# 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

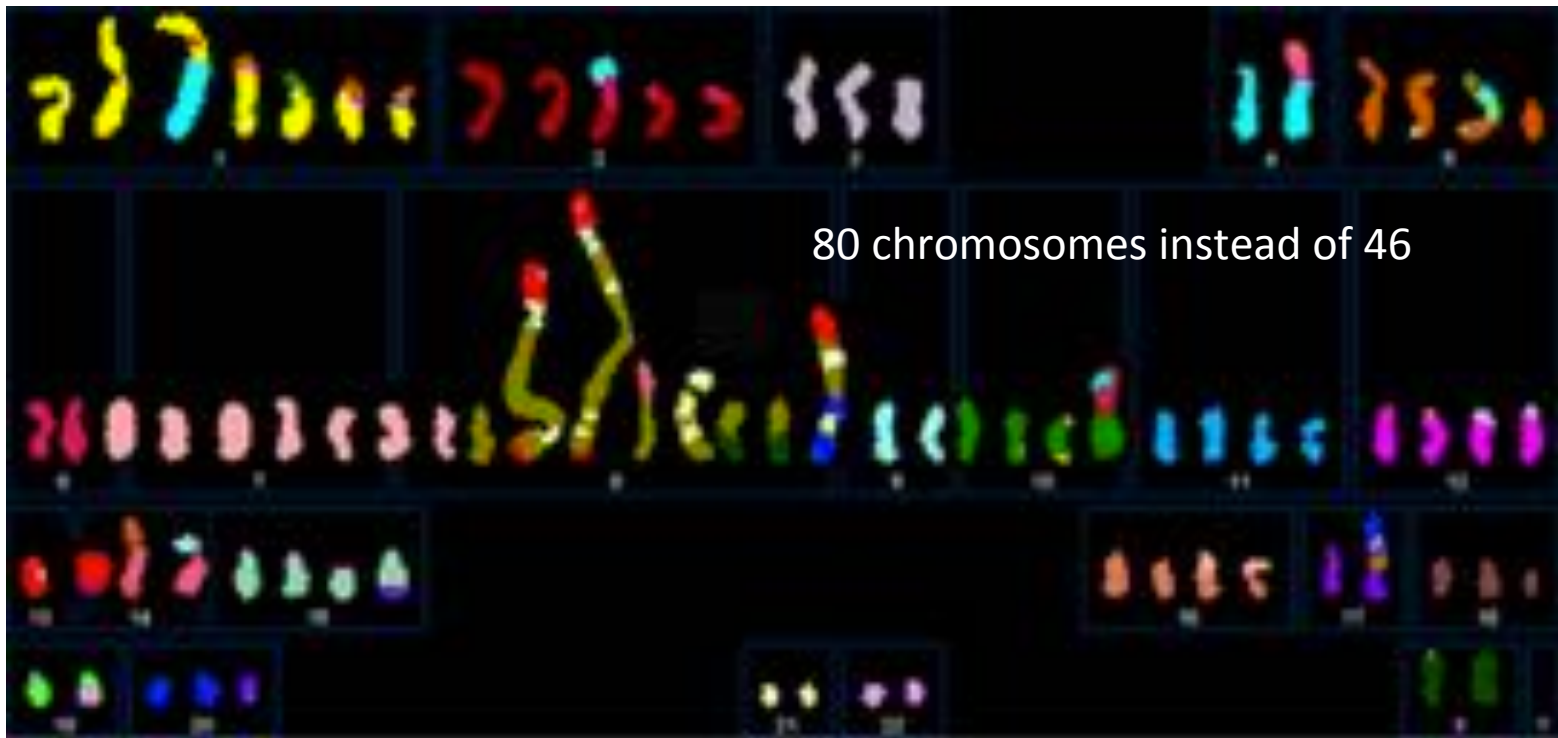McPherson + Beck at the Ontario Institute for Cancer Research

