# Comprehensive Genome and Transcriptome Structural Analysis of a Breast Cancer Cell Line using PacBio Long Read Sequencing

## Maria Nattestad

Schatz + McCombie + Hicks at Cold Spring Harbor Laboratory

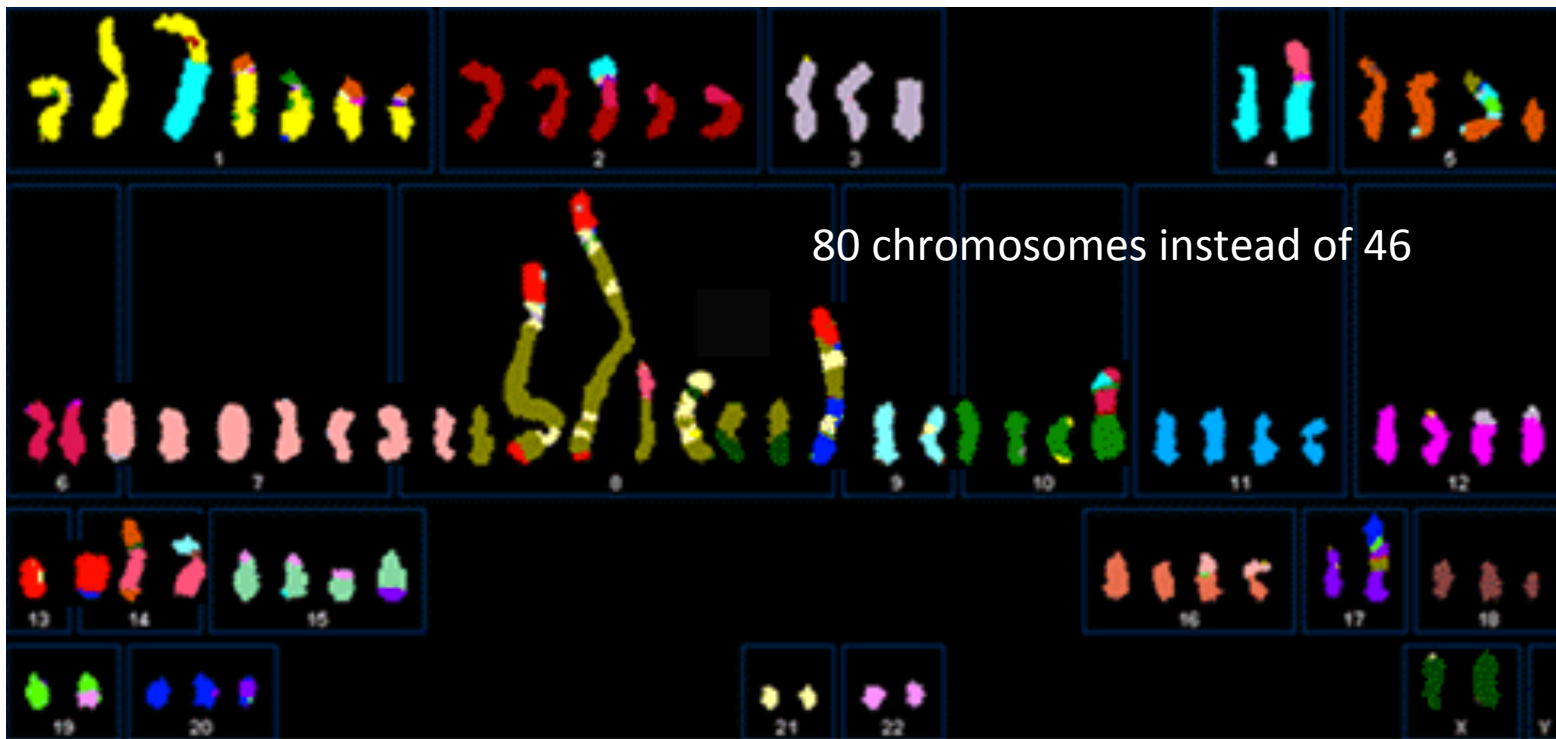McPherson + Beck at the Ontario Institute for Cancer Research
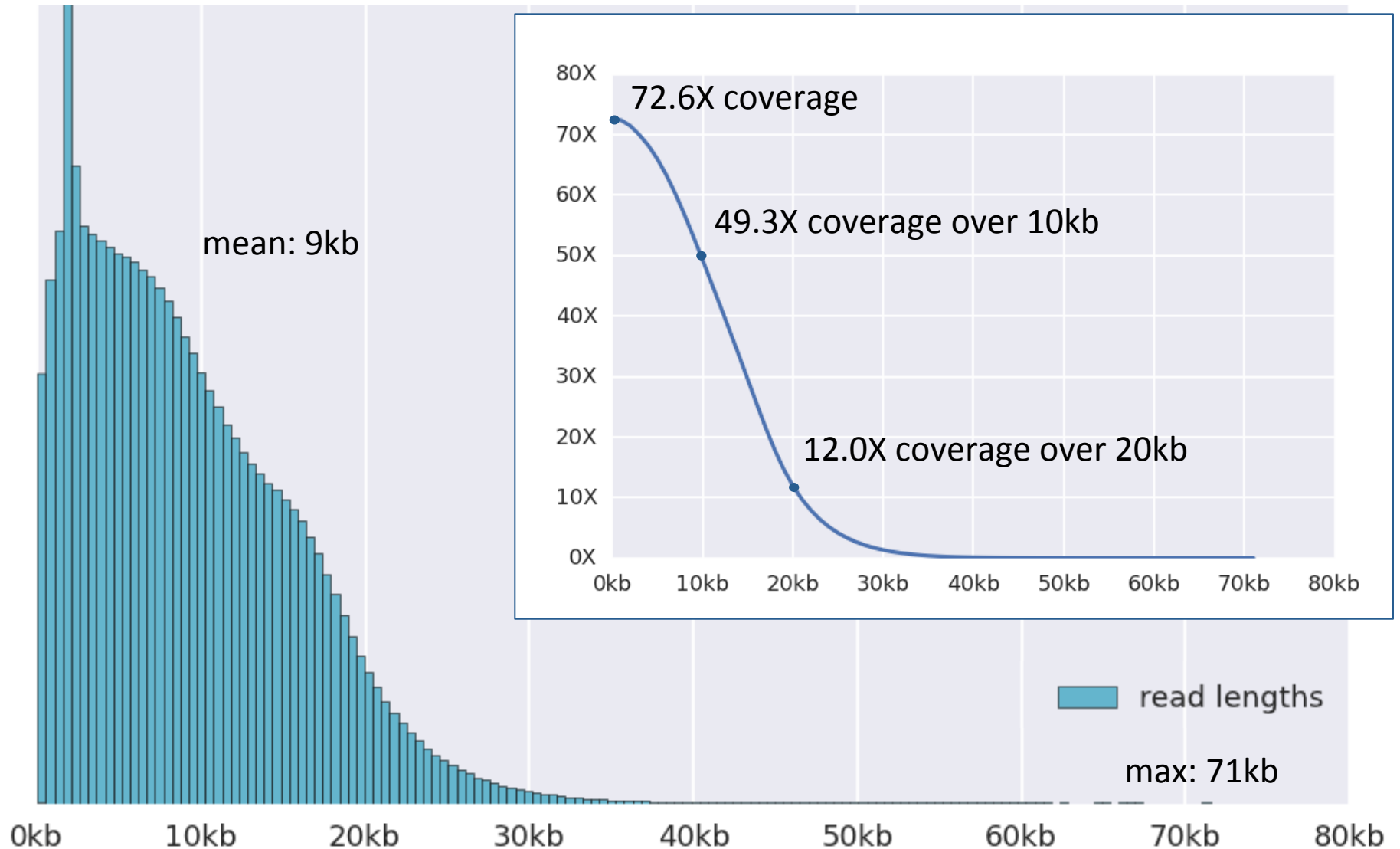
Pacific Biosciences

