# KSBi-BIML 2026

**Bioinformatics & Machine Learning(BIML) Workshop for Life Scientists**

생명정보학 & 머신러닝 워크샵 (온라인) ▶

# Single-cell multi-omics analysis to study tumor subclones

정효빈 _ 연세대학교

**KSBI** KOREAN SOCIETY FOR BIOINFORMATICS | 한국생명정보학회

# KSBi-BIML 2026

## Bioinformatics & Machine Learning (BIML) Workshop for Life Scientists

한국생명정보학회가 주최하는 BIML-2026 동계 Bioinformatics &Machine Learning 교육 워크숍에 여러분을 초대합니다.

BIML 워크숍은 생명정보학 연구자들이 최신 AI바이오 분야의 인공지능 기반 분석 기술과 바이오 데이터 분석 기법을 이론과 실습을 통해 체계적으로 배울 수 있는 전문 교육 프로그램입니다. 2015년에 시작된 BIML 워크숍은 올해로 12년 차를 맞이하며, 국내 생명정보학 분야의 최초이자 최고 수준의 교육 프로그램으로 자리 잡았습니다. 이번 워크숍은 크게 인공지능바이오(AI바이오) 분야와 디지털바이오 분야, 두 분야로 구성됩니다.

AI바이오 분야에서는 생명정보 분석에 폭넓게 응용되고 있는 다양한 인공지능 기반 자료 모델링 기법을 다룰 예정입니다. 특히, 인공지능 심층학습을 활용한 단백질 구조 예측, 유전체 분석, 신약 개발에 대한 이론 및 실습 강의가 진행됩니다.

또한 디지털바이오 분야에서는 단일세포오믹스, 공간오믹스, 멀티오믹스, 메타오믹스에 대한 강의도 마련되어 있어, 연구자들의 분석 역량 강화에 실질적인 도움을 줄 것으로 기대됩니다.

또한 2024년부터 추가된 의료정보 자료 분석을 다루는 강의를 올해도 지속해서 운영하고자 합니다. 이는 최근 의료정보 자료 분석에 관한 연구 수요 증가를 반영한 것으로, 관련 연구를 수행하는 의과학자 및 의료정보 연구자들에게 유용한 지침을 제공할 것입니다.

또한, 올해도 생명정보학 기술의 다양화에 발맞춰 온라인 강좌를 대폭 확대했습니다. 올해는 무료 강좌 10개를 포함한 총 40개 이상의 강좌가 개설되며, 연구 주제에 맞는 강좌 추천과 강연료 할인 혜택도 제공합니다.

BIML-2026는 국내 주요 연구 중심 대학의 전임 교수 및 각 분야 최고 전문가들의 강의로 구성되어 있으며, 기초 이론부터 최신 연구 동향까지 아우르는 심도 있는 교육의 장이 될 것으로 확신합니다.

여러분의 많은 관심과 참여를 기대합니다!

2026년 2월

한국생명정보학회장 류 성 호

# Single-cell multi-omics analysis
# to study tumor subclones

암의 종양 내 이질성 (intra-tumor heterogeneity)는 암 조직 내에 다양한 유전체적, 또는 후성 유전체적 특성을 가지는 세포들이 존재하면서 암의 진행을 가속화하고 항암제 내성을 심화시키는 현상을 의미한다. 특히 암의 진화 과정에서 축적되는 유전체 돌연변이와 구조변이들은 새로운 서브클론을 발생시키고, 이러한 서브클론들 각각의 특성을 파악하는 것이 암을 이해하고 치료 전략을 제시하는 데 필요하다. 그렇다면 암에서 이러한 서브클론들을 동정하기 위해 어떤 싱글셀 오믹스 기법들이 개발되어 있을까? 이러한 싱글셀 오믹스 데이터를 분석하기 위해 어떤 생명 정보학적인 도구들을 사용할 수 있을까? 서브클론의 동정 뿐 아니라 그 기능적 특성을 파악하기 위해서는 유전체와 전사체 또는 후성유전체 데이터를 함께 분석하는 싱글셀 멀티 오믹스 분석이 필요하다. 이를 구현하기 위한 생명 정보학적인 방법에는 어떤 것들이 있을까?

본 강의에서는 암에서 서브클론을 동정하기 위해 최근까지 개발되어 있는 다양한 싱글셀 오믹스 기법들에 대해 소개하고, 이들 중 scDNA-seq (Strand-seq, SDR-seq 등)을 이용하는 경우와, scRNA-seq을 이용하는 경우의 데이터 분석을 소개한다. 또한, 서브클론을 동정한 이후에 각각의 기능적인 특성들을 파악할 수 있는 싱글셀 멀티 오믹스를 위해 개발되어 있는 생명정보학 도구들을 소개한다. 이로써, 암의 종양 내 이질성을 심도적으로 탐구하고 의학적 연구에 응용할 수 있는 싱글셀 바이오 데이터 분석 역량을 갖출 수 있도록 하는 것이 최종 목표이다.

강의는 다음의 내용을 포함한다:

- 암에서 서브클론을 동정하기 위한 싱글셀 오믹스 기법들에 대한 소개
- scDNA-seq 기법 중 Strand-seq 데이터에서 서브클론을 동정하는 방법 소개
- Targeted scDNA-seq 기법 중 SDR-seq 소개 및 서브클론을 동정하는 방법 소개
- scRNA-seq 으로부터 서브클론을 유추하기 위한 데이터 분석 방법 소개
- 서브클론을 동정한 후 functional analysis를 위한 싱글셀 멀티오믹스 접근법과 최신 생명정보학 도구 소개

\* 교육생준비물: 노트북, R (또는 R studio)

\* 강의 난이도: 초급

\* 강의: 정효빈 교수 (연세대학교 시스템생물학과)

# Curriculum Vitae

## Speaker Name: Hyobin Jeong, Ph.D.

▶ **Personal Info**

| | |
|---|---|
| Name | Hyobin Jeong |
| Title | Assistant Professor |
| Affiliation | Yonsei University, Department of Systems Biology |

▶ **Contact Information**

| | |
|---|---|
| Address | 50 Yonsei-ro, Seodaemun-gu, Seoul, South Korea |
| Email | hyobinjeong@yonsei.ac.kr |

---

**Research Interest**

Computational biology, Single-cell Multi-omics, Intra-tumor heterogeneity

**Educational Experience**

2016.01-2017.11   Postdoc Fellow, Institute of Molecular Biology (IMB), Germany
2015.02-2015.12   Postdoc Fellow, Institute of Basic Science (IBS), Korea
2011.03-2015.02   POSTECH (Ph.D, Interdisciplinary Bioscience and Biotechnology, I-Bio)
2007.03-2011.02   POSTECH (Undergraduate study, Chemical Engineering)

**Professional Experience**

2024.03-Present   Assistant Professor, Department of Systems Biology, Yonsei University
2022.09-2024.02   Research Professor, Hanyang Institute of Bioscience and Biotechnology (HY-IBB)
2017.12-2022.08   Postdoc Fellow, EMBL, Germany

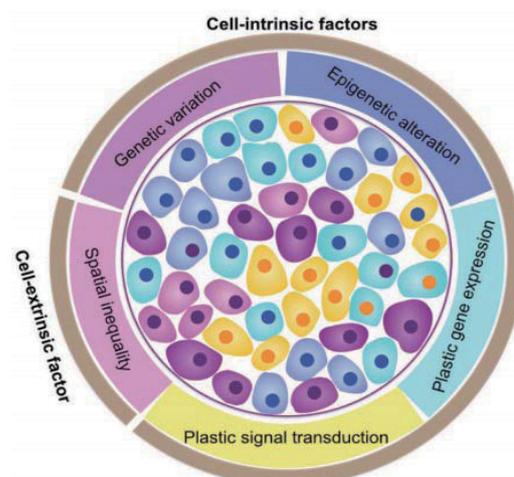**Selected Publications (3 maximum)**

1. 4. Leppä A-M*, Grimes K*, **Jeong H**\*, Huang F, Andreades A, Waclawiczek A, Boch T, Jauch A, Renders S, Stelmach P, Müller-Tidow C, Karpova D, Sohn M, Grünschläger F, Hasenfeld P, Garagorri E, Thiel V, Dolnik A, Rodriguez-Martin B, Bullinger L, Mrózek K, Eisfeld A-K, Krämer A, Sanders AD, Korbel JO#, Trumpp A#, (2024.12) "Single-cell multiomics analysis reveals dynamic clonal evolution and targetable phenotypes in acute myeloid leukemia with complex karyotype", **Nature Genetics [\*co-first]**

2. Grimes K*, **Jeong H**\*, Amoah A, Niemann J, Raeder B, Hasenfeld P, Benito E, Jann J-C, Nowak D, Ho A, Geiger H, Shui, S, Rausch T, Sanders AD#, Korbel JO#, (2024.05) "Cell type-specific consequences of mosaic structural variants in hematopoietic stem and progenitor cells", **Nature Genetics [\*co-first]**

3. **Hyobin Jeong**\*, Karen Grimes*, Kerstin K. Rauwolf, Peter-Martin Bruch, Tobias Rausch, Patrick Hasenfeld, Eva Benito Garagorri, Tobias Roider, Radhakrishnan Sabarinathan, David Porubsky, Sophie A. Herbst, Büşra Erarslan-Uysal, Johann-Christoph Jann, Tobias Marschall, Daniel Nowak, Jean-Pierre Bourquin, Andreas E. Kulozik, Sascha Dietrich, Beat Bornhauser, Ashley D. Sanders#, Jan O. Korbel#, (2023.06) "Functional analysis of structural variants in single cells using Strand-seq", **Nature Biotechnology [\*co-first]**

# KSBi-BIML 2026

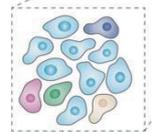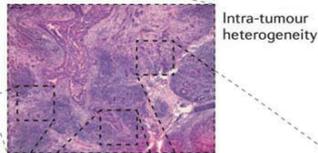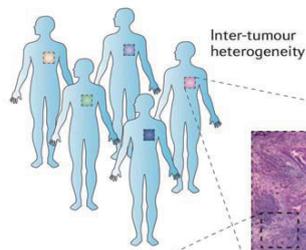## Single-cell multi-omics analysis to study tumor subclones

Hyobin Jeong, Assistant professor,
Department of Systems Biology, Yonsei University

---

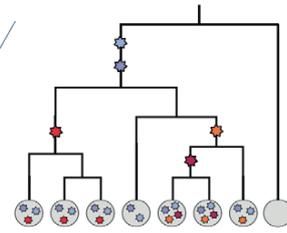## Tumor is composed of multiple subclones that makes intra-tumor heterogeneity
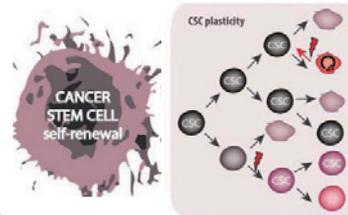


*Acta Pharmacologica Sinca (2015)*

2

# Multi-layered heterogeneity contributes to therapy failure and cancer progression



Inter-tumour heterogeneity

Intra-tumour heterogeneity

Dominance of clone 1

Dominance of clone 2

Mixed dominance

*Nature review cancer, 2012*

*Genetic variation*

*Genome Biology, 2016*

*Epigenetic (functional) alteration*

*Molecular cancer, 2017*

---

# How can we tackle the issues with intra-tumor heterogeneity?

- 암에서 이러한 서브클론들을 동정하기 위해 어떤 싱글셀 오믹스 기법들이 개발되어 있을까?

- 이러한 싱글셀 오믹스 데이터를 분석하기 위해 어떤 생명 정보학적인 도구들을 사용할 수 있을까?

- 서브클론의 동정 뿐 아니라 그 기능적 특성을 파악하기 위해서는 유전체와 전사체 또는 후성유전체 데이터를 함께 분석하는 싱글셀 멀티 오믹스 분석이 필요하다. 이를 구현하기 위한 생명 정보학적인 방법에는 어떤 것들이 있을까?
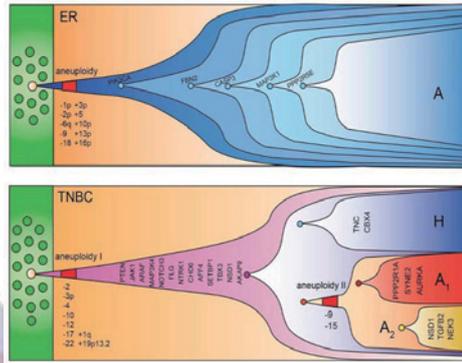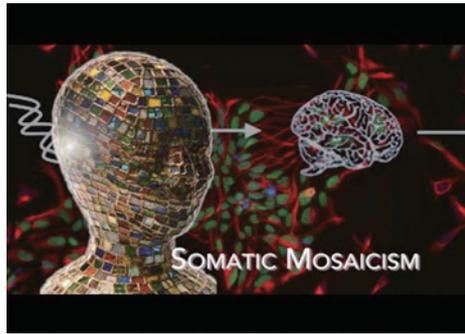
# Part1. 암에서 서브클론을 동정하기 위한 싱글셀 오믹스 기법들에 대한 소개

## Single-cell multi-omics analysis to study tumor subclones

---

## Single-cell technologies to explore cellular heterogeneity

Lineage

State

Trajectory

Cell surface proteins
- CITE-seq[20]
- REAP-seq[21]
- FACS[41,42]

Intracellular protein
- PEA[49,50]

Spatial position
- MERFISH[107,108,109]
- smFISH[102]
- STARmap[31]

DNA methylation
- scBS-seq[17]
- snmC-seq[16]
- sci-MET[19]

- scGESTALT[32]
- ScarTrace[33]
- LINNAEUS[34]
- MEMOIR[27]

Genome sequence
- SNS[9]
- SCI-seq[10]

Chromatin accessibility
- scATAC-seq[13]
- sciATAC-seq[14]
- scTHS-seq[15]
- 10X Genomics

Histone modifications
- scChIP-seq[23,24]

mRNA
- Drop-seq[4]
- InDrop[5]
- Smart-seq2[38]
- MARS-seq[3]
- 10X Genomics[6]
- SPLiT-seq[8]
- sci-RNA-seq[7]

Pseudotime
- Monocle[71,73]
- Wishbone[74]
- Velocyto[70]
- Diffusion[72]

*Tim Stuart & Rahul Satija*
*Nature review genetics, 2019*

6

# Genetic changes can happen in nucleotide level and also the form of larger rearrangement



*Single nucleotide variation*

*Structural Variation (SVs) – Strand-seq*

# Straucturual variation (SV) is a genomic rearrangement larger than 50bp



Oncogenic fusion

Inversion
(also Deletion, Translocation)

Oncogene amplification

Duplication

Enhancer hijacking

Translocation
(also Deletion, Inversion)

Tumour-suppressor deletion

Deletion

Genomic instability

Inversion
Translocation
Deletion
Amplification

*Macintyre et al. 2016*

# Structural variation (SV) is a key mutational process in cancer

Article | Open Access | Published: 05 February 2020

## Patterns of somatic structural variation in human cancer genomes

Yilong Li, Nicola D. Roberts, Jeremiah A. Wala, Ofer Shapira, Steven E. Schumacher, Kiran Kumar, Ekta Khurana, Sebastian Waszak, Jan O. Korbel, James E. Haber, Marcin Imielinski, PCAWG Structural Variation Working Group, Joachim Weischenfeldt ✉, Rameen Beroukhim ✉, Peter J. Campbell ✉ & PCAWG Consortium

