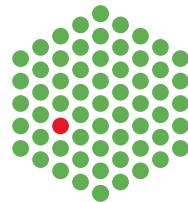
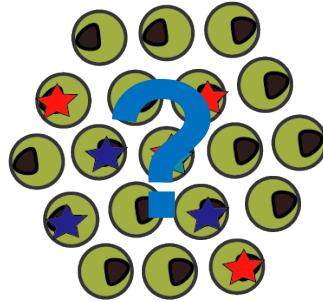


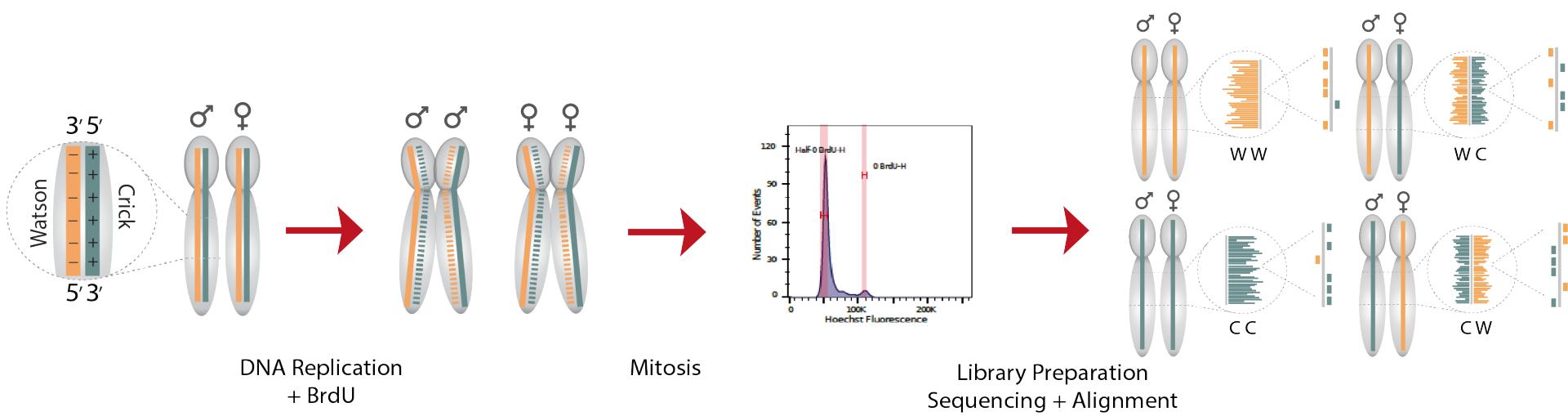
# Detection of structural variations at the single-cell level using Strand-seq



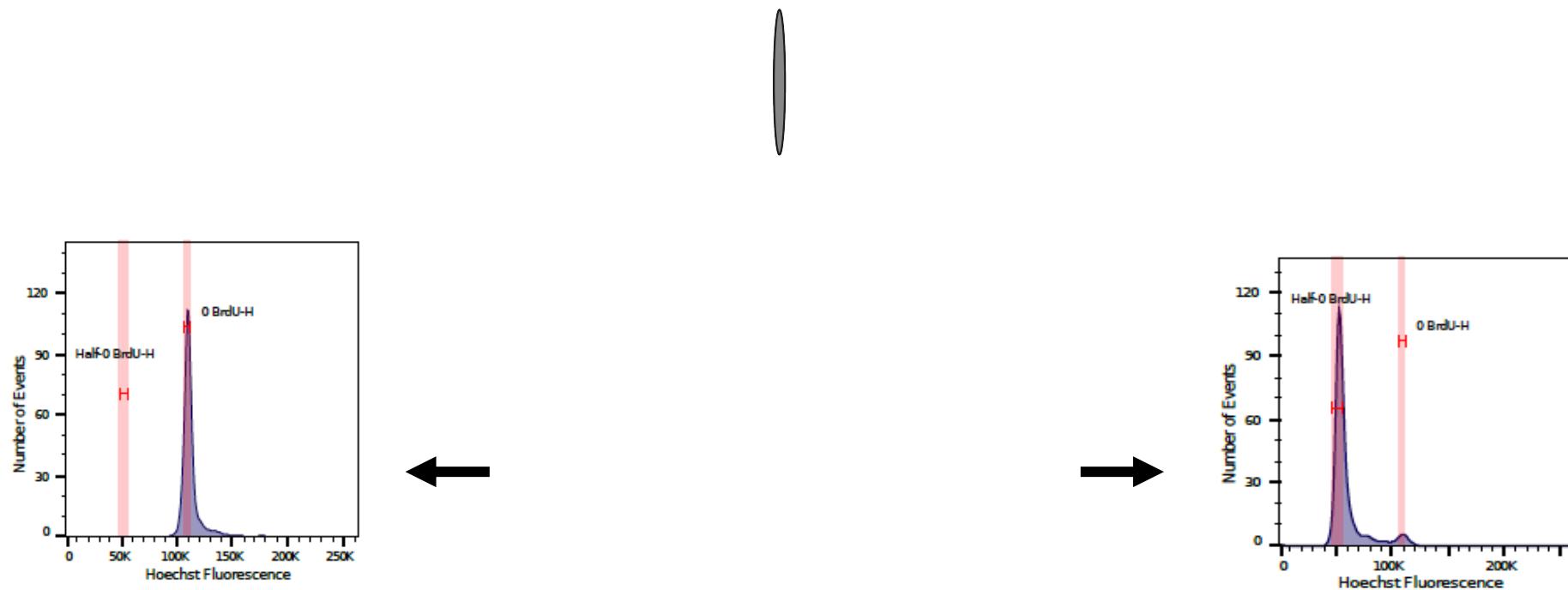


1. Structural variant detection
2. Single-cell resolution

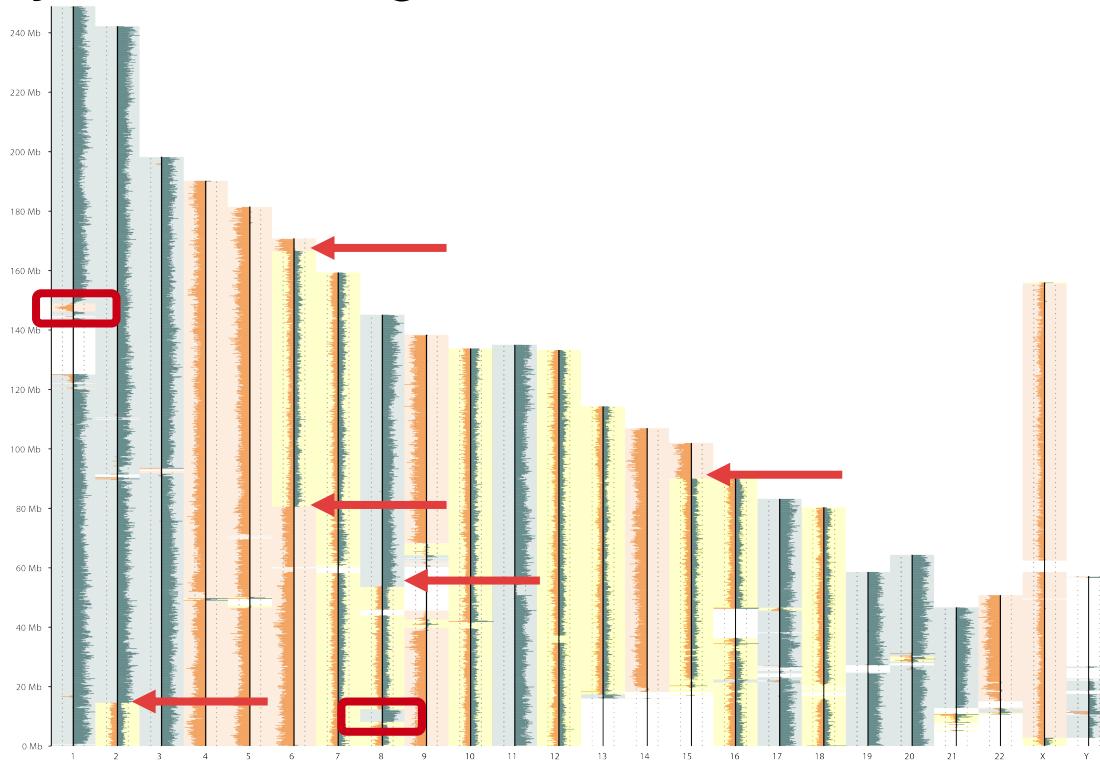
# Strand-seq selectively sequences DNA template from single cells to preserve DNA directionality



