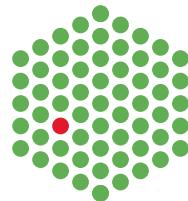
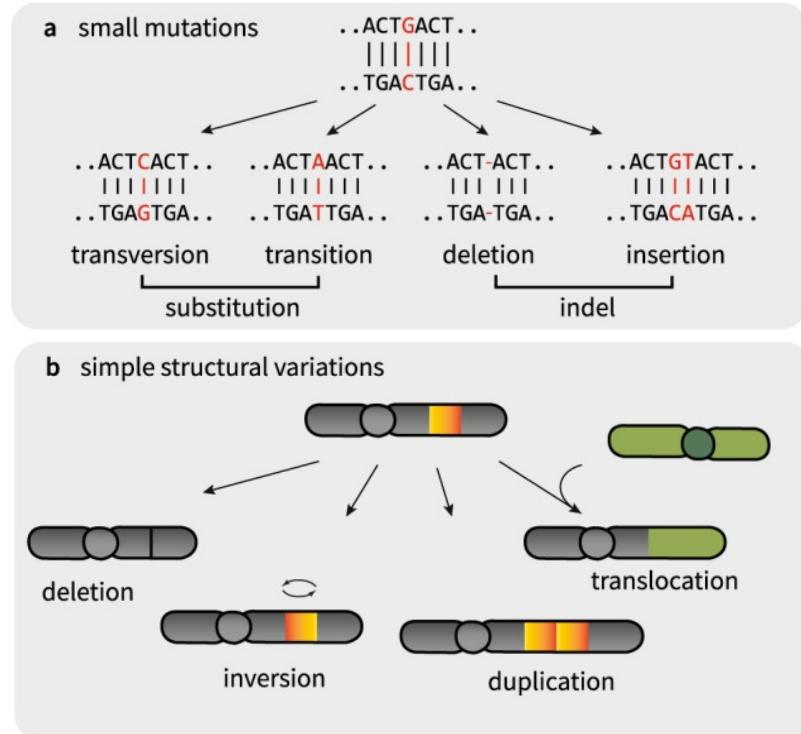


# Detection of structural variations using bulk whole-genome sequencing and Delly



# Structural variants are diverse, and occur commonly in the human germline

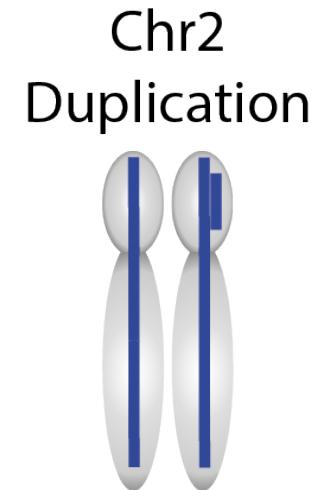
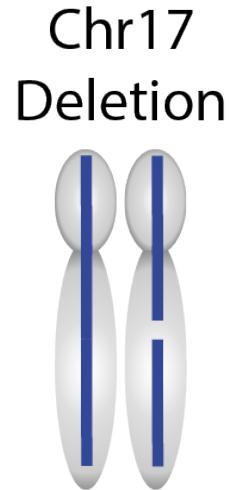
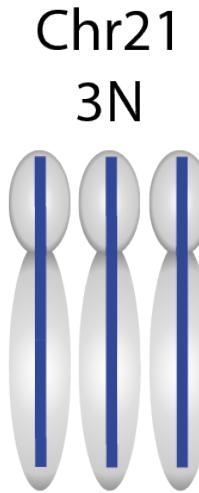