Pacific Biosciences

DNAnexus

# SK-BR-3

Most commonly used Her2-amplified breast cancer cell line
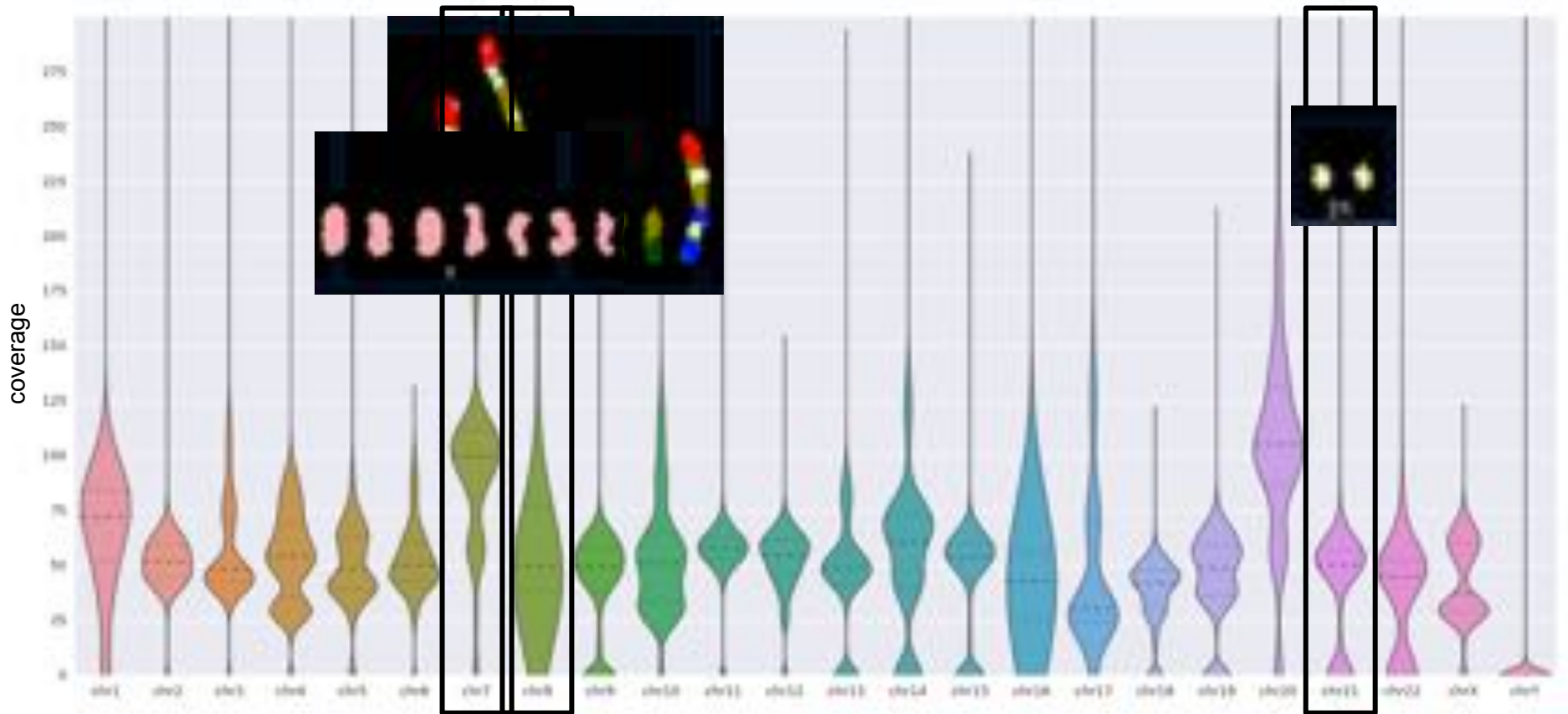


80 chromosomes instead of 46

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

(Davidson et al, 2000)

# PacBio long-read DNA sequencing

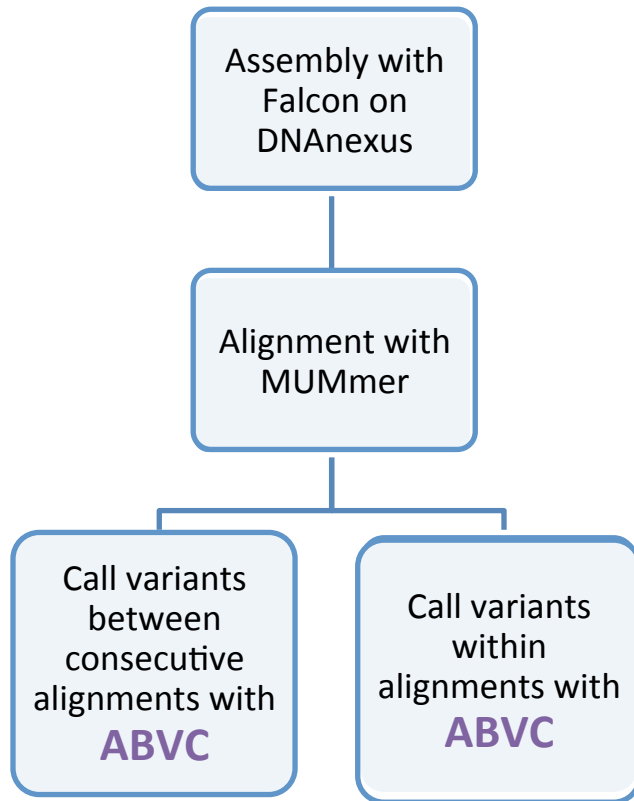mean read length: 9 kb
max read length: 71 kb