DNAnexus

CSH

# SK-BR-3

Most commonly used Her2-amplified breast cancer cell line
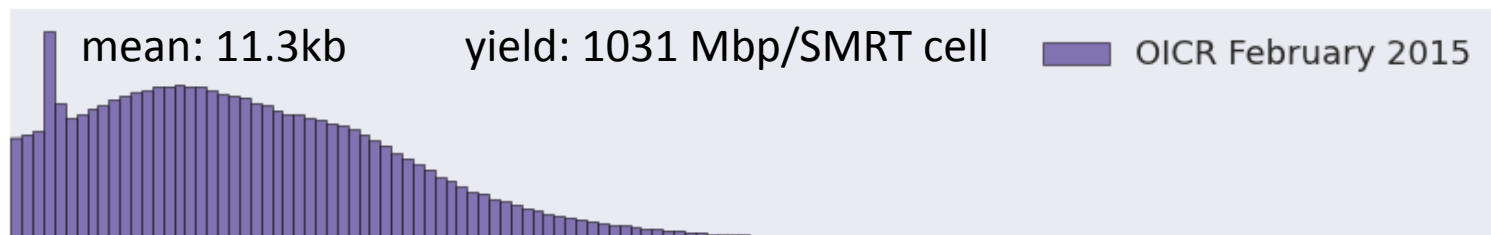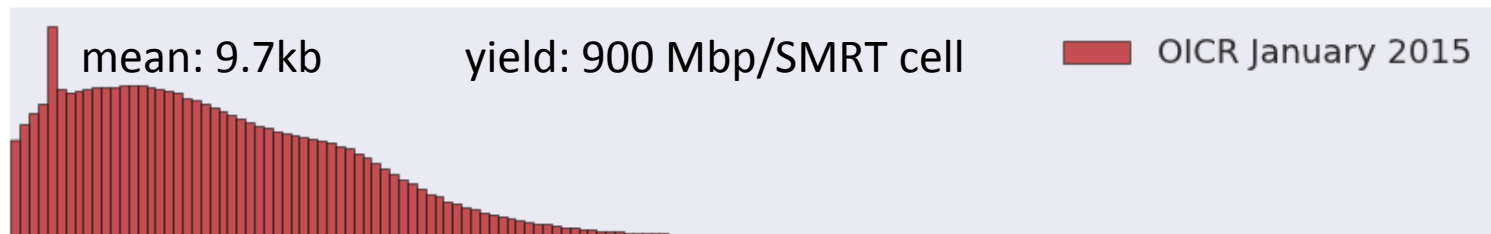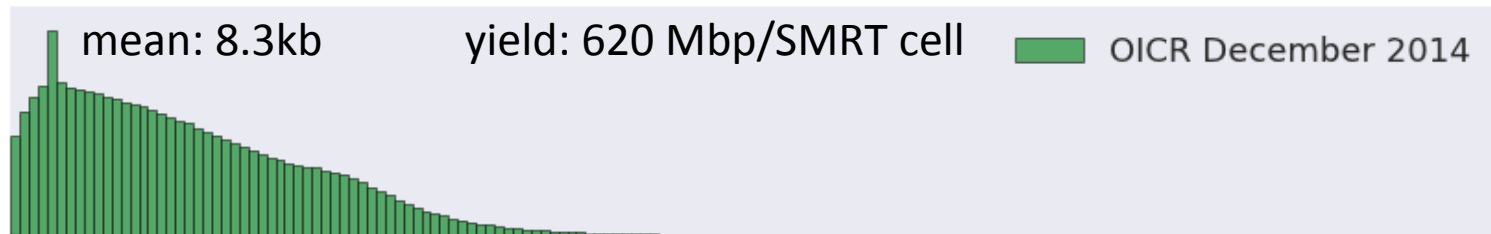
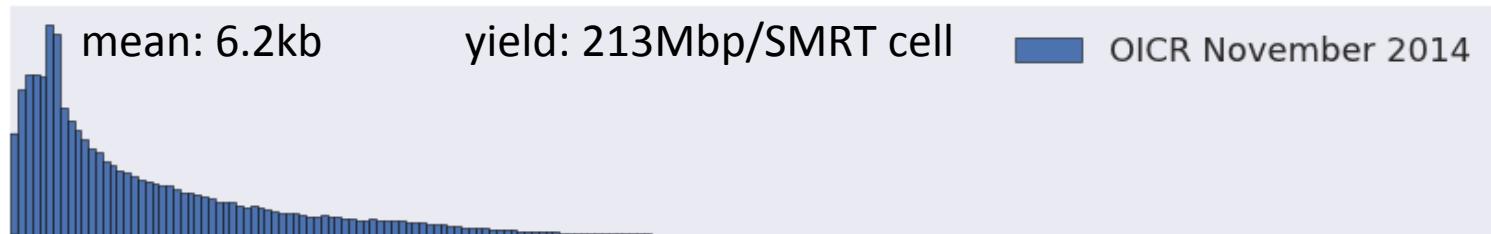80 chromosomes instead of 46

Often used for pre-clinical research on Her2-targeting therapeutics such as Herceptin (Trastuzumab) and resistance to these therapies.

(Davidson et al, 2000)

