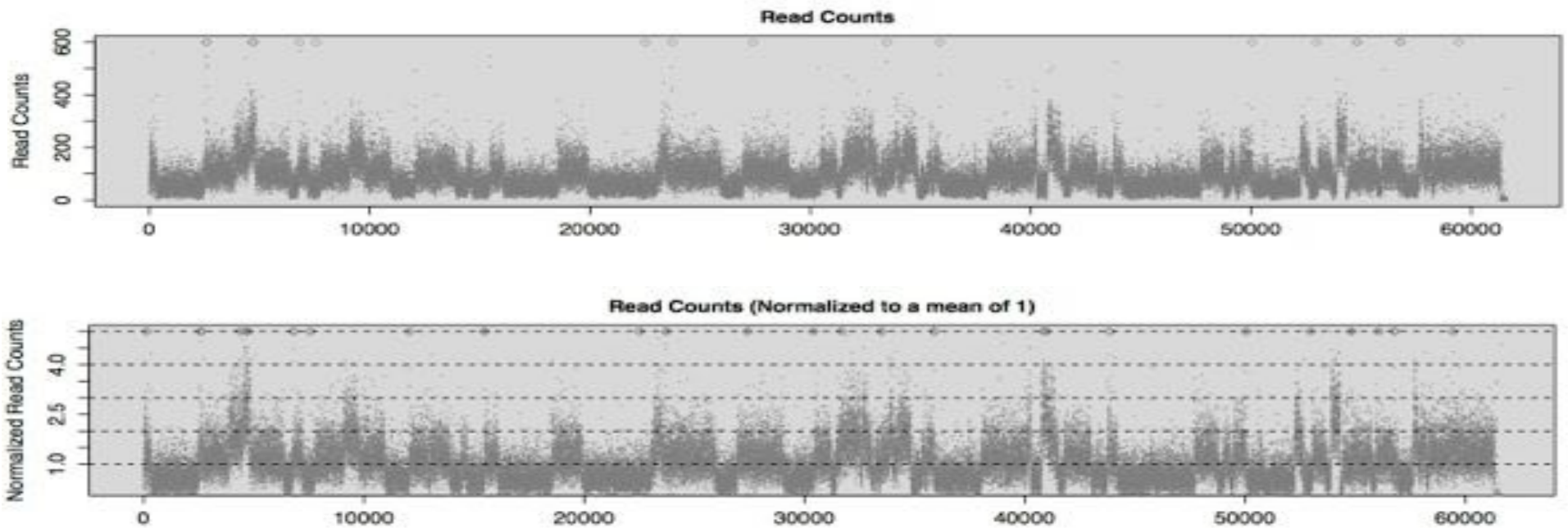


Single Cell CNV analysis

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

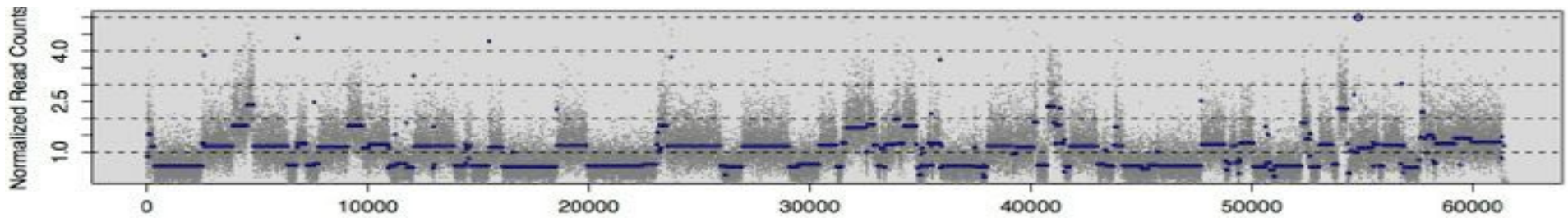
Use uniquely mappable bases to establish bins

