# The resurgence of reference quality genomes

Michael Schatz
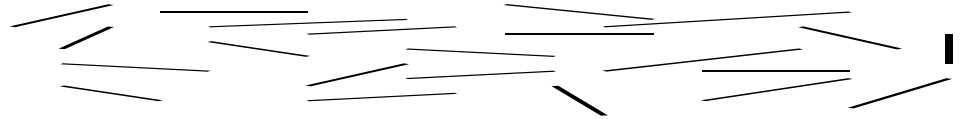
May 22, 2015
NYU Genomics Symposium

# Sequence Assembly Problem

1. Shear & Sequence DNA
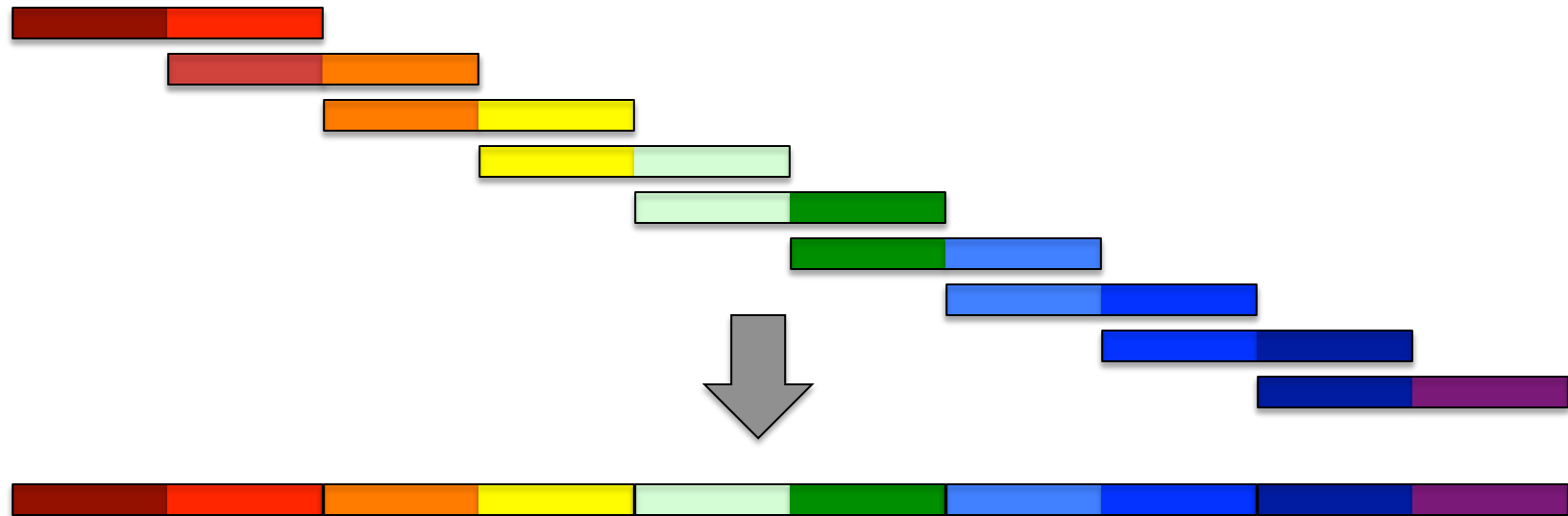
2. Construct assembly graph from overlapping reads

   ...AGCCTAG GGATGCGCGACACGT

      GGATGCGCGACACGT CGCATATCCGGTTTGGT CAACCTCGGACGGAC

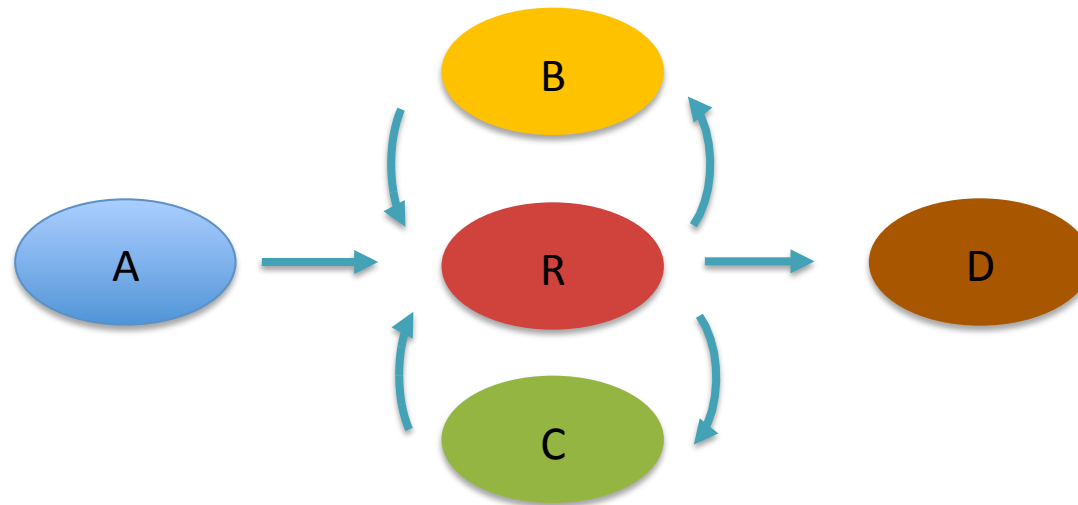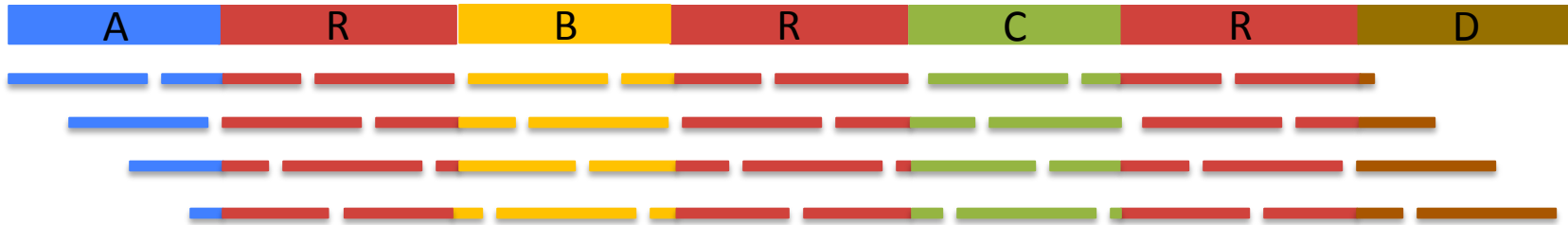                                              CAACCTCGGACGGAC CTCAGCGAA...
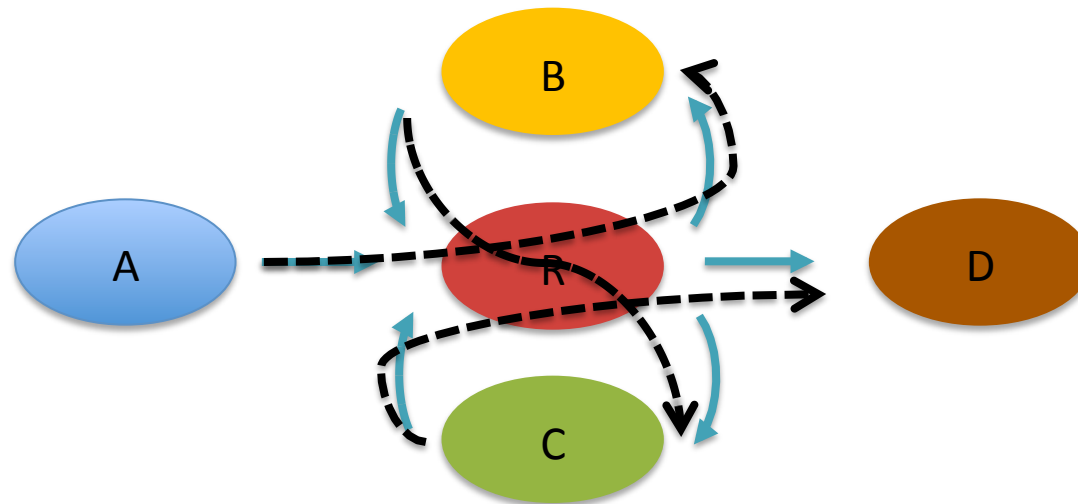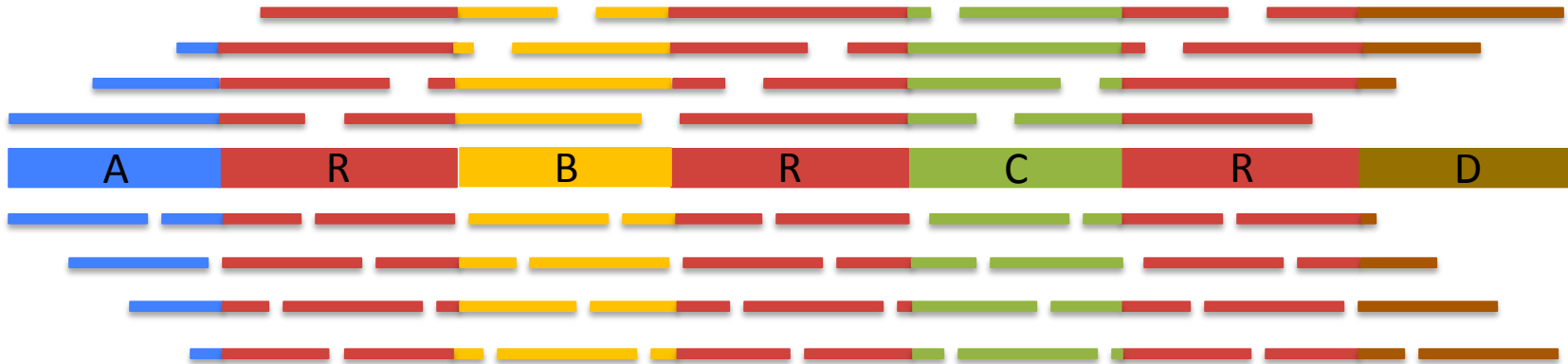
3. Simplify assembly graph