# Sequencing SK-BR-3:
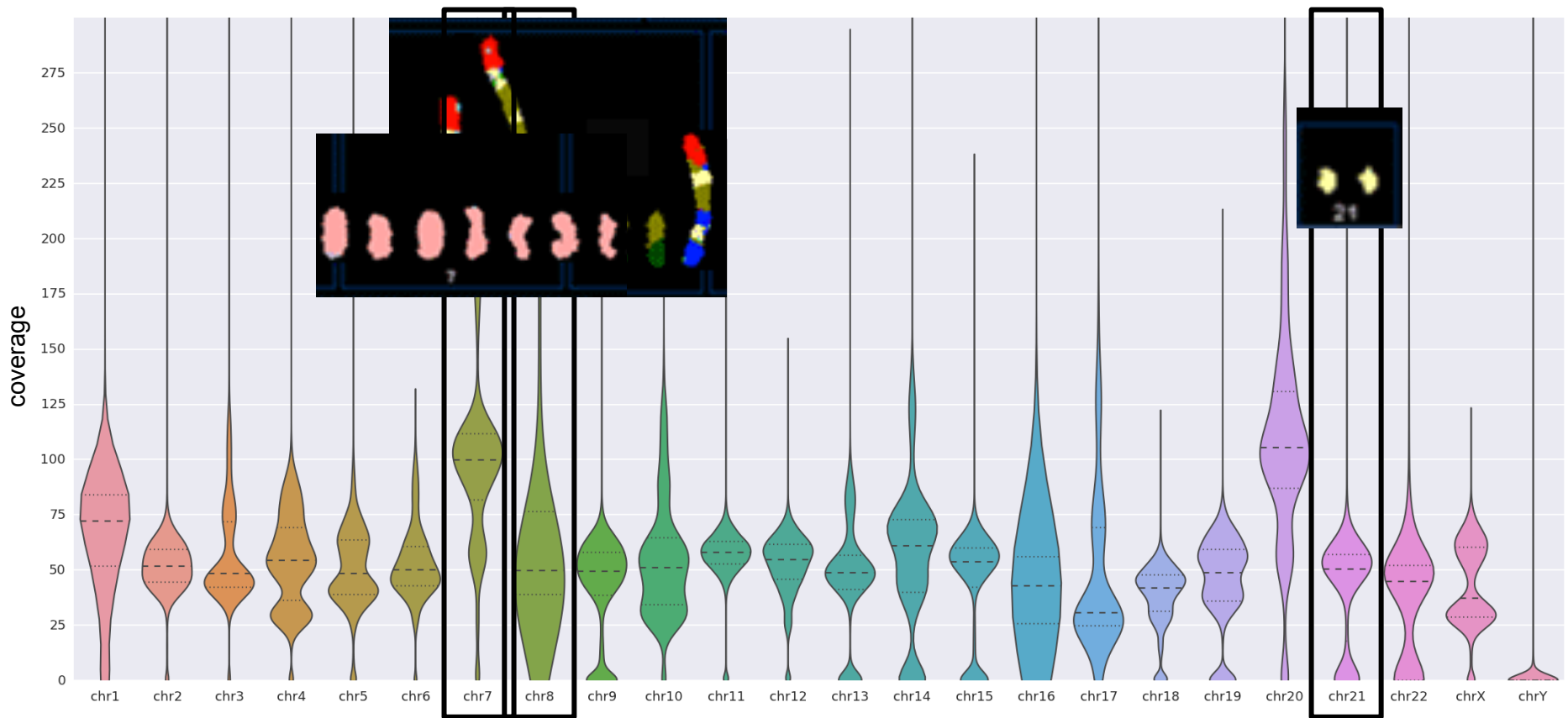# PacBio read length distribution

# Dramatic changes just by experimenting with library preparation



mean: 6.2kb    yield: 213Mbp/SMRT cell    OICR November 2014

mean: 8.3kb    yield: 620 Mbp/SMRT cell    OICR December 2014

mean: 9.7kb    yield: 900 Mbp/SMRT cell    OICR January 2015

mean: 11.3kb    yield: 1031 Mbp/SMRT cell    OICR February 2015

0kb    10kb    20kb    30kb    40kb    50kb    60kb    70kb

# Copy-number analysis is consistent with karyotype results



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

# We could call SNPs if we wanted to

We recovered a known missense mutation in p53: **R175H**

Arg

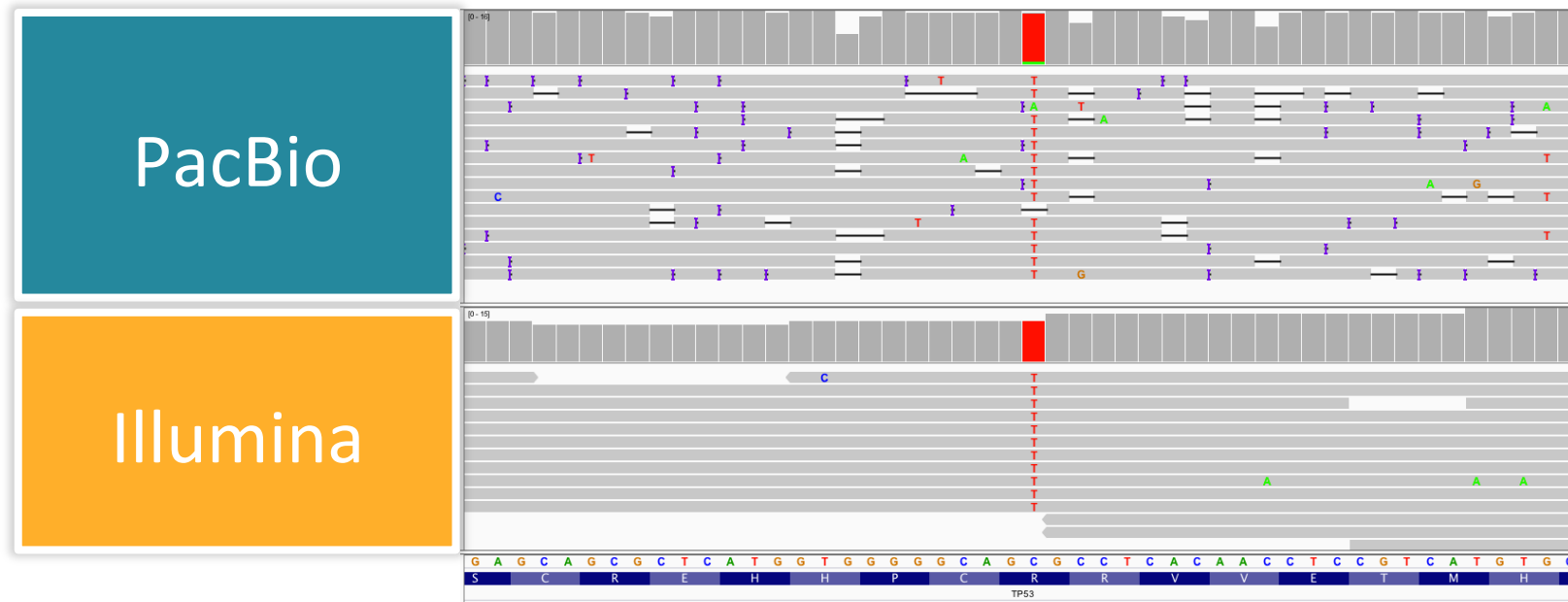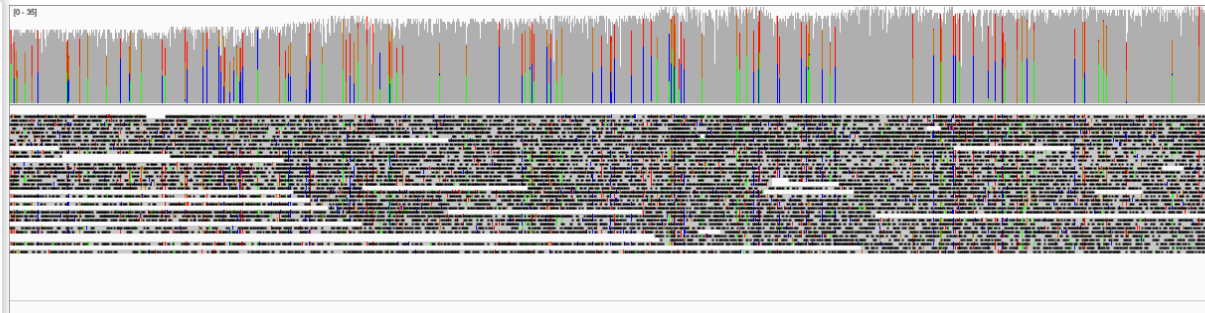| | |
|---|---|
| Reference | ATCTGAGCAGCGCTCATGGTGGGGGCAG**C**GCCTCACAACCTCCGTCATGTGCTGTGACTGCTT |
| Illumina | ATCTGAGCAGCGCTCATGGTGGGGGCAG**T**GCCTCACAACCTCCGTCATGTGCTGTGACTGCTT |
| PacBio | ATCTGAGCAGCGCTCATGGTGGGGGCAG**T**GCCTCACAACCTCCGTCATGTGCTGTGACTGCTT |

His



**PacBio**

**Illumina**

**TP53 gene**

| | |
|---|---|
| Insertion rate | 11.5% |
| Deletion rate | 3.4% |
| Mismatch rate | 1.4% |

# PacBio reads are longer and less susceptible to mapping issues



**HLA-A gene**

# Genome structural analysis

**Assembly-based**

Assembly with Falcon on DNAnexus

Alignment with MUMmer

Call variants between consecutive alignments with **ABVC**

Call variants within alignments with **ABVC**

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

**Alignment-based**

Alignment with BWA-MEM

Copy number analysis

SV-calling from split reads with **Sniffles**

Validations

**SplitThreader**

Detailed analysis of Her2 amplifications

661 long-range variants
(>10kb distance)

# Genome assembly with FALCON

Find overlaps between PacBio reads

Error-correct reads against each other
• High enough coverage to self-correct so we don't need Illumina reads

Filtering to avoid confounding by repeats

Constructing a graph from the overlaps

Constructing contigs from graph

# Iterations of Falcon assembly on DNAnexus

| Standard FALCON parameters | ⇒ | Higher allowance for coverage differences | ⇒ | Require only 2X error-corrected coverage |
|---|---|---|---|---|

Repeat filters in FALCON also remove aneuploid regions

Diploid regions have ~25X coverage --> low after error correction, especially if heterozygous

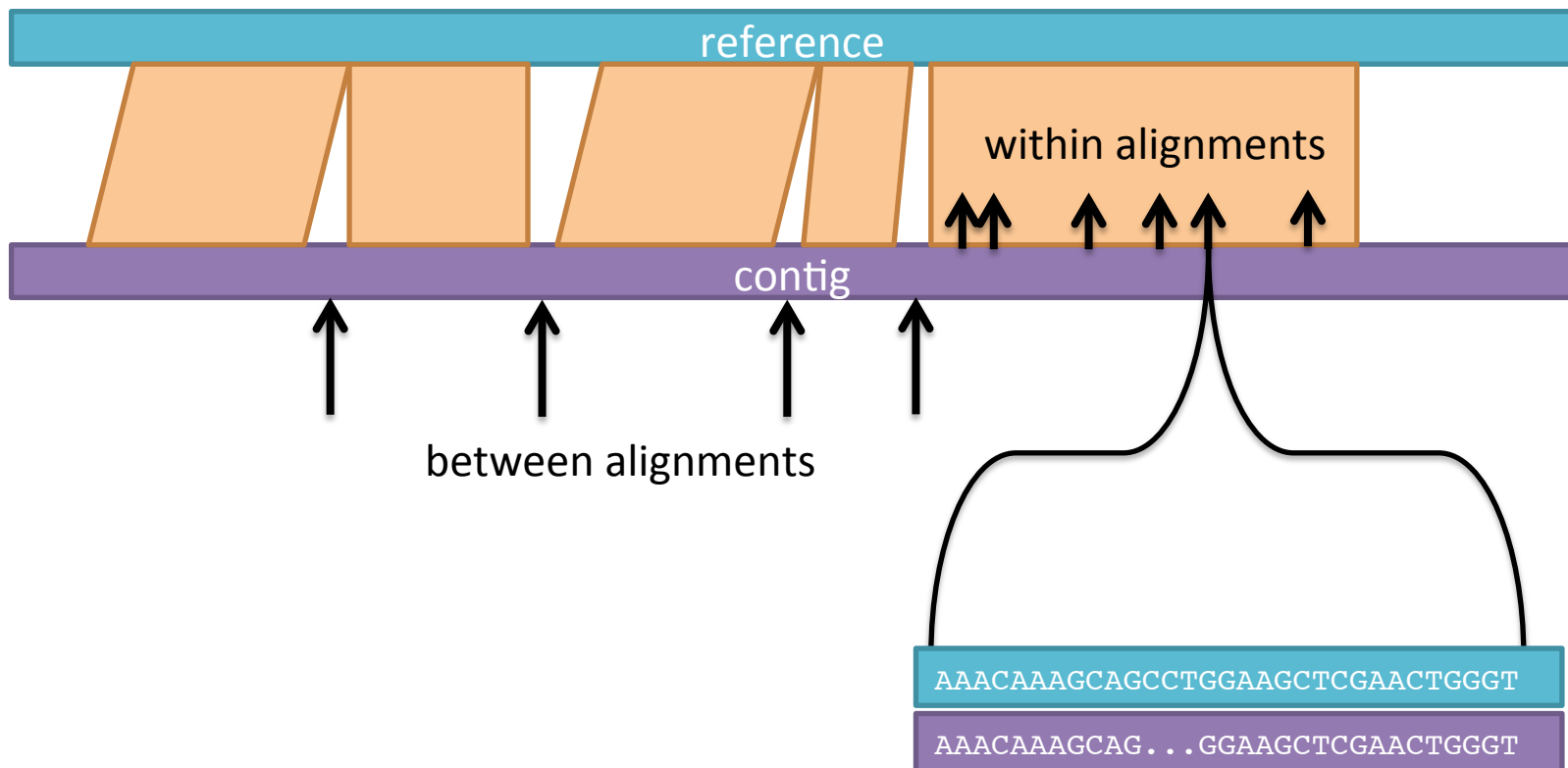# Assembly using PacBio yields far better contiguity