72X coverage



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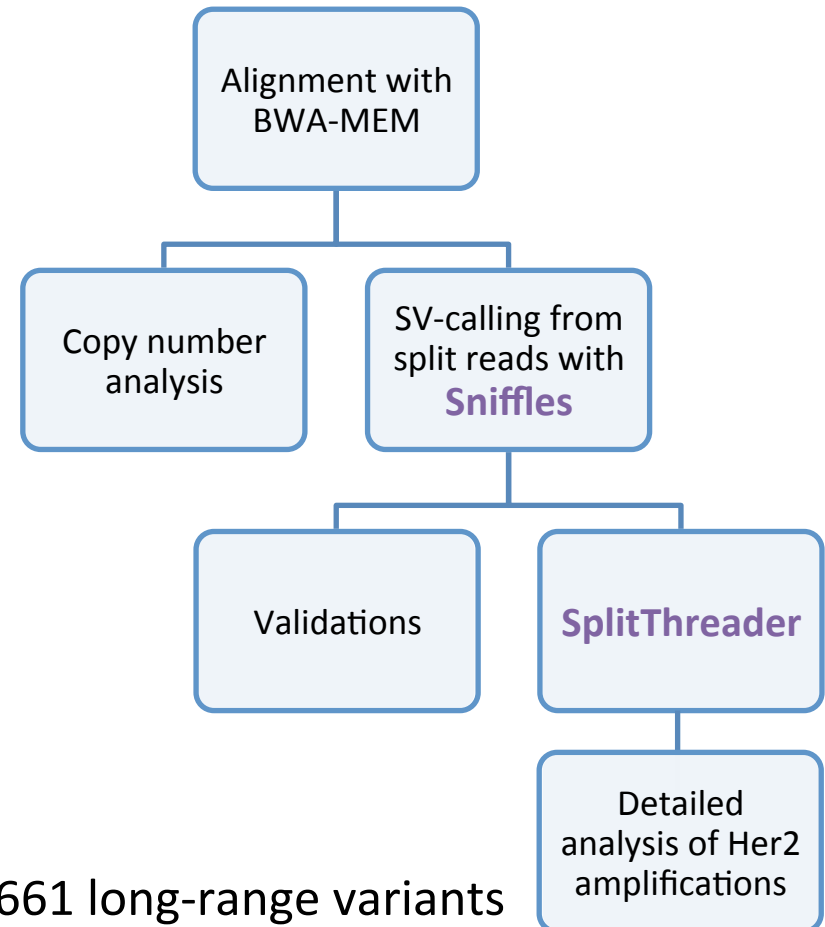Coverage per chromosome varies greatly as expected from previous karyotyping results

# Genome structural analysis

## Assembly-based

```
Assembly with
Falcon on
DNAnexus
```
│
```
Alignment with
MUMmer
```
│
```
Call variants          Call variants
between                within
consecutive            alignments with
alignments with        ABVC
ABVC
```

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

## Alignment-based

```
Alignment with
BWA-MEM
```
│
```
Copy number        SV-calling from
analysis           split reads with
                   Sniffles
```
│
```
Validations        SplitThreader
```
│
```
Detailed
analysis of Her2
amplifications
```

661 long-range variants
(>10kb distance)