# Nuclei from cells hemi-substituted with BrdU can be identified based on the Hoechst profile

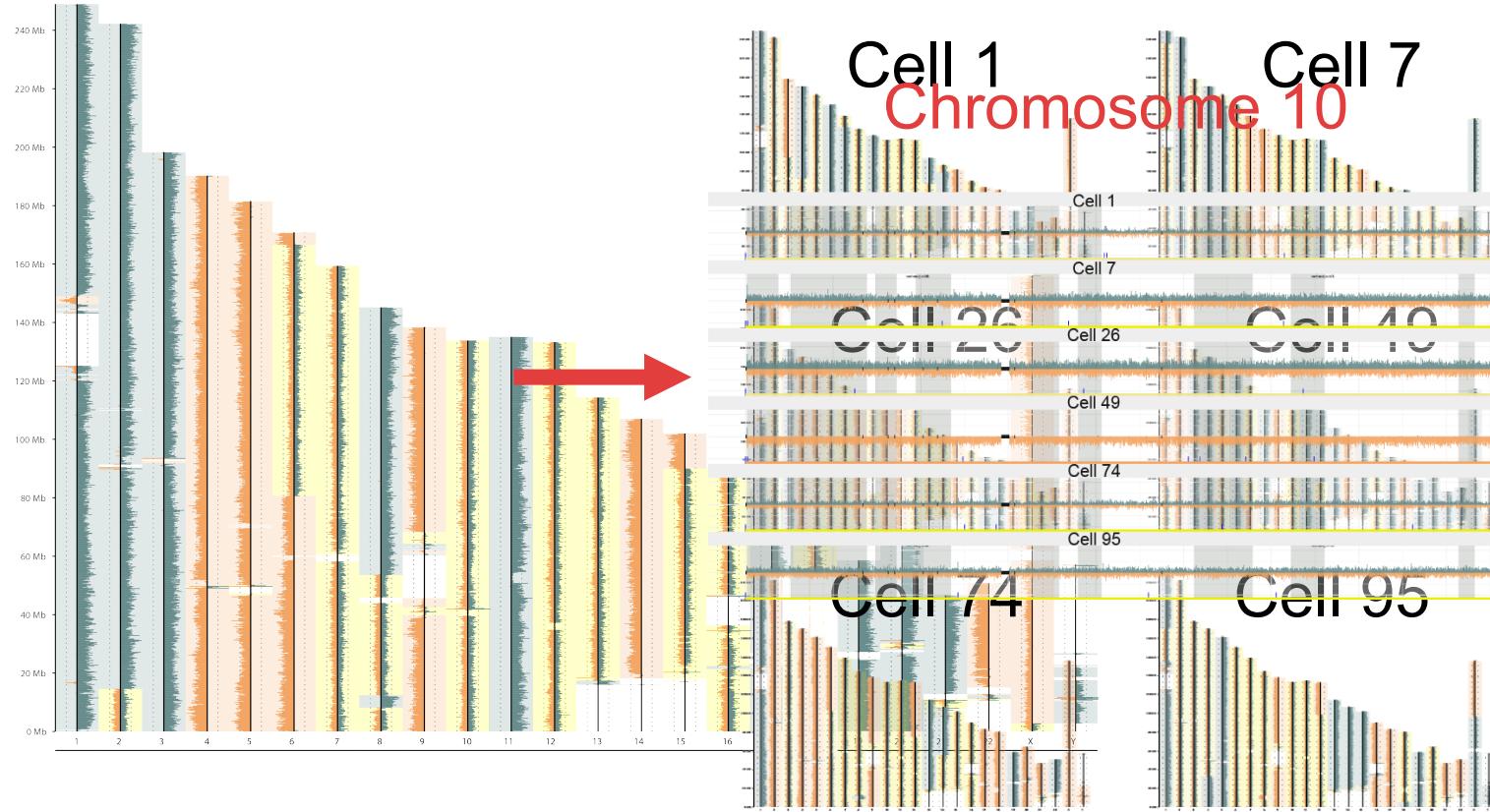


# Strand-seq allows for structural variant detection in single cells at very low coverage



Median uniquely mapped reads: ~400,000  
Median coverage: 0.03 X

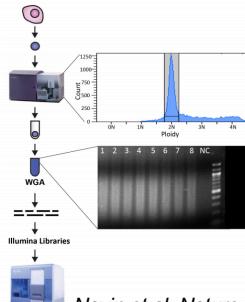
# Strand-seq allows for structural variant detection in single cells at very low coverage



# Overview of the single-cell genome data analysis – scWGS

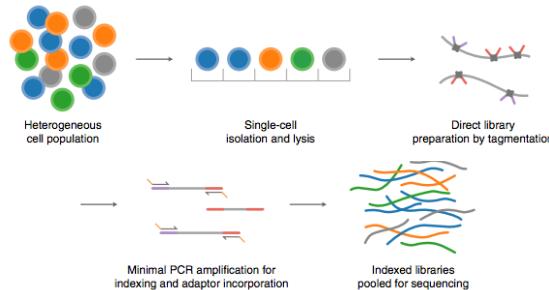
## Single nucleotide variation (SNVs) - scWGS

### Single-nucleus sequencing (SNS)



Navin et al. *Nature*, 2011

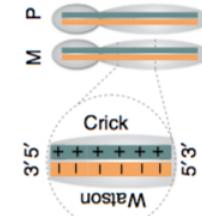
### Direct library preparation (DLP)



Zahn et al. *Nat Methods*, 2017

## Structural Variation

### Strand-seq



Sanders et al. *Nat protocol*, 2017

Step1. Alignment - Finding a correct position of reads: Bowtie2, BWA

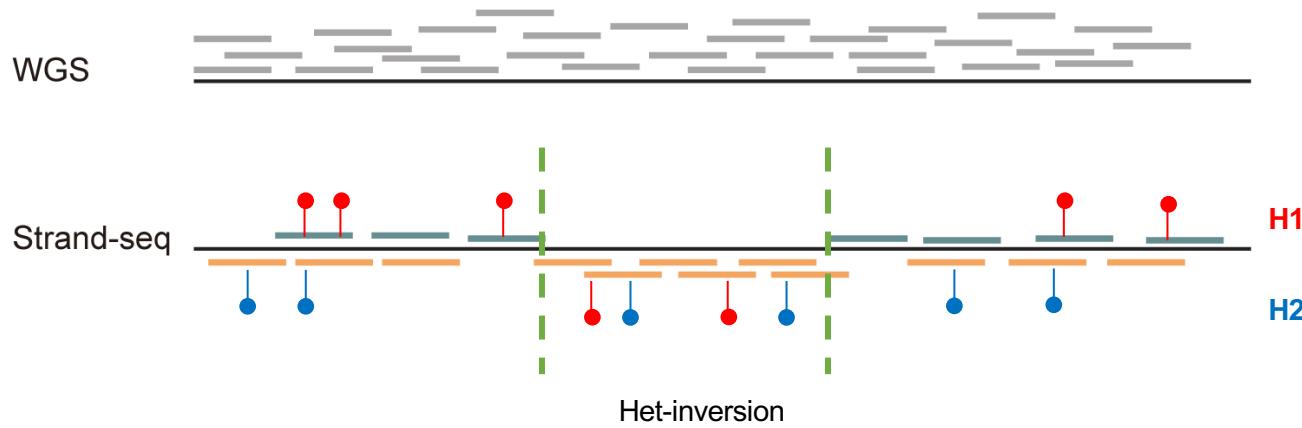
Step2. Remove PCR duplicate: Picard mark duplicate, Biobambam