**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*
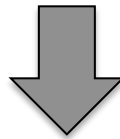
# Assembly Complexity

# 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* | *Pacific Biosciences* | *Oxford Nanopore* |
|---|---|---|
|  |  |  |
| (Kuleshov et al. 2014) | (Berlin et al, 2014) | (Quick et al, 2014) |

**Long Span Sequencing: Chromosome Scaffolding, SV analysis, phasing**

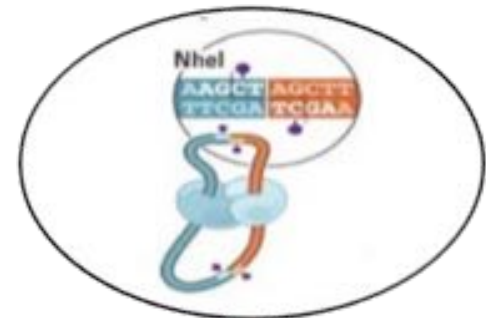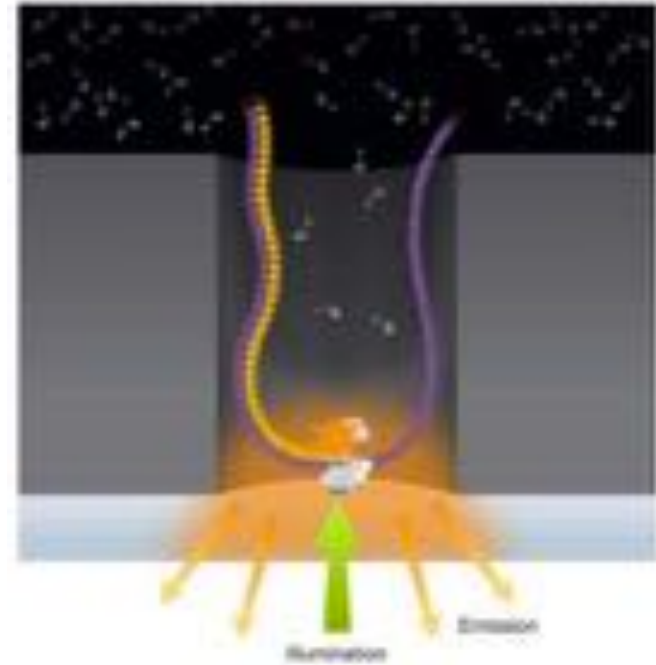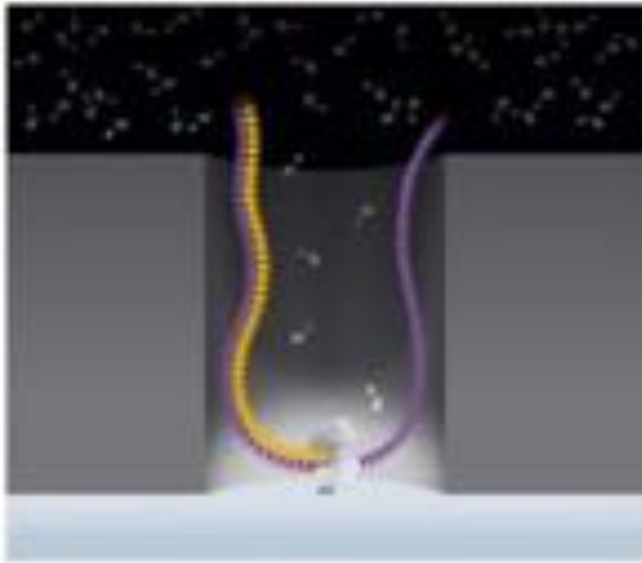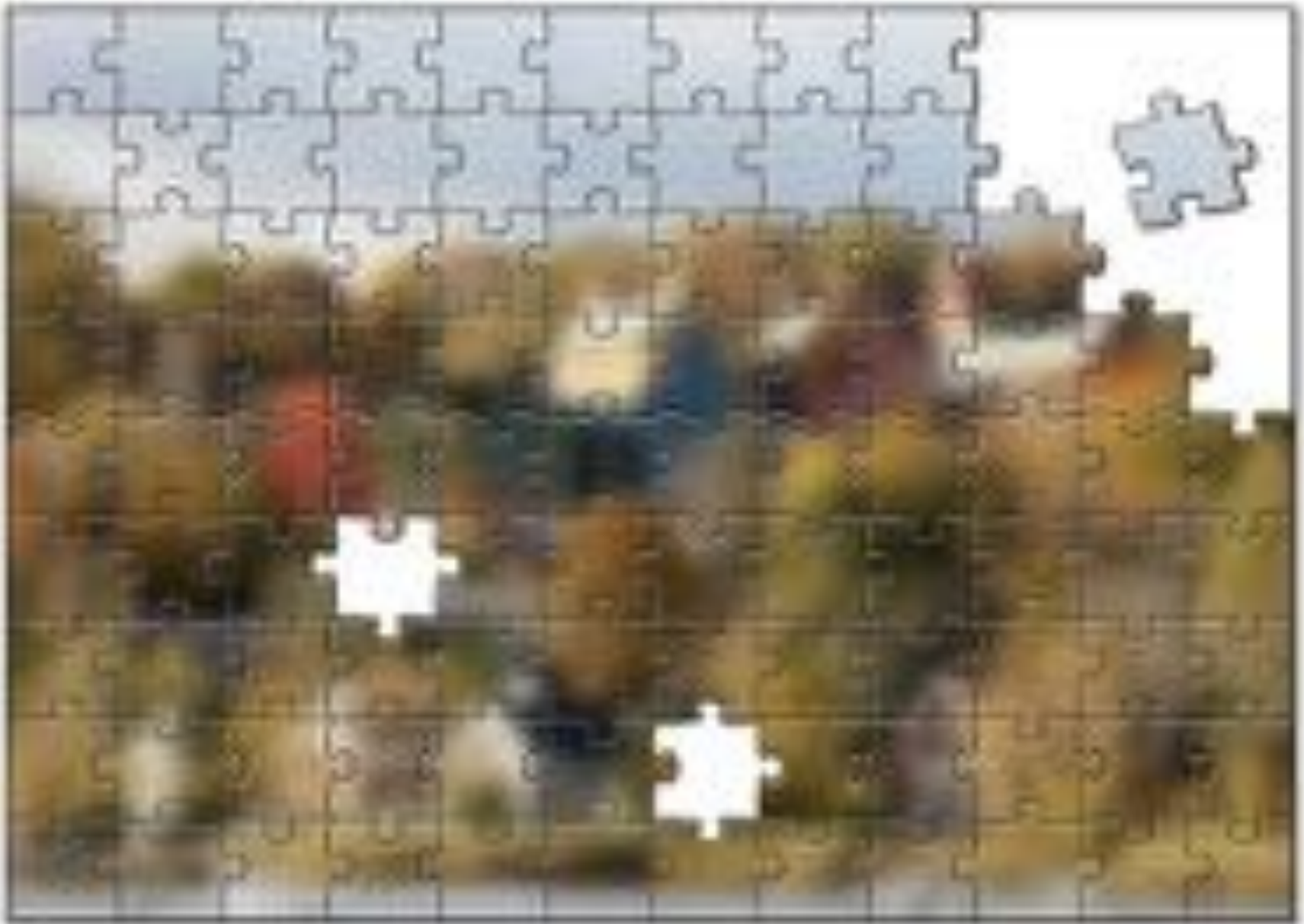| *Molecular Barcoding* | *Optical Mapping* | *Chromatin Assays* |
|---|---|---|
|  |  |  |
| (10Xgenomics.com) | (Cao et al, 2014) | (Putnam et al, 2015) |