2) Normalization

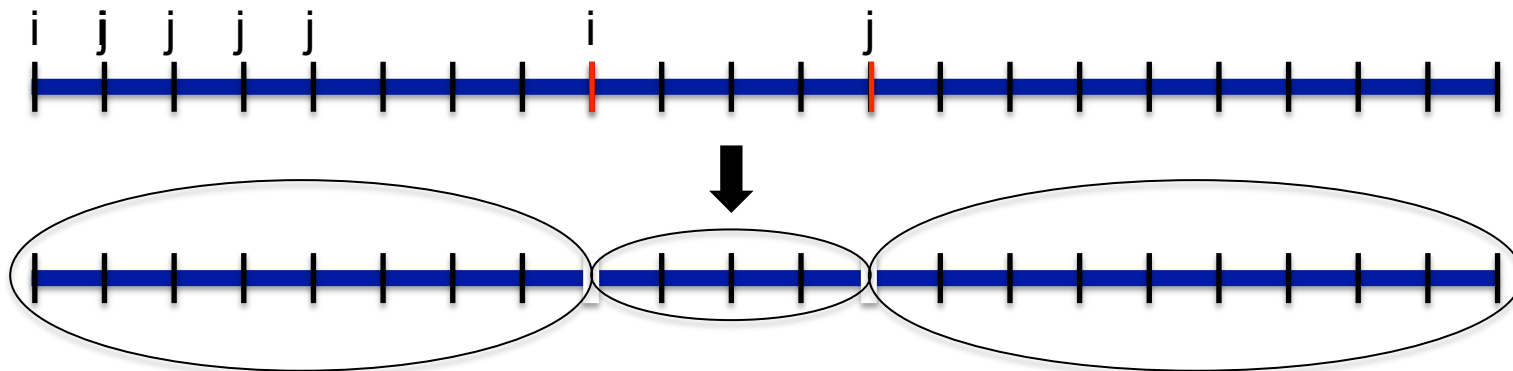


Also correct for mappability, GC content, amplification biases

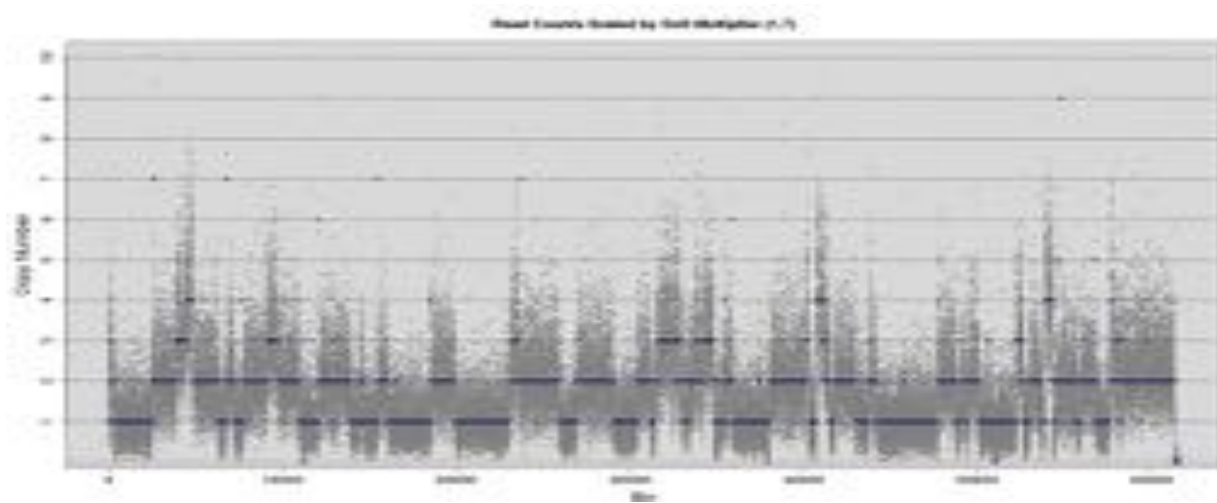
3) Segmentation



Circular Binary Segmentation (CBS)