Heritable sequence differences between individuals

0.1%

0.5-1%

# Germline structural variation has been linked to various syndromes



Down Syndrome

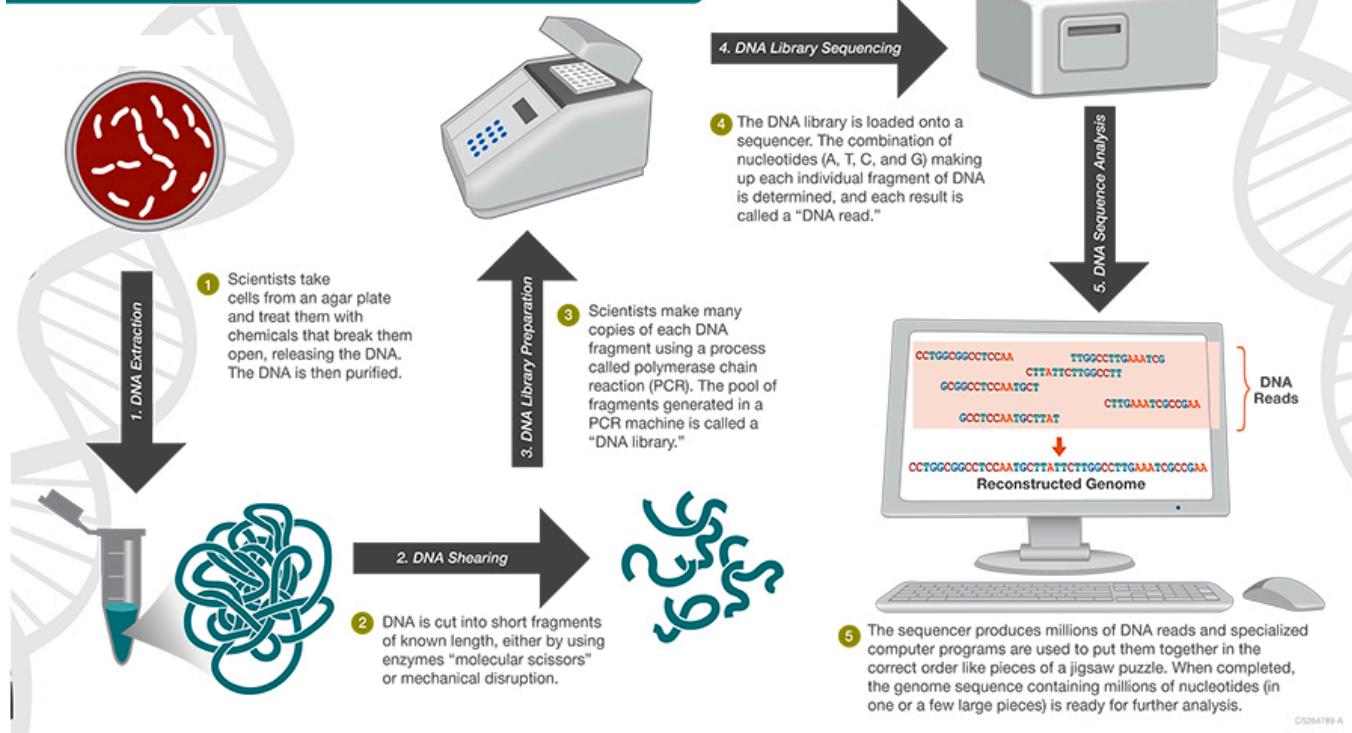
Smith-Magenis Syndrome

Lynch Syndrome

# How do we get the genome sequence information?

## The Whole Genome Sequencing (WGS) Process

WGS is a laboratory procedure that determines the order of bases in the genome of an organism in one process. WGS provides a very precise DNA fingerprint that can help link cases to one another allowing an outbreak to be detected and solved sooner.



# Fastq files are the starting point of genomics data analysis



Paired-end

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

ATAC TTT

AAAG TAT

CTGT AAA

TTT AGAG

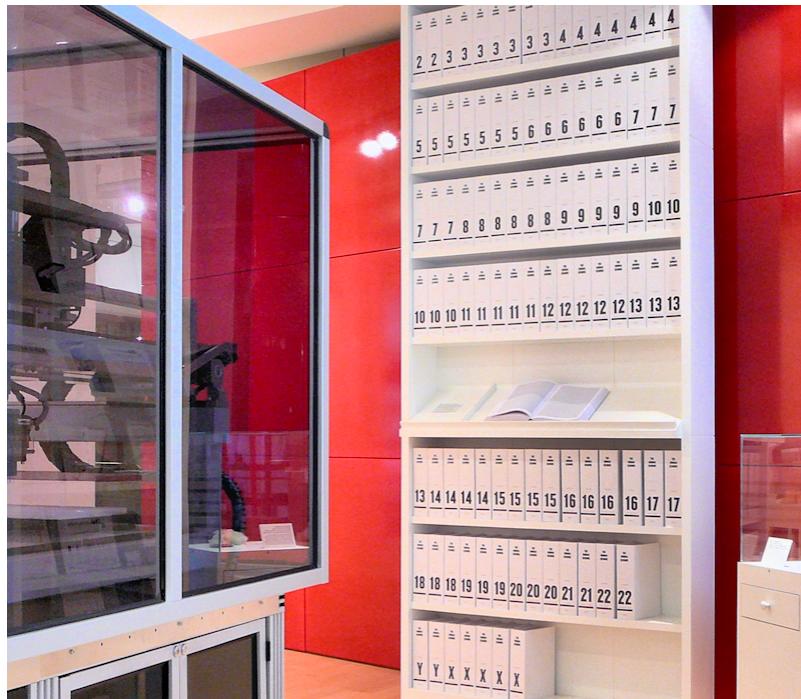
Read 1



1

Read 2

# Searching for the origin of DNA fragment using reference genome



Wellcome Collection, London

The first printout of the human  
reference genome



A T A C T T T

A A A G T A T

BWA



bam file

A T A C T T T . Chr3:1000

A A A G T A T . Chr11:35082

# Quality checking of bam file after the alignment

- Checkpoint1: How much portion of the reads were aligned to the reference genome?
- Checkpoint2: How much portion of the reads were from PCR duplicates?
- If the bam file pass the QC, we are ready to detect structural variations!