*Nature* **578**, 112–121 (2020) | Cite this article

79k Accesses | 267 Citations | 175 Altmetric | Metrics



---

# Single-cell technologies to explore *genetic* heterogeneity



*Single-cell WGS*

Navin et al. Nature, 2011

Zahn et al. Nat Methods, 2017

*Tapestri (SDR-seq)*

Linderhofer et al. Nat Methods, 2025

*Strand-seq*

Sanders et al. Nat protocol, 2017

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

Step2. Remove PCR duplicate: Picard mark duplicate, Biobambam

Step3. Genotyping: Freebayes, GATK

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

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

10

# Single-cell technologies to explore *genetic* heterogeneity (Tapestri)

Missionbio Tapestri platform and Mosaic



https://missionbio.github.io/mosaic/manual/install.html

---

# Single-cell technologies to explore *genetic* heterogeneity (SDR-seq)
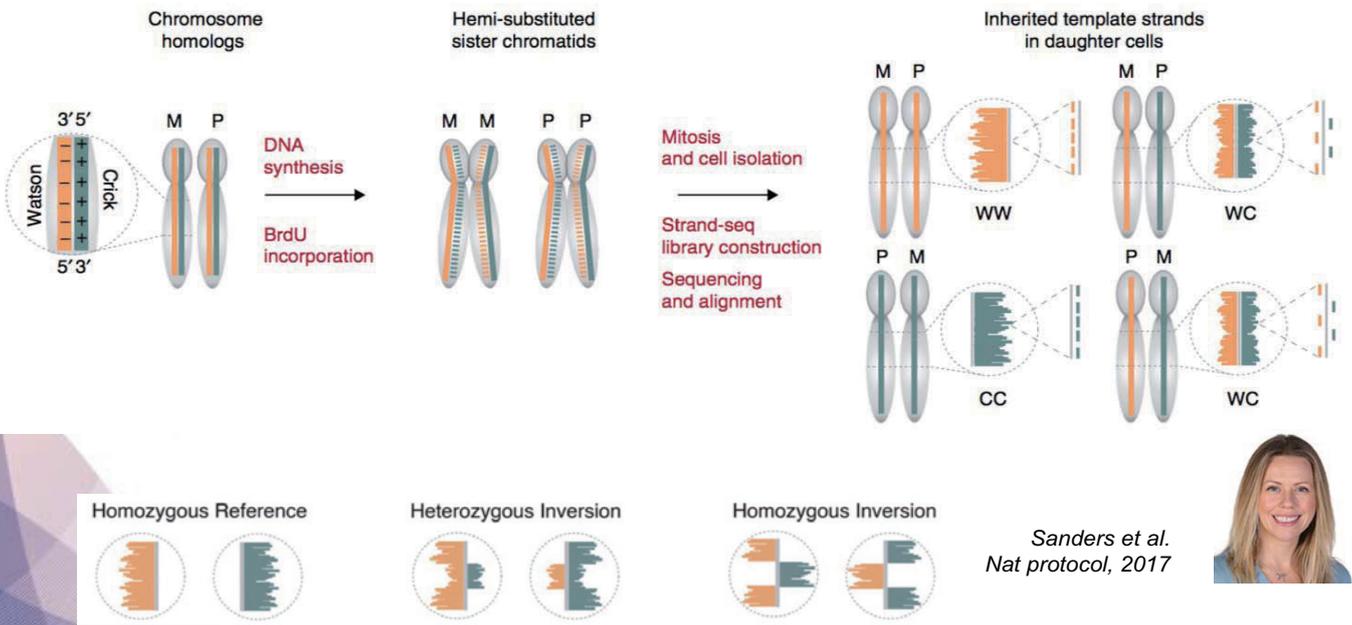
SDR-seq (Nature Methods, 2025) and SDRranger



https://github.com/hawkjo/SDRranger

# Single-cell technologies to explore *genetic* heterogeneity (Strand-seq)

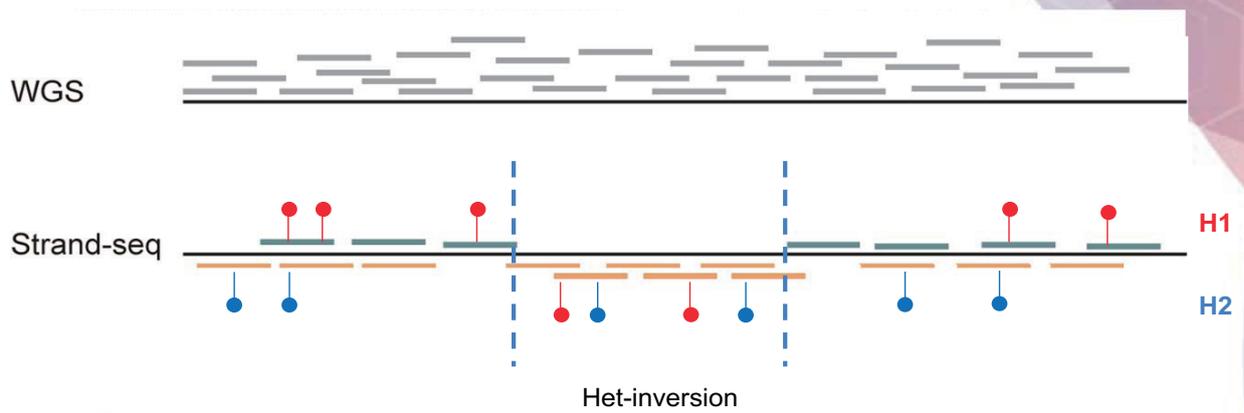Strand-seq (Nat protocol, 2017) and mosaicatcher (Nat biotech, 2020),
scNOVA (Nat biotech, 2023)



*Sanders et al.
Nat protocol, 2017*

---

# Part2. scDNA-seq 기법 중 Strand-seq 데이터에서 서브클론을 동정하는 방법 소개

## Single-cell multi-omics analysis to study tumor subclones

# Specialties of the Strand-seq data analysis



Het-inversion

- **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**

15

# Challenges of the Strand-seq data analysis

**Strand sequencing**



Less than 0.1x coverage

**Conventional sequencing (short-read)**



~ 30x coverage

16

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

**Single nucleotide variation (SNVs) - scWGS**
**Structural Variation**



*Single-nucleus sequencing (SNS)*

*Direct library*

*Strand-seq*

Navin et al. Nature, 2011

Zahn et al.
Nat Methods, 2017

Sanders et al.
Nat protocol, 2017

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

Step2. Remove PCR duplicate: Biobambam          **Quality checking!**

Step3. Genotyping, Haplotyping, Segmentation: StrandPhaseR, breakpointR

Step4. Structural variation calling: MosaiCatcher

Step5. Single-cell clustering and Phylogenetics

---

# Why the orientation of the reads are important?



*Porubsky et al. Nat comm, 2017*

Homozygous Reference          Heterozygous Inversion          Homozygous Inversion
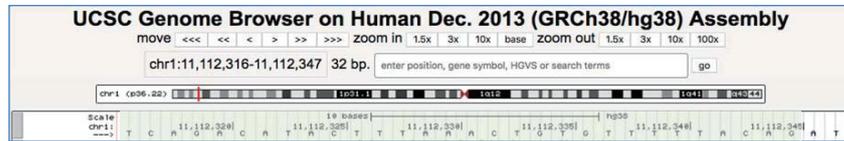
*Sanders et al. Genome Res, 2016*

- 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

A A A G T A T

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

**A A A G T A T**   Reverse (-) ☐ Watson (SAMFLAG 16)

19

---

**Decoding SAM flags**

This utility makes it easy to identify what are the properties of a read base
be for a given combination of properties.

To decode a given SAM flag value, just enter the number in the field below.

SAM Flag: 99    [Explain]

[Switch to mate]  Toggle first in pair / second in pair

**Find SAM flag by property:**

To find out what the SAM flag value would be for a given combination of properties, tick the boxes
for those that you'd like to include. The flag value will be shown in the SAM Flag field above.
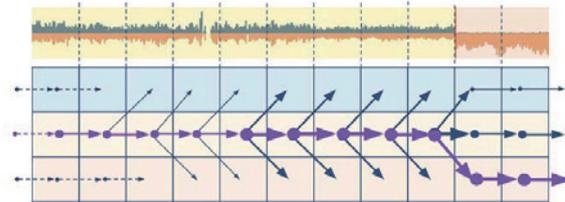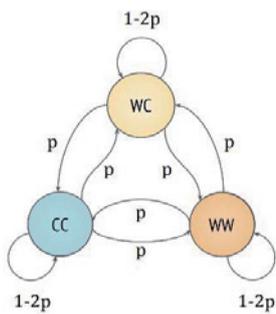
- ☑ read paired
- ☑ read mapped in proper pair
- ☐ read unmapped
- ☐ mate unmapped
- ☐ read reverse strand
- ☑ mate reverse strand
- ☑ first in pair
- ☐ second in pair
- ☐ not primary alignment
- ☐ read fails platform/vendor quality checks
- ☐ read is PCR or optical duplicate
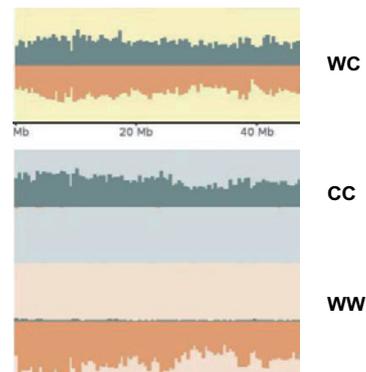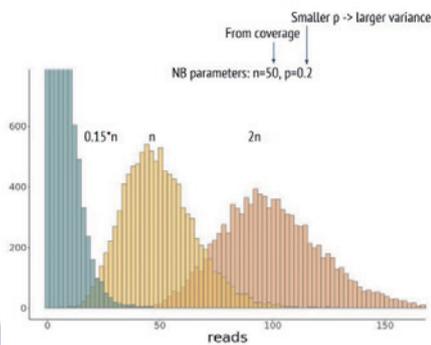- ☐ supplementary alignment

*https://broadinstitute.github.io/picard/explain-flags.html*

20

# Count the Watson and Crick reads using genomic windows



Genomic windows
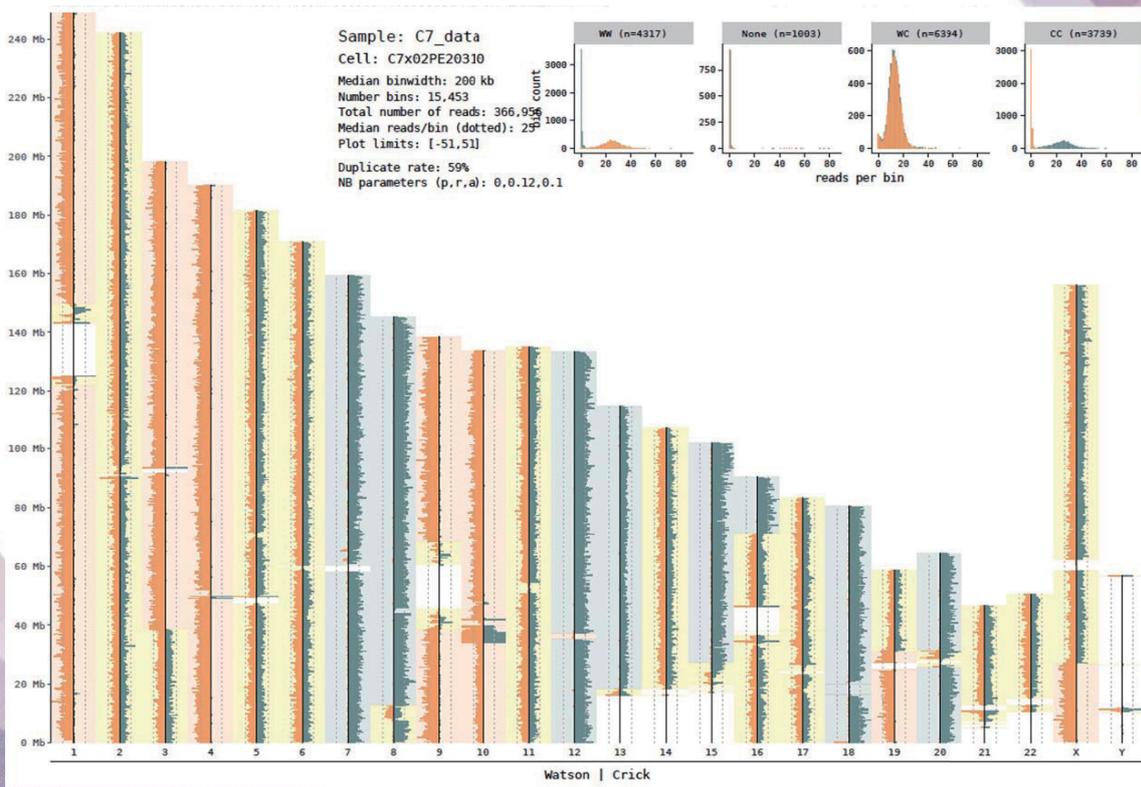20000 bp
50000 bp
100000 bp
200000 bp
500000 bp

# 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

Smaller p -> larger variance

From coverage

NB parameters: n=50, p=0.2

wc

cc

ww

# Strand-seq result of example single-cell

# Strand-seq result of T-ALL (leukemia) sample



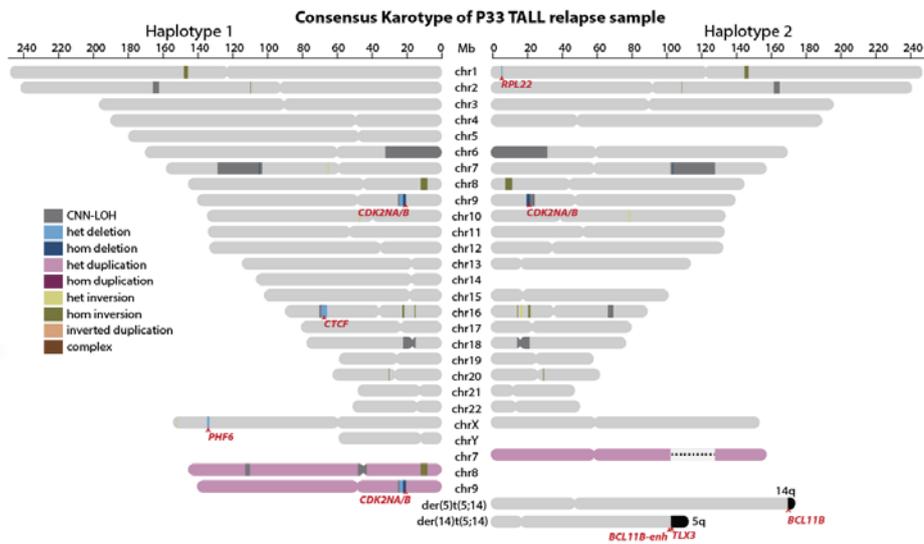*Figure from scTRIP manuscript,*
*Sanders et al. 2020*

# Mosaicatcher towards the automatic single-cell SV calling and clustering

https://github.com/friendsofstrandseq/mosaicatcher-pipeline



**v 1**

**Korbel group, EMBL**

*Ashley Sanders*  *Sasha Meiers*

**Marschall group, MPI informatics**

*David Porubsky*  *Maryam Ghareghani*

**v 2**

*Thomas Weber*

*Sanders et al. 2020*
*Weber et al. 2022 (ongoing)*

---

# Mosaicatcher calls single-cell SV using Bayesian framework



**Bayesian framework**

WT   Deletion   Inversion

# Mosaicatcher calls single-cell SV using Bayesian framework



*Figure from scTRIP manuscript, Sanders et al.*

- Input: single-cell BAM files
- Workflow management: Snakemake

- Binned read counting (100kb) and normalization
- Assign strand-specific read data into genomic bins
- Detects and haplotype-phases heterozygous SNPs
- Segments the single cell sequence data
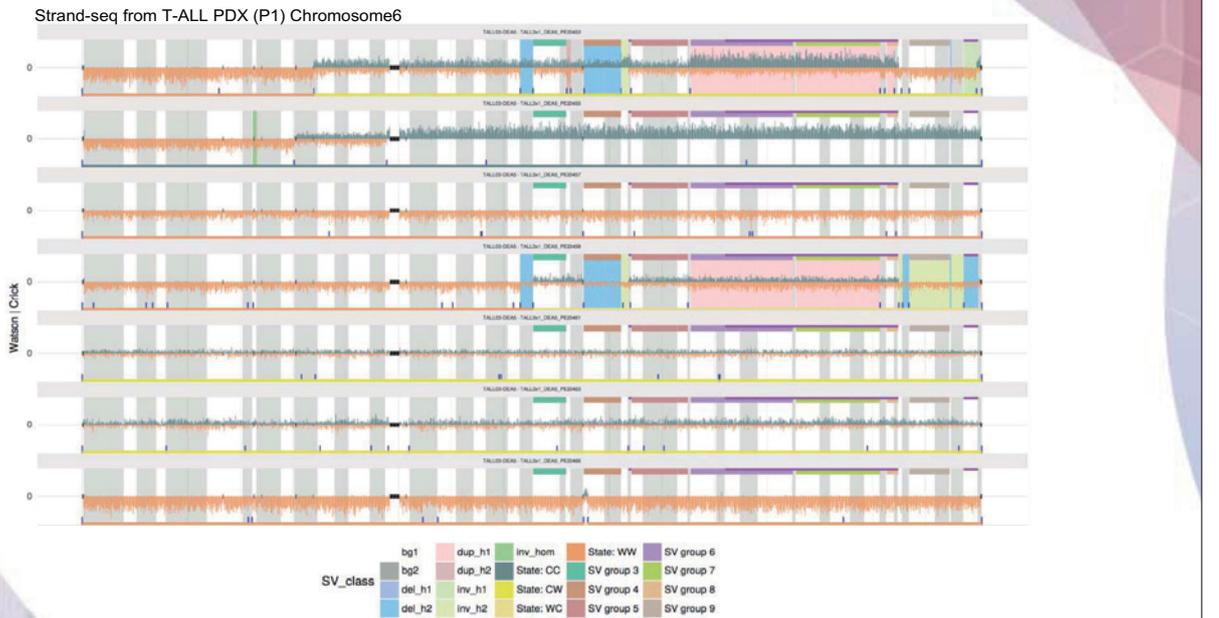- Calculates genotype likelihoods for each segment and single cell using Bayesian framework