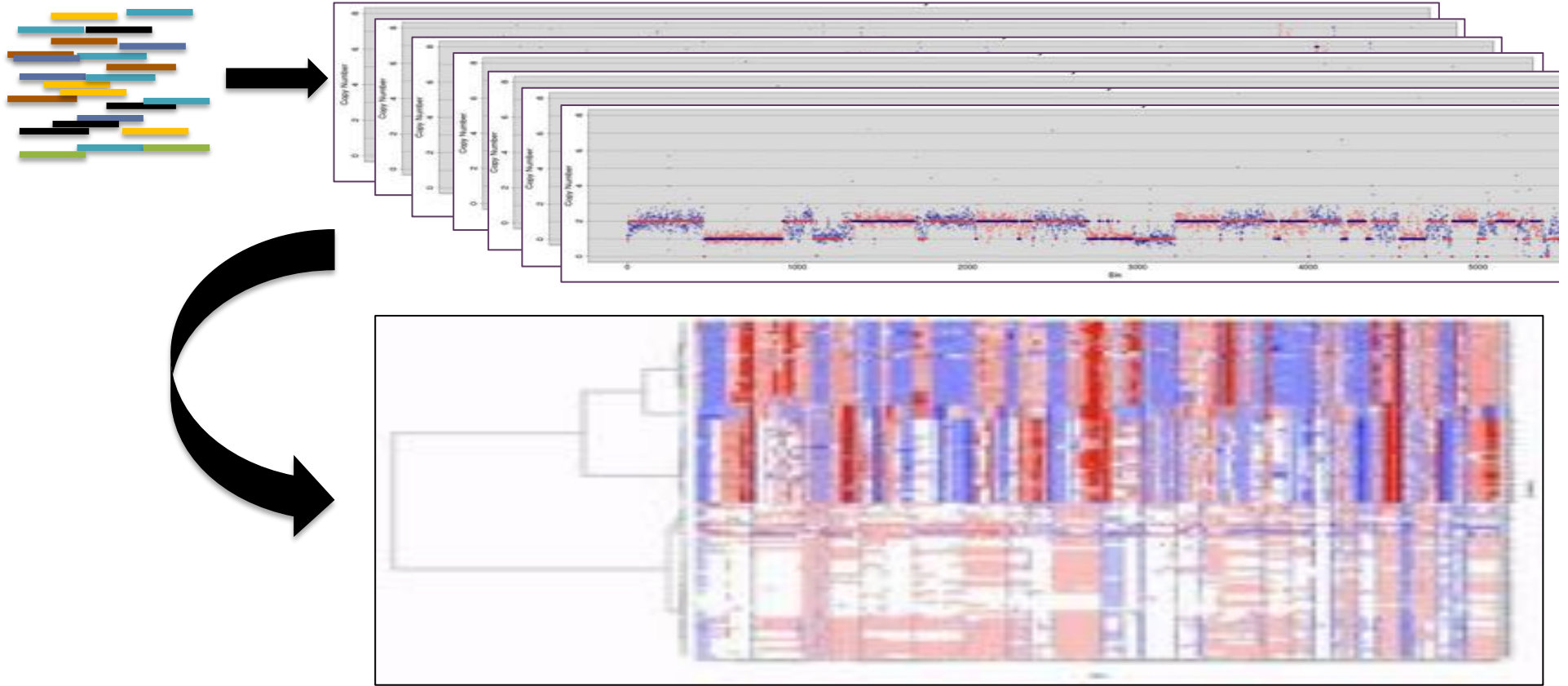


4) Estimating Copy Number



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

5) Cells to Populations



Ginkgo

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

Interactive Single Cell CNV analysis & clustering

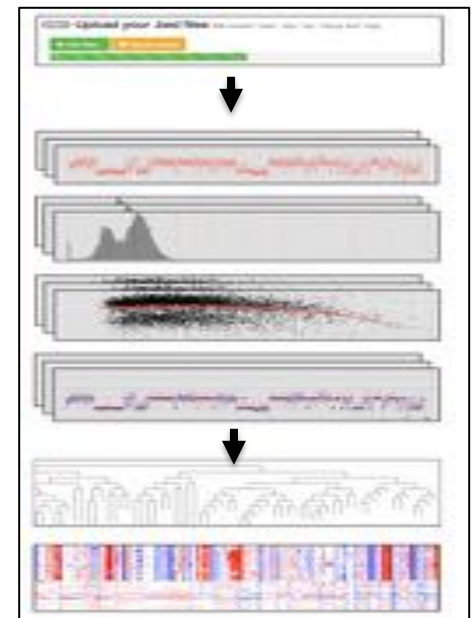
- Easy-to-use, web interface, parameterized for binning, segmentation, clustering, etc
- Per cell through project-wide analysis in any species