Number of sequences: 13,532
Total sequence length: 2.97Gb
Mean: 266 kb
Max: 19.9 Mb
**N50: 2.46 Mb**

Relative to a genome size of 3 Gb

Number of sequences: 748,955
Total sequence length: 2.07 Gb
Mean: 2.8 kb
Max: 61 kb
**N50: 3.3 kb**

# Variant detection from a genome assembly

reference

within alignments

contig

between alignments

AAACAAAGCAGCCTGGAAGCTCGAACTGGGT

AAACAAAGCAG...GGAAGCTCGAACTGGGT

# ABVC: Variants within alignments

Insertion

| reference | AAACAAAGCAG...CCTGGGAAGCTCGAACTGGGT |
|---|---|
| contig | AAACAAAGCAGTACCCTGGGAAGCTCGAACTGGGT |

Deletion

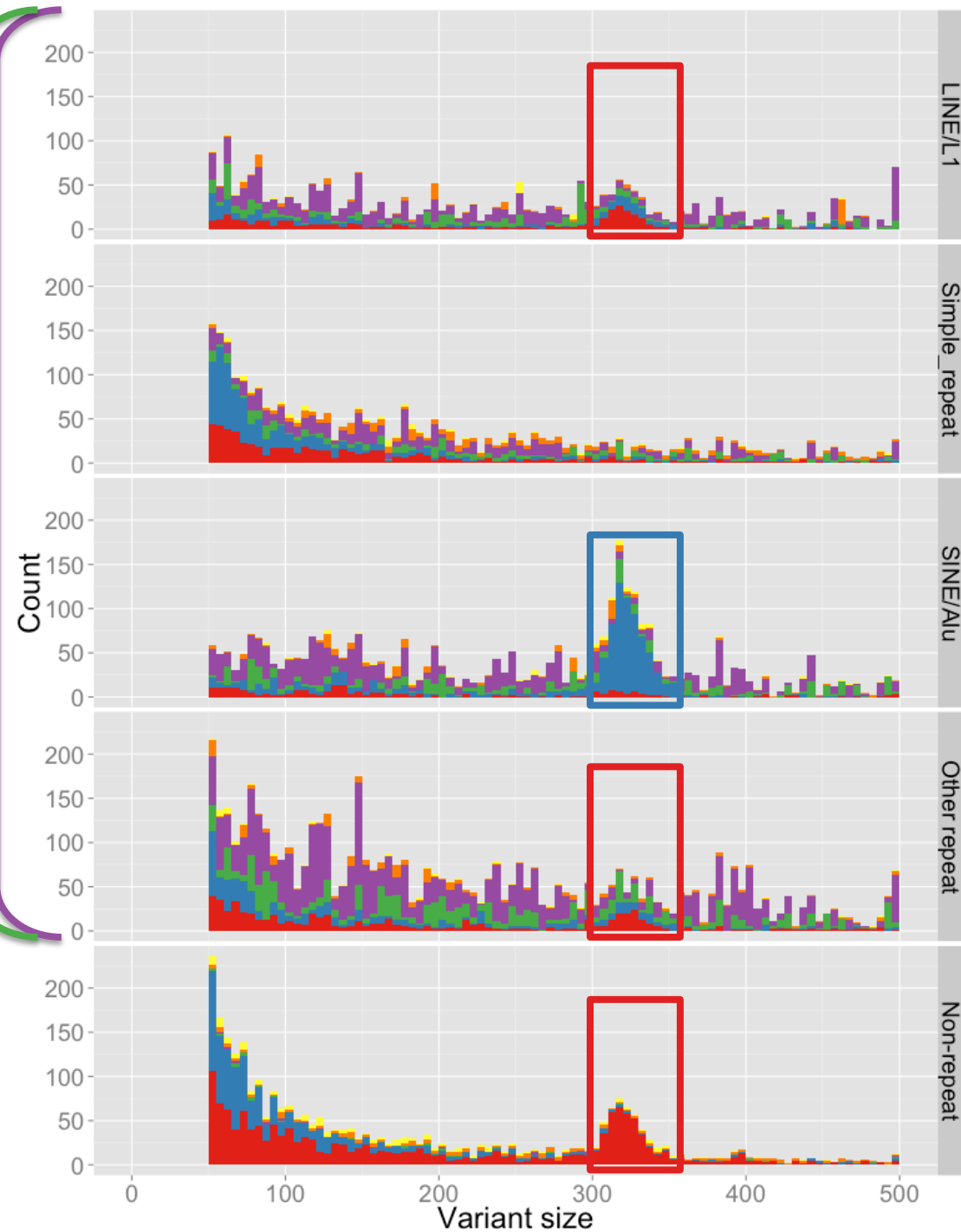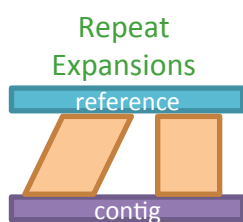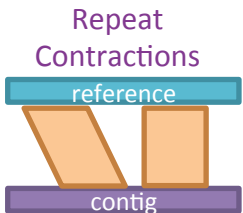| reference | AAACAAAGCAGCCTGGAAGCTCGAACTGGGT |
|---|---|
| contig | AAACAAAGCAG...GGAAGCTCGAACTGGGT |

# ABVC: Unique length filtering is needed to prevent false positives due to repetitive elements

All alignments

reference

contig

Repeat maps in multiple locations

Traditional filtering by MUMmer
- choose the best alignment for each query
- random choice if multiple alignments with the same score

reference

contig

May be falsely called as a translocation

# ABVC: Unique length filtering is needed to prevent false positives due to repetitive elements

All alignments

reference

contig

Repeat maps in multiple locations

Filtering based on requiring a certain length of unique alignment

reference

contig

All repeats with less than 10kb of unique sequence are filtered out

# Types of variants detected by ABVC



**Defined point**

Insertion

Deletion

**Overlapping alignments suggest tandem repeat**

Tandem Expansions

Tandem Contractions

**Gap where sequences do not align uniquely suggests a repeat**

Repeat Expansions

Repeat Contractions

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

Repeat Contractions

reference

contig

Repeat Expansions

reference

contig

BLASTed 515 insertions: 427 (83%) of them matched Alu elements

**Variant type**
- Insertion
- Deletion
- Repeat expansion
- Repeat contraction
- Tandem expansion
- Tandem contraction

LINE/L1

Simple_repeat

SINE/Alu

Other repeat

Non-repeat

Count

Variant size

# Why assembly doesn't capture long-range variants

Normal chromosome

A → B → C → D

Translocation

X → Y

Turn assembly graph into a fasta file of contigs

Other chromosome, now connected

A → B

C → D

X → Y

Major SVs are not contained within any contigs

# Genome structural analysis

## Assembly-based

Assembly with Falcon on DNAnexus

Alignment with MUMmer

Call variants between consecutive alignments with **ABVC**

Call variants within alignments with **ABVC**

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

## Alignment-based

Alignment with BWA-MEM

Copy number analysis

SV-calling from split reads with **Sniffles**

Validations

**SplitThreader**
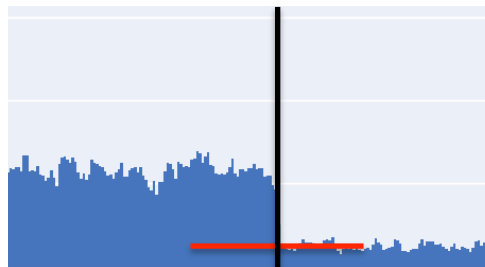
Detailed analysis of Her2 amplifications

661 long-range variants
(>10kb distance)

# Variant-calling from split-read alignment



Software: Sniffles by Fritz Sedlazeck

# Long-range structural variants found by Sniffles

Her2 oncogene



661 long-range variants
(>10kb distance)

# Chromosome 8 has the most intra- and inter- chromosomal long-range variants

# Long-range structural variants found by Sniffles



Her2 oncogene

661 long-range variants
(>10kb distance)

# SplitThreader:
# Graphical threading to retrace complex history of rearrangements in cancer genomes

800
600 Her2
GSDMB
400 RARA
200

36 Mb                                                                                    41 Mb

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

# Bridging the gap



variant-calling from sequencing

?