# PacBio SMRT Sequencing

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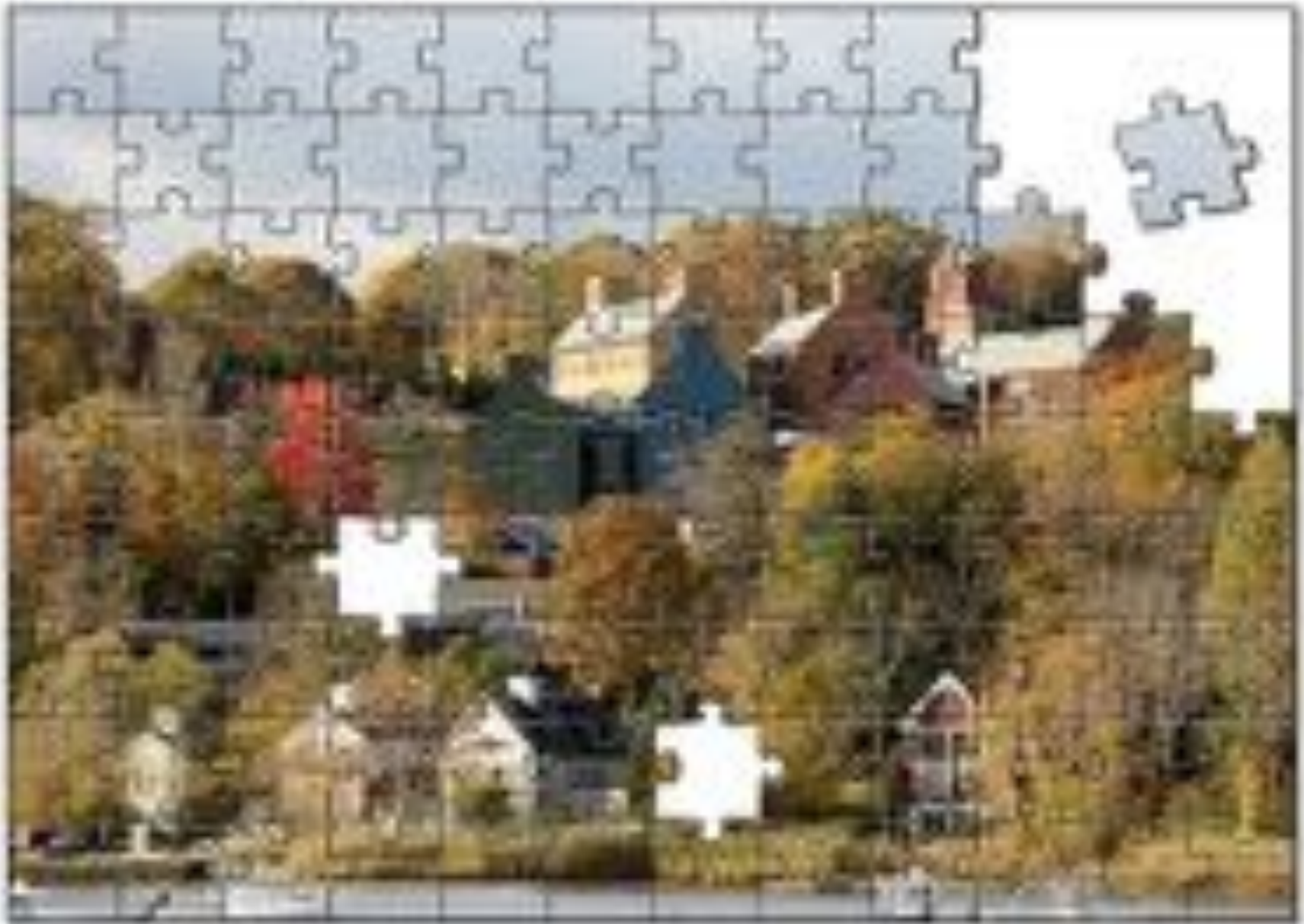


http://www.pacificbiosciences.com/assets/files/pacbio_technology_backgrounder.pdf

# Single Molecule Sequences

# "Corrective Lens" for Sequencing

# PacBio Assembly Algorithms

| PBJelly | PacBioToCA & ECTools | HGAP/MHAP & Quiver |
|---|---|---|
|  |  |  |
| **Gap Filling and Assembly Upgrade** | **Hybrid/PB-only Error Correction** | **PB-only Correction & Polishing** |
| English *et al* (2012) *PLOS One.* 7(11): e47768 | Koren, Schatz, *et al* (2012) *Nature Biotechnology.* 30:693–700 | Chin *et al* (2013) *Nature Methods.* 10:563–569 |