Compare MDA, DOP-PCR, and MALBAC

- DOP-PCR shows superior resolution and consistency

Available for collaboration

- Extending clustering methods, prototyping scRNA

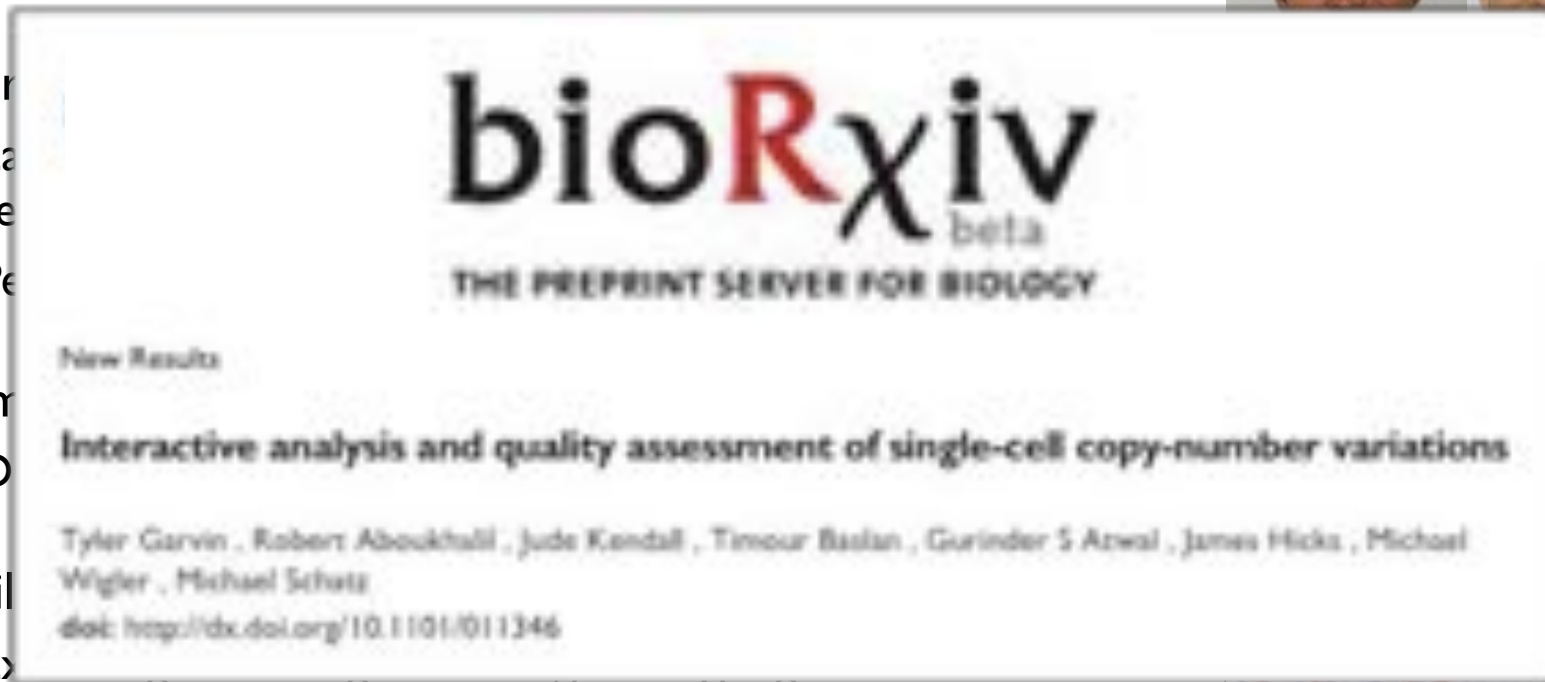


Ginkgo

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



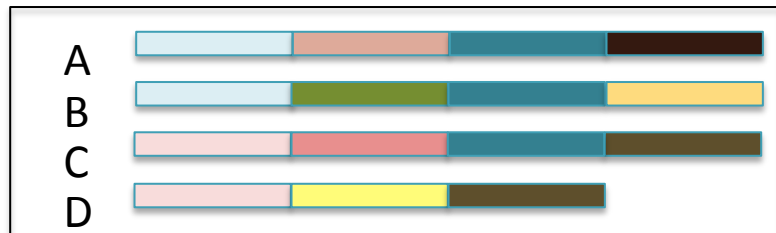
Inter
– Ea
se
– Pe
Com
– D
Avail
– Ex



The screenshot shows the bioRxiv beta logo at the top, with the tagline 'THE PREPRINT SERVER FOR BIOLOGY'. Below the logo, it says 'New Results' and features the title 'Interactive analysis and quality assessment of single-cell copy-number variations'. The authors listed are Tyler Garvin, Robert Aboukhalil, Jude Kendall, Timour Bastan, Gurinder S Arwal, James Hicks, Michael Wigler, and Michael Schatz. The DOI is provided as <http://dx.doi.org/10.1101/011346>.

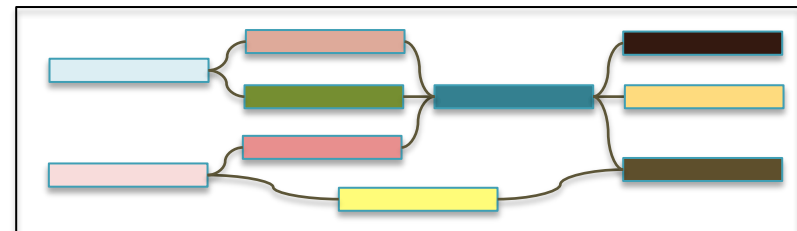