big-picture view from karyotyping

Context

Resolution

# Threading through the whole-genome graph to produce a synthetic karyotype

# Synthetic karyotype with SplitThreader

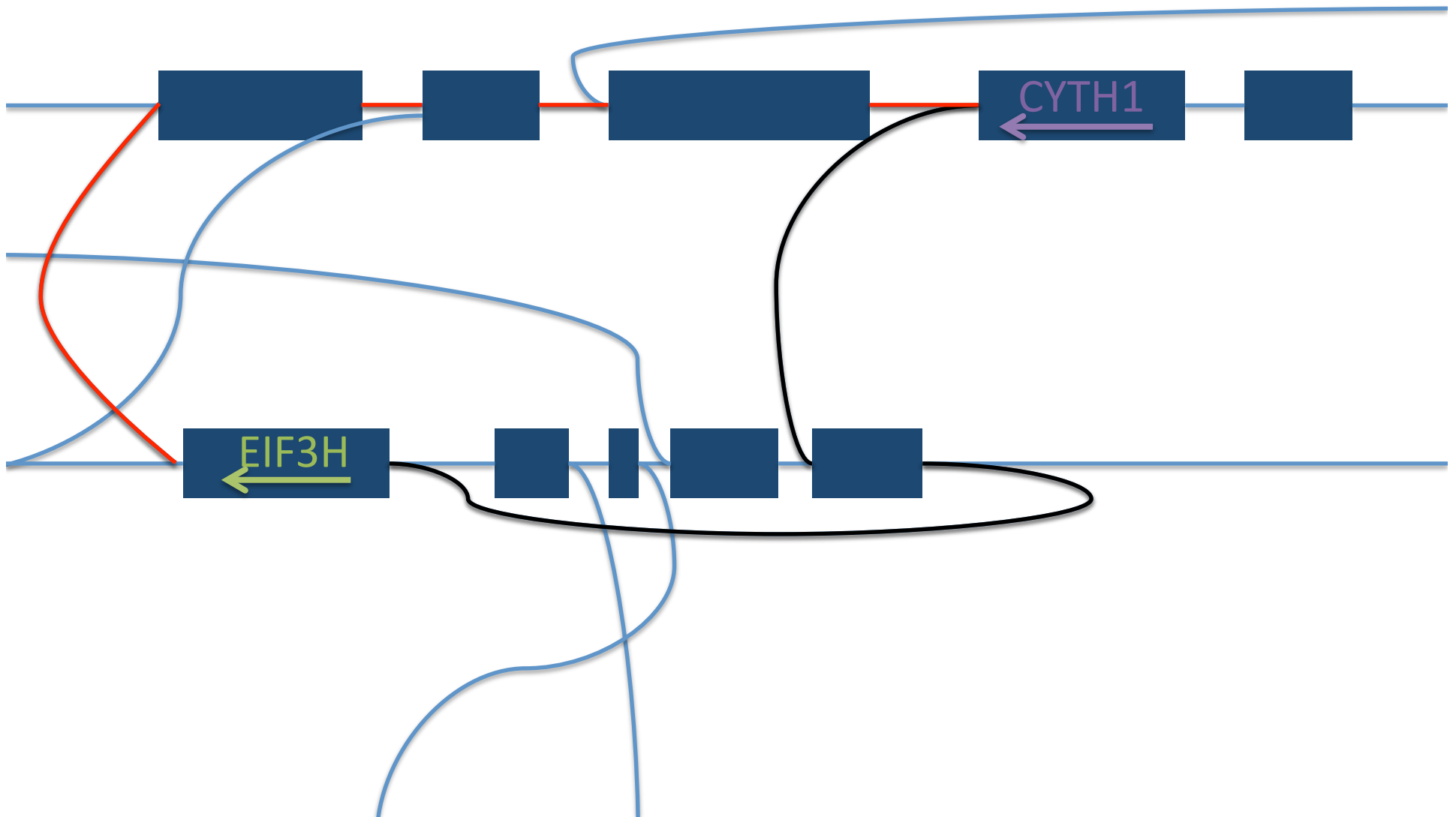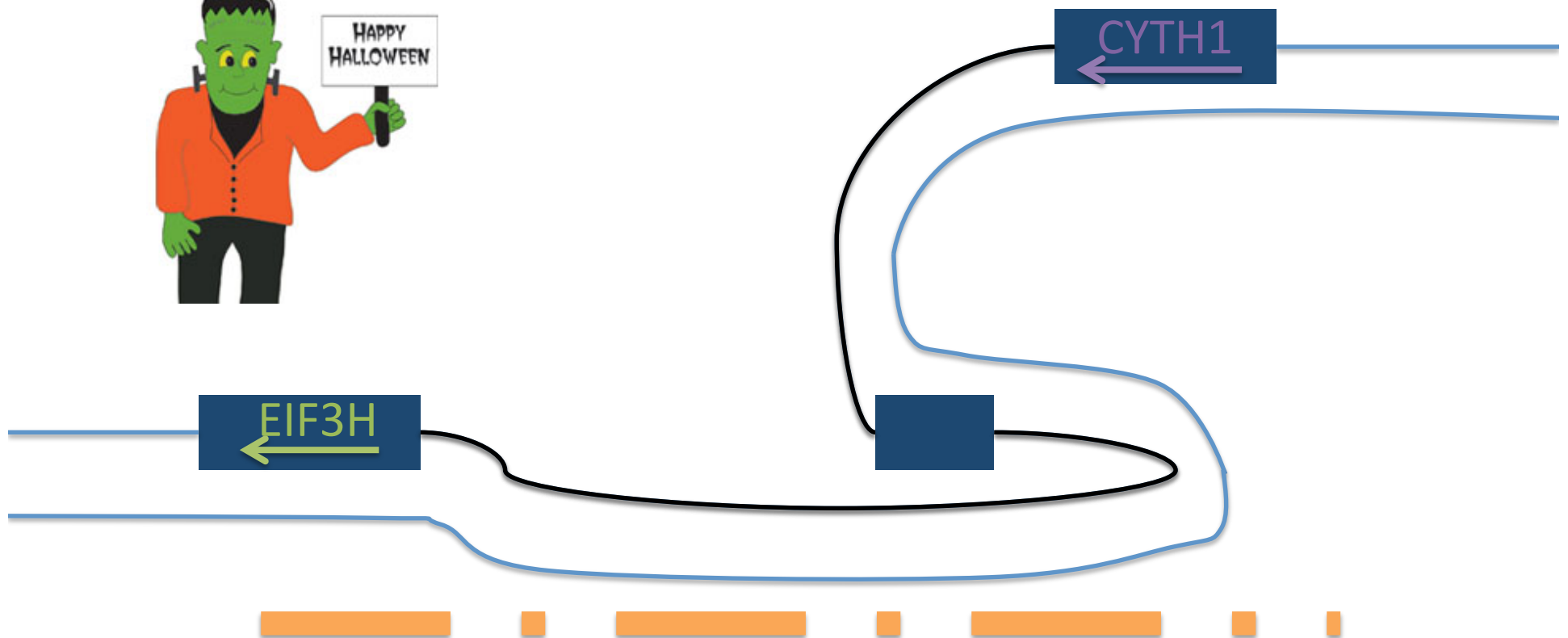Preliminary SplitThreader synthetic karyotype