< 5x                    PacBio Coverage                    > 50x
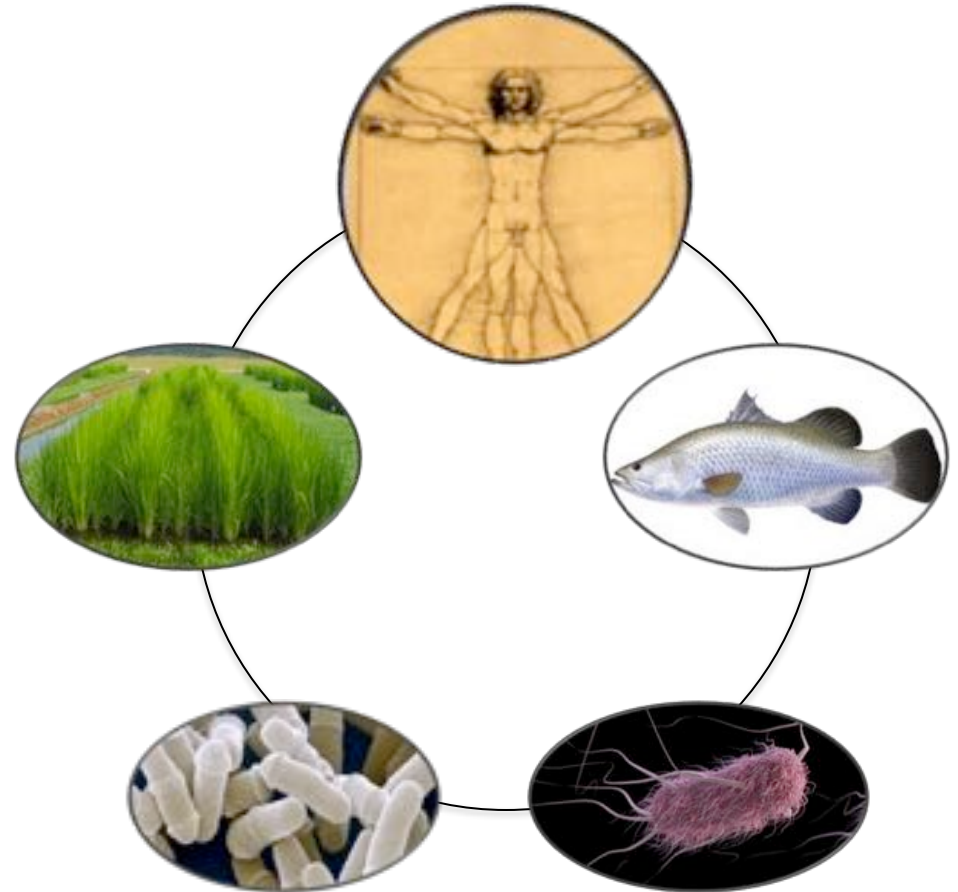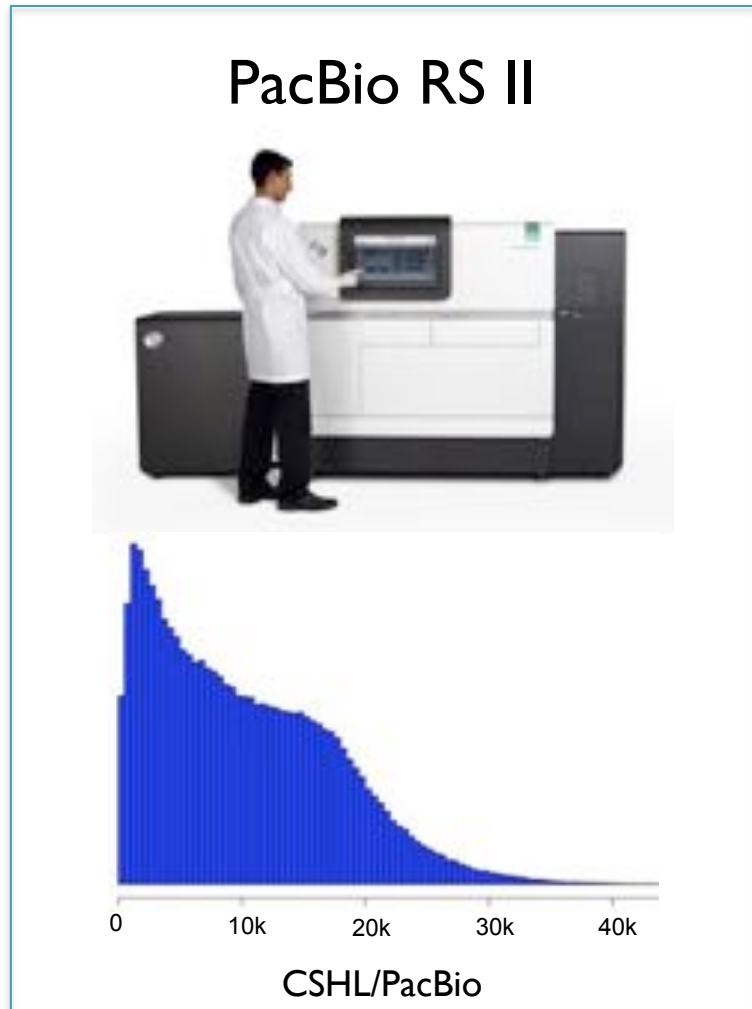
# 3rd Gen Long Read Sequencing
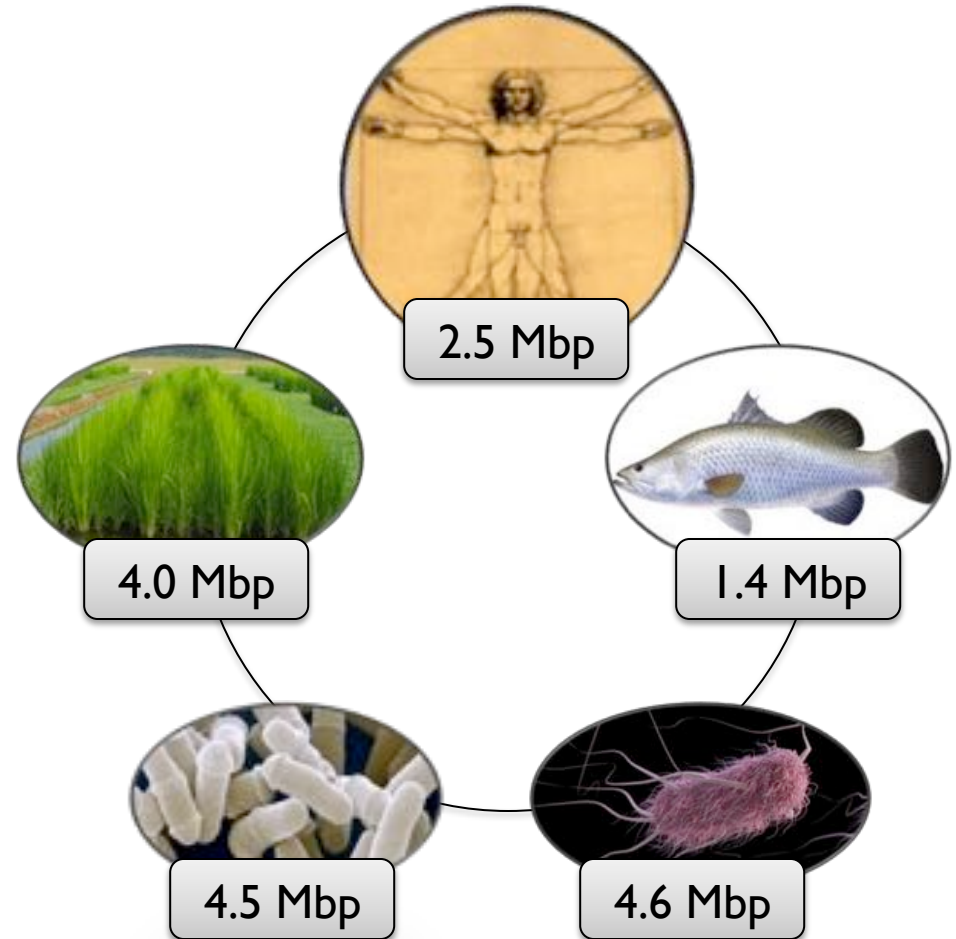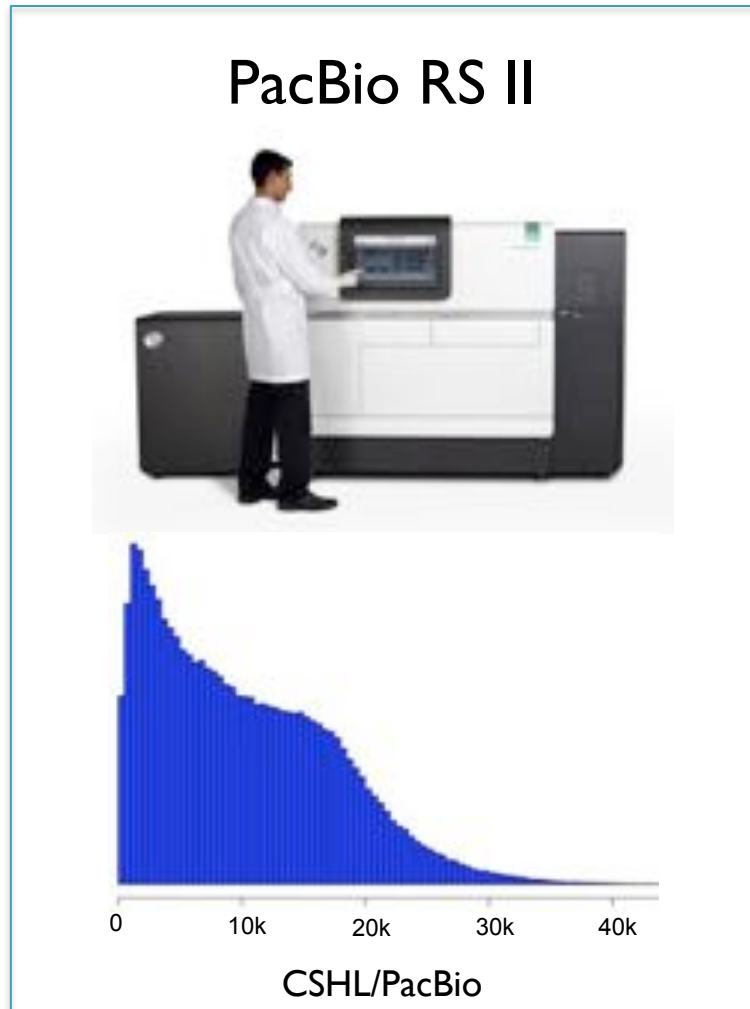