# 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)

# 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
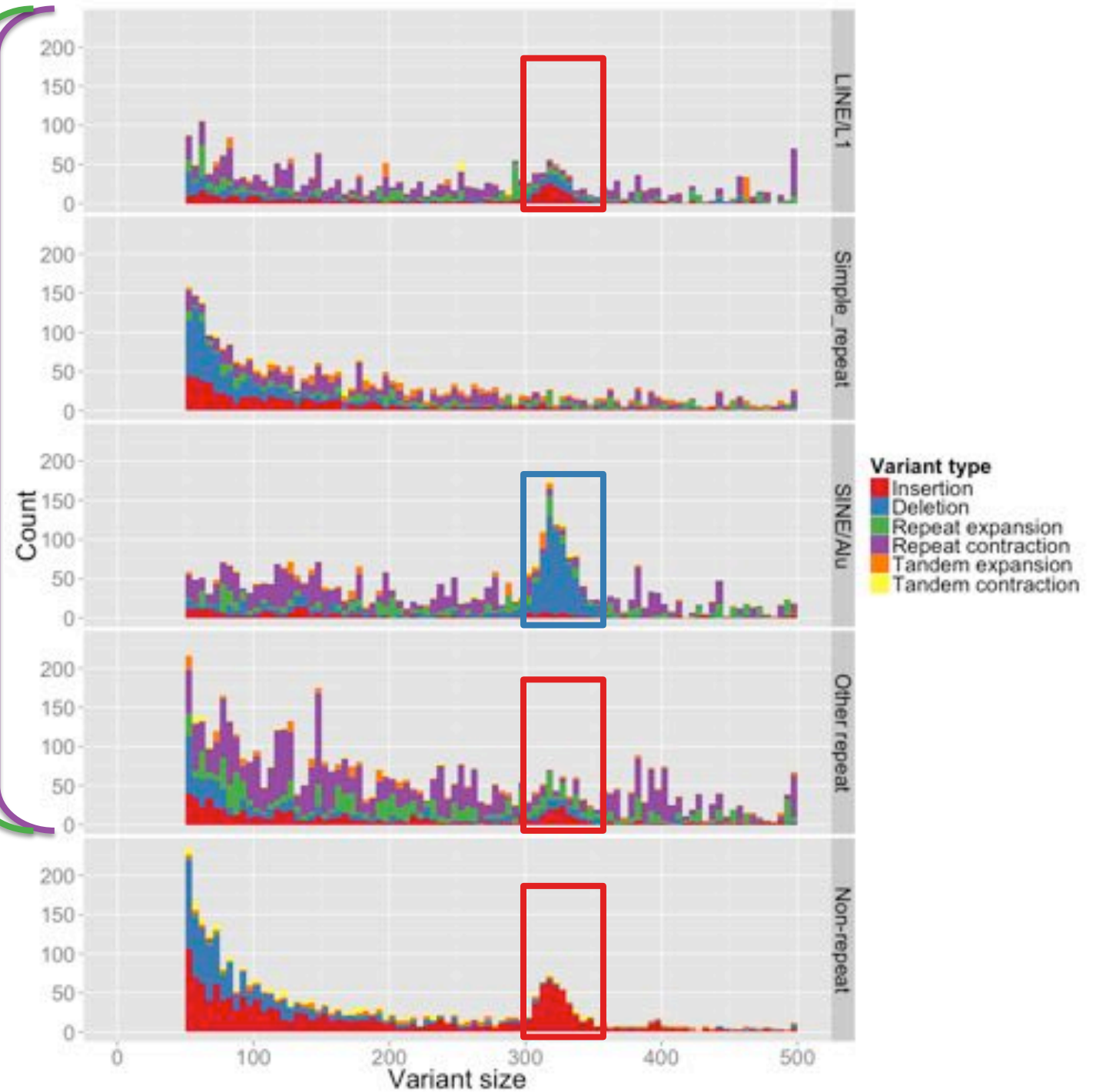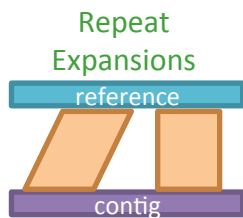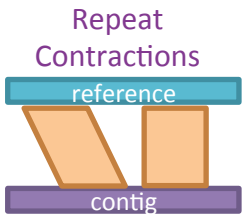Total sequence length: 2.07 Gb
Mean: 2.8 kb
Max: 61 kb
## N50: 3.3 kb