Real karyotype

# Transcriptome analysis with IsoSeq: Long-read RNA sequencing

- Full-length transcripts
- Found 17 gene fusions with both DNA and RNA evidence
  - 13 seen in previous RNA-seq literature
  - 4 novel fusions
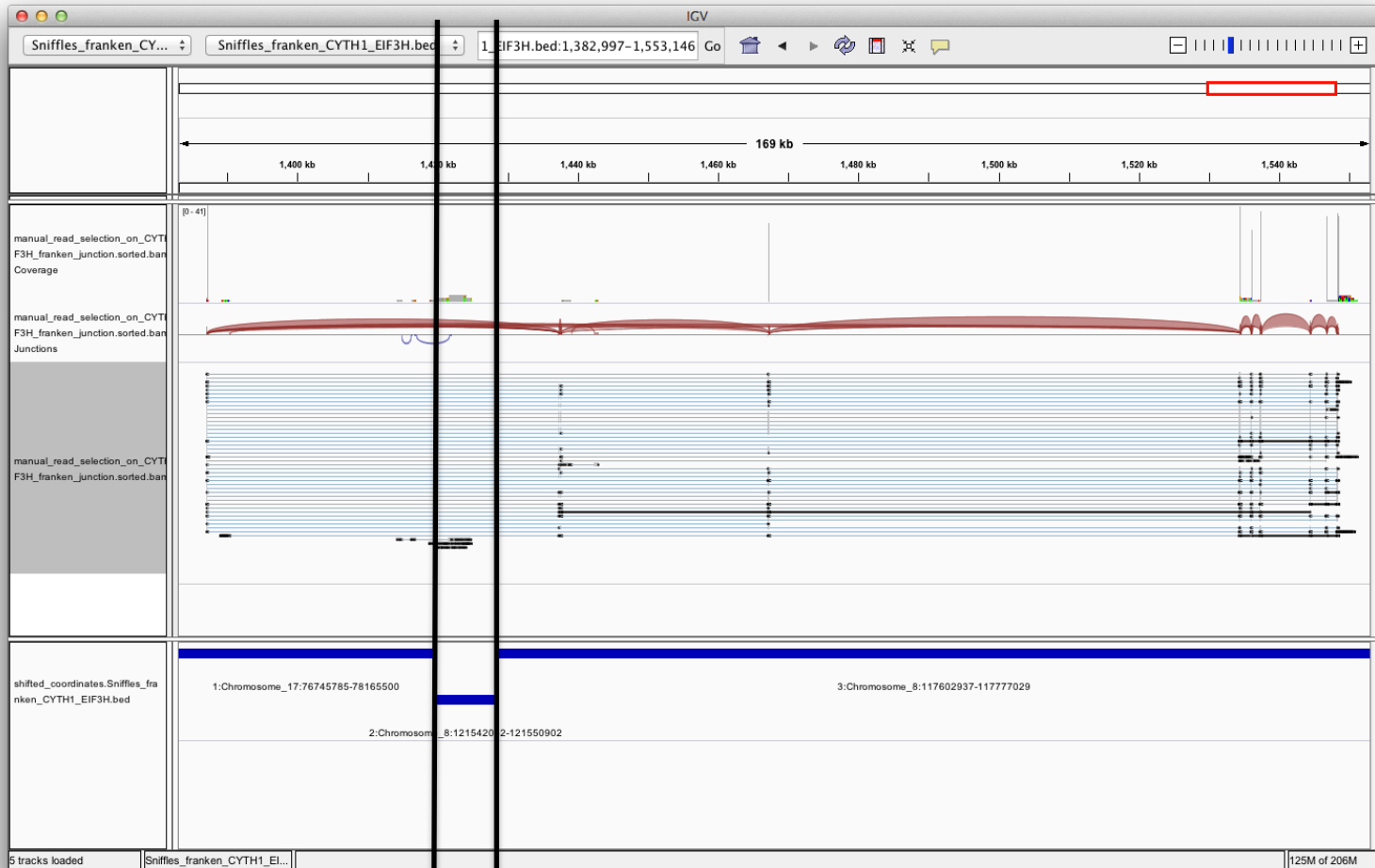- 2 previously observed fusions had RNA evidence but no direct link in the DNA
  - Confirmed using SplitThreader
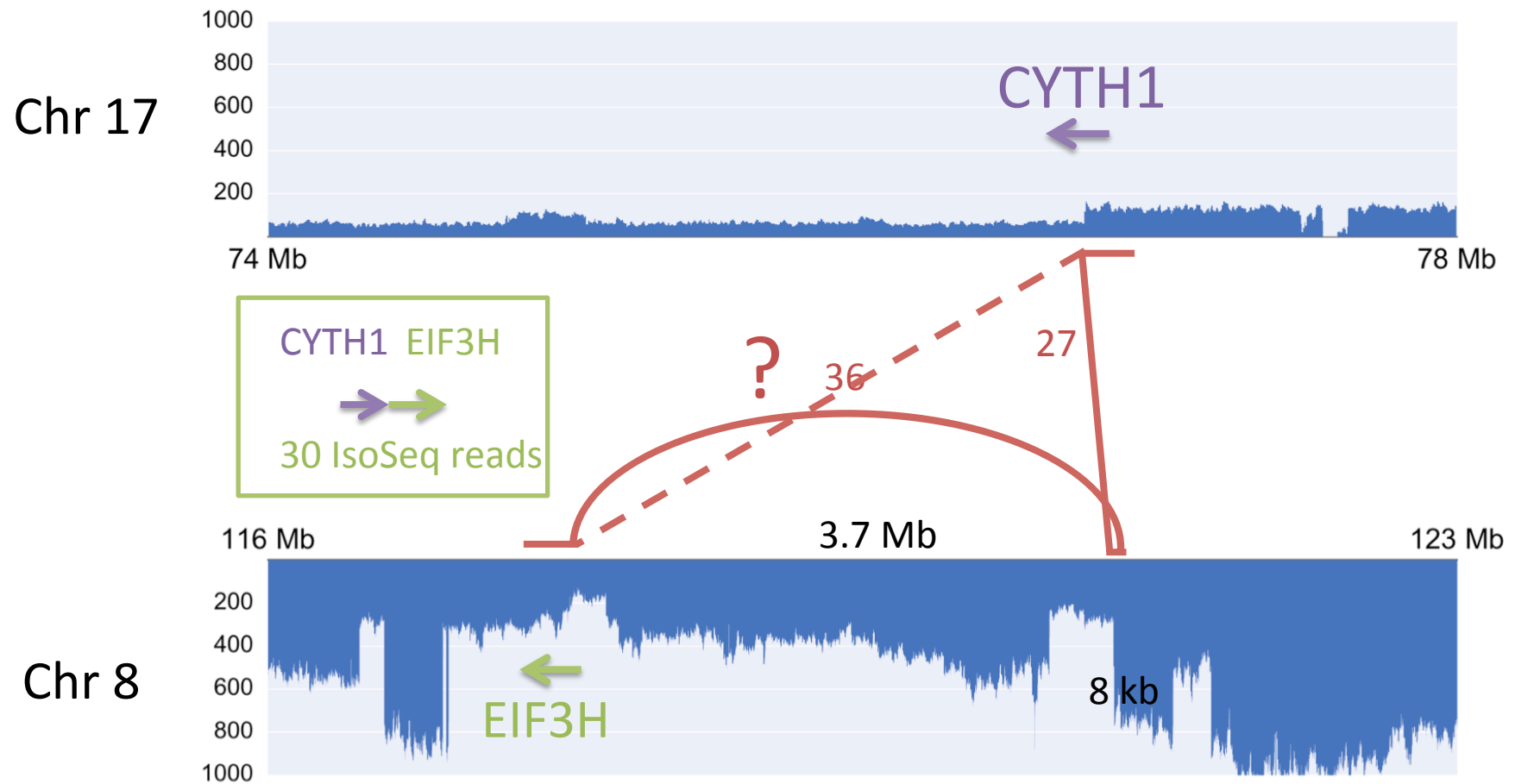
# CYTH1-EIF3H gene fusion in the SplitThreader graph

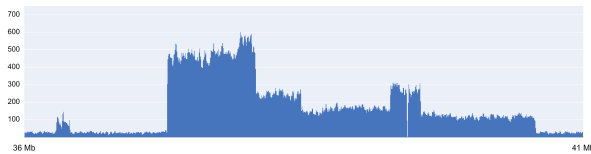# CYTH1-EIF3H gene fusion in the SplitThreader graph

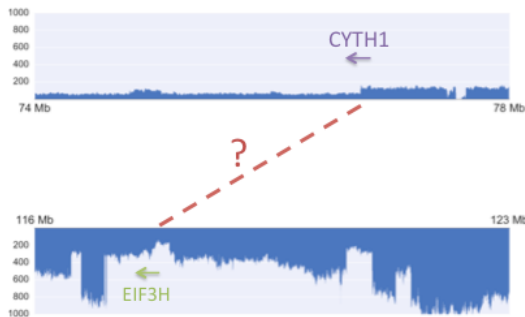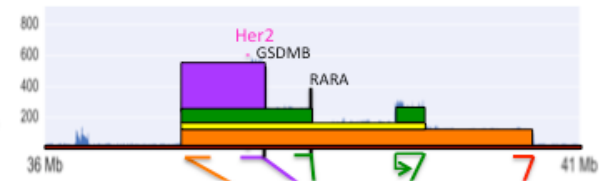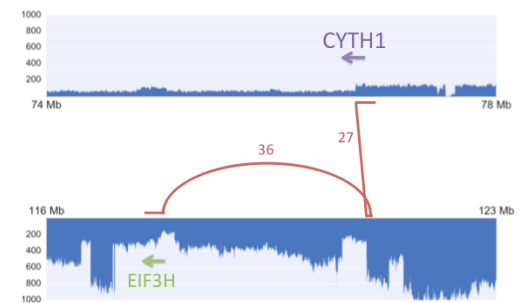# Frankensteining the CYTH1-EIF3H gene fusion