# Mosaicatcher calls single-cell SV using Bayesian framework

# Chromosome plot with SVs called by MosaiCatcher framework

Strand-seq from T-ALL PDX (P1) Chromosome6

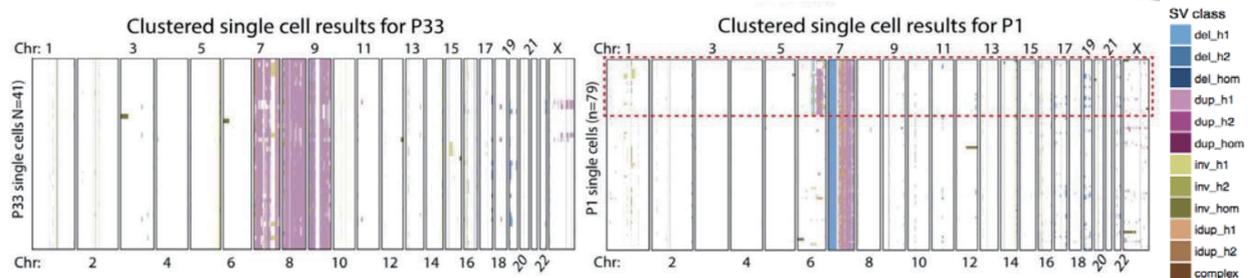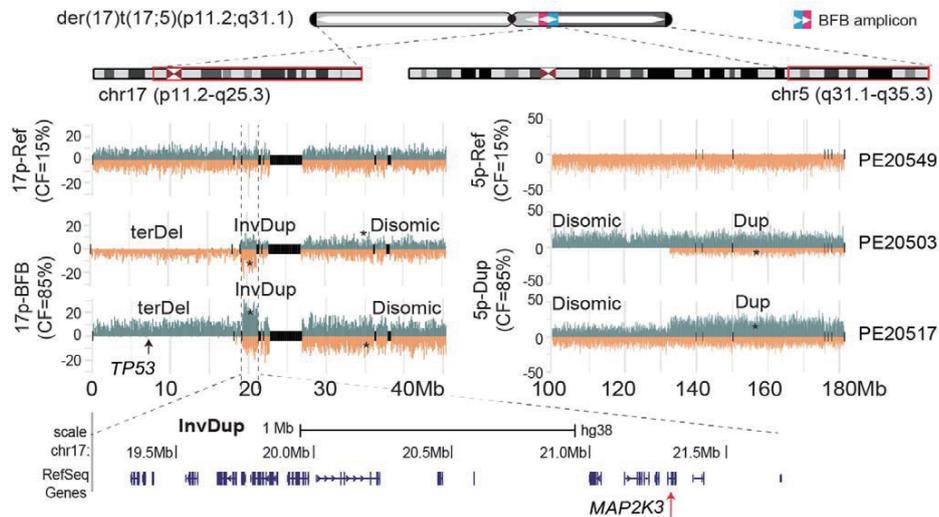# Heatmap of single-cells based on SVs called by MosaiCatcher framework



*Figure from scTRIP manuscript, Sanders et al. 2019*

- This heatmap was arranged using Ward's method for hierarchical clustering of SVs genotype likelihoods in two PDX samples
- P33 shows single dominant clone but P1 shows subclonal population in the sample represented by 23 cells
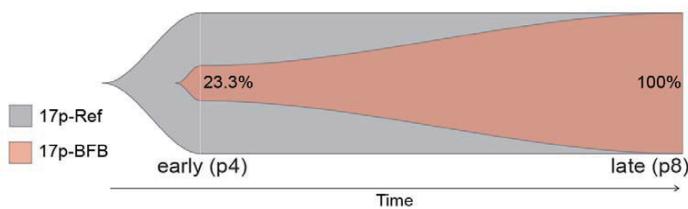
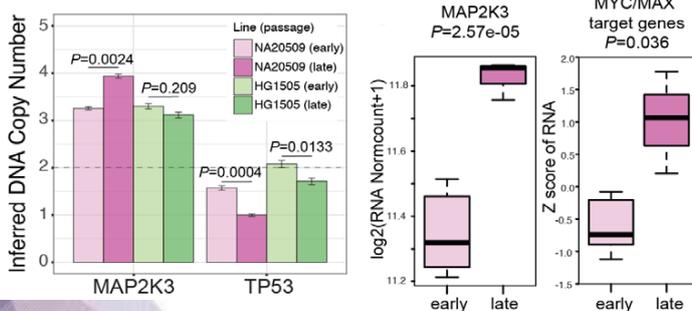## Subclones identified from Strand-seq and MosaiCatcher (Lymphoblastoid cell line, GM20509)



## Subclonal evolution can be analyzed using Strand-seq



- NA20509 (=GM20509) cell line was in culture for passage 4 (early) and passage 8 (late)

- MAP2K3, and MYC/MAX target genes were increased in late passage
- MYC expression was not changed

# Practical session – how to run Mosaicatcher

Genome analysis

## MosaiCatcher v2: a single-cell structural variations detection and analysis reference framework based on Strand-seq

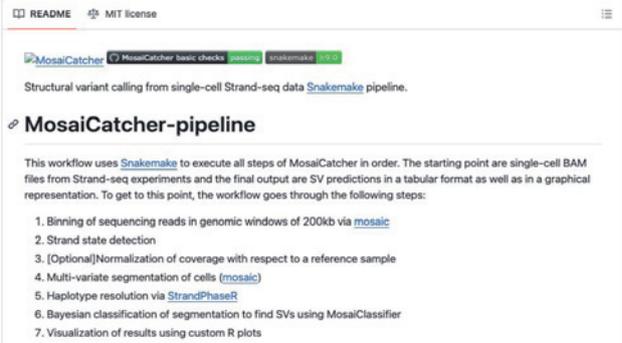Thomas Weber [1], Marco Raffaele Cosenza[1], Jan Korbel[1,2,*]

[1]European Molecular Biology Laboratory, Genome Biology Unit, Heidelberg, Germany
[2]Bridging Research Division on Mechanisms of Genomic Variation and Data Science, German Cancer Research Center (DKFZ), Heidelberg, Germany

*Corresponding author. Genome Biology Unit, EMBL Heidelberg, Meyerhofstraße 1, Heidelberg 69117, Germany. Tel: +49 6221 387-8822, fax: +49 6221 387-8518.
E-mail: jan.korbel@embl.de
Associate Editor: Can Alkan

**Abstract**
**Summary:** Single-cell DNA template strand sequencing (Strand-seq) allows a range of various genomic analysis including chromosome length haplotype phasing and structural variation (SV) calling in individual cells. Here, we present MosaiCatcher v2, a standardized workflow and reference framework for single-cell SV detection using Strand-seq. This framework introduces a range of functionalities, including: an automated upstream Quality Control (QC) and assembly sub-workflow that relies on multiple genome assemblies and incorporates a multistep normalization module, integration of the single-cell nucleosome occupancy and genetic variation analysis SV functional characterization and of the ArbiGent SV genotyping modules, platform portability, as well as a user-friendly and shareable web report. These new features of MosaiCatcher v2 enable reproducible computational processing of Strand-seq data, which are increasingly used in human genetics and single-cell genomics, toward production environments. MosaiCatcher v2 is compatible with both container and conda environments, ensuring reproducibility and robustness and positioning the framework as a cornerstone in computational processing of Strand-seq data.
**Availability and implementation:** MosaiCatcher v2 is a standardized workflow, implemented using the Snakemake workflow management system. The pipeline is available on GitHub: https://github.com/friendsofstrandseq/mosaicatcher-pipeline/ and on the snakemake-workflow-catalog: https://snakemake.github.io/snakemake-workflow-catalog/?usage=friendsofstrandseq/mosaicatcher-pipeline. Strand-seq example input data used in the publication can be found in the Data availability statement. Additionally, a lightweight dataset for test purposes can be found on the GitHub repository.

https://github.com/friendsofstrandseq/mosaicatcher-pipeline/

*Based on the tutorial written by Chiwon Chung*

---

# Data directory tree



Data directory must follow a certain format
- Pipeline requires you to have either
  1. BAM
  2. Fastq

Placed in like the image on the left.
- Fastq files (for ASHLEYS) must follow the following syntax:
→ [SAMPLE_NAME].x.fastq.gz
→ X is the strand number (coding/non-coding)

# Run_mosaicatcher.sh

```
5 #!/bin/bash
4 docker run --rm -it \
3     -v ./data:/pipeline/data \
2     -v ./out:/pipeline/out \
1     -e USER_ID="$(id -u)" \
6     mosaicatcher
1     /bin/bash
~
~
~
~
```

To start the pipeline
1. Copy the run_mosaicatcher file in /home/shared to directory you desire
2. Run the following command:
   /path/to/run_mosaicatcher.sh

./data: where your test data (and outputs) will be stored

./output: where your ASHLEYS result (and other checkpoint files) will be stored

You can change these (red boxes) to whatever you want.

Ex) You have data folder at ~/my_data:

-v ~/my_data:/pipeline/data

Do NOT change the rest of the script!

# Within the Docker file

```
(base) root@fb59a8d00a65:/pipeline# ll
total 56
drwxr-xr-x  1 root root   31 Jul 15 03:50 ./
drwxr-xr-x  1 root root    6 Jul 17 02:49 ../
drwxr-xr-x  9 root root  178 Jul 15 03:49 .git/
drwxr-xr-x  4 root root   86 Jul 15 03:49 .github/
-rw-r--r--  1 root root 5076 Jul 15 03:49 .gitignore
-rw-r--r--  1 root root  403 Jul 15 03:49 .gitmodules
-rw-r--r--  1 root root 1697 Jul 15 03:49 .gitpod.Dockerfile
-rw-r--r--  1 root root 2238 Jul 15 03:49 .gitpod.yml
-rw-r--r--  1 root root  547 Jul 15 03:49 .pre-commit-config.yaml
-rw-r--r--  1 root root  118 Jul 15 03:49 .snakemake-workflow-catalog.yml
drwxr-xr-x  5 root root  154 Jul 15 03:49 .tests/
-rw-r--r--  1 root root 1472 Jul 15 03:49 CHANGELOG.md
-rw-r--r--  1 root root 1092 Jul 15 03:49 LICENSE
-rw-r--r--  1 root root 4226 Jul 15 03:49 README.md
drwxr-xr-x  5 root root 4096 Jul 15 03:49 afac/
drwxr-xr-x  2 root root   70 Jul 15 03:49 config/
drwxr-xr-x  5 1001 1001   60 Jul 15 03:14 data/
drwxr-xr-x  3 root root  116 Jul 15 03:49 docs/
drwxr-xr-x  2 root root   25 Jul 15 03:49 envs/
drwxr-xr-x  3 root root 4096 Jul 15 03:49 github-actions-runner/
drwxr-xr-x  3 1001 1001  176 Jul 15 11:07 out/
-rwxrwxr-x  1 root root  546 Jul 15 01:01 run_pipelinev8.sh*
drwxr-xr-x 10 root root  153 Jul 15 03:49 workflow/
```

If the command is run correctly, you should now be inside the docker container. Using the ll command as seen in the image, we can see there is a shell script:

→ run_pipelinev8.sh

This is the script we will use to run the actual pipeline.

The data/ and out/ folders should be connected to the data and out folders you set in the run_mosaicatcher.sh file.

# Within the Docker file

```bash
#!/bin/bash
snakemake \
    --cores 10 \
    --configfile /pipeline/config/config.yaml \
    --config \
        data_location=/pipeline/data \
        ashleys_pipeline=True \
        ashleys_pipeline_only=False \
        scNOVA=False \
        scNOVA_manual_cell_selection=False \
        chromosomes_to_exclude=[] \
        mosaicatcher_pipeline=True \
        use_light_data=False \
        publishdir=/pipeline/out \
        user=${USER_ID} \
    --profile workflow/snakemake_profiles/mosaicatcher-pipeline/v8/local/conda/ \
    --forceall \
    --dry-run \
~
~
```

Use vim run_mosaicatcherv8.sh to access the script. Please refer to online tutorials for vim controls, as it is out of scope with this manual.

You only need to change the lines accentuated by red boxes.

The pipeline will be run TWICE:
1. First run is for ASHLEYS and mosaicatcher pipeline
2. Second is for scNOVA pipeline

The dry run command is for testing purposes only, and will be removed when the pipeline is run.

# Running the pipeline

```
total                                  882

Reasons:
    (check individual jobs above for details)
    forced:
        aggregate_summary_statistics, all, ashleys_ashleys_final_r
le_for_mosaic_count, ashleys_generate_features, ashleys_mark_dupli
gative_control_bypass, ashleys_predict, ashleys_publishdir_outputs
ns, change_ownership, check_single_paired_end, combine_strandphase
aplotag_likelihoods, download_hg38_reference, estimate_ploidy, fil
ser_output, merge_blacklist_bins_for_norm, merge_haplotag_tables,
plot_clustering, plot_clustering_chromosome_dev, plot_clustering_p

This was a dry-run (flag -n). The order of jobs does not reflect t
The run involves checkpoint jobs, which will result in alteration
(base) root@fb59a8d00a65:/pipeline#
```

First, run the mosaicatcher_v8.sh file as is; it should do a test run.

If all is without issues, you will see no errors and the CLI will show the total number of rules (which will vary depending on your dataset size)

In case you see errors, the most likely reason is that your files are not set up correctly. Make sure your fastq/bam files are placed like it's shown in slide 2.

# Running the pipeline

```bash
#/bin/bash
snakemake \
    --cores 10 \
    --configfile /pipeline/config/config.yaml \
    --config \
        data_location=/pipeline/data \
        ashleys_pipeline=True \
        ashleys_pipeline_only=False \
        scNOVA=False \
        scNOVA_manual_cell_selection=False \
        chromosomes_to_exclude=[] \
        mosaicatcher_pipeline=True \
        use_light_data=False \
        publishdir=/pipeline/out \
        user=${USER_ID} \
    --profile workflow/snakemake_profiles/mosaicatcher-pipeline/v8/local/conda/ \
    --forceall \
~
```

If you see no errors, remove the dry run option that was originally in the run_mosaicatcherv8.sh pipeline, then run it again. The first step of the pipeline should start running.

The duration will vary but expect it to take at least a couple of hours. Depending on server CPU usage, it may even take a full day.

# Running the pipeline

```bash
#/bin/bash
snakemake \
    --cores 10 \
    --configfile /pipeline/config/config.yaml \
    --config \
        data_location=/pipeline/data \
        ashleys_pipeline=False \
        ashleys_pipeline_only=False \
        scNOVA=True \
        scNOVA_manual_cell_selection=False \
        chromosomes_to_exclude=["chrY"] \
        mosaicatcher_pipeline=True \
        use_light_data=False \
        publishdir=/pipeline/out \
        user=${USER_ID} \
    --profile workflow/snakemake_profiles/mosaicatcher-pipeline/v8/local/conda/ \
    --forceall \
```

```
├── bam
├── fastq
│   ├── SAMPLE1.1.fastq.gz
│   ├── SAMPLE1.2.fastq.gz
│   ├── SAMPLE2.1.fastq.gz
│   ├── SAMPLE2.2.fastq.gz
│   ├── SAMPLE3.1.fastq.gz
│   ├── SAMPLE3.2.fastq.gz
└── scNOVA_input_user
    └── input_subclonality.txt
```

```
Filename        Subclonality
TALL3x01PE20406 clone2
TALL3x01PE20414 clone2
TALL3x01PE20415 clone1
TALL3x01PE20416 clone1
TALL3x01PE20417 clone1
TALL3x01PE20418 clone1
TALL3x01PE20419 clone1
TALL3x01PE20421 clone1
TALL3x01PE20422 clone1
TALL3x01PE20424 clone2
TALL3x01PE20427 clone1
TALL3x01PE20430 clone1
TALL3x01PE20433 clone1
TALL3x01PE20435 clone2
```

Once the pipeline is finished, create a new directory scNOVA_input_user and add your subclonality file. It MUST be named input_subclonality.txt, and it MUST be a tsv file with the correct header names. Otherwise, the pipeline will throw errors.