# ABVC: Assembly-Based Variant-Caller

Defined point

### Insertion

reference

contig

### Deletion

reference

contig

Overlapping
alignments suggest
tandem repeat

### Tandem Expansions

reference

contig

### Tandem Contractions

reference

contig

Gap where sequences
do not align uniquely
suggests a repeat

### Repeat Expansions

reference

contig

### Repeat Contractions

reference

contig

~ 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

# 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)

# Split-read variant calling with Sniffles to capture the long-range variants



See Fritz at Poster 183

# Long-range structural variants found by Sniffles

Her2 oncogene



661 long-range variants
(>10kb distance)
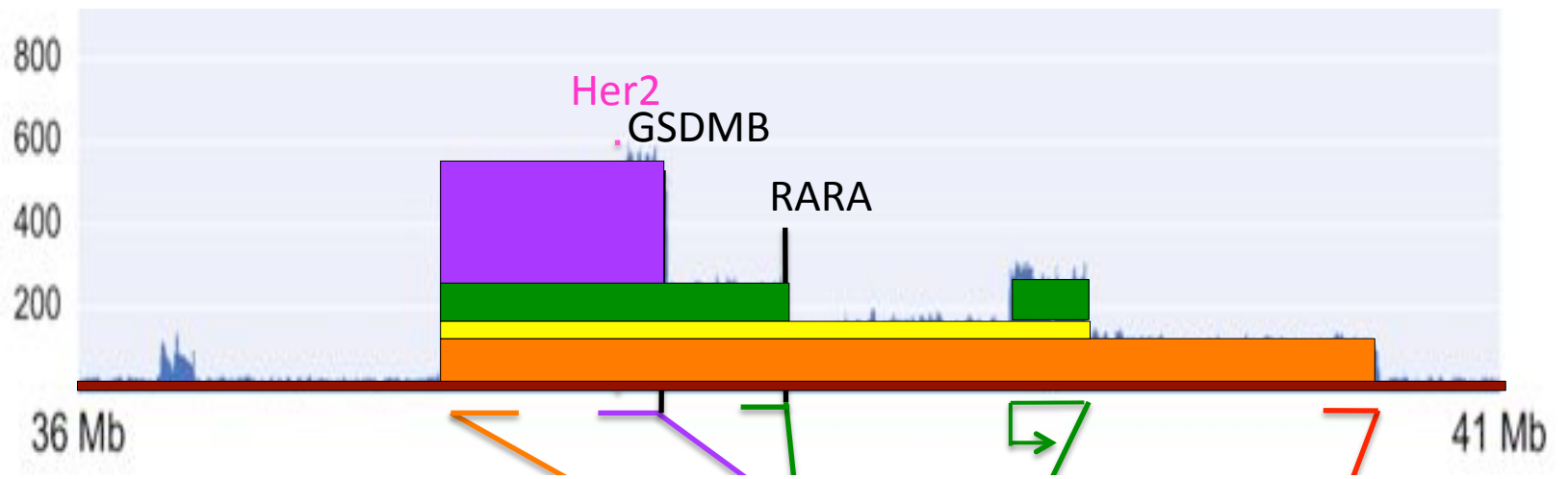
# SplitThreader:
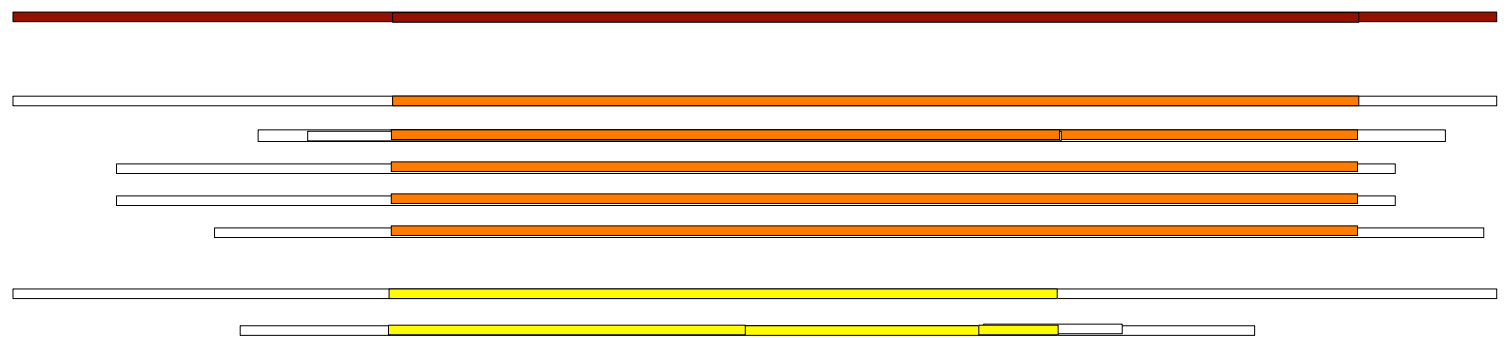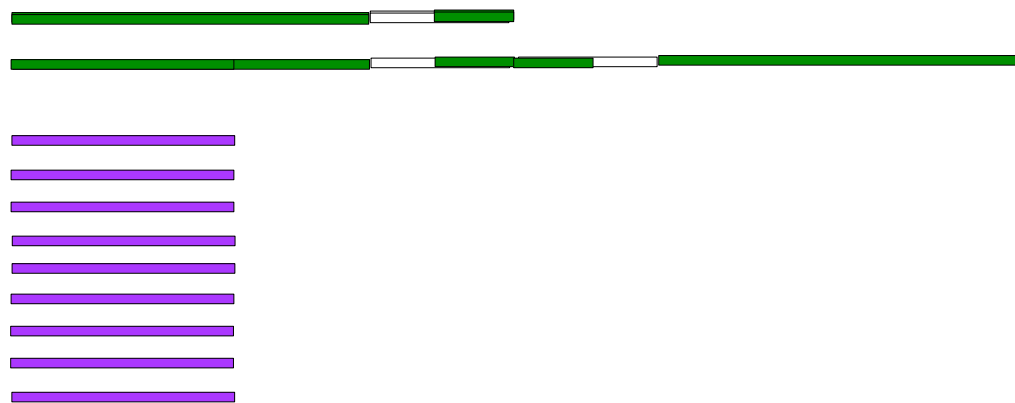# Graphical threading to retrace complex history of rearrangements in cancer genomes