# CYTH1-EIF3H gene fusion

# The genome informs the transcriptome



Explain amplifications

Trace gene fusions

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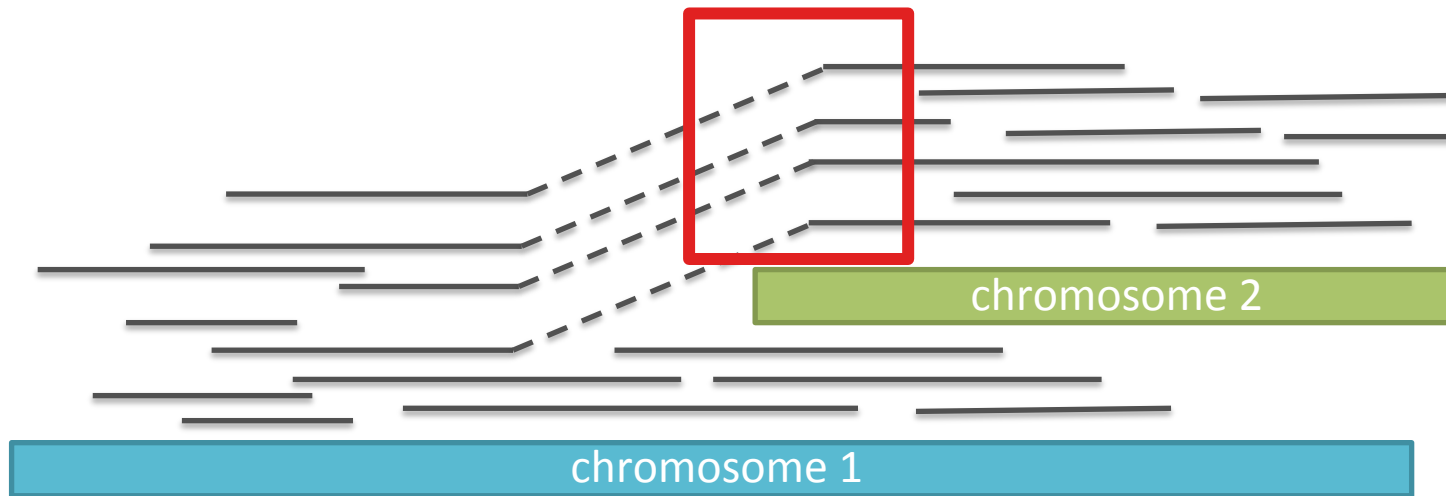
PacBio

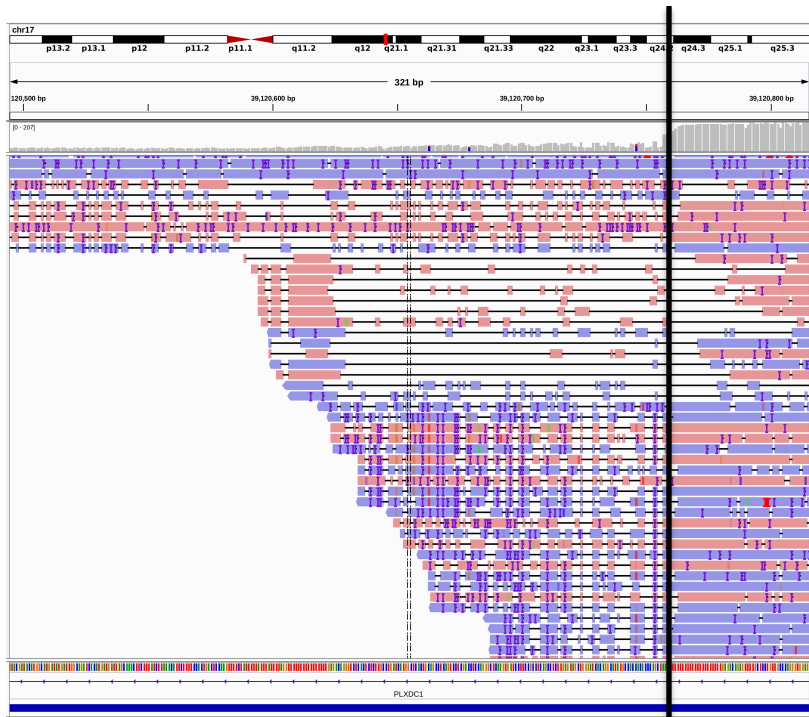Illumina

# Zooming in on the breakpoint

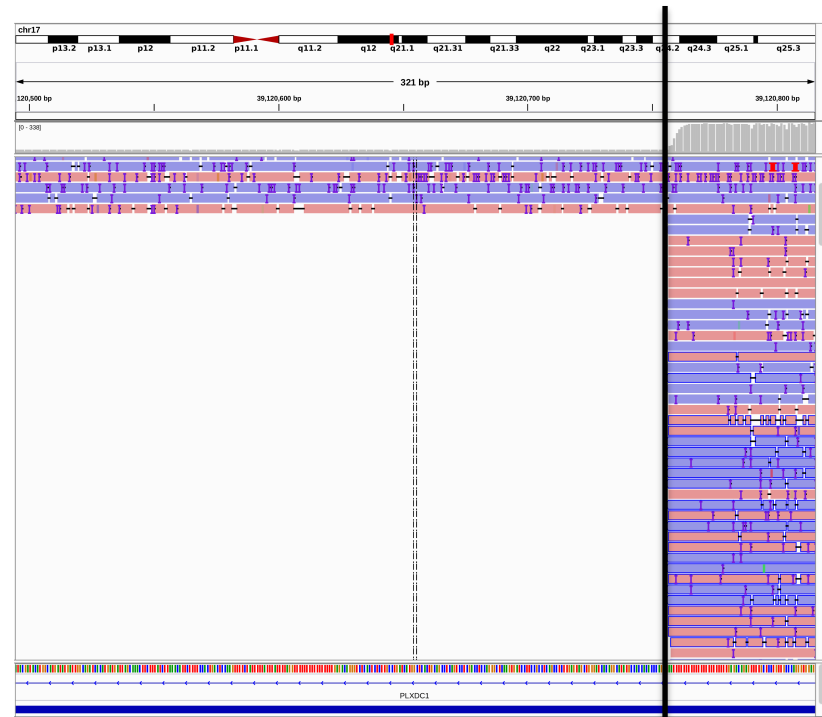# New aligner NGM-LR narrows down the breakpoint to base-pair resolution
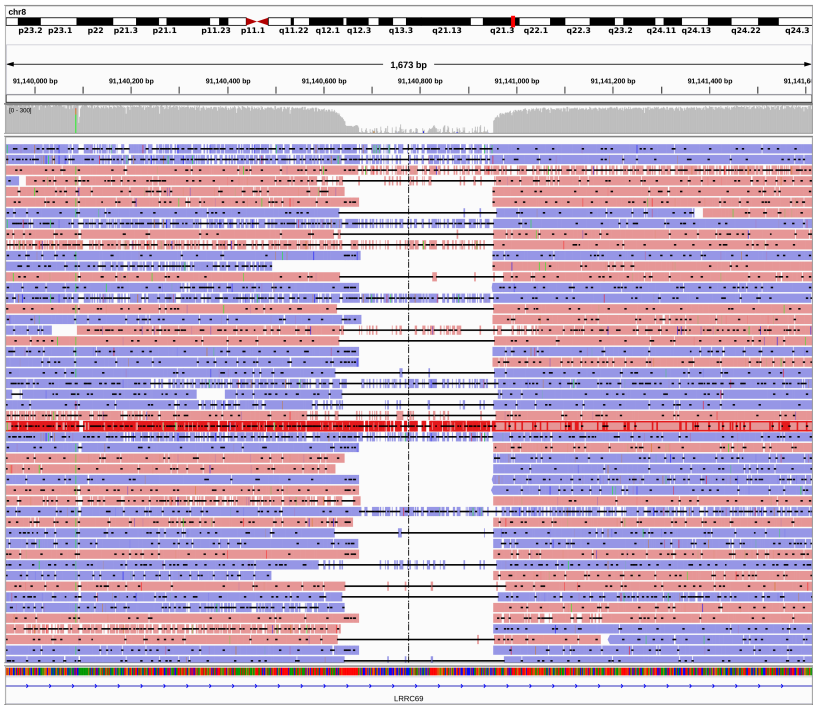
Philipp Rescheneder