To run the scNOVA pipeline, please change the run_mosaicatcherv8.sh as shown on the left. Note that scNOVA currently does not support chromosome Y, and must be place in the chromosomes_to_exclude list.
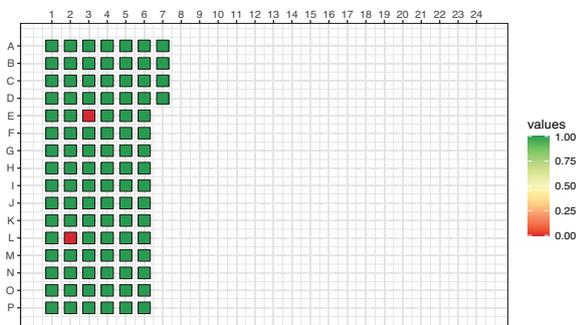
# Output – QC result based on ASHLEY algorithm

out/data/RPE_mixture/cell_selection/labels_positive_control_corrected.tsv

```
cell       prediction      probability      sample
BM510x3PE20401.sort.mdup.bam    1       0.8672   RPE_mixture
BM510x3PE20402.sort.mdup.bam    1       0.9472   RPE_mixture
BM510x3PE20403.sort.mdup.bam    1       0.905    RPE_mixture
BM510x3PE20406.sort.mdup.bam    1       0.9052   RPE_mixture
BM510x3PE20407.sort.mdup.bam    1       0.8948   RPE_mixture
BM510x3PE20408.sort.mdup.bam    1       0.8764   RPE_mixture
BM510x3PE20410.sort.mdup.bam    1       0.9333   RPE_mixture
BM510x3PE20411.sort.mdup.bam    1       0.9145   RPE_mixture
BM510x3PE20414.sort.mdup.bam    1       0.8525   RPE_mixture
BM510x3PE20415.sort.mdup.bam    1       0.8981   RPE_mixture
BM510x3PE20416.sort.mdup.bam    1       0.8384   RPE_mixture
BM510x3PE20417.sort.mdup.bam    1       0.9306   RPE_mixture
BM510x3PE20418.sort.mdup.bam    1       0.841    RPE_mixture
BM510x3PE20419.sort.mdup.bam    1       0.9133   RPE_mixture
BM510x3PE20421.sort.mdup.bam    1       0.7389   RPE_mixture
BM510x3PE20422.sort.mdup.bam    1       0.8608   RPE_mixture
BM510x3PE20423.sort.mdup.bam    1       0.9013   RPE_mixture
BM510x3PE20424.sort.mdup.bam    1       0.8732   RPE_mixture
BM510x3PE20425.sort.mdup.bam    1       0.8763   RPE_mixture
BM510x3PE20426.sort.mdup.bam    1       0.9577   RPE_mixture
RPE1WTPE20401.sort.mdup.bam     1       0.9291   RPE_mixture
RPE1WTPE20402.sort.mdup.bam     1       0.8426   RPE_mixture
RPE1WTPE20403.sort.mdup.bam     1       0.9345   RPE_mixture
RPE1WTPE20404.sort.mdup.bam     1       0.8848   RPE_mixture
RPE1WTPE20405.sort.mdup.bam     1       0.8993   RPE_mixture
RPE1WTPE20406.sort.mdup.bam     1       0.9239   RPE_mixture
RPE1WTPE20407.sort.mdup.bam     1       0.9367   RPE_mixture
RPE1WTPE20409.sort.mdup.bam     0       0.0018   RPE_mixture
RPE1WTPE20410.sort.mdup.bam     1       0.9106   RPE_mixture
RPE1WTPE20411.sort.mdup.bam     1       0.9282   RPE_mixture
RPE1WTPE20412.sort.mdup.bam     1       0.9132   RPE_mixture
RPE1WTPE20413.sort.mdup.bam     1       0.9341   RPE_mixture
RPE1WTPE20414.sort.mdup.bam     1       0.9612   RPE_mixture
RPE1WTPE20415.sort.mdup.bam     1       0.9289   RPE_mixture
```

---

# Output – QC result based on ASHLEY algorithm

out/data/RPE_mixture/plots/plate/ashleys_plate_predictions.pdf
out/data/RPE_mixture/plots/plate/ashleys_plate_probabilities.pdf



Sample: RPE_mixture l ASHLEYS binary predictions (cutoff=0.5)

Sample: RPE_mixture l ASHLEYS probabilities

**반올림해서 0.5가 되면 통과되는 결과임 (0.4759 → PASS)

# Output – Plotting pipeline

out/data/RPE_mixture/plots/counts/CountComplete.raw.pdf



---

# Output – Plotting pipeline

Plotting pipeline – basic information for each single-cell libraries

out/data/RPE_mixture/counts/RPE_mixture.info_raw

# Output – Plotting pipeline

Plotting pipeline – single-cell count matrix using 200kb bins

out/data/RPE_mixture/counts/RPE_mixture.txt.raw.gz

```
chrom   start    end      sample        cell             c       w       class
chr1    0        200000   RPE_mixture   BM510x3PE20401   4       6       WW
chr1    200000   400000   RPE_mixture   BM510x3PE20401   0       0       WW
chr1    400000   600000   RPE_mixture   BM510x3PE20401   2       0       WW
chr1    600000   800000   RPE_mixture   BM510x3PE20401   1       8       WW
chr1    800000   1000000  RPE_mixture   BM510x3PE20401   0       30      WW
chr1    1000000  1200000  RPE_mixture   BM510x3PE20401   1       35      WW
chr1    1200000  1400000  RPE_mixture   BM510x3PE20401   0       18      WW
chr1    1400000  1600000  RPE_mixture   BM510x3PE20401   1       22      WW
chr1    1600000  1800000  RPE_mixture   BM510x3PE20401   0       23      WW
chr1    1800000  2000000  RPE_mixture   BM510x3PE20401   0       22      WW
chr1    2000000  2200000  RPE_mixture   BM510x3PE20401   0       38      WW
chr1    2200000  2400000  RPE_mixture   BM510x3PE20401   0       30      WW
chr1    2400000  2600000  RPE_mixture   BM510x3PE20401   0       33      WW
chr1    2600000  2800000  RPE_mixture   BM510x3PE20401   0       37      WW
chr1    2800000  3000000  RPE_mixture   BM510x3PE20401   0       45      WW
chr1    3000000  3200000  RPE_mixture   BM510x3PE20401   0       50      WW
chr1    3200000  3400000  RPE_mixture   BM510x3PE20401   0       36      WW
chr1    3400000  3600000  RPE_mixture   BM510x3PE20401   0       31      WW
chr1    3600000  3800000  RPE_mixture   BM510x3PE20401   2       42      WW
chr1    3800000  4000000  RPE_mixture   BM510x3PE20401   0       45      WW
chr1    4000000  4200000  RPE_mixture   BM510x3PE20401   0       36      WW
chr1    4200000  4400000  RPE_mixture   BM510x3PE20401   1       46      WW
chr1    4400000  4600000  RPE_mixture   BM510x3PE20401   0       37      WW
chr1    4600000  4800000  RPE_mixture   BM510x3PE20401   0       37      WW
chr1    4800000  5000000  RPE_mixture   BM510x3PE20401   0       51      WW
chr1    5000000  5200000  RPE_mixture   BM510x3PE20401   0       39      WW
chr1    5200000  5400000  RPE_mixture   BM510x3PE20401   0       51      WW
chr1    5400000  5600000  RPE_mixture   BM510x3PE20401   0       44      WW
chr1    5600000  5800000  RPE_mixture   BM510x3PE20401   0       38      WW
chr1    5800000  6000000  RPE_mixture   BM510x3PE20401   0       36      WW
chr1    6000000  6200000  RPE_mixture   BM510x3PE20401   0       44      WW
chr1    6200000  6400000  RPE_mixture   BM510x3PE20401   0       39      WW
chr1    6400000  6600000  RPE_mixture   BM510x3PE20401   0       28      WW
chr1    6600000  6800000  RPE_mixture   BM510x3PE20401   0       33      WW
chr1    6800000  7000000  RPE_mixture   BM510x3PE20401   0       40      WW
chr1    7000000  7200000  RPE_mixture   BM510x3PE20401   0       42      WW
```

# Output – haplotype phasing result

data/RPE_mixture/strandphaser/StrandPhaseR_final_output.txt

```
chrom   start   end          sample        cell             class
chr1    0       1800000      RPE_mixture   RPE1WTPE20430    WW
chr1    0       4800000      RPE_mixture   RPE1WTPE20441    CW
chr1    0       13800000     RPE_mixture   RPE1WTPE20432    WC
chr1    0       22000000     RPE_mixture   RPE1WTPE20457    WC
chr1    0       22600000     RPE_mixture   RPE1WTPE20452    WW
chr1    0       34800000     RPE_mixture   RPE1WTPE20426    CW
chr1    0       43800000     RPE_mixture   BM510x3PE20414   WW
chr1    0       59000000     RPE_mixture   RPE1WTPE20423    CW
chr1    0       67400000     RPE_mixture   RPE1WTPE20448    WW
chr1    0       78000000     RPE_mixture   RPE1WTPE20438    CW
chr1    0       93200000     RPE_mixture   RPE1WTPE20465    CC
chr1    0       119000000    RPE_mixture   RPE1WTPE20418    CC
chr1    0       119600000    RPE_mixture   RPE1WTPE20490    WW
chr1    0       119800000    RPE_mixture   BM510x3PE20419   CC
chr1    0       119800000    RPE_mixture   RPE1WTPE20428    WW
chr1    0       119800000    RPE_mixture   RPE1WTPE20459    CC
chr1    0       119800000    RPE_mixture   RPE1WTPE20462    WW
chr1    0       119800000    RPE_mixture   RPE1WTPE20480    WW
chr1    0       122400000    RPE_mixture   RPE1WTPE20402    CW
chr1    0       122400000    RPE_mixture   RPE1WTPE20483    WW
chr1    0       150800000    RPE_mixture   RPE1WTPE20440    WC
chr1    0       194200000    RPE_mixture   RPE1WTPE20424    WW
chr1    0       201000000    RPE_mixture   BM510x3PE20408   WC
chr1    0       209200000    RPE_mixture   RPE1WTPE20439    WW
chr1    0       213800000    RPE_mixture   RPE1WTPE20429    CW
chr1    0       219600000    RPE_mixture   RPE1WTPE20494    CC
chr1    0       224200000    RPE_mixture   BM510x3PE20418   WW
chr1    0       248956422    RPE_mixture   BM510x3PE20401   WW
chr1    0       248956422    RPE_mixture   BM510x3PE20402   WW
chr1    0       248956422    RPE_mixture   BM510x3PE20403   CW
chr1    0       248956422    RPE_mixture   BM510x3PE20406   WC
chr1    0       248956422    RPE_mixture   BM510x3PE20407   CC
chr1    0       248956422    RPE_mixture   BM510x3PE20410   CC
```
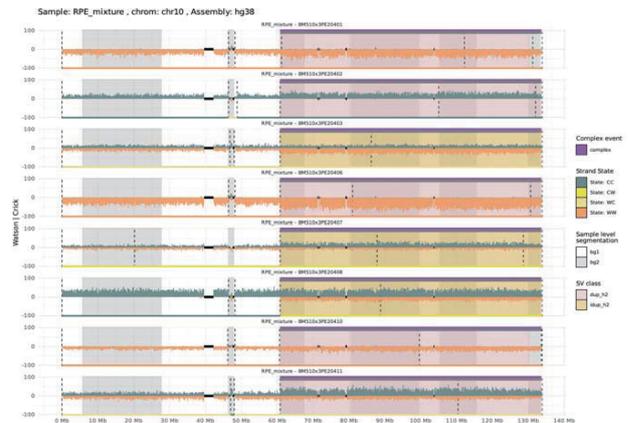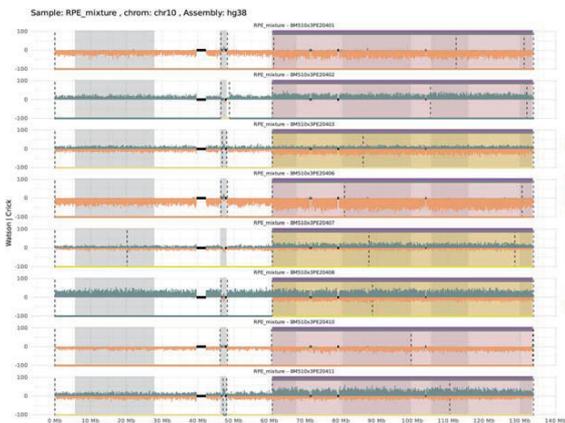
# Output – SV calling result