Her2
GSDMB
RARA

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
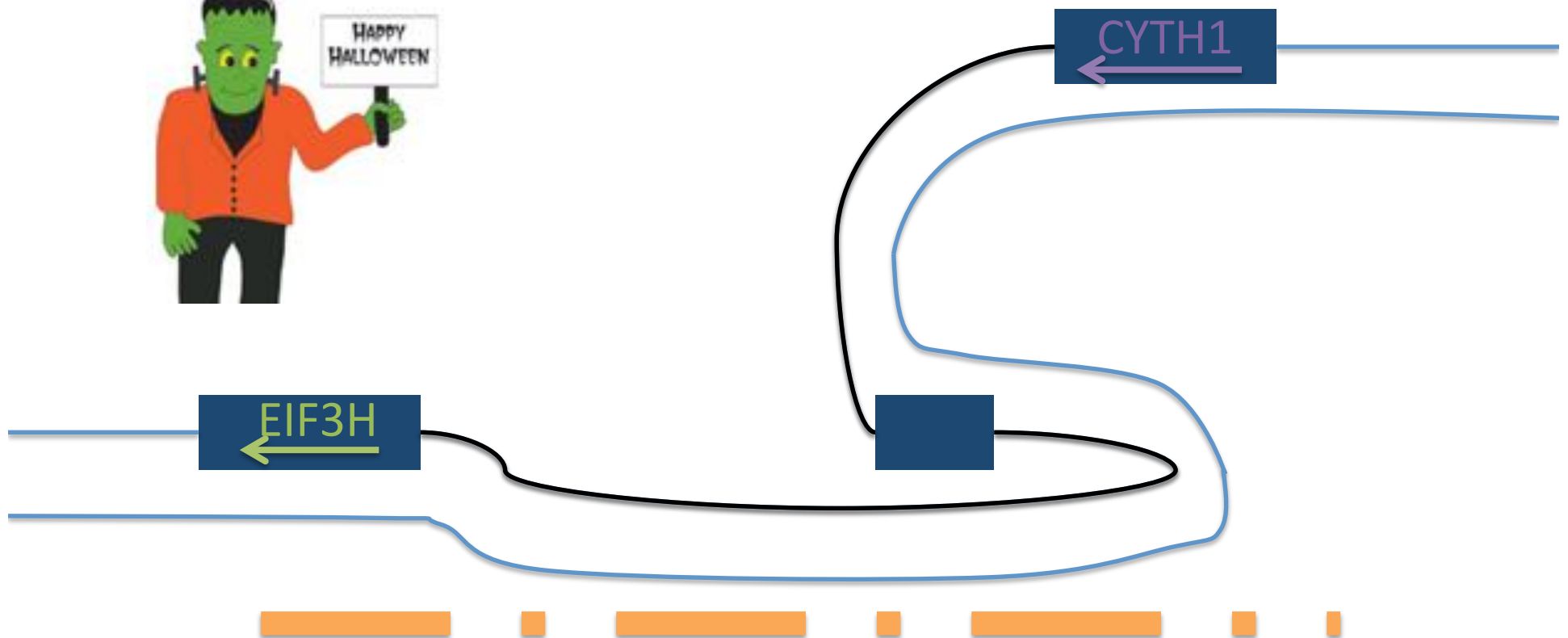5. Translocation within chromosome 8