PacBio RS II

CSHL/PacBio

# 3<sup>rd</sup> Gen Long Read Sequencing



PacBio RS II

CSHL/PacBio

# 3rd Gen Long Read Sequencing



PacBio RS II

CSHL/PacBio

2.5 Mbp

1.4 Mbp

4.0 Mbp

4.5 Mbp

4.6 Mbp

# SK-BR-3



Maria Nattestad

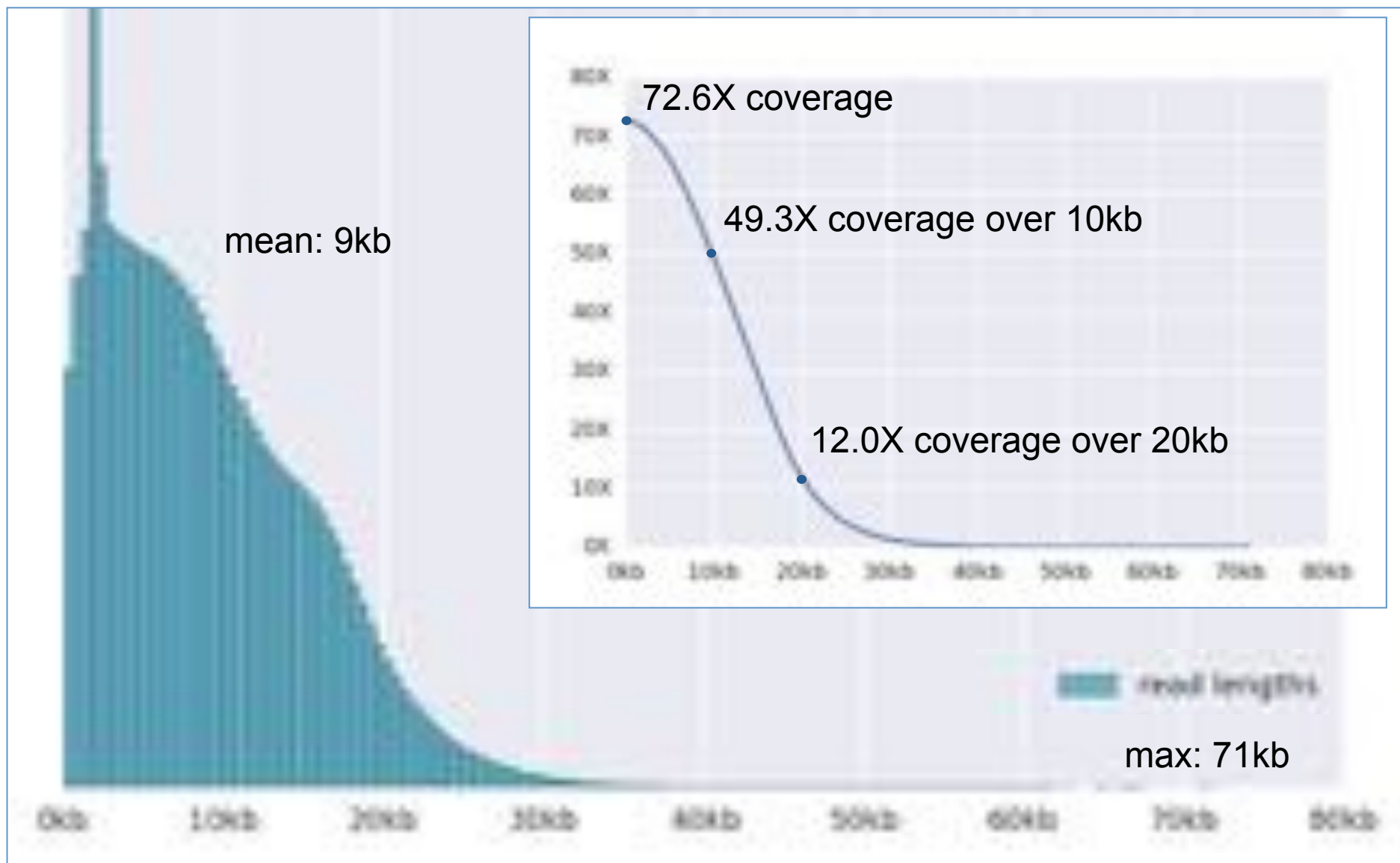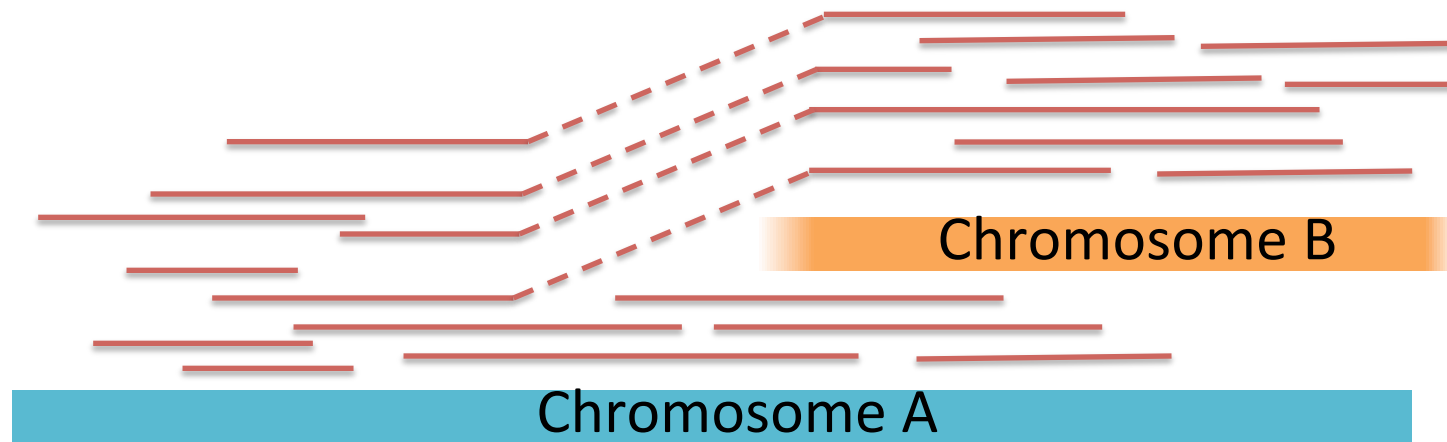Most commonly used Her2-amplified breast cancer

(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



mean: 9kb

72.6X coverage

49.3X coverage over 10kb

12.0X coverage over 20kb

read lengths

max: 71kb

# 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***
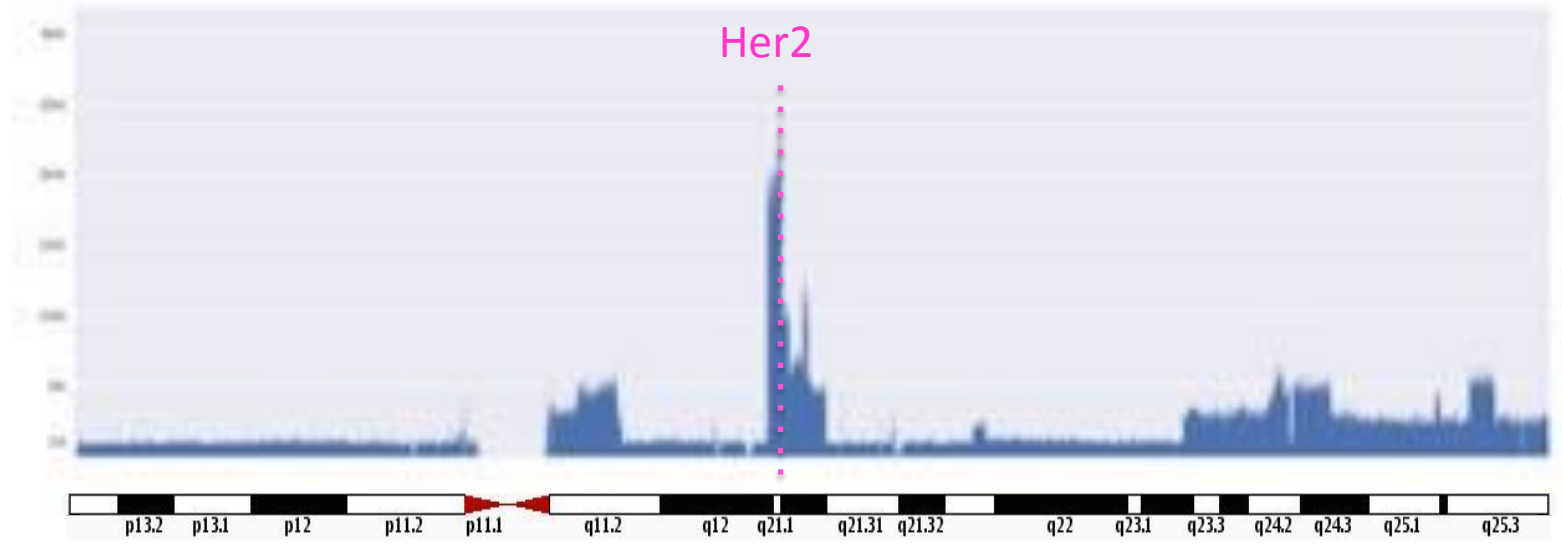   DNAnexus-enabled version of Falcon

   **Total Assembly: 2.64Gbp          Contig N50: 2.56 Mbp          Max Contig: 23.5Mbp**
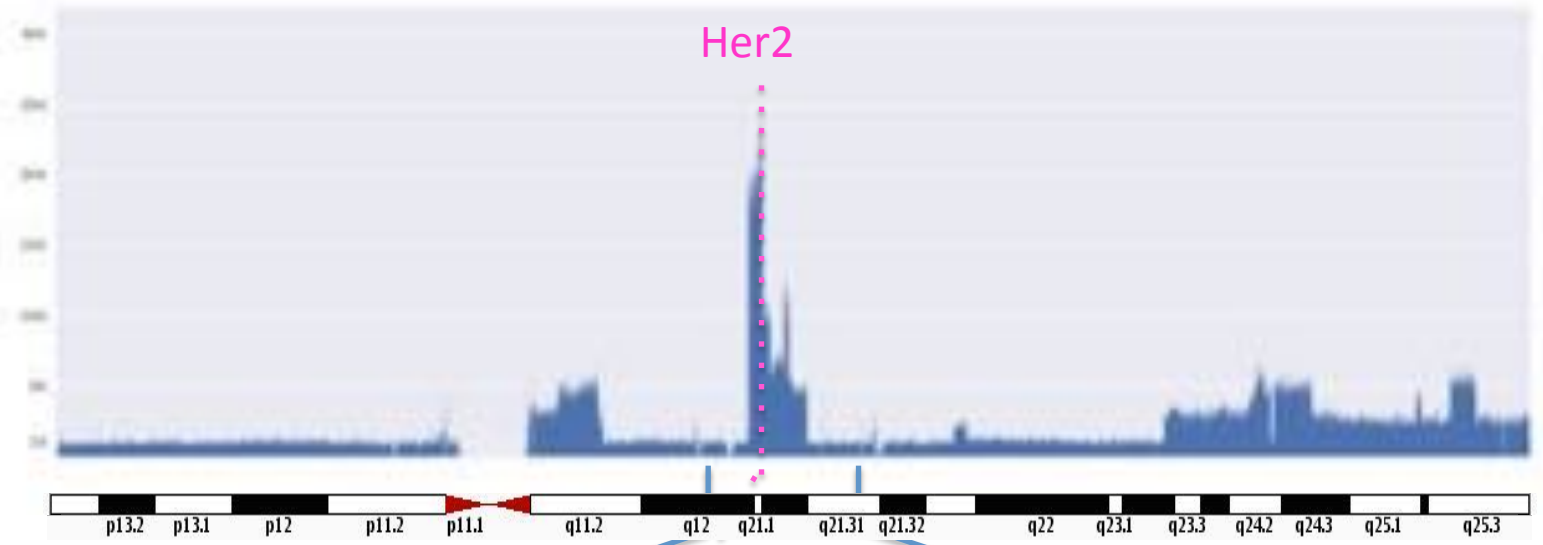
PacBio