data/RPE_mixture/mosaiclassifier/sv_calls/stringent_filterTRUE.tsv

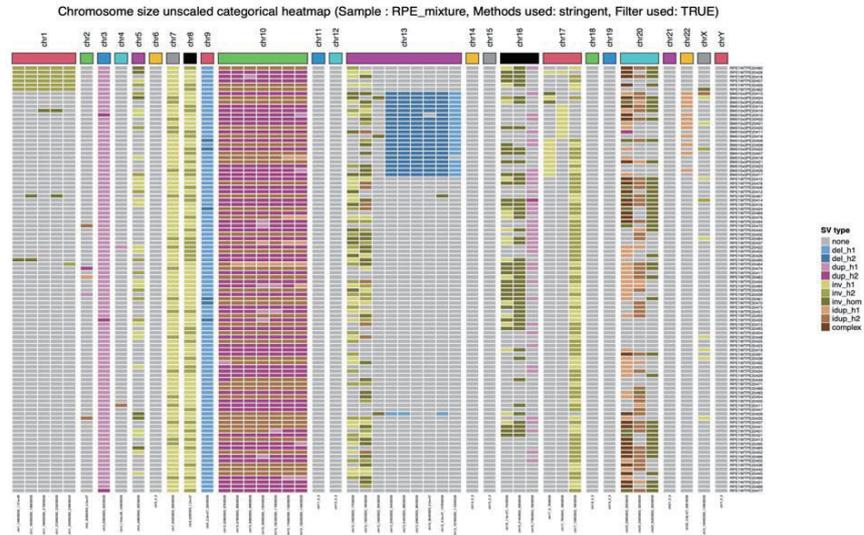| chrom | start | end | sample | cell | class | scalar | num_bins | sv_call_name | sv_call_haplotype | sv_call_name_2nd | sv_call_haplotype_2nd | llr_to_ref | llr_to_2nd | af |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| chr1 | 149600000 | 157000000 | RPE_mixture | RPE1WTPE20418 | CW | 1 | 37 | inv_h2 1001 | ref_hom 1010 | 78.2183473084405 | 78.2183473084405 | 0.06 | | |
| chr1 | 149600000 | 157000000 | RPE_mixture | RPE1WTPE20428 | WC | 1 | 37 | inv_h2 1001 | inv_hom 101 | 117.162866864643 | 114.552001025681 | 0.06 | | |
| chr1 | 149600000 | 157000000 | RPE_mixture | RPE1WTPE20457 | WC | 1 | 37 | inv_h2 1001 | ref_hom 1010 | 102.184205496873 | 102.184205496873 | 0.06 | | |
| chr1 | 149600000 | 157000000 | RPE_mixture | RPE1WTPE20459 | CW | 1 | 37 | inv_h2 1001 | inv_hom 101 | 93.1970086748642 | 63.6245523739165 | 0.06 | | |
| chr1 | 149600000 | 157000000 | RPE_mixture | RPE1WTPE20462 | WC | 1 | 37 | inv_h2 1001 | inv_hom 101 | 87.205544088843 | 75.6074814292192 | 0.06 | | |
| chr1 | 149600000 | 157000000 | RPE_mixture | RPE1WTPE20480 | WC | 1 | 37 | inv_h2 1001 | ref_hom 1010 | 117.162866864643 | 117.162866864643 | 0.06 | | |
| chr1 | 149600000 | 157000000 | RPE_mixture | RPE1WTPE20490 | CW | 1 | 37 | inv_hom 101 | ref_hom 1010 | 26.5767240273936 | 26.5767240273936 | 0.01 | | |
| chr1 | 185800000 | 198600000 | RPE_mixture | RPE1WTPE20418 | CW | 1 | 64 | inv_h2 1001 | ref_hom 1010 | 101.586495522311 | 101.586495522311 | 0.06 | | |
| chr1 | 185800000 | 198600000 | RPE_mixture | RPE1WTPE20428 | WC | 1 | 64 | inv_h2 1001 | ref_hom 1010 | 134.539550531405 | 134.539550531405 | 0.06 | | |
| chr1 | 185800000 | 198600000 | RPE_mixture | RPE1WTPE20457 | WC | 1 | 64 | inv_h2 1001 | inv_hom 101 | 170.488337814053 | 164.580823397578 | 0.06 | | |
| chr1 | 185800000 | 198600000 | RPE_mixture | RPE1WTPE20462 | WC | 1 | 64 | inv_h2 1001 | ref_hom 1010 | 140.531015078513 | 140.531015078513 | 0.06 | | |
| chr1 | 185800000 | 198600000 | RPE_mixture | RPE1WTPE20467 | CW | 1 | 64 | inv_h2 1001 | inv_hom 101 | 119.560889163617 | 95.6789811058175 | 0.06 | | |
| chr1 | 185800000 | 198600000 | RPE_mixture | RPE1WTPE20480 | WC | 1 | 64 | inv_h2 1001 | ref_hom 1010 | 65.8221598875551 | 65.8221598875551 | 0.02 | | |
| chr1 | 185800000 | 198600000 | RPE_mixture | RPE1WTPE20490 | CW | 1 | 64 | inv_hom 101 | ref_hom 1010 | 332.257880585969 | 332.257880585969 | 0.06 | | |
| chr1 | 185800000 | 198600000 | RPE_mixture | BM510x3PE20419 | WC | 1 | 75 | inv_hom 101 | ref_hom 1010 | 29.873372604907 | 29.873372604907 | 0.02 | | |
| chr1 | 198600000 | 213800000 | RPE_mixture | RPE1WTPE20418 | CW | 1 | 75 | inv_h2 1001 | ref_hom 1010 | 59.6809873811256 | 59.6809873811256 | 0.01 | | |
| chr1 | 198600000 | 213800000 | RPE_mixture | RPE1WTPE20428 | WC | 1 | 75 | inv_h2 1001 | ref_hom 1010 | 251.46362341753 | 251.46362341753 | 0.06 | | |
| chr1 | 198600000 | 213800000 | RPE_mixture | RPE1WTPE20457 | WC | 1 | 75 | inv_h2 1001 | inv_hom 101 | 320.365465709272 | 302.624730157902 | 0.06 | | |
| chr1 | 198600000 | 213800000 | RPE_mixture | RPE1WTPE20459 | CW | 1 | 75 | inv_h2 1001 | inv_hom 101 | 245.472158870422 | 209.757029677729 | 0.06 | | |
| chr1 | 198600000 | 213800000 | RPE_mixture | RPE1WTPE20462 | WC | 1 | 75 | inv_h2 1001 | inv_hom 101 | 200.536174767112 | 173.80824239508 | 0.06 | | |
| chr1 | 198600000 | 213800000 | RPE_mixture | RPE1WTPE20480 | WC | 1 | 75 | inv_h2 1001 | ref_hom 1010 | 524.075260310944 | 524.075260310944 | 0.06 | | |
| chr1 | 213800000 | 224200000 | RPE_mixture | BM510x3PE20419 | WC | 1 | 52 | inv_hom 101 | ref_hom 1010 | 17.7454792146946 | 17.7454792146946 | 0.02 | | |
| chr1 | 213800000 | 224200000 | RPE_mixture | RPE1WTPE20418 | CW | 1 | 52 | inv_h2 1001 | ref_hom 1010 | 153.000317013216 | 153.000317013216 | 0.06 | | |
| chr1 | 213800000 | 224200000 | RPE_mixture | RPE1WTPE20428 | WC | 1 | 52 | inv_h2 1001 | ref_hom 1010 | 224.897891578511 | 224.897891578511 | 0.06 | | |
| chr1 | 213800000 | 224200000 | RPE_mixture | RPE1WTPE20457 | WC | 1 | 52 | inv_h2 1001 | inv_hom 101 | 185.953372022309 | 180.190821901831 | 0.06 | | |
| chr1 | 213800000 | 224200000 | RPE_mixture | RPE1WTPE20459 | CW | 1 | 52 | inv_h2 1001 | inv_hom 101 | 153.000317013214 | 93.3145859687631 | 0.06 | | |
| chr1 | 213800000 | 224200000 | RPE_mixture | RPE1WTPE20462 | WC | 1 | 52 | inv_h2 1001 | inv_hom 101 | 123.04299427766 | 99.306050515857 | 0.06 | | |
| chr1 | 213800000 | 224200000 | RPE_mixture | RPE1WTPE20467 | CW | 1 | 52 | inv_hom 101 | ref_hom 1010 | 14.7497469411407 | 14.7497469411407 | 0.02 | | |
| chr1 | 213800000 | 224200000 | RPE_mixture | RPE1WTPE20480 | WC | 1 | 52 | inv_h2 1001 | ref_hom 1010 | 344.727182520671 | 344.727182520671 | 0.06 | | |
| chr1 | 243600000 | 248956422 | RPE_mixture | RPE1WTPE20418 | CW | 1 | 27 | inv_h2 1001 | ref_hom 1010 | 15.9733231459201 | 15.9733231459201 | 0.07 | | |
| chr1 | 243600000 | 248956422 | RPE_mixture | RPE1WTPE20428 | WC | 1 | 27 | inv_h2 1001 | ref_hom 1010 | 15.9733231459225 | 15.9733231459225 | 0.07 | | |
| chr1 | 243600000 | 248956422 | RPE_mixture | RPE1WTPE20457 | CC | 1 | 27 | inv_h2 1001 | ref_hom 1010 | 15.2900282760901 | 15.2900282760901 | 0.07 | | |
| chr1 | 243600000 | 248956422 | RPE_mixture | RPE1WTPE20459 | CW | 1 | 27 | inv_h2 1001 | ref_hom 1010 | 15.9733231255162 | 15.9733231255162 | 0.07 | | |
| chr1 | 243600000 | 248956422 | RPE_mixture | RPE1WTPE20462 | WC | 1 | 27 | inv_h2 1001 | ref_hom 1010 | 15.9733221189223 | 15.9733221189223 | 0.07 | | |
| chr1 | 243600000 | 248956422 | RPE_mixture | RPE1WTPE20464 | WW | 1 | 27 | inv_h2 1001 | ref_hom 1010 | 15.2900282760917 | 15.2900282760917 | 0.07 | | |
| chr1 | 243600000 | 248956422 | RPE_mixture | RPE1WTPE20480 | WC | 1 | 27 | inv_h2 1001 | ref_hom 1010 | 15.9733231459229 | 15.9733231459229 | 0.07 | | |
| chr10 | 60800000 | 67800000 | RPE_mixture | BM510x3PE20401 | WW | 1 | 35 | dup_h2 1020 | dup_h1 2010 | inf | 22.2185094368695 | 0.57 | | |
| chr10 | 60800000 | 67800000 | RPE_mixture | BM510x3PE20402 | CC | 1 | 35 | dup_h2 1020 | dup_h1 2010 | inf | 11.8213017284703 | 0.57 | | |

# Output – SV calling result

data/RPE_mixture/plots/sv_calls_dev/stringent_filterTRUE/*.pdf



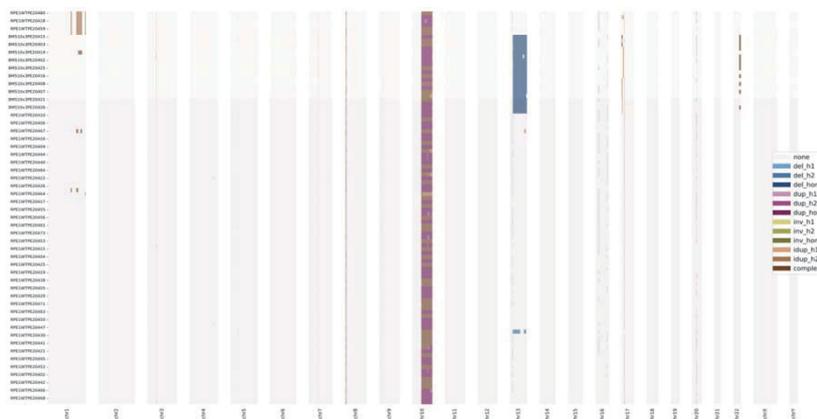data/RPE_mixture/plots/sv_calls_dev/lenient_filterFALSE/*.pdf

# Output – SV calling result

data/RPE_mixture/plots/sv_clustering_dev/stringent-filterTRUE-position.pdf



Chromosome size unscaled categorical heatmap (Sample : RPE_mixture, Methods used: stringent, Filter used: TRUE)
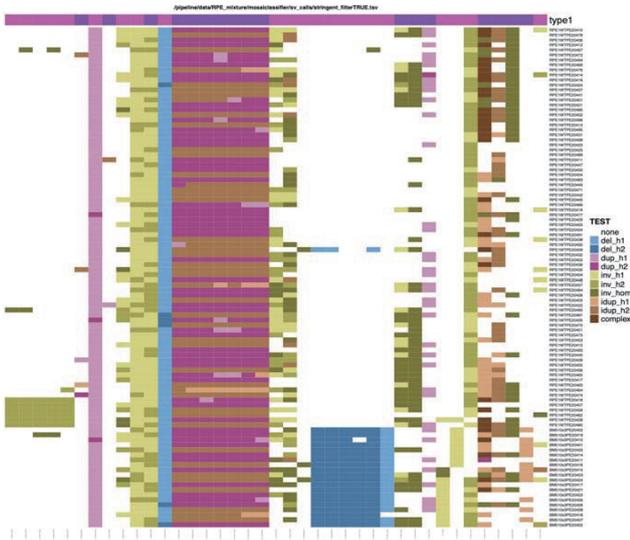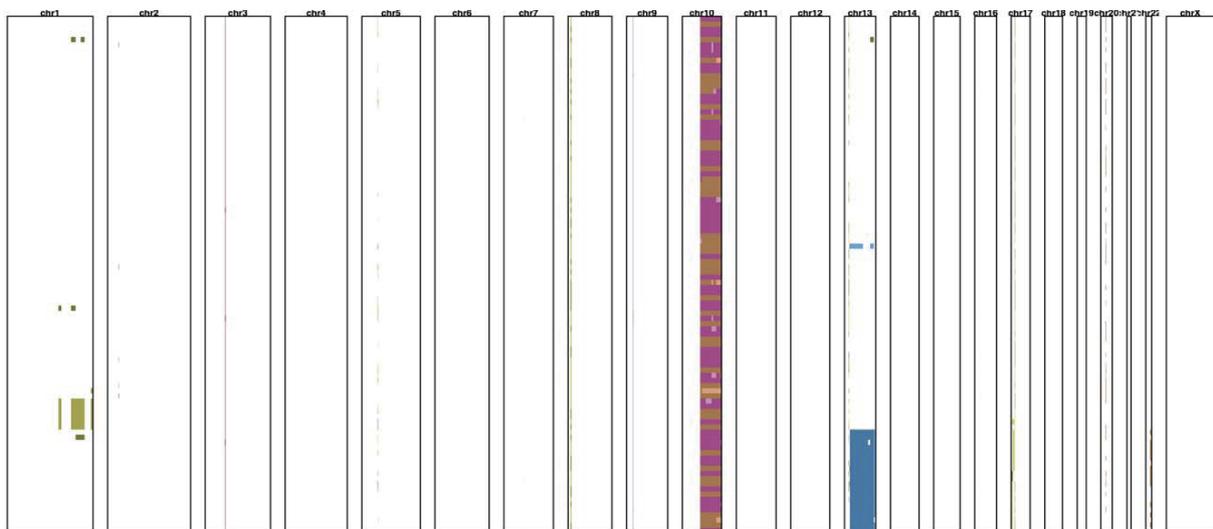
# Output – SV calling result

data/RPE_mixture/plots/sv_clustering_dev/stringent-filterTRUE-chromosome.pdf

# Output – SV calling result

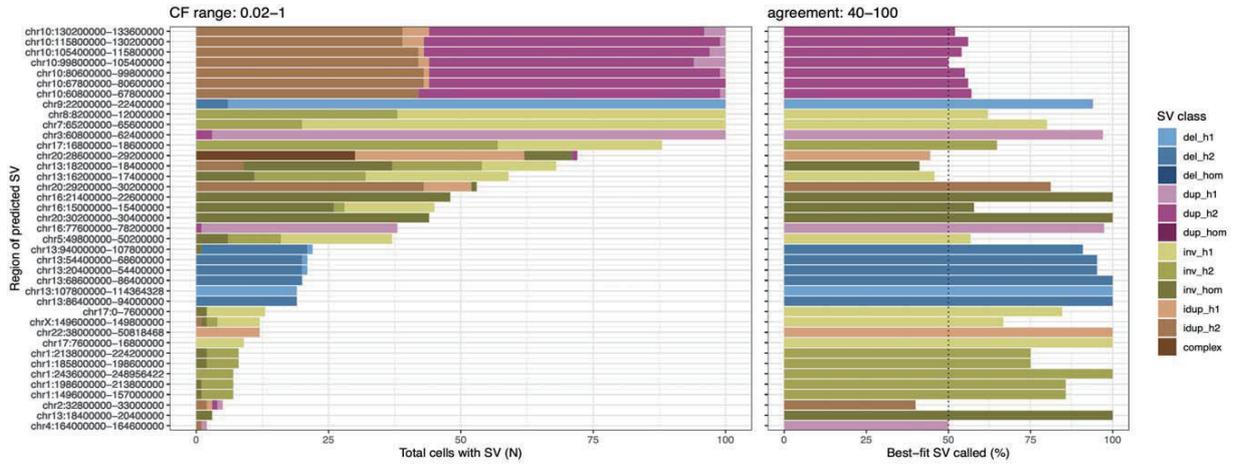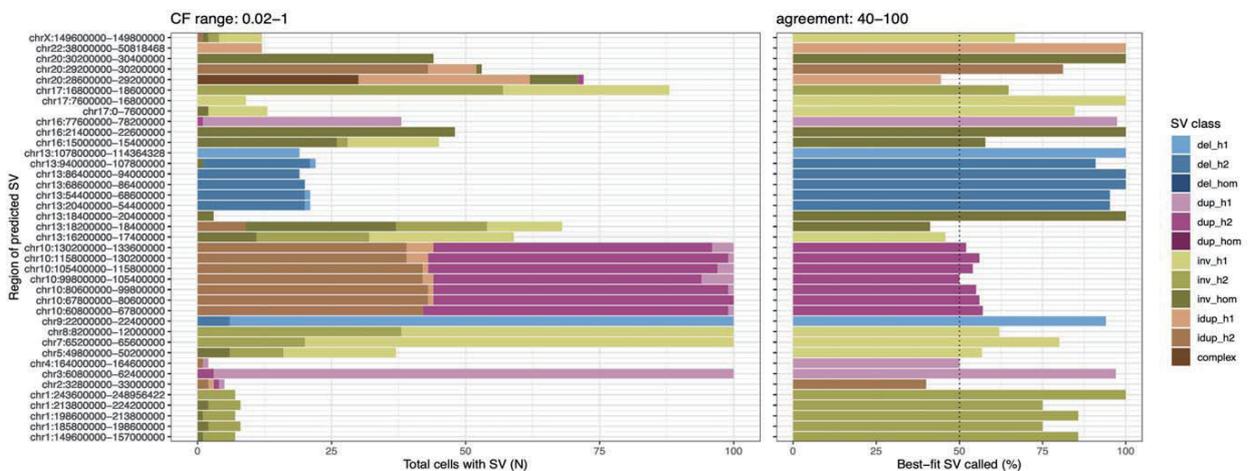data/RPE_mixture/plots/sv_clustering/stringent-filterTRUE-position.pdf
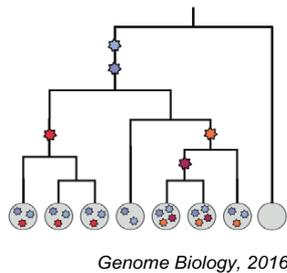


# Output – SV calling result

data/RPE_mixture/plots/sv_clustering/stringent-filterTRUE-chromosome.pdf

# Output – SV calling result

data/RPE_mixture/plots/sv_consistency/stringent_filterTRUE.consistency-barplot-byaf.pdf



# Output – SV calling result

data/RPE_mixture/plots/sv_consistency/stringent_filterTRUE.consistency-barplot-bypos.pdf

# Part3. scNOVA – Strand-seq 에서 동정한 서브클론의 기능적 분석을 위한 멀티오믹스 기법

Single-cell multi-omics analysis to study tumor subclones

---

## How can we measure functional consequence of somatic structural variants in different subclones?



*Genome Biology, 2016*

**Genetic variation**

**Epigenetic (functional) alteration**

?