Pan-Genome Alignment & Assembly



Time to start focusing on problems studying populations of complete genomes

- Available today for many microbial species, near future for higher eukaryotes



Pan-genome colored de Bruijn graph

- Encodes all the sequence relationships between the genomes
- How well conserved is a given sequence?

SplitMEM: A graphical algorithm for pan-genome analysis with suffix skips

Marcus, S, Lee, H, Schatz MC (2014) *Bioinformatics*. doi: 10.1093/bioinformatics/btu756

Extending reference assembly models

Church, D. et al. (2015) *Genome Biology*. doi:10.1186/s13059-015-0587-3

Understanding Genome Structure & Function

Single Molecule Sequencing

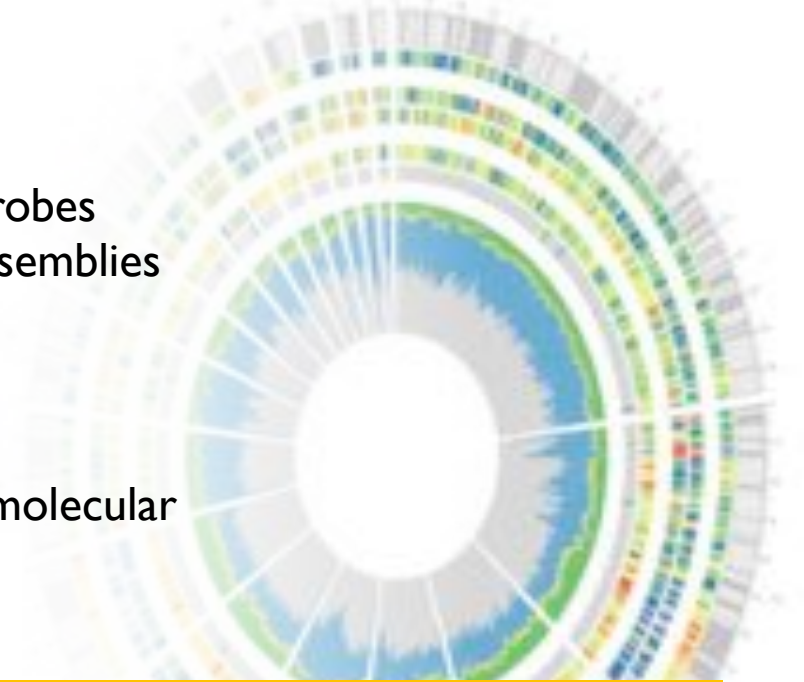
Now have the ability to **perfectly assemble** microbes and many small eukaryotes, **reference quality** assemblies of larger eukaryotes