Her2

Chr 17: 83 Mb

PacBio

Her2

p13.2 p13.1 p12 p11.2 p11.1 q11.2 q12 q21.1 q21.31 q21.32 q22 q23.1 q23.3 q24.2 q24.3 q25.1 q25.3

8 Mb

PacBio
chr17

Her2

PacBio

Her2

p13.2 p13.1 p12 p11.2 p11.1 q11.2 q12 q21.1 q21.31 q21.32 q22 q23.1 q23.3 q24.2 q24.3 q25.1 q25.3

8 Mb

Her2

PacBio
chr17

PacBio chr17

PacBio chr8

GSDMB
Her2
RARA
PKIA
TATDN1

50 Mb

Confirmed both known gene fusions in this region

PacBio chr17

GSDMB
Her2
RARA

PacBio chr8

PKIA
TATDN1

50 Mb

1.6 Mb

Confirmed both known gene fusions in this region

PacBio
chr17

Her2

RARA

chr8

PKIA

1.6 Mb

Joint coverage and breakpoint analysis to discover underlying events

# Cancer lesion Reconstruction



PacBio

chr17

Her2

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

# Cancer lesion Reconstruction

**Available *today* under the Toronto Agreement:**
- Fastq & BAM files of aligned reads
- Interactive Coverage Analysis with BAM.IOBIO
- Whole genome assembly & alignment

**Available soon**
- Whole genome methylation analysis
- Full length cDNA transciptome analysis
- Comparison to single cell analysis of >100 individual cells

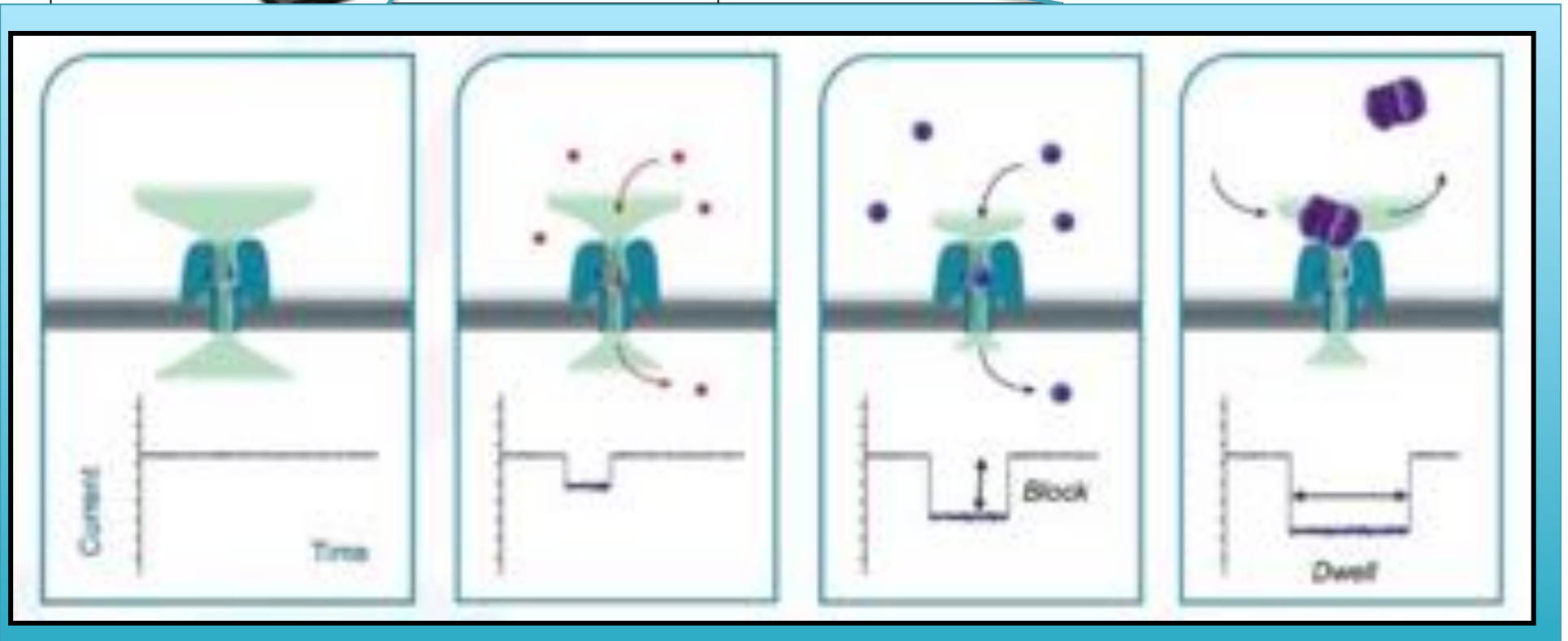***http://schatzlab.cshl.edu/data/skbr3/***

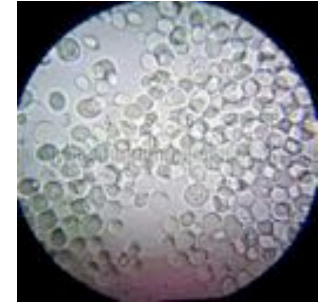4. Final duplication from within chromosome 8

# Oxford Nanopore MinION



- Thumb drive sized sequencer powered over USB

- Capacity for 512 reads at once

- Senses DNA by measuring changes to ion flow

# Nanopore Accuracy
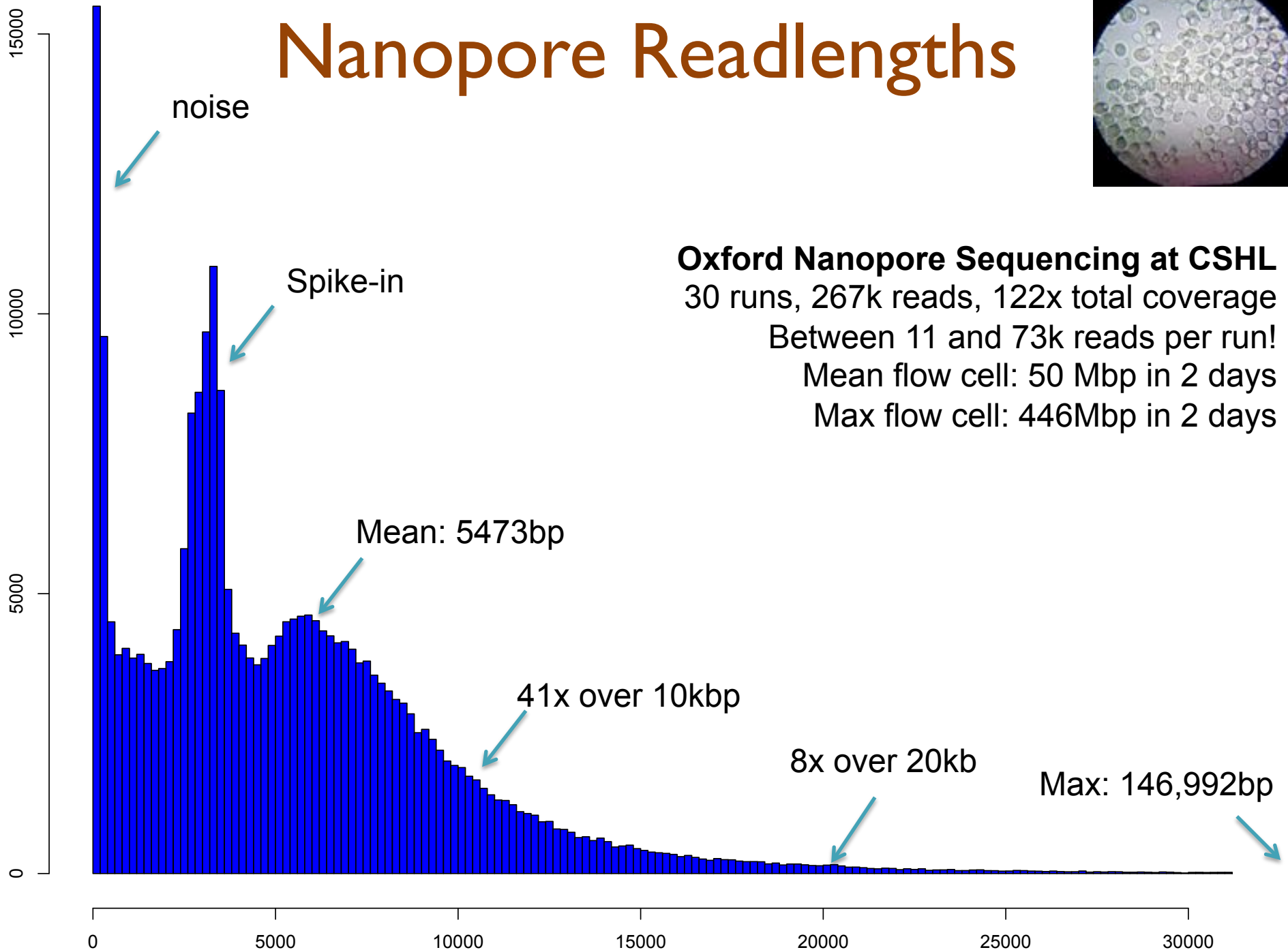
**Alignment Quality (BLASTN)**
Of reads that align, average ~64% identity
"2D base-calling" improves to ~70% identity