Step3. Genotyping: Freebayes, GATK

Step4. Somatic mutation and CNV calling: SCcaller, Monovar, Aneufinder

Step5. Single-cell clustering and Phylogenetics: SCIPhi, TimeScape

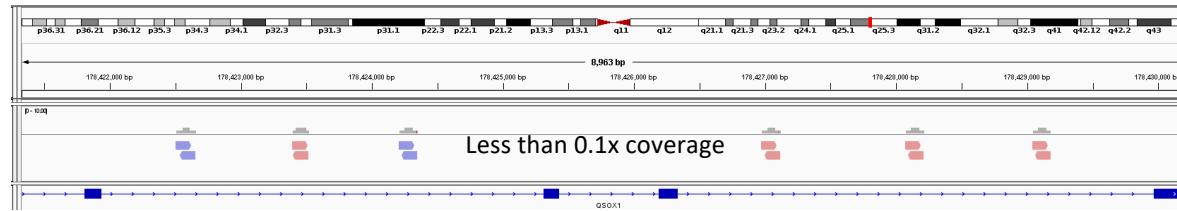
# Specialties of the Strand-seq data analysis



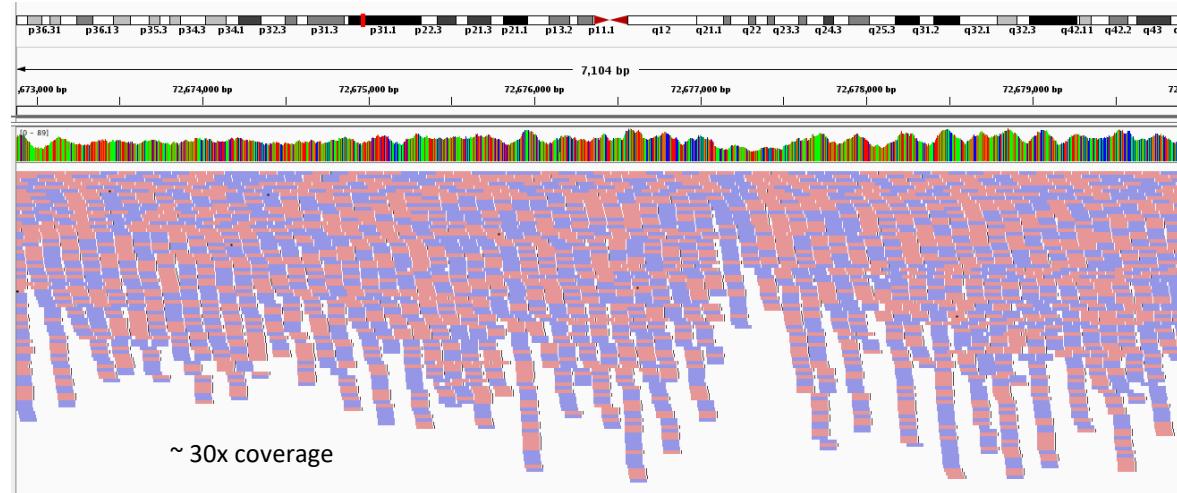
- Sequence orientation is important (Crick or Watson)
- Breakpoint needs to be detected
- Strand state and haplotypes can be assigned
- Multiple types of structural variations need to be classified

# Challenges of the Strand-seq data analysis

## Strand sequencing



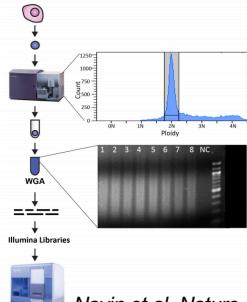
## Conventional sequencing (short-read)



# Overview of the single-cell genome data analysis – Strand-seq

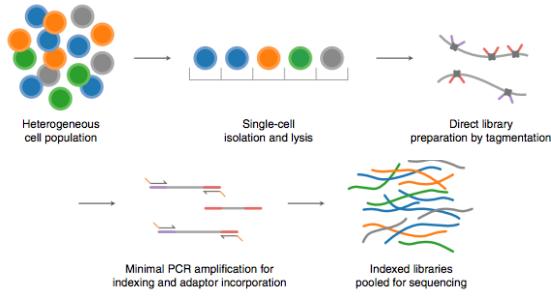
## Single nucleotide variation (SNVs) - scWGS

### Single-nucleus sequencing (SNS)



Navin et al. Nature, 2011

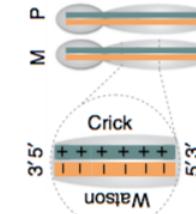
### Direct library preparation (DLP)



Zahn et al. Nat Methods, 2017

## Structural Variation

### Strand-seq



Sanders et al. Nat protocol, 2017

Step1. Alignment - Finding a correct position of reads: **BWA, sequence orientation**

Step2. Remove PCR duplicate: Biobambam

Step3. Genotyping, Haplotyping, Segmentation: **BreakpointR**

Step4. Structural variation calling: **MosaiCatcher**

Step5. Single-cell clustering and Phylogenetics

**Quality  
checking!**

# Fastq files are the starting point of Strand-seq data analysis



Paired-end

\* \_1\_sequence.txt.gz \* \_2\_sequence.txt.gz

ATAC TTT

AAAG TAT

CTGT AAA

TTT AGAG

Read 1



Read 2

# Why the orientation of the reads are important?



Porubsky et al. *Nat comm*, 2017

Homozygous Reference



Heterozygous Inversion



Homozygous Inversion

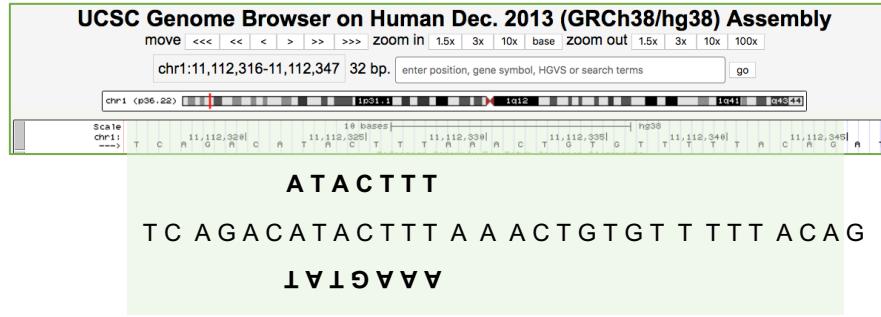


Sanders et al. *Genome Res*, 2016

# How can we assign sequencing reads into Crick and Watson?



- Crick (C) aligns to the plus (forward) strand of the reference assembly
- Watson (W) aligns to the minus (reverse) strand