Alfred: BAM alignment statistics, feature counting and feature annotation

Alfred is available as a [Bioconda package](#), as a statically linked binary from the [GitHub release page](#) or as a minimal [Docker container](#). Please have a look at [Alfred's documentation](#) for any installation or usage questions.



*Tobias Rausch  
Senior Bioinformatician*

# Structural Variant Calling (germline & somatic)

BIOINFORMATICS

Vol. 28 ECCB 2012, pages i333–i339  
doi:10.1093/bioinformatics/bts378

## DELLY: structural variant discovery by integrated paired-end and split-read analysis

Tobias Rausch,<sup>1,3,\*</sup>, Thomas Zichner<sup>1</sup>, Andreas Schlattl<sup>1</sup>, Adrian M. Stütz<sup>1</sup>, Vladimir Benes<sup>3</sup> and Jan O. Korbel<sup>1,2</sup>

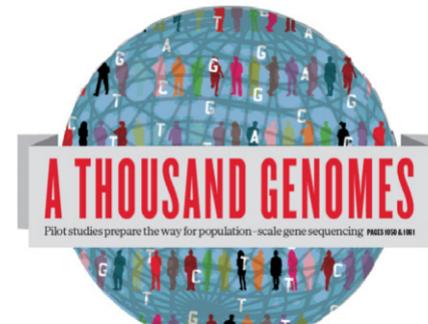
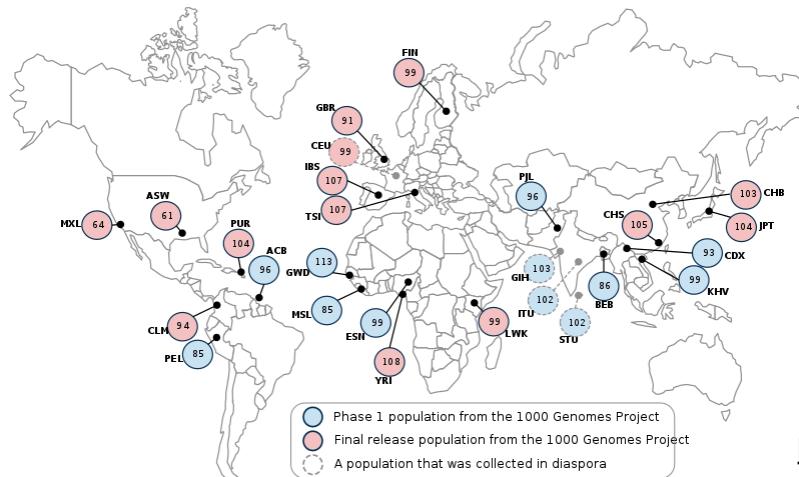
<sup>1</sup>European Molecular Biology Laboratory (EMBL), Genome Biology, Meyerhofstr. 1, 69117 Heidelberg, Germany and

<sup>2</sup>EMBL European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK

and <sup>3</sup>European Molecular Biology Laboratory (EMBL), Core Facilities and Services, Meyerhofstr. 1, 69117 Heidelberg, Germany