*Molecular cancer, 2017*

Oncogenic fusion — Inversion (also Deletion, Translocation)

Oncogene amplification — Duplication

Enhancer hijacking — Translocation (also Deletion, Inversion)

Tumour-suppressor deletion — Deletion

Genomic instability — Inversion Translocation Deletion Amplification

*Macintyre et al. 2016*

56

# Single-cell technologies to explore _functional_ heterogeneity



DNA accessibility
scNOME-seq
scATAC-seq
scDNAse-seq

Chromosome organization
scHIC

Transcription
scRNA-seq

DNA modifications
scBS-seq
scAba-seq
CLEVER-seq

Histone modifications
scChIP-seq

_Kelsey et al. Science, 2017_

_scMNase-seq, Lai et al. 2018_

https://doi.org/10.1038/s41586-018-0567-3

## LETTER

### Principles of nucleosome organization revealed by single-cell micrococcal nuclease sequencing

Binbin Lai[1], Weiwu Gao[1,2], Kairong Cui[1], Wanli Xie[1,3], Qingsong Tang[1], Wenfei Jin[1], Gangqing Hu[1], Bing Ni[2] & Keji Zhao[1]*

**_Can we use Nucleosome Occupancy to study functional consequence of SVs ?_**

---

# Nucleosomes are the basic unit of chromatin which slide along DNA



- Nucleosome is composed of two copies of four core histones together with 146~147bp of DNA

- Human diploid genomes have 30 million nucleosomes

- Transcriptionally active gene promoters exhibit a prominent nucleosome-depleted region (NDR) directly upstream of the TSS

_Nat Rev Mol Cell Biol, 2017_

# Nucleosomes pattern is informative for the gene expression and cell type of origin

Cell free DNA protected by nucleosome is secreted to the blood



# Nucleosomes pattern is informative for the gene expression and cell type of origin



Inferring expressed genes by whole-genome sequencing of plasma DNA

*Nat Genet, 2016*

Cell-free DNA Comprises an In Vivo Nucleosome Footprint that Informs Its Tissues-Of-Origin

*Cell, 2016*

# Nucleosome dynamics can be measured by genomic assays



**MNase-seq**

MNase digestion

TF

Endonuclease activity of MNase

Exonuclease activity of MNase

DNA purification

Library preparation and sequencing

Sequencing adapter

Composite signal

*Nat Rev Genet, 2014*

- MNase is a secreted glycoprotein with a preference for single-stranded DNA and RNA

- It cleave one strand of DNA when the helix 'breathes' and subsequently cleave the other strand to generate a double-strand break

- It then 'nibbles' the exposed DNA end until it reaches an obstruction, such as a nucleosome



FAIRE
DNase
MNase
ATAC

MNase
DNase
FAIRE
ATAC

*Epigenetics & Chromatin, 2014*

---

# Strand-seq protocol involves MNase treatment



**Strand-seq protocol**

Steps 1–11
Isolation of hemi-substituted cells

Nuclei preparation

Single-nuclei sort

*Sanders et al. Nature protocol, 2017*

Steps 12–19 MNase digestion

96-well plate with single nuclei

Nuclei isolation

MNase

MNase digestion

*Trends in genetics, 2017*

- Strand-seq is a single-cell based DNA sequencing method which gives haplotype-resolved structural variation information.

- Question: Does Strand-seq profile reflects nucleosome occupancy?

- Question: If then, can Strand-seq additionally provides information of gene expression and cell identity?

# Nucleosome position and occupancy can be detected from Strand-seq data

**Han Chinese (CHS) trio (Lymphoblastoid cell line) chr12:34,346,260-34,349,260**

# Nucleosome occupancy is negatively correlated with gene expression level

# Nucleosome occupancy in the genebody is informative for differential expression

## Input data (Strand-seq)

### RPE-1 (182 cells)

| | cell1 | cell2 | ... | cell N |
|---|---|---|---|---|
| Gene1 | 10 | 30 | ... | 5 |
| Gene2 | 3 | 2 | ... | 0 |
| ... | ... | ... | ... | ... |
| Gene N | 30 | 50 | ... | 80 |

19770 genes

### LCL (224 cells)

| | cell1 | cell2 | ... | cell N |
|---|---|---|---|---|
| Gene1 | 1 | 2 | ... | 1 |
| Gene2 | 8 | 4 | ... | 5 |
| ... | ... | ... | ... | ... |
| Gene N | 14 | 25 | ... | 10 |

19770 genes

## Approach (DESeq of nucleosome occupancy)

*Anders et al. 2010, Love et al. 2014*

RPE1 up-regulated DEGs    LCL up-regulated DEGs

AUC=0.885    AUC=0.772

65

---

# Nucleosome occupancy can be used to classify cell-type



66

## scNOVA : Coupling genome-epigenome using Strand-seq technology

**Strand-seq**
Nuclei isolation with BrdU labeled chromosome

MNase digestion

Sequencing

Haplotype-aware SV

1. Depth
2. Strand
3. Phase

[scTRIP]

Haplotype-aware NO

cell1
cell2
cell3

Gene1    Gene2

CRE1  CRE2

*Jeong\* and Grimes\* et al…. Sanders and Korbel Nature Biotech, 2022*

**Haplotype specific NO (local/cis-effect of SVs)**

H1    SV (e.g. Inversion)
H2
                    CRE        Gene
CRE occupancy          Differential gene activity
(NO at CREs)            (NO at gene bodies)

**Clone specific NO (global/trans-effect of SVs)**

clone1        clone2

67

---

## Computational pipeline of scNOVA

**Pre-processing: Haplotype-resolved SV discovery in single cells (scTRIP)**

scNOVA input : bam files of single-cell libraries, SV calls, and subclone assignment

| Clone1 Plate1, N cells | Clone1 Plate2, N cells | Clone2 Plate1, N cells | Clone2 Plate2, N cells |

**Module1 : Compute NO at gene bodies and cis-regulatory elements (CREs)**

Output result both from haplotype aware/unaware manner

*Module2 : Infer haplotype-specific NO (local effect of SVs)*

**Haplotype resolved NO**

Haplotype comparison at Gene body

CREs grouped by nearest gene

Sliding window (300kb)

*Module3 : Infer altered gene activity (global effect of SVs)*

**Model : CNN**

Train : RPE-1
Features : NO, GC contents, CpG%, RT, single-cell variance

**Single-cell DE analysis**

Gene body region, Generalized linear model DESeq2 (*Love et al.*)

**Filtering Expressed genes**

| | Clone1 | Clone2 |
|---|---|---|
| Gene1 | expressed | expressed |
| Gene2 | expressed | expressed |
| Gene3 | none | none |

**Combine result**

Expressed & DE

**Link somatic SVs to single-cell functional readout**

*How can it be helpful to understand the global effect of SV?*

https://github.com/jeongdo801/scNOVA

68

- 34 -

# How subclonal SVs alter the epigenome and phenotype?

| System | SVs |
|---|---|

**T-ALL P1**

*Andreas Kulozik group,
Beat Bornhauser,
Jean-Pierre Bourquin
Uni Zurich*



*Major clone (68%)*

*Subclone (32%) (only appeared in relapse)*
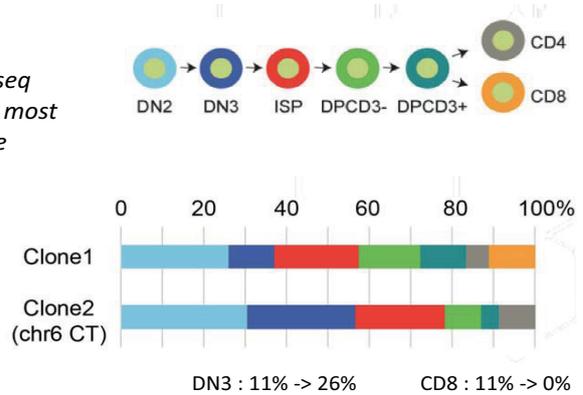
*Sanders et al. 2020*

# SV subclone in P1 shows increase of premature stages in the cellular hierarchy

*ATAC-seq signature matrix (2020 peaks)*


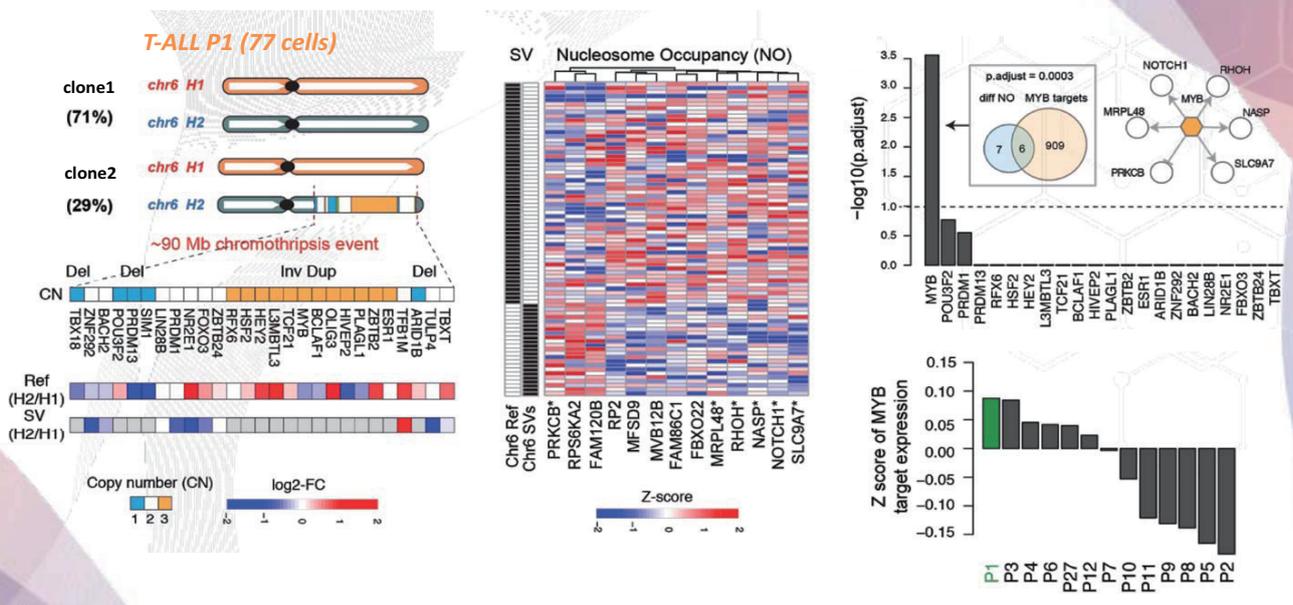
*Project Strand-seq single-cell data to most likely cell type*

*T cell differentiation stages*

DN3 : 11% -> 26%        CD8 : 11% -> 0%

*ATAC-seq and signature matrix from Erarslan-Uysal et al. EMBO Mol Med, 2020*

# SV subclone in T-ALL P1 shows altered MYB target genes including NOTCH1



*T-ALL P1 (77 cells)*

*How the cell type (state) composition different in clone1 and clone2?*

---

# Notch signaling and MYB has been reported in T-ALL oncogenesis

## Validation of increased dosage of MYB expression in rearranged haplotype



**Bulk RNA-seq**

**MYB H2/H1 relapse**
log2-FC = 0.45
(1.37 fold increase)
p-value = 0.0317

*Single-cell experiment is needed to confirm subclonal level transcriptome changes*

---

## Practical session – how to run scNOVA

OXFORD

Genome analysis

**MosaiCatcher v2: a single-cell structural variations detection and analysis reference framework based on Strand-seq**

Thomas Weber[1], Marco Raffaele Cosenza[1], Jan Korbel[1,2,*]

[1]European Molecular Biology Laboratory, Genome Biology Unit, Heidelberg, Germany
[2]Bridging Research Division on Mechanisms of Genomic Variation and Data Science, German Cancer Research Center (DKFZ), Heidelberg, Germany

*Corresponding author. Genome Biology Unit, EMBL Heidelberg, Meyerhofstraße 1, Heidelberg 69117, Germany. Tel: +49 6221 387-8822, fax: +49 6221 387-8518. E-mail: jan.korbel@embl.de
Associate Editor: Can Alkan

**Abstract**
**Summary:** Single-cell DNA template strand sequencing (Strand-seq) allows a range of various genomic analysis including chromosome length haplotype phasing and structural variation (SV) calling in individual cells. Here, we present MosaiCatcher v2, a standardized workflow and reference framework for single-cell SV detection using Strand-seq. This framework introduces a range of functionalities, including: an automated upstream Quality Control (QC) and assembly sub-workflow that relies on multiple genome assemblies and incorporates a multistep normalization module, integration of the single-cell nucleosome occupancy and genetic variation analysis SV functional characterization and of the ArbiGent SV genotyping modules, platform portability, as well as a user-friendly and shareable web report. These new features of MosaiCatcher v2 enable reproducible computational processing of Strand-seq data, which are increasingly used in human genetics and single-cell genomics, toward production environments. MosaiCatcher v2 is compatible with both container and conda environments, ensuring reproducibility and robustness and positioning the framework as a cornerstone in computational processing of Strand-seq data.
**Availability and implementation:** MosaiCatcher v2 is a standardized workflow, implemented using the Snakemake workflow management system. The pipeline is available on GitHub: https://github.com/friendsofstrandseq/mosaicatcher-pipeline/ and on the snakemake-workflow-catalog: https://snakemake.github.io/snakemake-workflow-catalog/?usage=friendsofstrandseq/mosaicatcher-pipeline. Strand-seq example input data used in the publication can be found in the Data availability statement. Additionally, a lightweight dataset for test purposes can be found on the GitHub repository.

📖 README  ⚖ MIT license

MosaiCatcher  MosaiCatcher basic checks passing  snakemake

Structural variant calling from single-cell Strand-seq data Snakemake pipeline.

🔗 **MosaiCatcher-pipeline**

This workflow uses Snakemake to execute all steps of MosaiCatcher in order. The starting point are single-cell BAM files from Strand-seq experiments and the final output are SV predictions in a tabular format as well as in a graphical representation. To get to this point, the workflow goes through the following steps:

1. Binning of sequencing reads in genomic windows of 200kb via mosaic
2. Strand state detection
3. [Optional]Normalization of coverage with respect to a reference sample
4. Multi-variate segmentation of cells (mosaic)
5. Haplotype resolution via StrandPhaseR
6. Bayesian classification of segmentation to find SVs using MosaiClassifier
7. Visualization of results using custom R plots

https://github.com/friendsofstrandseq/ mosaicatcher-pipeline/

*Based on the tutorial written by Chiwon Chung*

# Practical session – running the pipeline

```bash
#!/bin/bash
snakemake \
    --cores 10 \
    --configfile /pipeline/config/config.yaml \
    --config \
        data_location=/pipeline/data \
        ashleys_pipeline=False \
        ashleys_pipeline_only=False \
        scNOVA=True \
        scNOVA_manual_cell_selection=False \
        chromosomes_to_exclude=["chrY"] \
        mosaicatcher_pipeline=True \
        use_light_data=False \
        publishdir=/pipeline/out \
        user=${USER_ID} \
    --profile workflow/snakemake_profiles/mosaicatcher-pipeline/v8/local/conda/ \
    --forceall \
```

Once the pipeline is finished, create a new directory scNOVA_input_user and add your subclonality file. It MUST be named input_subclonality.txt, and it MUST be a tsv file with the correct header names. Otherwise, the pipeline will throw errors.