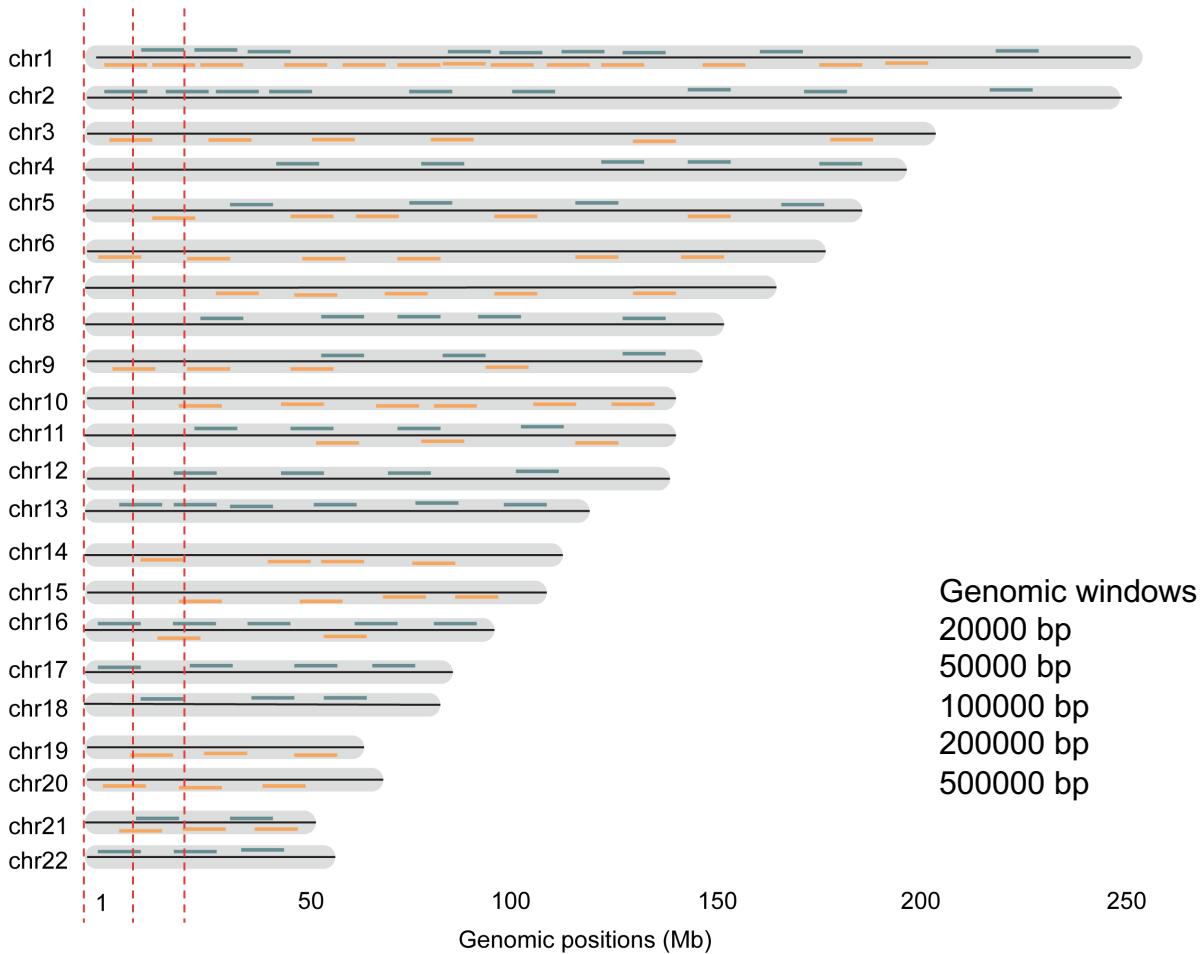
**A T A C T T T** Forward (+) → Crick (SAMFLAG 0)

**AAA G T A T** Reverse (-) → Watson (SAMFLAG 16)

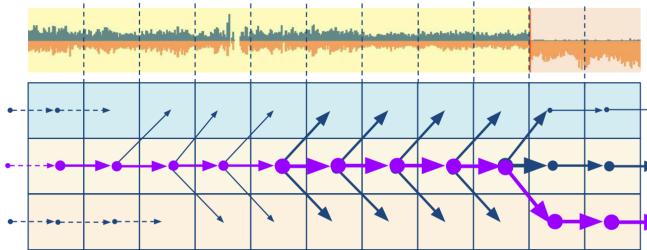
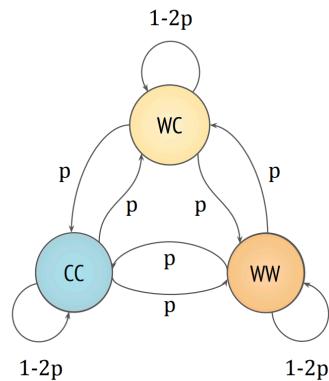
Crick = A1 read matches + or A2 read matches -

Watson = A1 read matches - or A2 read matches +

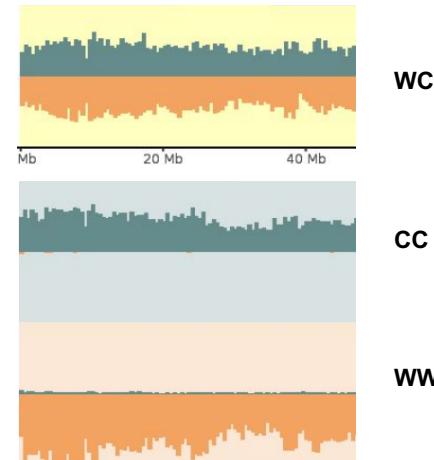
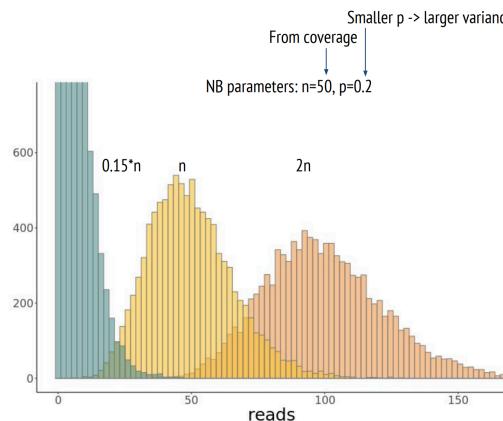
## Count the Watson and Crick reads using genomic windows