# NanoCorr: Nanopore-Illumina Hybrid Error Correction

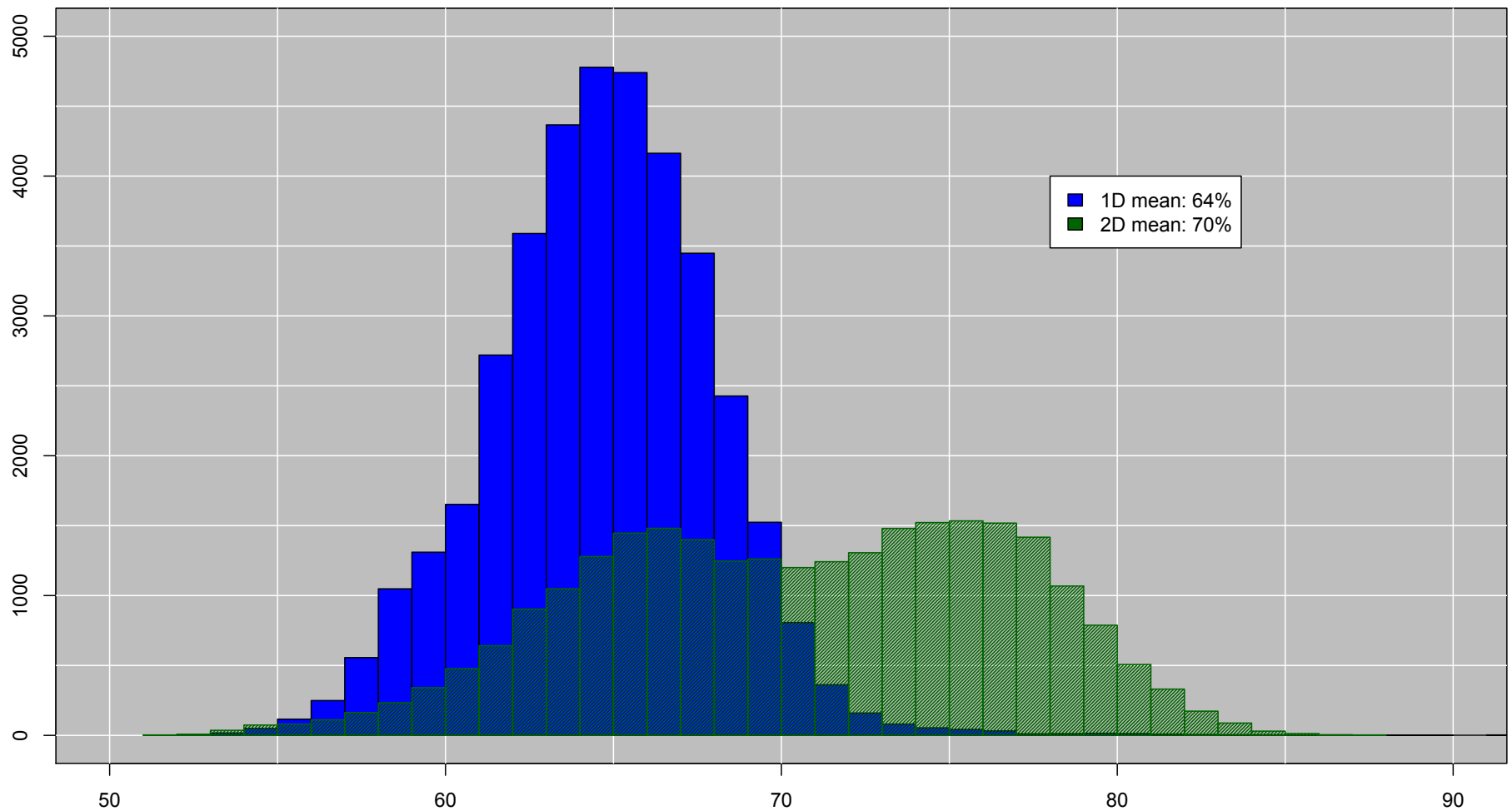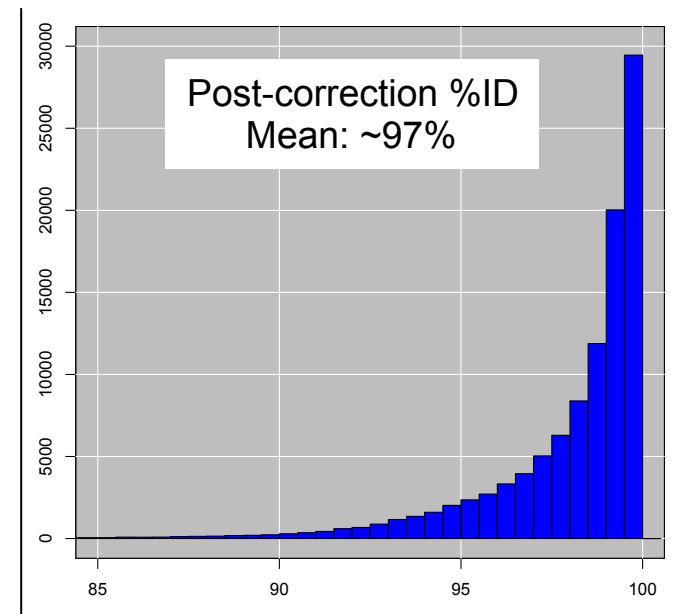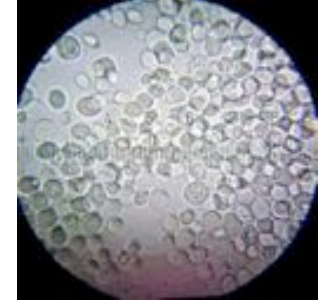https://github.com/jgurtowski/nanocorr

1. BLAST Miseq reads to all raw Oxford Nanopore reads

2. Select non-repetitive alignments
   - First pass scans to remove "contained" alignments
   - Second pass uses Dynamic Programming (LIS) to select set of high-identity alignments with minimal overlaps

3. Compute consensus of each Oxford Nanopore read
   - State machine of most commonly observed base at each position in read



Post-correction %ID
Mean: ~97%

**Oxford Nanopore Sequencing and de novo Assembly of a Eukaryotic Genome**
Goodwin, S, Gurtowski, J *et al.* (2015) bioRxiv doi: http://dx.doi.org/10.1101/013490