# Transcriptome analysis with IsoSeq: 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
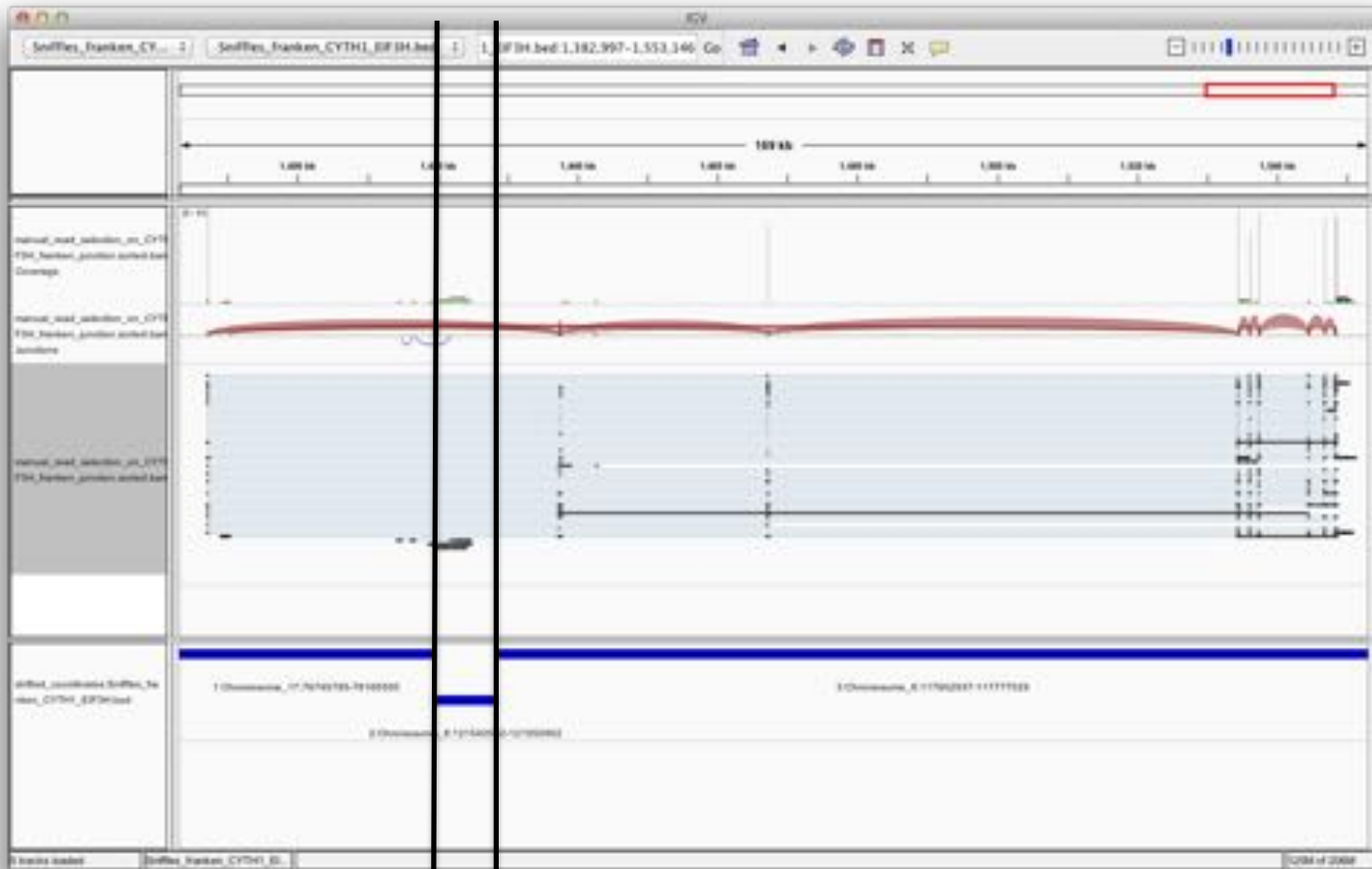  - Confirmed using SplitThreader