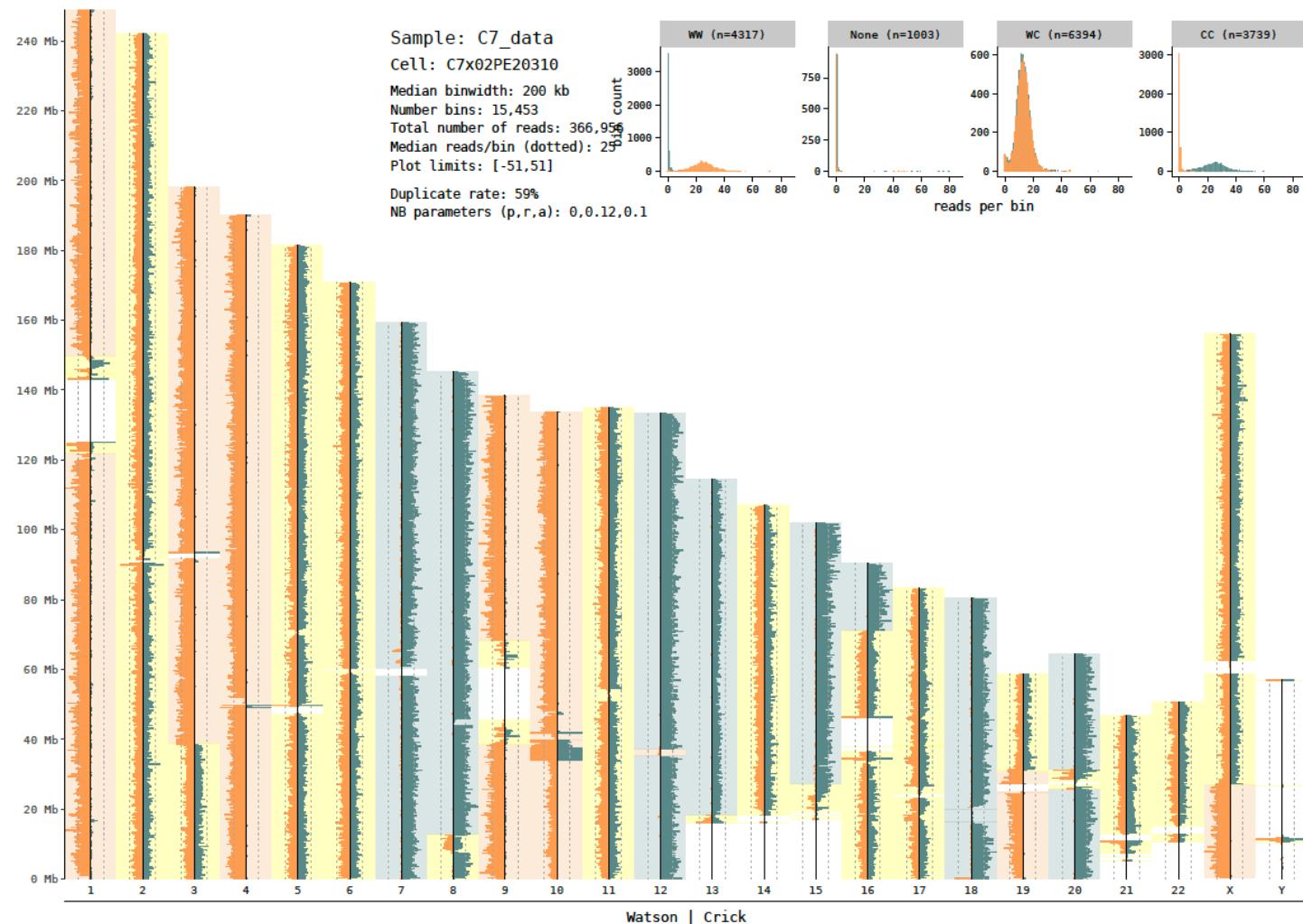


# Strand states are called using a Hidden Markov Model



Arrows show the most probable sequence of state transitions  
Thickness of line = probability of the path from start  
Purple path is the most probable path in the end





# Processing the Strand-seq library after sequencing

Input: Paired-end fastq files : \*\_1\_sequence.txt.gz \*\_2\_sequence.txt.gz



Step1. Align reads, mark duplicates, add read groups

```
bwa mem -t {threads} \ -R '@RG\tID:{params.name}\tSM:{SAMPLE_NAME}\tLB:{wildcards.cell}' \
{params.genome} \ {input} \ | samtools view -bT {params.genome} - \ > {output}
bammarkduplicates markthreads={threads} I={input} O={output} index=1 rmdup=0
```



Step2. Filter & count bins in fixed windows

```
Rscript {params.qc_plot} {input.counts} {input.info} {output}
```



Step3. Plots and statistics

```
Rscript {params.sv_plot} {input} plot_chromosomes/{SAMPLE_NAME}/window_{wildcards.window}
```



***Don't worry, we can do it with one touch  
using Snakemake!***

Sascha Meiers



# Workflow management with Snakemake

- Snakemake is a Python-based workflow management tool
- It can run workflows consisting of multiple inter-dependent tasks. These can be computed either locally or on the cluster.
- Intermediate results are stored, so the pipeline only executes what is left to do.
- See <http://snakemake.readthedocs.io> or <https://bitbucket.org/snakeyed/snakeyed/overview>



```
rule bwa_map:
    input:
        "data/genome.fa",
        "data/samples/{sample}.fastq"
    output:
        "mapped_reads/{sample}.bam"
    shell:
        "bwa mem {input} | samtools view -Sb - > {output}"
```

```
rule samtools_sort:
    input:
        "mapped_reads/{sample}.bam"
    output:
        "sorted_reads/{sample}.bam"
    shell:
        "samtools sort -T sorted_reads/{wildcards.sample} "
        "-O bam {input} > {output}"
```

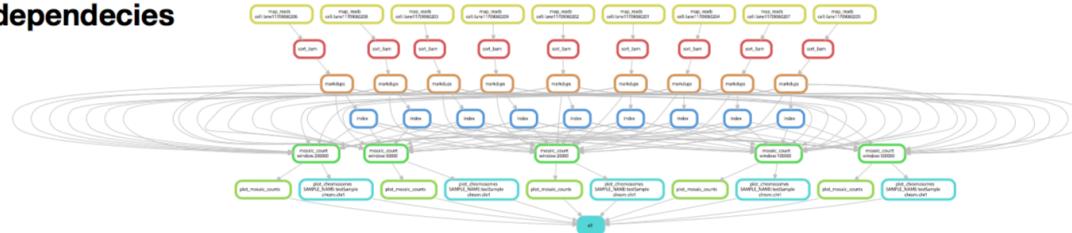
# What will be happened in the pipeline?

- Align to GRCh38 (recommended version GCA\_000001405.15\_GRCh38\_no\_alt\_analysis\_set.fna)
- During that step, add read groups (1 SM tag per sample, but different (possibly multiple) IDs per cell)
- Mark duplicates
- Count reads of all cells in bins (20,50,100, and 200kb)
- Plot overviews of all cells
- Plot chromosome by chromosome

## Job overview

count	1	all
9	index	index
9	map_reads	map_reads cell_line112906206
9	markdups	markdups
5	mosaic_count	mosaic_count mosaic 2000
5	plot_chromosomes	plot_chromosomes 20 chromosomes
5	plot_mosaic_counts	plot_mosaic_counts sort_mosaic
9	sort_bam	sort_bam
52		

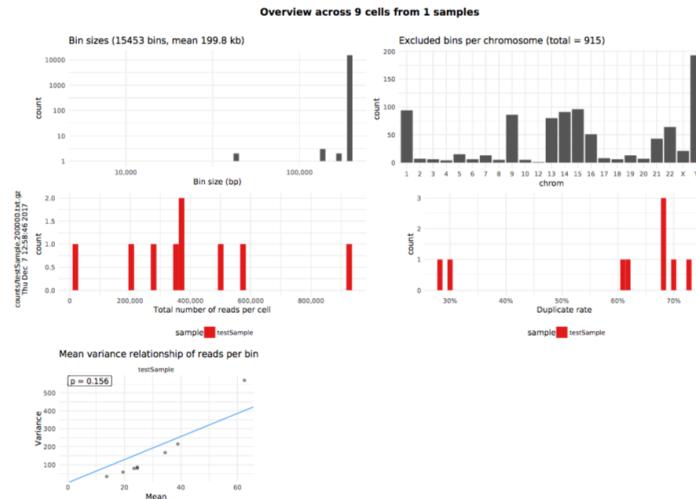
## Job dependencies



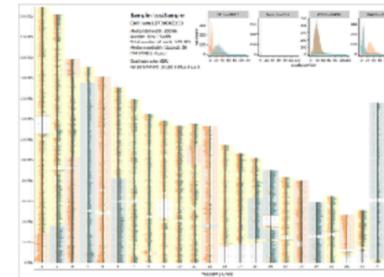
# Output of the pipeline

- Three different plots for quality checking and first look for the SVs
  - Count matrix for read depth and ploidy diagnosis
  - Bam files from this pipeline will be used as an input of downstream advanced analysis