Single Cell Sequencing

Exciting technologies to probe the genetic and molecular **composition of complex environments**

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

Largely limited by our quantitative power to make comparisons and find patterns



Acknowledgements

Schatz Lab

Rahul Amin
Eric Biggers
Han Fang
Tyler Gavin
James Gurtowski
Ke Jiang
Hayan Lee
Zak Lemmon
Shoshana Marcus
Giuseppe Narzisi
Maria Nattestad
Aspyn Palatnick
Srividya Ramakrishnan
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

OICR

Karen Ng
Timothy Beck
Yogi Sundaravadanam
John McPherson

NBACC

Adam Phillippy
Serge Koren

Pacific Biosciences
Oxford Nanopore



National Human
Genome Research
Institute



U.S. DEPARTMENT OF
ENERGY

SFARI

SIMONS FOUNDATION
AUTISM RESEARCH INITIATIVE

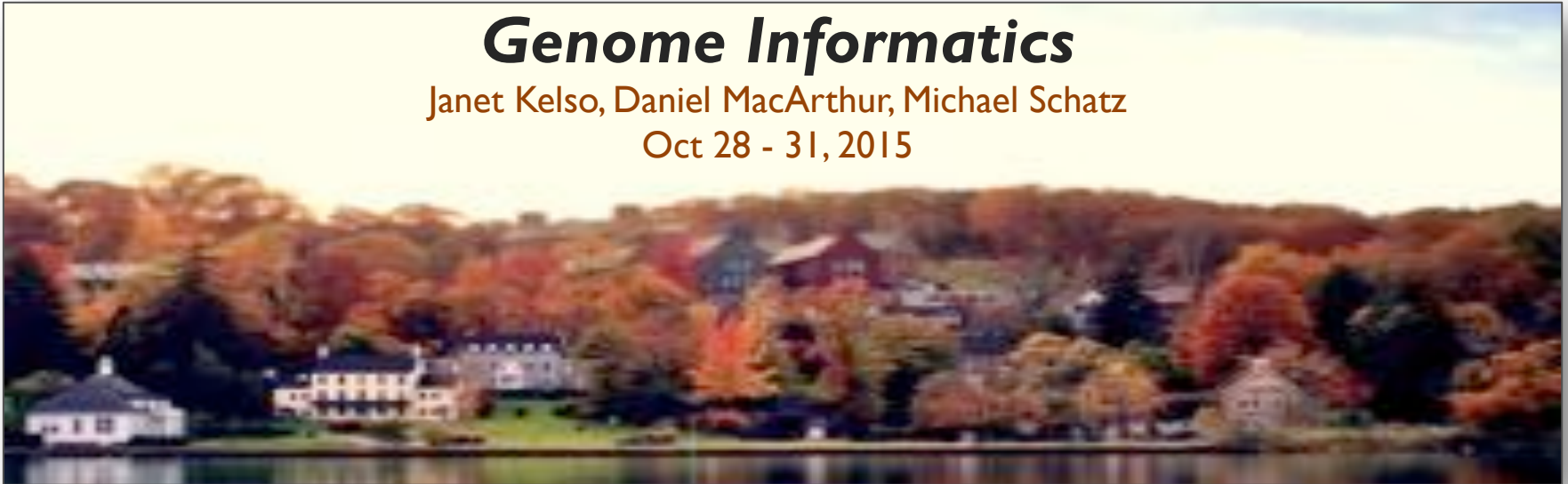


ALFRED P. SLOAN
FOUNDATION

Genome Informatics

Janet Kelso, Daniel MacArthur, Michael Schatz

Oct 28 - 31, 2015



Thank you

<http://schatzlab.cshl.edu>

@mike_schatz