To run the scNOVA pipeline, please change the run_mosaicatcherv8.sh as shown on the left. Note that scNOVA currently does not support chromosome Y, and must be place in the chromosomes_to_exclude list.

```
├── bam
├── fastq
│   ├── SAMPLE1.1.fastq.gz
│   ├── SAMPLE1.2.fastq.gz
│   ├── SAMPLE2.1.fastq.gz
│   ├── SAMPLE2.2.fastq.gz
│   ├── SAMPLE3.1.fastq.gz
│   ├── SAMPLE3.2.fastq.gz
└── scNOVA_input_user
    └── input_subclonality.txt
```

```
Filename        Subclonality
TALL3x01PE20406 clone2
TALL3x01PE20414 clone2
TALL3x01PE20415 clone1
TALL3x01PE20416 clone1
TALL3x01PE20417 clone1
TALL3x01PE20418 clone1
TALL3x01PE20419 clone1
TALL3x01PE20421 clone1
TALL3x01PE20422 clone1
TALL3x01PE20424 clone2
TALL3x01PE20427 clone1
TALL3x01PE20430 clone1
TALL3x01PE20433 clone1
TALL3x01PE20435 clone2
```

# Output – folder structure

```
(base) [hyobinjeong@node01 Project_mosaicatcher_RPE_mixture_Tutorial]$ ls -lh data/RPE_mixture/
total 260K
drwxr-xr-x  2 hyobinjeong hyobinjeong  20K Oct  5 22:10 bam
drwxr-xr-x  2 hyobinjeong hyobinjeong  12K Sep 18 14:09 bam_ashleys
drwxr-xr-x  2 hyobinjeong hyobinjeong 4.0K Oct  5 22:13 cell_selection
drwxr-xr-x  2 hyobinjeong hyobinjeong 4.0K Oct  5 22:10 checks
drwxr-xr-x  3 hyobinjeong hyobinjeong 4.0K Oct  6 01:02 config
drwxr-xr-x  3 hyobinjeong hyobinjeong 4.0K Oct  5 23:18 counts
drwxrwxrwx  2 hyobinjeong hyobinjeong 8.0K Sep 18 09:56 fastq
drwxr-xr-x  5 hyobinjeong hyobinjeong   57 Sep 18 15:01 haplotag
drwxr-xr-x 48 hyobinjeong hyobinjeong  12K Oct  6 02:18 log
drwxr-xr-x  7 hyobinjeong hyobinjeong  135 Sep 18 15:13 mosaiclassifier
drwxr-xr-x  3 hyobinjeong hyobinjeong   26 Sep 18 10:10 normalizations
drwxr-xr-x  2 hyobinjeong hyobinjeong   10 Oct  6 01:26 nucleosome_sampleA
drwxr-xr-x  2 hyobinjeong hyobinjeong   10 Oct  6 01:26 nucleosome_sampleB
drwxr-xr-x  2 hyobinjeong hyobinjeong  103 Oct  6 01:02 ploidy
drwxr-xr-x 10 hyobinjeong hyobinjeong  192 Sep 18 15:15 plots
drwxr-xr-x  2 hyobinjeong hyobinjeong  126 Sep 18 14:09 predictions
drwxr-xr-x  4 hyobinjeong hyobinjeong  166 Oct  5 23:49 scNOVA_bam_merge
drwxr-xr-x 48 hyobinjeong hyobinjeong 100K Oct  6 01:32 scNOVA_bam_modified
drwxr-xr-x  2 hyobinjeong hyobinjeong 4.0K Oct  6 02:04 scNOVA_input_user
drwxrwxr-x  4 hyobinjeong hyobinjeong   70 Oct  6 01:26 scNOVA_nucleosomes_bam
drwxr-xr-x  7 hyobinjeong hyobinjeong 4.0K Oct  6 03:02 scNOVA_result
drwxr-xr-x 25 hyobinjeong hyobinjeong 4.0K Oct  6 02:12 scNOVA_result_CNN
drwxr-xr-x  2 hyobinjeong hyobinjeong  170 Oct  6 02:18 scNOVA_result_haplo
drwxr-xr-x  2 hyobinjeong hyobinjeong  122 Oct  6 02:13 scNOVA_result_plots
drwxr-xr-x  3 hyobinjeong hyobinjeong 4.0K Oct  5 23:50 segmentation
drwxr-xr-x  2 hyobinjeong hyobinjeong  12K Oct  5 22:16 selected
drwxr-xr-x  2 hyobinjeong hyobinjeong 4.0K Oct  6 01:12 snv_calls
drwxr-xr-x  2 hyobinjeong hyobinjeong  138 Oct  6 01:42 stats
drwxr-xr-x 26 hyobinjeong hyobinjeong 4.0K Oct  6 01:26 strandphaser


(base) [hyobinjeong@node01 Project_mosaicatcher_RPE_mixture_Tutorial]$ ls -lh data/RPE_mixture/scNOVA_bam_merge
total 4.2G
drwxr-xr-x 2 hyobinjeong hyobinjeong  12K Oct  5 23:45 clone1
-rw-r--r-- 1 hyobinjeong hyobinjeong 3.3G Oct  5 23:49 clone1.merge.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 2.3M Oct  5 23:49 clone1.merge.bam.bai
drwxr-xr-x 2 hyobinjeong hyobinjeong 4.0K Oct  5 23:47 clone2
-rw-r--r-- 1 hyobinjeong hyobinjeong 889M Oct  5 23:48 clone2.merge.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 2.1M Oct  5 23:48 clone2.merge.bam.bai


(base) [hyobinjeong@node01 Project_mosaicatcher_RPE_mixture_Tutorial]$ ls -lh data/RPE_mixture/scNOVA_nucleosomes_bam/nucleosome_sampleA
total 1.1G
-rw-r--r-- 1 hyobinjeong hyobinjeong 1.1G Oct  6 01:34 result.H1.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 1.9M Oct  6 01:34 result.H1.bam.bai
(base) [hyobinjeong@node01 Project_mosaicatcher_RPE_mixture_Tutorial]$ ls -lh data/RPE_mixture/scNOVA_nucleosomes_bam/nucleosome_sampleB
total 1.1G
-rw-r--r-- 1 hyobinjeong hyobinjeong 1.1G Oct  6 01:34 result.H2.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 2.0M Oct  6 01:34 result.H2.bam.bai
```

# Output – Processed single-cell bam files

data/RPE_mixture/scNOVA_bam_modified

```
-rw-r--r-- 1 hyobinjeong hyobinjeong  75M Oct  5 22:10 RPE1WTPE20493.sc_pre_mono.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 1.7K Oct  5 22:15 RPE1WTPE20493.sc_pre_mono.metrix_dup.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong  75M Oct  5 22:12 RPE1WTPE20493.sc_pre_mono_sort_for_mark.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  34M Oct  5 22:15 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 4.3M Oct  5 23:14 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam.bai
-rw-r--r-- 1 hyobinjeong hyobinjeong  11M Oct  5 22:23 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam.C1.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  11M Oct  5 22:23 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam.C2.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  21M Oct  5 22:23 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam.C.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 3.4M Oct  5 22:25 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam.C.bam.bai
-rw-r--r-- 1 hyobinjeong hyobinjeong 6.8M Oct  5 22:23 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam.W1.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 6.8M Oct  5 22:23 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam.W2.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  13M Oct  5 22:23 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam.W.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 2.8M Oct  5 22:25 RPE1WTPE20493.sc_pre_mono_sort_for_mark_uniq.bam.W.bam.bai
-rw-r--r-- 1 hyobinjeong hyobinjeong 7.2K Oct  5 22:10 RPE1WTPE20494.header_test.sam
-rw-r--r-- 1 hyobinjeong hyobinjeong 8.1K Oct  5 22:36 RPE1WTPE20494.header_WC.sam
-rw-r--r-- 1 hyobinjeong hyobinjeong 134M Oct  5 22:10 RPE1WTPE20494.sc_pre_mono.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 1.7K Oct  5 22:15 RPE1WTPE20494.sc_pre_mono.metrix_dup.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 134M Oct  5 22:13 RPE1WTPE20494.sc_pre_mono_sort_for_mark.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  57M Oct  5 22:15 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 4.6M Oct  5 22:15 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam.bai
-rw-r--r-- 1 hyobinjeong hyobinjeong  16M Oct  5 22:36 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam.C1.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  16M Oct  5 22:36 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam.C2.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  29M Oct  5 22:43 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam.C.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 3.7M Oct  5 22:43 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam.C.bam.bai
-rw-r--r-- 1 hyobinjeong hyobinjeong  15M Oct  5 22:36 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam.W1.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  15M Oct  5 22:36 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam.W2.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  27M Oct  5 22:43 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam.W.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 3.7M Oct  5 22:43 RPE1WTPE20494.sc_pre_mono_sort_for_mark_uniq.bam.W.bam.bai
-rw-r--r-- 1 hyobinjeong hyobinjeong 7.2K Oct  5 22:10 RPE1WTPE20495.header_test.sam
-rw-r--r-- 1 hyobinjeong hyobinjeong 8.1K Oct  5 22:16 RPE1WTPE20495.header_WC.sam
-rw-r--r-- 1 hyobinjeong hyobinjeong  82M Oct  5 22:10 RPE1WTPE20495.sc_pre_mono.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 1.7K Oct  5 22:15 RPE1WTPE20495.sc_pre_mono.metrix_dup.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong  82M Oct  5 22:11 RPE1WTPE20495.sc_pre_mono_sort_for_mark.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  35M Oct  5 22:15 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 4.4M Oct  5 23:35 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam.bai
-rw-r--r-- 1 hyobinjeong hyobinjeong 8.3M Oct  5 22:16 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam.C1.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 8.2M Oct  5 22:16 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam.C2.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  16M Oct  5 22:16 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam.C.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 3.1M Oct  5 22:22 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam.C.bam.bai
-rw-r--r-- 1 hyobinjeong hyobinjeong  11M Oct  5 22:16 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam.W1.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  11M Oct  5 22:16 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam.W2.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong  20M Oct  5 22:16 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam.W.bam
-rw-r--r-- 1 hyobinjeong hyobinjeong 3.4M Oct  5 22:22 RPE1WTPE20495.sc_pre_mono_sort_for_mark_uniq.bam.W.bam.bai
(base) [hyobinjeong@node01 Project_mosaicatcher_RPE_mixture_Tutorial]$ ls -lh data/RPE_mixture/scNOVA_bam_modified
```

# Output – Result from scNOVA pipeline

data/RPE_mixture/scNOVA_result

```
(base) [hyobinjeong@node01 Project_mosaicatcher_RPE_mixture_Tutorial]$ ls -lh data/RPE_mixture/scNOVA_result
total 4.3G
drwxr-xr-x 2 hyobinjeong hyobinjeong   94 Oct  5 23:49 count_reads_chr_length
drwxr-xr-x 2 hyobinjeong hyobinjeong 4.0K Oct  5 23:48 count_reads_CREs
drwxr-xr-x 2 hyobinjeong hyobinjeong  100 Oct  5 23:50 count_reads_for_DNN
drwxr-xr-x 2 hyobinjeong hyobinjeong 8.0K Oct  5 23:49 count_reads_for_DNN_sc
drwxr-xr-x 2 hyobinjeong hyobinjeong 4.0K Oct  5 23:50 count_reads_split
-rw-r--r-- 1 hyobinjeong hyobinjeong  880 Oct  5 23:50 Deeptool_chr_length_RPE_mixture.tab
-rw-r--r-- 1 hyobinjeong hyobinjeong 649M Oct  6 00:00 Deeptool_Genes_for_CNN_RPE_mixture_sc_sort_lab_final.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 649M Oct  6 00:00 Deeptool_Genes_for_CNN_RPE_mixture_sc_sort_lab.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 628M Oct  5 23:59 Deeptool_Genes_for_CNN_RPE_mixture_sc_sort.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 628M Oct  5 23:59 Deeptool_Genes_for_CNN_RPE_mixture_sc.tab
-rw-r--r-- 1 hyobinjeong hyobinjeong 102M Oct  6 01:08 Deeptool_Genes_for_CNN_RPE_mixture_sort_lab.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong  80M Oct  6 01:08 Deeptool_Genes_for_CNN_RPE_mixture_sort.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong  80M Oct  6 01:02 Deeptool_Genes_for_CNN_RPE_mixture.tab
-rw-r--r-- 1 hyobinjeong hyobinjeong  39K Oct  6 02:08 Expression_all_clone1.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong  39K Oct  6 02:08 Expression_all_clone2.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 167M Oct  6 02:08 Features_reshape_all_orientation_norm_var_GC_CpG_RT_T_comb3_clone1.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 154M Oct  6 02:08 Features_reshape_all_orientation_norm_var_GC_CpG_RT_T_comb3_clone2.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 838K Oct  6 02:08 Features_reshape_all_TSS_matrix_woM_all_RT_clone1.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 838K Oct  6 02:08 Features_reshape_all_TSS_matrix_woM_all_RT_clone2.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 5.7M Oct  6 02:04 Features_reshape_clone1_orientation_CN_correct0.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong  53M Oct  6 02:07 Features_reshape_clone1_orientation_norm.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong   19 Oct  6 02:06 Features_reshape_clone1_Resid_orientation_IQR.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong  46M Oct  6 02:08 Features_reshape_clone1_Resid_orientation.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 5.7M Oct  6 02:04 Features_reshape_clone2_orientation_CN_correct0.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong  53M Oct  6 02:06 Features_reshape_clone2_orientation_norm.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong   19 Oct  6 02:06 Features_reshape_clone2_Resid_orientation_IQR.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong  32M Oct  6 02:07 Features_reshape_clone2_Resid_orientation.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 118M Oct  6 02:07 Features_reshape_RPE_mixture_clone1_orientation_norm_qc.pdf
-rw-r--r-- 1 hyobinjeong hyobinjeong  86M Oct  6 02:08 Features_reshape_RPE_mixture_clone1_Resid_orientation_qc.pdf
-rw-r--r-- 1 hyobinjeong hyobinjeong 107M Oct  6 02:06 Features_reshape_RPE_mixture_clone2_orientation_norm_qc.pdf
-rw-r--r-- 1 hyobinjeong hyobinjeong  33M Oct  6 02:07 Features_reshape_RPE_mixture_clone2_Resid_orientation_qc.pdf
-rw-r--r-- 1 hyobinjeong hyobinjeong 2.1M Oct  6 03:02 result_PLSDA_RPE_mixture.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 2.9M Oct  6 02:13 Result_scNOVA_infer_expression_table.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 165M Oct  5 23:51 RPE_mixture_CREs_2kb_sort_num_sort_for_chromVAR.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 165M Oct  5 23:51 RPE_mixture_CREs_2kb_sort_num.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 170M Oct  5 23:50 RPE_mixture_CREs_2kb_sort.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 170M Oct  5 23:50 RPE_mixture_CREs_2kb.tab
-rw-r--r-- 1 hyobinjeong hyobinjeong 5.1M Oct  5 23:51 RPE_mixture_sort.txt
-rw-r--r-- 1 hyobinjeong hyobinjeong 5.1M Oct  5 23:51 RPE_mixture.tab
```

# Output – Single-cell Nucleosome count matrix

Single-cell Nucleosome count matrix (gene X cell)

data/RPE_mixture/scNOVA_result/RPE_mixture_sort.txt



# Output – Single-cell Nucleosome count matrix

Single-cell Nucleosome count matrix (CRE X cell)

data/RPE_mixture/scNOVA_result/RPE_mixture_CREs_2kb_sort.txt

# Output – altered gene activity between clones (DESeq2)

data/RPE_mixture/scNOVA_result/Result_scNOVA_infer_expression_table.txt



# Output – altered gene activity between clones (DESeq2) plots

**Result of infer altered gene activity between clones (DESeq2) - plots**

data/RPE_mixture/scNOVA_result_plots/Result_scNOVA_plots_RPE_mixture.pdf



Clone1
(RPE1_WT)
Clone2 (RPE1_BM510)

# Output – altered gene activity between clones (PLSDA)

data/RPE_mixture/scNOVA_result/result_PLSDA_RPE_mixture.txt



# Output – altered gene activity between clones (PLSDA) plots

data/RPE_mixture/scNOVA_result_plots/Result_scNOVA_plots_RPE_mixture_alternative_PLSDA.pdf



Clone1 (RPE1_WT)
Clone2 (RPE1_BM510)

# Part4. scRNA-seq에서 서브클론을 유추하고 기능적으로 분석하는 멀티오믹스 기법 소개

Single-cell multi-omics analysis to
study tumor subclones

---

## Recent strategy to study genome and functional readout from single-cell RNA-seq

Patient tumor

*Infer CNV from single-cell RNA-seq*



*Neftel et al. Cell, 2019*

# Recent strategy to study genome and functional readout from single-cell RNA-seq

## SCNA inference methods based on transcriptome



| | | Method detail | | | |
|---|---|---|---|---|---|
| | Method | SV class | Require pre-defined SV breakpoint | Size resolution in the paper | Chr6 SV detection |
| Discovery | InferCNV (Science, 2014) | CNV only | N | entire chromosomes or large segments of chromosomes | N |
| Discovery | HoneyBADGER (Genome Res, 2018) | CNV only | N | 10Mb | N |
| Discovery | CONICSmat 'discovery mode' (Bioinformatics, 2018) | CNV only | N | 100 expressed genes (by default) | N |
| Genotyping | CONICSmat 'genotype mode' (Bioinformatics, 2018) User provide candidate SCNA | CNV only | Y | 100 expressed genes (by default) | Y |

---

# Recent strategy to study genome and functional readout from single-cell RNA-seq (InferCNV)

## InferCNV: Inferring copy number alterations from tumor single cell RNA-Seq data



InferCNV is used to explore tumor single cell RNA-Seq data to identify evidence for somatic large-scale chromosomal copy number alterations, such as gains or deletions of entire chromosomes or large segments of chromosomes. This is done by exploring expression intensity of genes across positions of tumor genome in comparison to a set of reference 'normal' cells. A heatmap is generated illustrating the relative expression intensities across each chromosome, and it often becomes readily apparent as to which regions of the tumor genome are over-abundant or less-abundant as compared to that of normal cells.