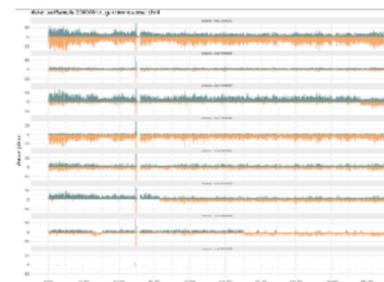
## Overview plot



## Each single cell

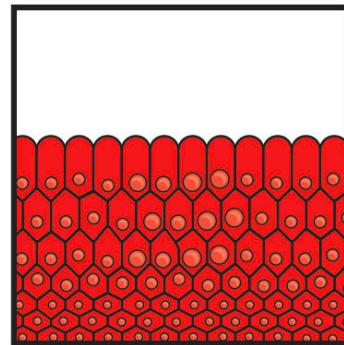


## Single chromosomes

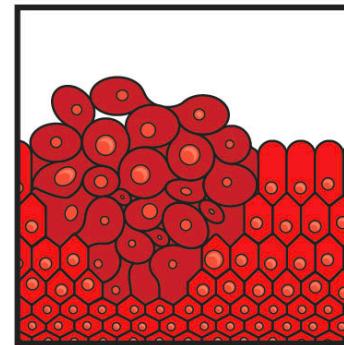


# Which data do we analyze today?

Real  
world

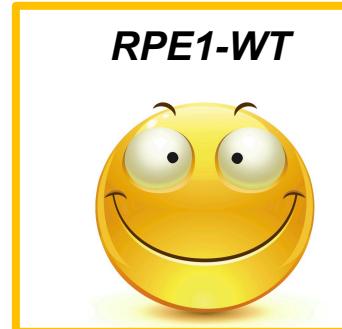