# NanoCorr Yeast Assembly

S288C Reference sequence
- 12.1Mbp; 16 chromo + mitochondria; N50: 924kbp
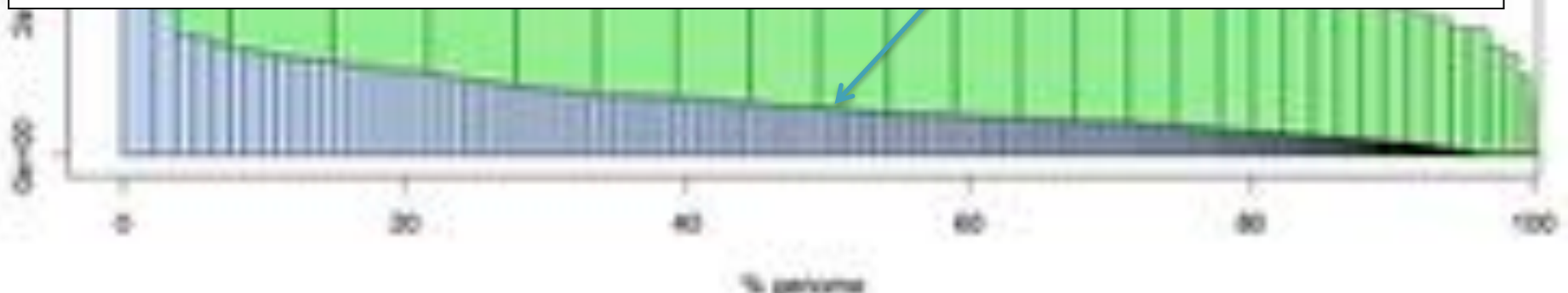


## bioRxiv beta
### THE PREPRINT SERVER FOR BIOLOGY

New Results

### Oxford Nanopore Sequencing and de novo Assembly of a Eukaryotic Genome

Sara Goodwin , James Gurtowski , Scott Ethe-Sayers , Panchajanya Deshpande , Michael Schatz , W Richard McCombie

doi: http://dx.doi.org/10.1101/013490

# Genomic Futures?

# Genomic Futures?

# iGenomics: Mobile Sequence Analysis

Aspyn Palatnick, Elodie Ghedin, Michael Schatz



***The worlds first genomics analysis app for iOS devices***

*BWT + Dynamic Programming + UI*

First application:
- Handheld diagnostics and therapeutic recommendations for influenza infections

- In the iOS AppStore now!

**Future applications**
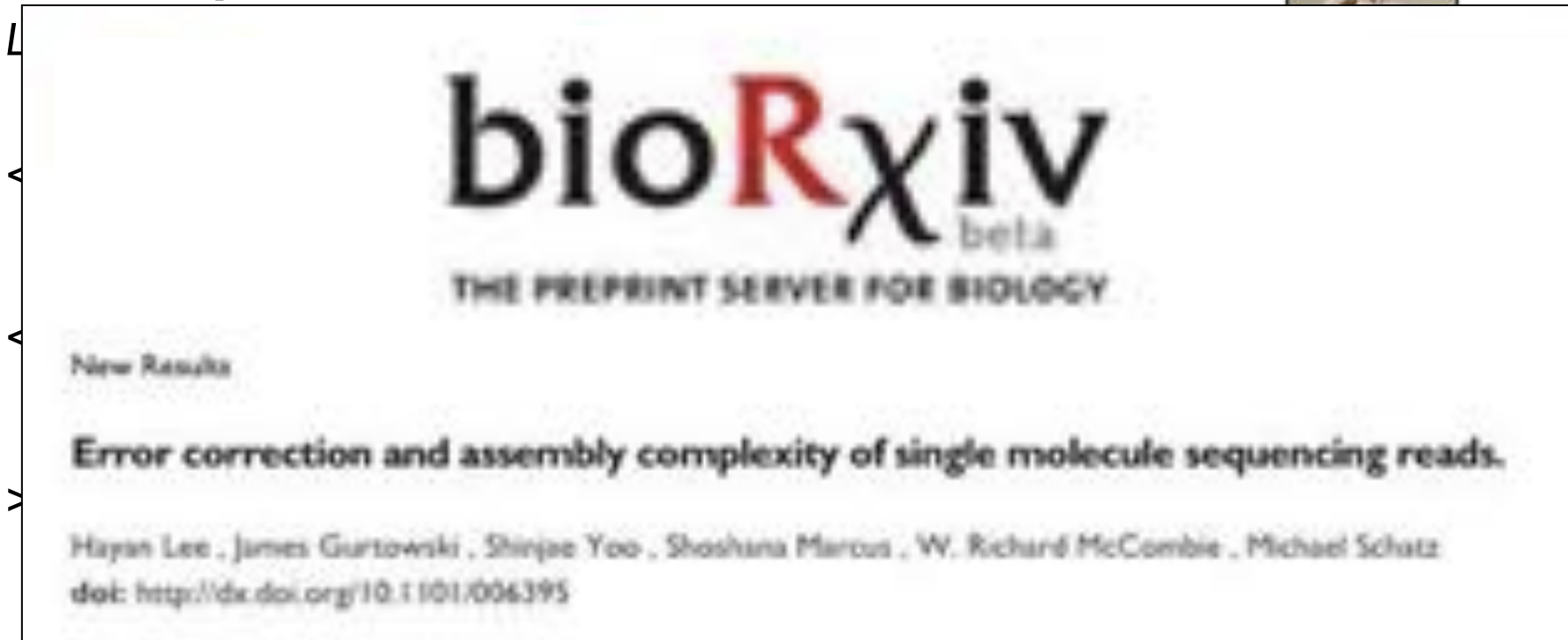- Pathogen detection
- Food safety
- Biomarkers
- etc..

***http://schatzlab.cshl.edu/iGenomics***

# What should we expect from an assembly?

*Summary & Recommendations*



The year 2015 will mark the return to
reference quality genome sequence

# 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
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**
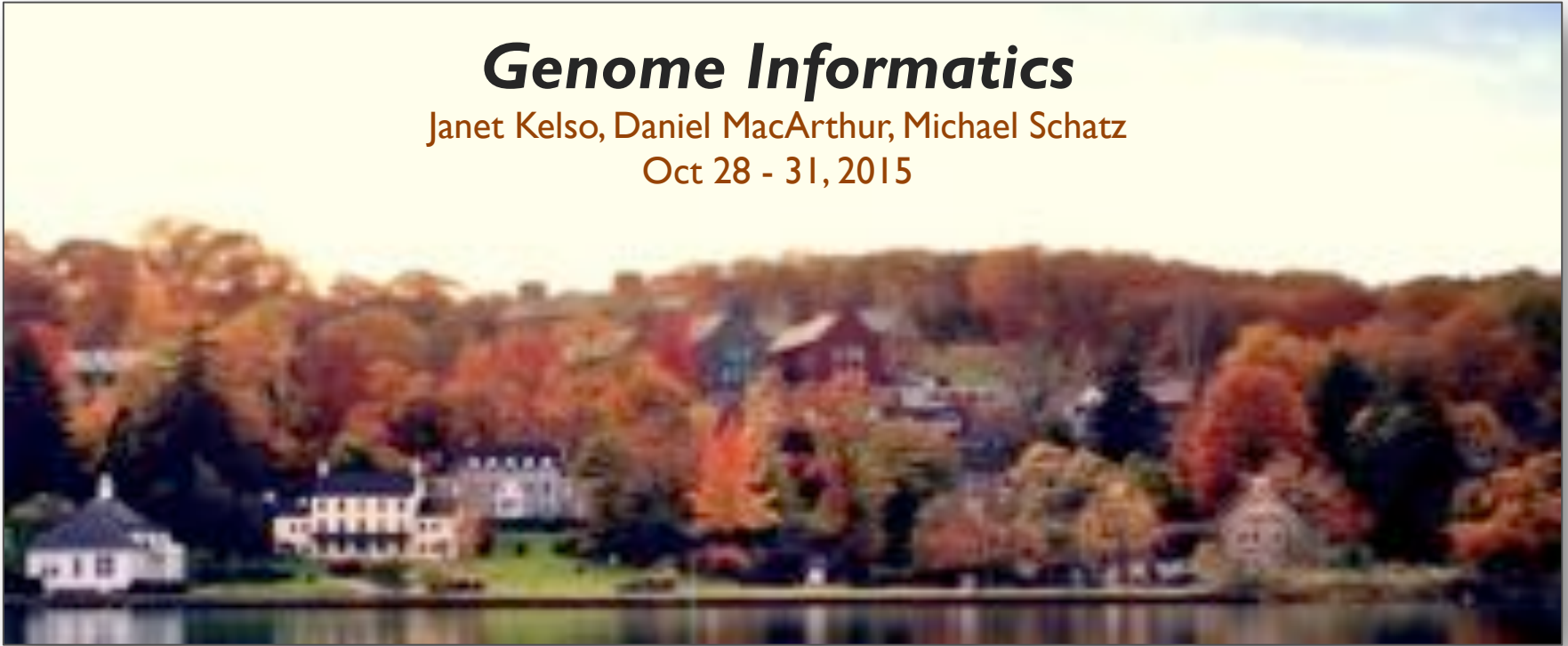John McPherson
Karen Ng
Timothy Beck
Yogi Sundaravadanam

**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