InferCNV provides access to several residual expression filters to explore minimizing noise and further revealing the signal supporting CNA. Additionally, inferCNV includes methods to predict CNA regions and define cell clusters according to patterns of heterogeneity.

InferCNV is one component of the TrinityCTAT toolkit focused on leveraging the use of RNA-Seq to better understand cancer transcriptomes. To find out more about Trinity CTAT please visit TrinityCTAT.

https://github.com/broadinstitute/inferCNV/wiki

# Recent strategy to study genome and functional readout from single-cell RNA-seq (HoneyBADGER)

# Recent strategy to study genome and functional readout from single-cell RNA-seq (NumBat)

# Recent strategy to study genome and functional readout from single-cell RNA-seq (CONICS)



## CONICS

*CONICS*: *CO*py-*N*umber analysis *I*n single-*C*ell RNA-*S*equencing

CONICS works with either full transcript (e.g. Fluidigm C1) or 5'/3' tagged (e.g. 10X Genomics) data!

The CONICS paper has been accepted for publication in Bioinformatics. Check it out here !

### Table of contents

- CONICSmat - Identifying CNVs fro
- Identifying CNVs from scRNA-seq u
- Integrating the minor-allele frequen
- Phylogenetic tree contruction
- Intra-clone co-expression networks
- Assessing the correlation of CNV st
- False discovery rate estimation: Cro
- False discovery rate estimation: Em

https://github.com/diazlab/CONICS

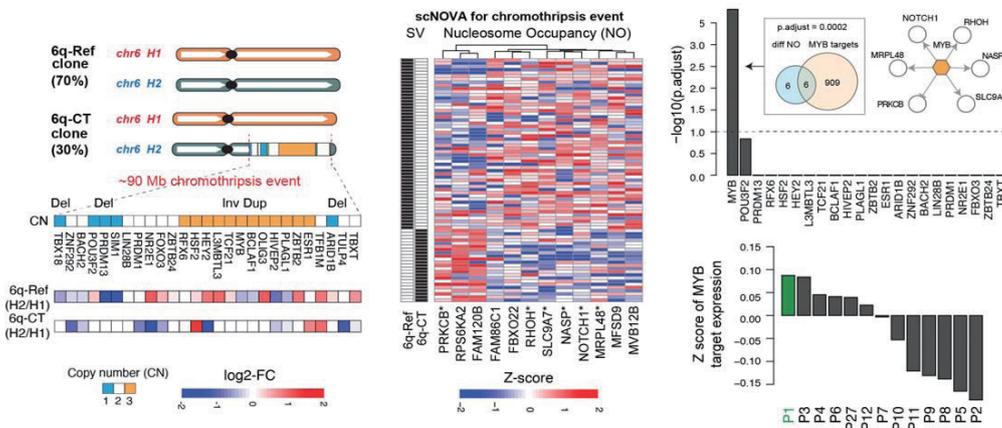### CONICSmat - Identifying CNVs from scRNA-seq using a count table

CONICSmat is an R package that can be used to identify CNVs in single cell RNA-seq data from a gene expression table, without the need of an explicit normal control dataset. CONICSmat works with either full transcript (e.g. Fluidigm C1) or 5'/3' tagged (e.g. 10X Genomics) data. A tutorial on how to use CONICSmat, and a Smart-Seq2 dataset, can be found on the CONICSmat Wiki page [CLICK here].

*Visualizations of scRNA-seq data from* Oligodendroglioma *(Tirosh et al., 2016) generated with CONICSmat.*

---

# Applying CNV inference of scRNA-seq to the T-ALL case study

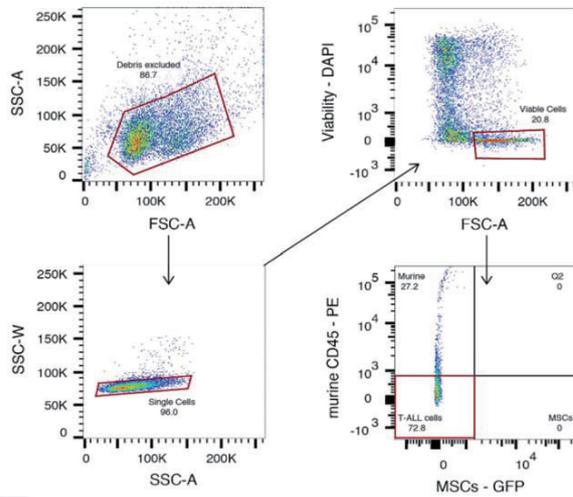Hypothesis : 6q-CT cells have MYB-Notch activation compared to 6-Ref cells



*Single-cell experiment is needed to confirm subclonal level transcriptome changes*
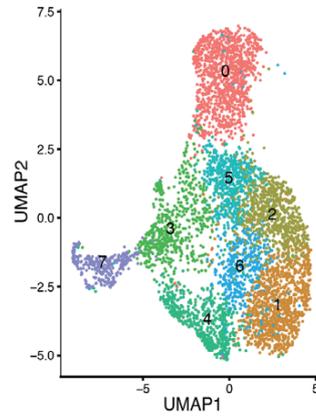
# Applying CONICSmat to the T-ALL case study

Gating strategy for single, viable T-ALL cell isolation from T-ALL sample T-ALL_P1 for scRNA-seq.
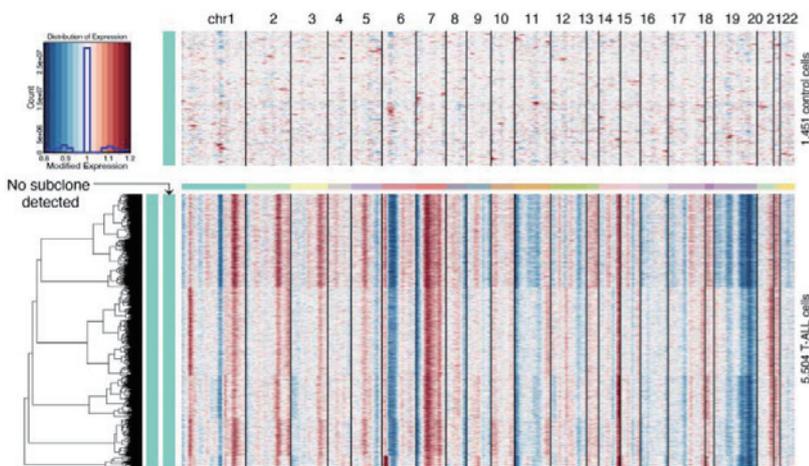


**Unsupervised analysis of transcriptome**

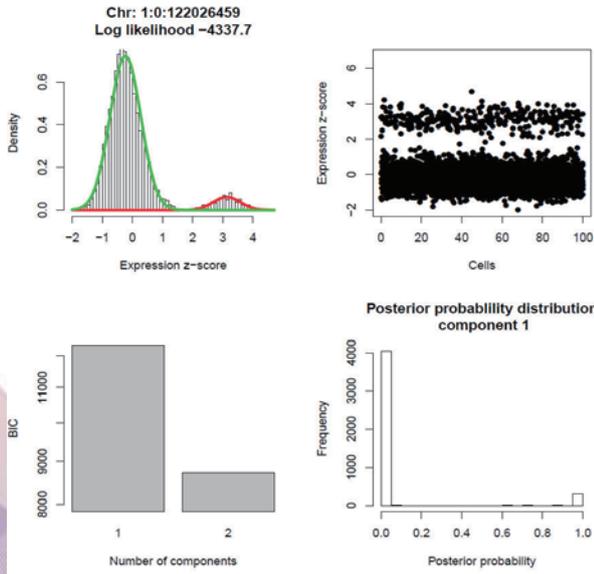# Applying InferCNV to the T-ALL case study



InferCNV analysis of 5,504 high quality T-ALL_P1 cells, and 1,451 control cells. Control cells were downloaded from PBMC data provided by 10X Genomics. This analysis did not discover subclones in 5,504 T-ALL cells.
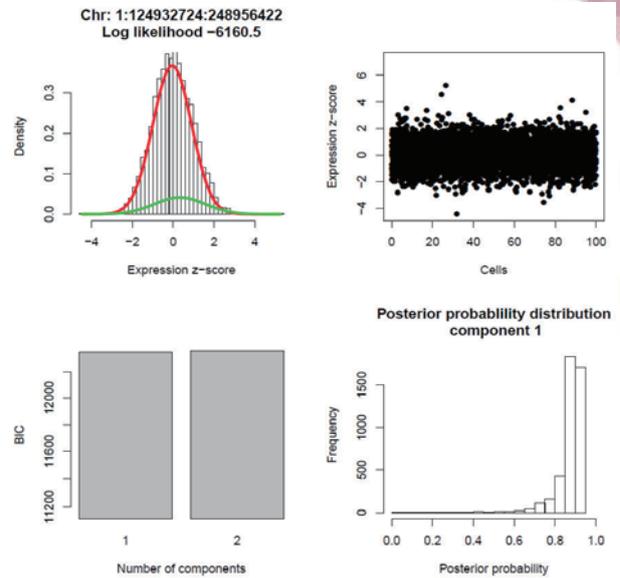
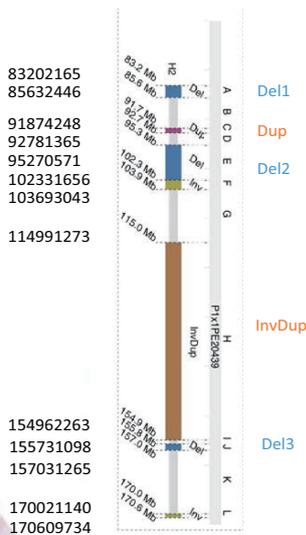Applying CONICSmat to the T-ALL scRNA-seq (Genotyping mode)
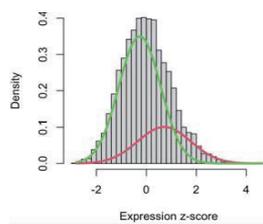
Presence of Subclonal CNA | Absence of Subclonal CNA

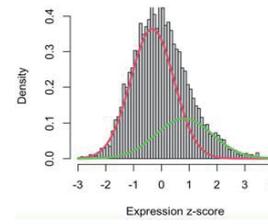CONICSmat analysis supports the presence of chr6 deletions and duplications

Sanders et al. 2020

10X transcriptome experiment from Karen Grimes

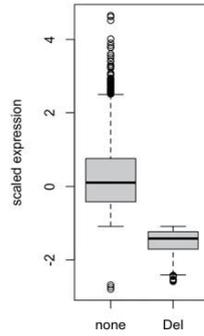| Genomic range | BIC 1 component | BIC 2 components | BIC difference | LRT adj. p-val | CNV call (CF%) |
|---|---|---|---|---|---|
| chr6_Del 35 genes | 15635.9018 | 15553.1101 | 82.7917307 | 0 | 729 cells (13.2%) |
| chr6_Dup 192 genes | 15635.9018 | 15481.839 | 154.062863 | 0 | 265 cells (4.8%) |

# Cluster3 and Cluster7 cells are highly enriched with deletion calls

**Deletion call (del1 + del2 + del3) probability>0.9**



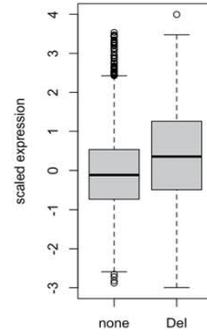| Row Labels | Del(0.9) | %Del(all) | p-value | adjustedP |
|---|---|---|---|---|
| Cluster0 | 166 | 14.901 | 0.039 | 0.122 |
| Cluster1 | 97 | 9.454 | 1.000 | 1.000 |
| Cluster2 | 118 | 15.013 | 0.066 | 0.131 |
| Cluster3 | 118 | 19.440 | 0.000 | 0.000 |
| Cluster4 | 40 | 6.981 | 1.000 | 1.000 |
| Cluster5 | 83 | 15.690 | 0.049 | 0.122 |
| Cluster6 | 32 | 6.695 | 1.000 | 1.000 |
| Cluster7 | 62 | 20.395 | 0.000 | 0.001 |
| Cluster8 | 6 | 10.526 | 0.785 | 1.000 |
| Cluster9 | 7 | 23.333 | 0.092 | 0.154 |
| Grand Total | 729 | 13.245 | - | - |

**Deleted region**
Ranksum test, t.test
p-value < 2.2e-16

**Duplicated region**
Ranksum test, t.test
p-value < 2.2e-16


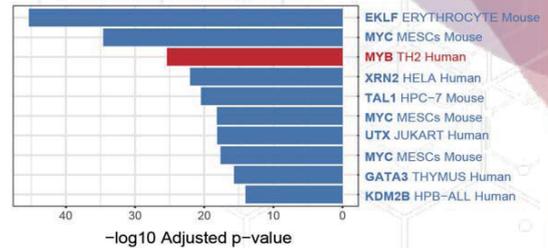
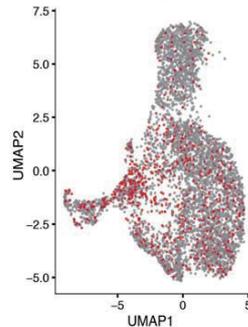| Type | Dup | none |
|---|---|---|
| Del | 96 | 633 |
| none | 169 | 4606 |

Fisher exact test
p-value < 2.2e-16

---

# SV subclone in P1 shows increase of MYB target expression and cells with premature stages in the cellular hierarchy



**Unsupervised analysis of transcriptome**

**SCNA inference from transcriptome**

**Integration & interpretation**

**TF-target Enrichment (Cluster 3)**

- EKLF ERYTHROCYTE Mouse
- MYC MESCs Mouse
- MYB TH2 Human
- XRN2 HELA Human
- TAL1 HPC-7 Mouse
- MYC MESCs Mouse
- UTX JUKART Human
- MYC MESCs Mouse
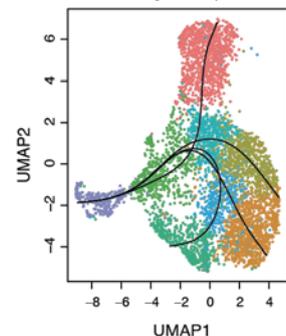- GATA3 THYMUS Human
- KDM2B HPB-ALL Human

−log10 Adjusted p−value

**Signature gene activity**

**Lineage analysis**

## Summary

- 암에서 이러한 서브클론들을 동정하기 위해 어떤 싱글셀 오믹스 기법들이 개발되어 있을까? → *single-cell WGS, SDR-seq, Strand-seq, etc*

- 이러한 싱글셀 오믹스 데이터를 분석하기 위해 어떤 생명 정보학적인 도구들을 사용할 수 있을까? → *MosaiCatcher for Strand-seq analysis*

- 서브클론의 동정 뿐 아니라 그 기능적 특성을 파악하기 위해서는 유전체와 전사체 또는 후성유전체 데이터를 함께 분석하는 싱글셀 멀티 오믹스 분석이 필요하다. 이를 구현하기 위한 생명 정보학적인 방법에는 어떤 것들이 있을까?

  → *scNOVA for Strand-seq analysis*

  → *Infer copy number alteration from scRNA-seq (InferCNV, HoneyBADGER, Numbat, CONICS etc.)*

---

SBi 한국생명정보학회
Korean Society for Bioinformatics

# 감사합니다.