## BWA-MEM
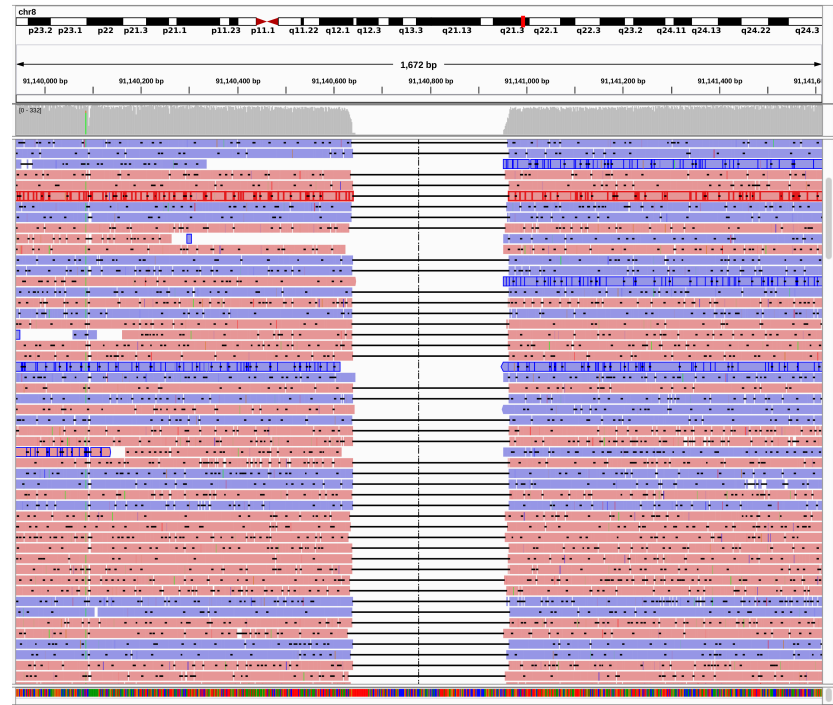
## NGM-LR

One side of an interchromosomal translocation

# NGM-LR also enables better small variant calling
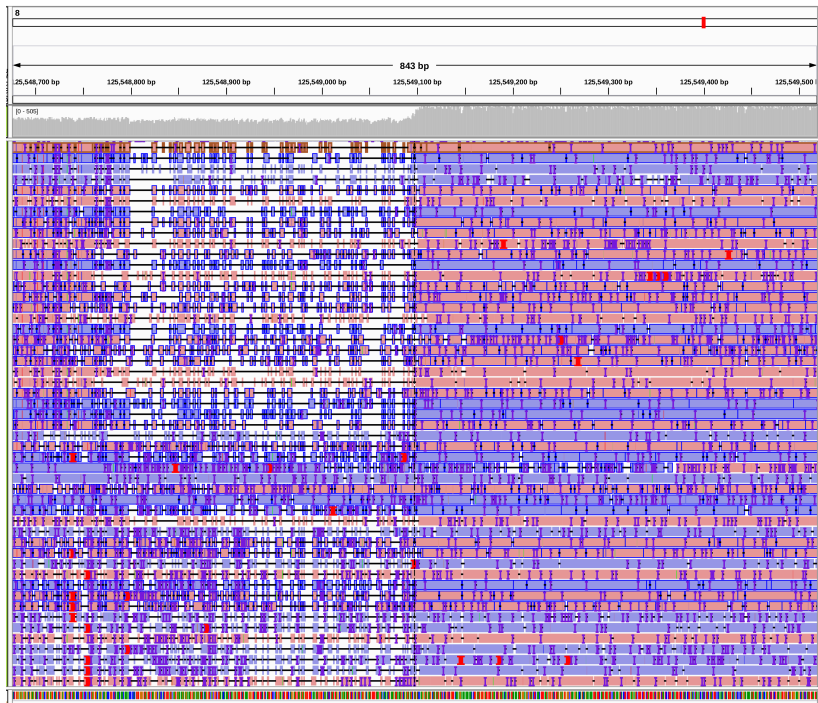
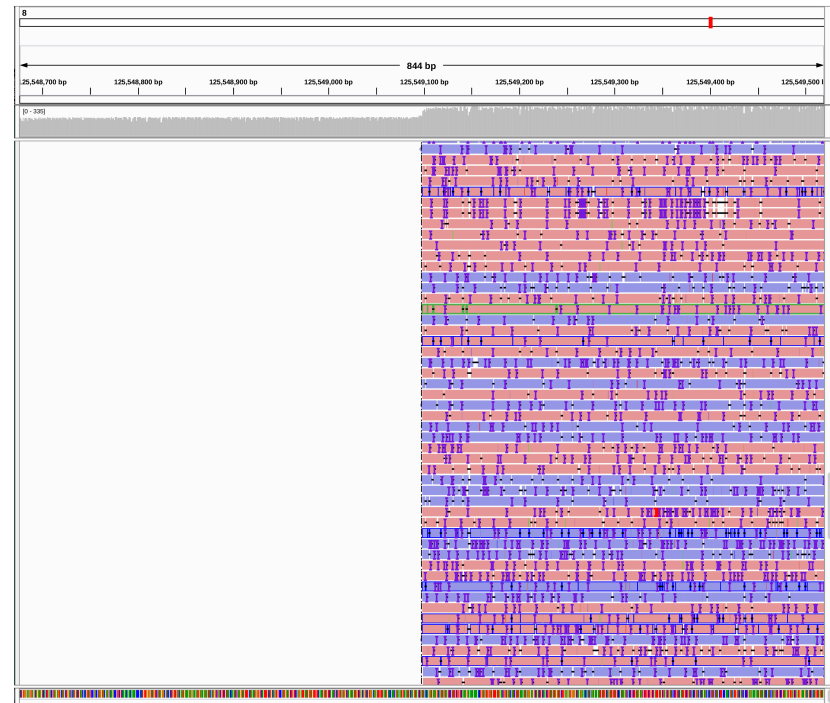## BWA-MEM

## NGM-LR



deletion

deletion

# Without NGM-LR, alignments can be smudged over hundreds of base-pairs away from the breakpoint
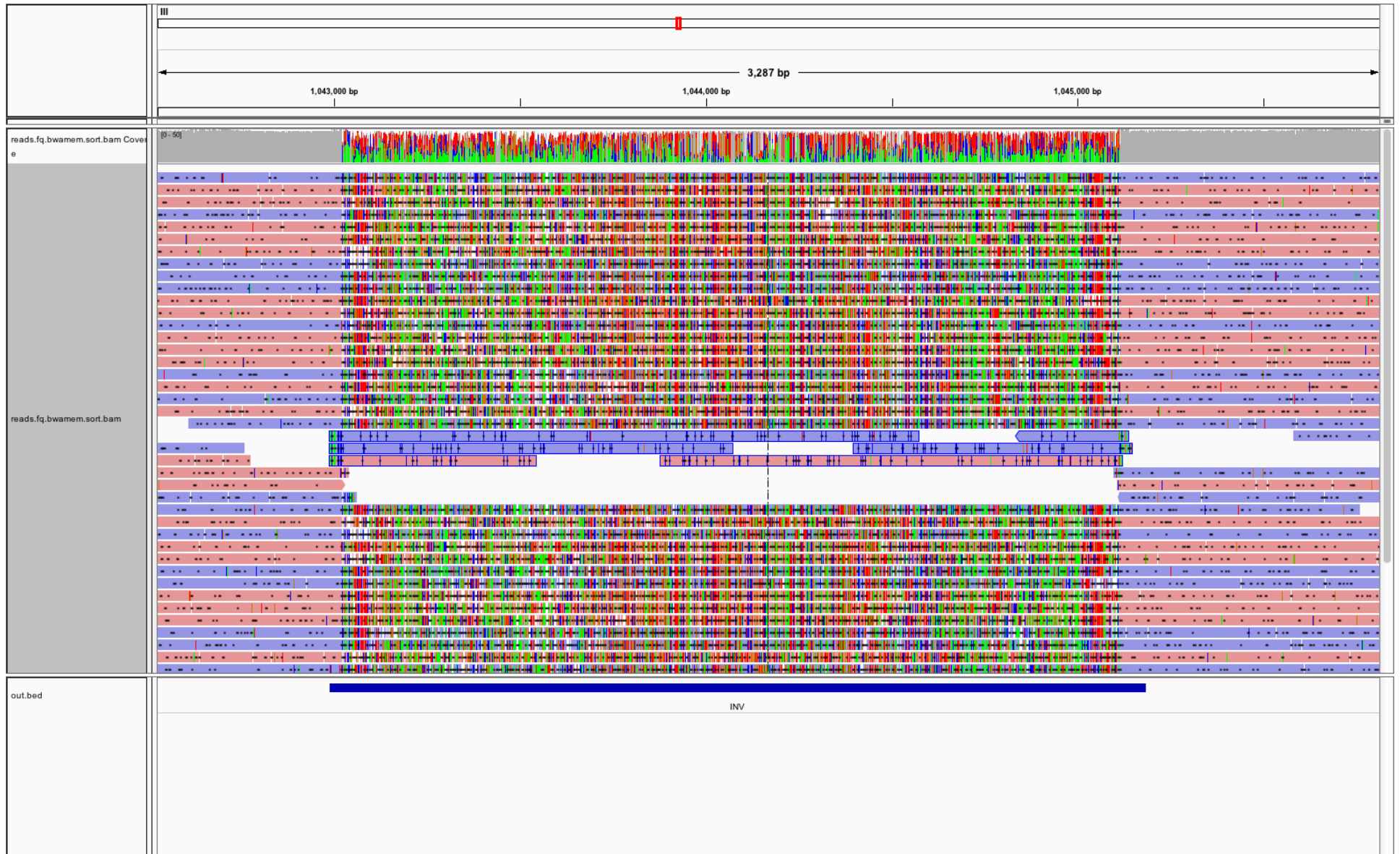
## BWA-MEM

## NGM-LR



translocation

translocation

# Inversion in BWA-MEM

# Acknowledgments

**Cold Spring Harbor Laboratory**

Sara Goodwin
**Fritz Sedlazeck = Sniffles**
**Philipp Rescheneder = NGM-LR**
Timour Baslan
Tyler Garvin
Han Fang
James Gurtowski
Elizabeth Hutton
Marley Alford
Melissa Kramer
Eric Antoniou
James Hicks
Michael Schatz
W. Richard McCombie

**OICR — Ontario Institute for Cancer Research**

Karen Ng
Timothy Beck
Yogi Sundaravadanam
John McPherson

**DNAnexus**

**PACIFIC BIOSCIENCES®**

Elizabeth Tseng
Jason Chin

**Genentech FOUNDATION**