Normal cells



Cells forming a tumour

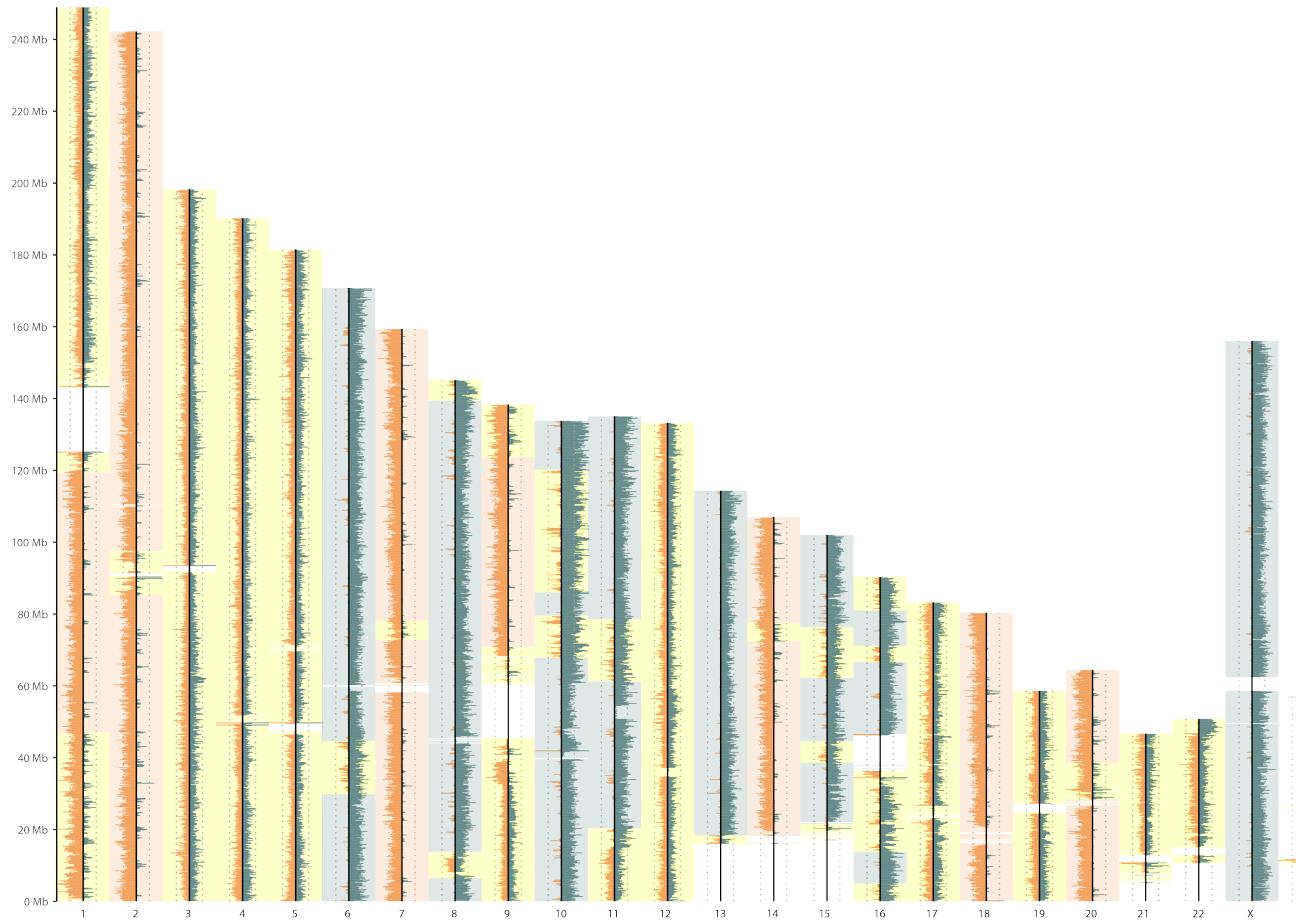
Model  
system



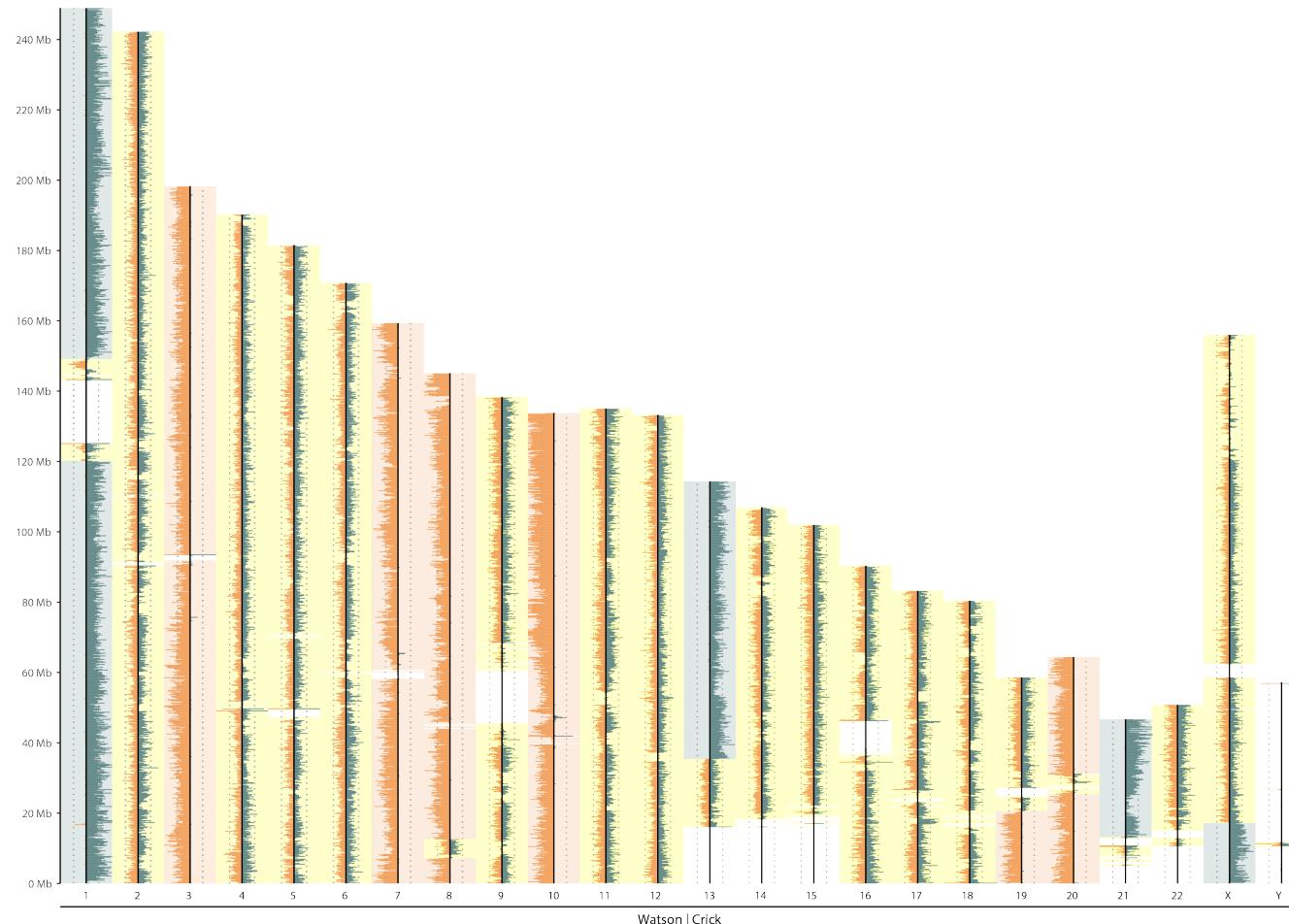
vs



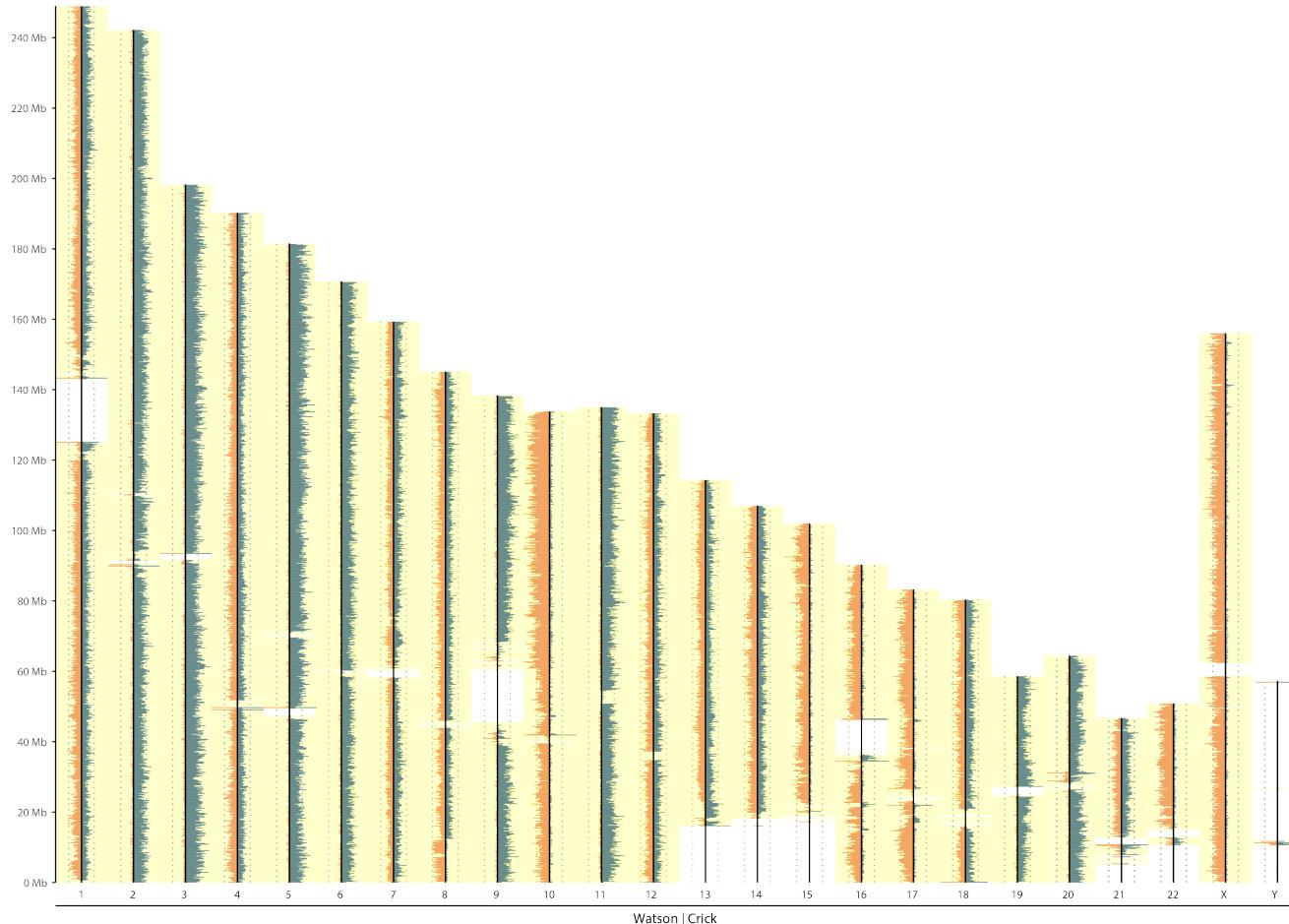
# Filtering Strand-seq Libraries: Under-BrdU incorporation