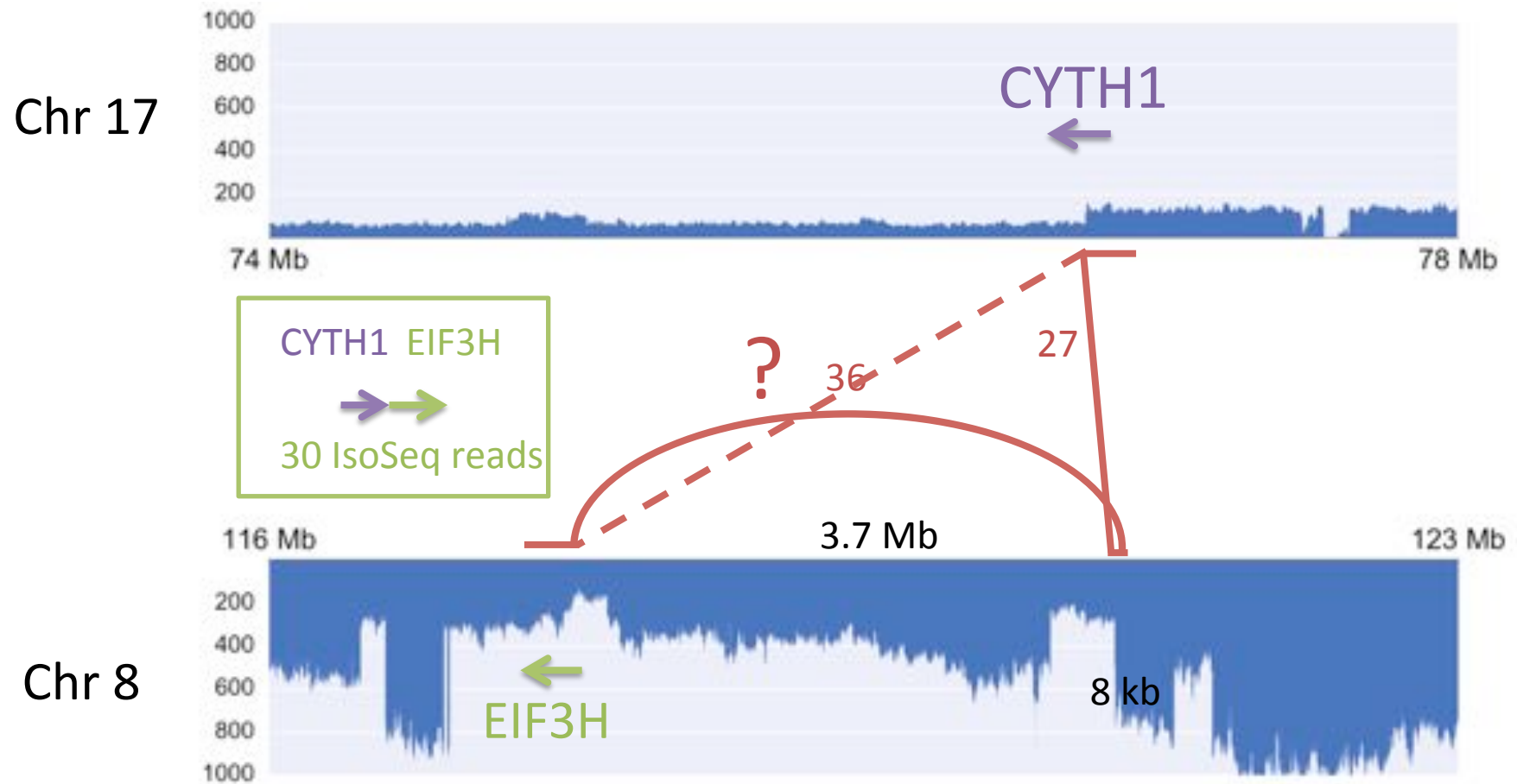# CYTH1-EIF3H gene fusion in the SplitThreader graph

# CYTH1-EIF3H gene fusion in the SplitThreader graph

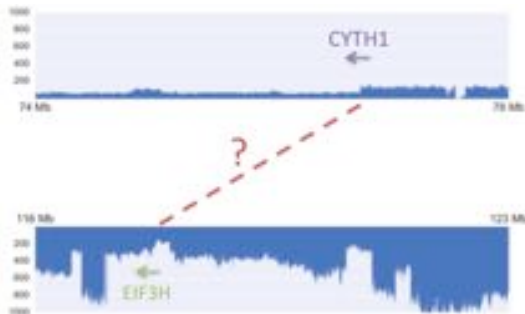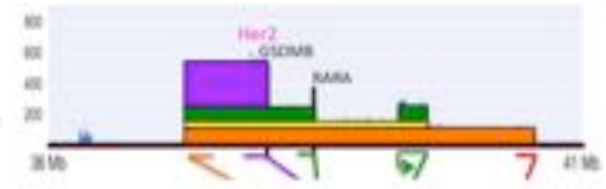# Frankensteining the CYTH1-EIF3H gene fusion

# CYTH1-EIF3H gene fusion
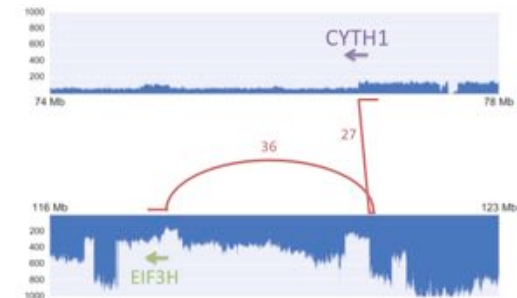
# The genome informs the transcriptome



Explain amplifications

Trace gene fusions

More genomes coming soon!

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

# Acknowledgments