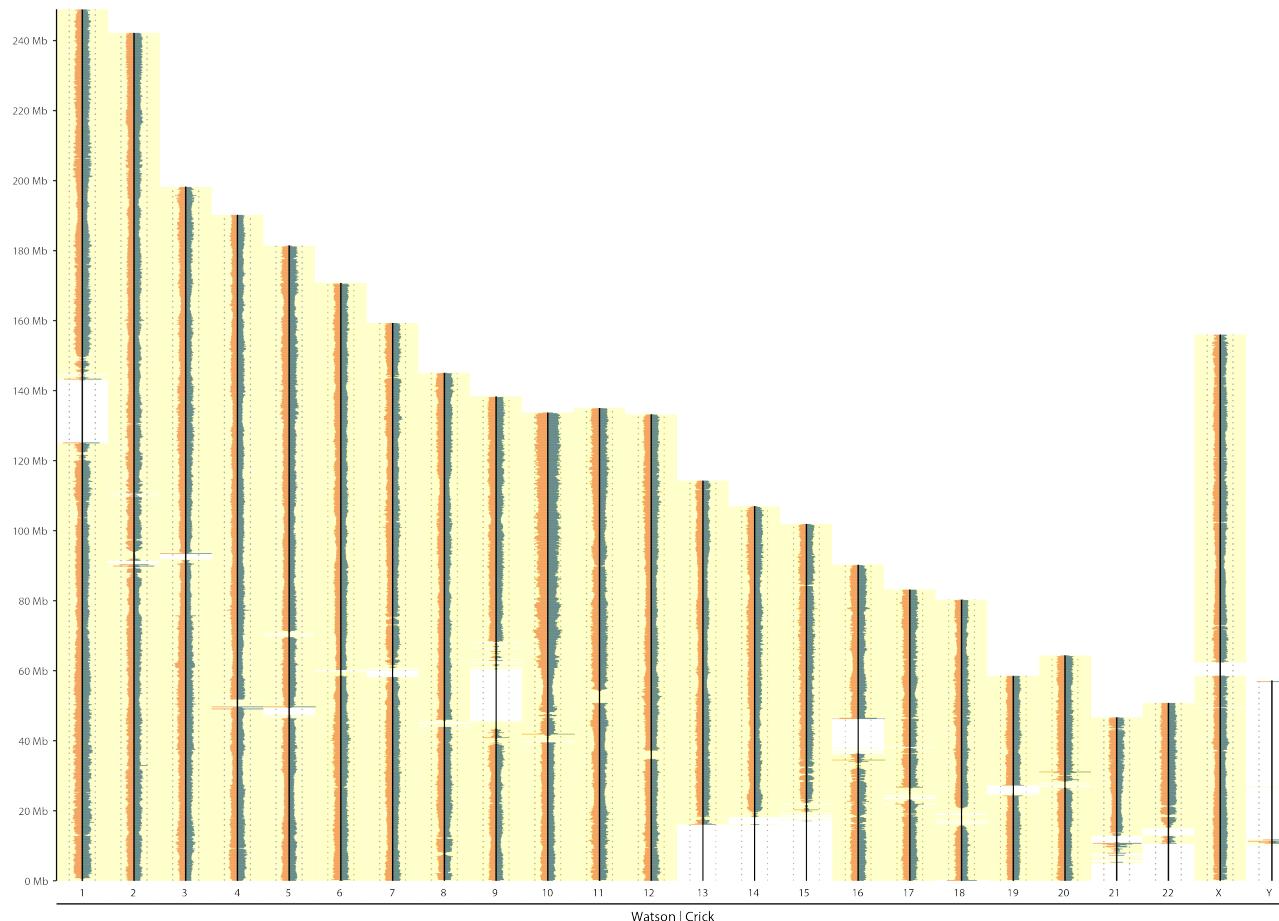
# Filtering Strand-seq Libraries: Over-BrdU incorporation



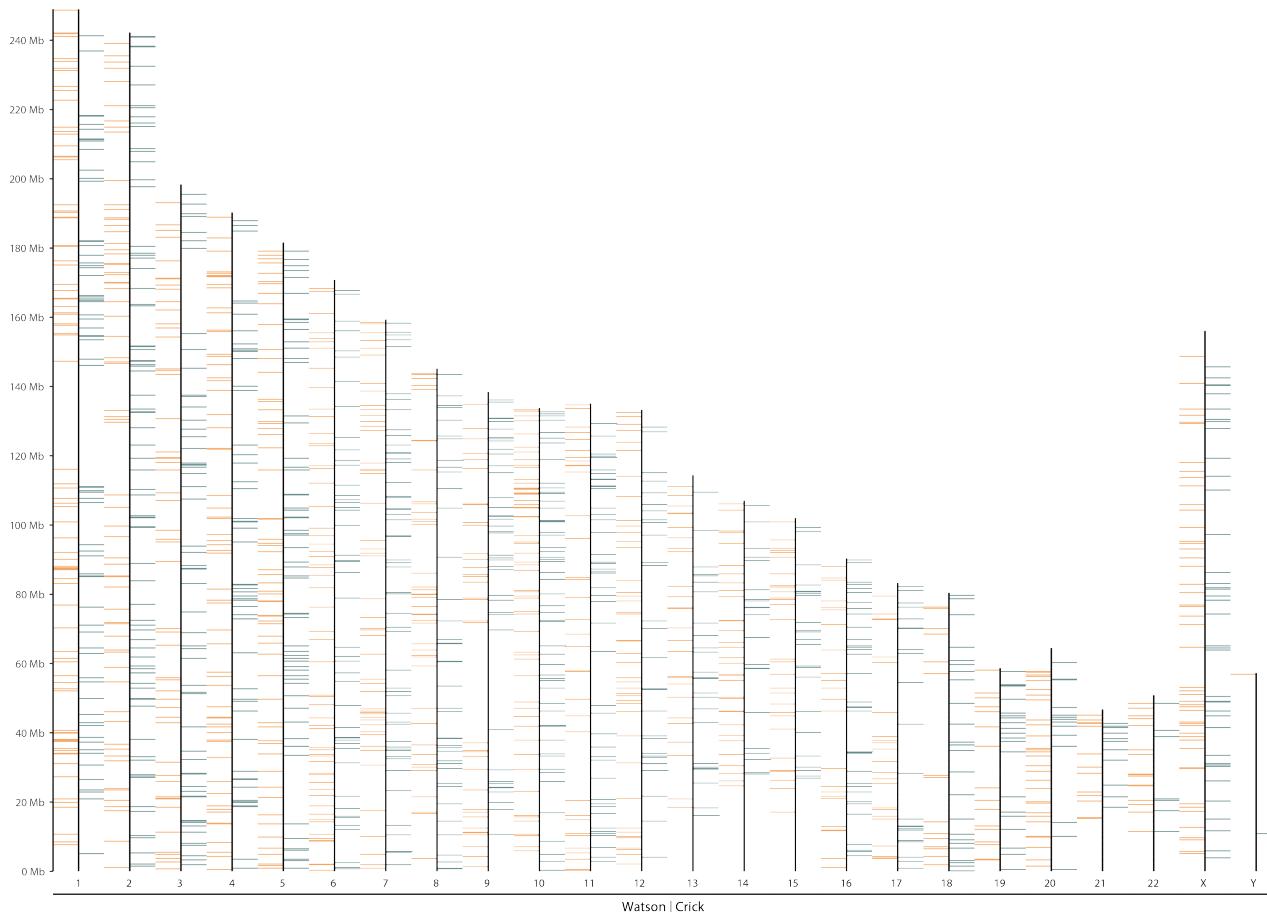
# Filtering Strand-seq Libraries: Background Contamination



# Filtering Strand-seq Libraries: Multi-cell (Positive Control)



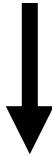
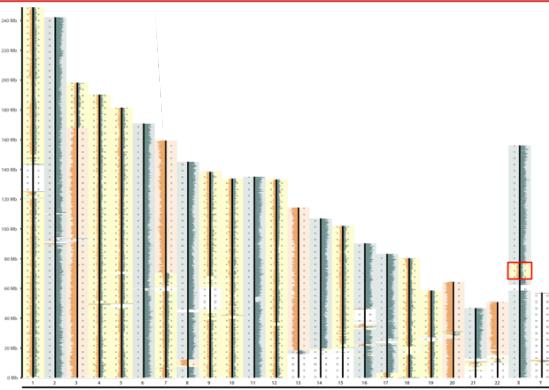
# Filtering Strand-seq Libraries: 0-cell (Negative Control)



# BreakPointR analysis allows refinement of breakpoints and alignment of Strand-seq data to public datasets on UCSC



27



chrX (q11.2-q21.1) p22.2 21.1 112 421.1 Xq23 q24 Xq25 Xq28

BreakpointSummary



**Input:**

- Curated Strand-seq bam files

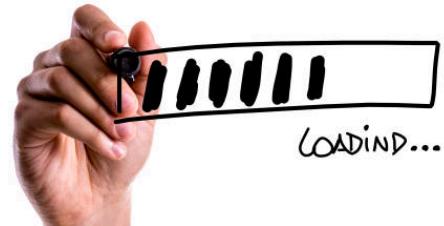
**Output:**

- Bed files
- Breakpoint summaries

*(genome.ucsc.edu)*

*Thank you for  
listening*

**COMING SOON**



**Day3**