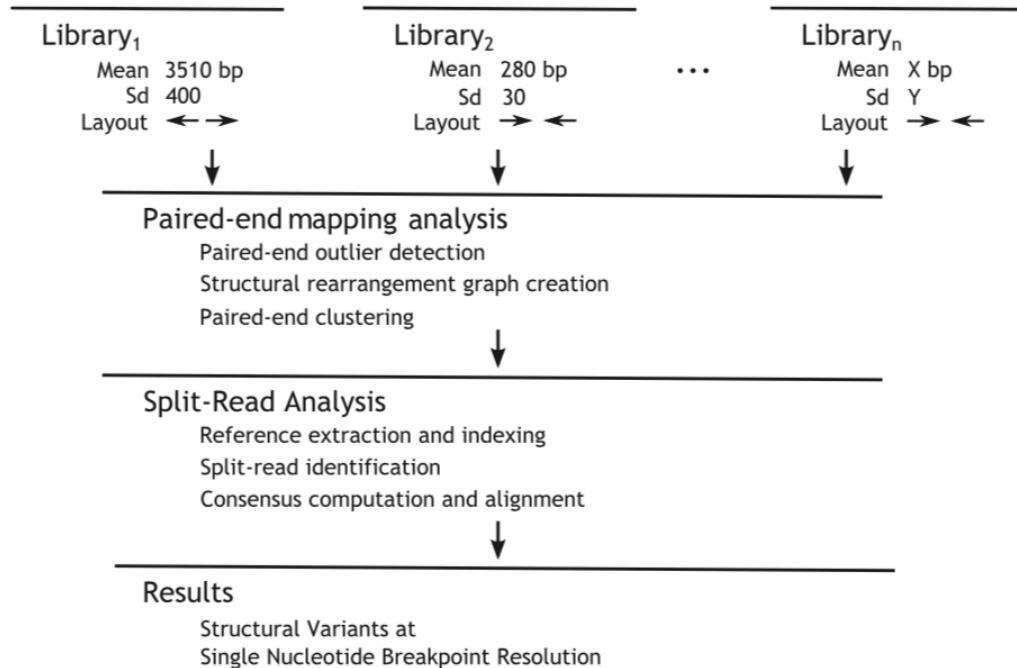
Rausch et al. 2012



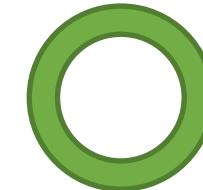
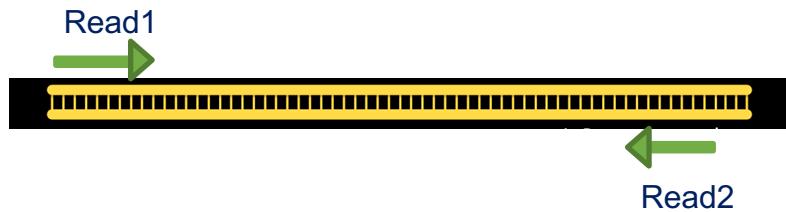
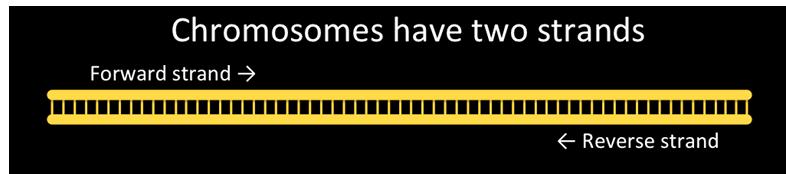
<https://www.nature.com/articles/nature15393>

# Two components of DELLY analysis

- Paired-end mapping analysis
- Split-read analysis



# Paired reads have different orientation by default



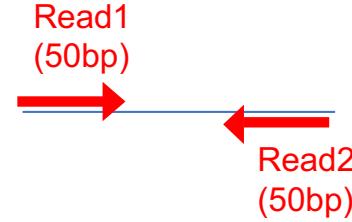
# Paired reads gives estimation of probable insert size of that library



Library 1  
Fragment size : 500bp  
Insert size : 400bp



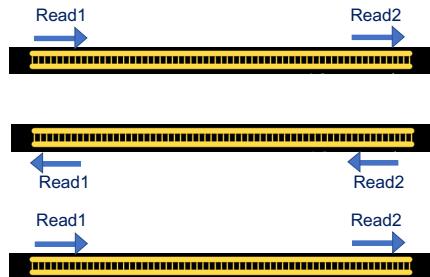
Library 2  
Fragment size : 350bp  
Insert size : 250bp



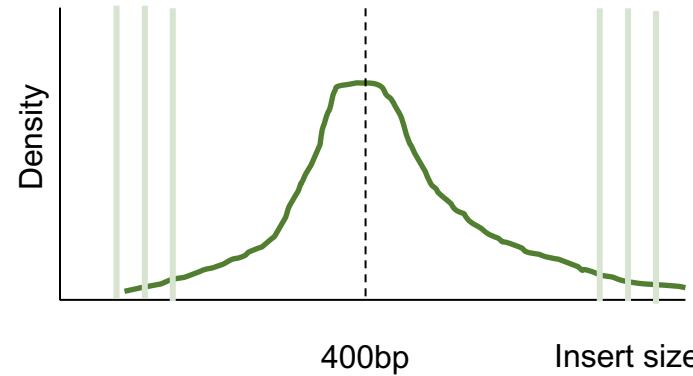
- Standard Illumina sequencing libraries tend to have a fragment size of 100-700bp
- Insert size is impacted by the process of cluster generation in which libraries are denatured, diluted and distributed on the two-dimensional surface of the flow cell and then amplified.

# Discordantly mapped read-pairs can be identified by insert size & orientation

Abnormal orientation

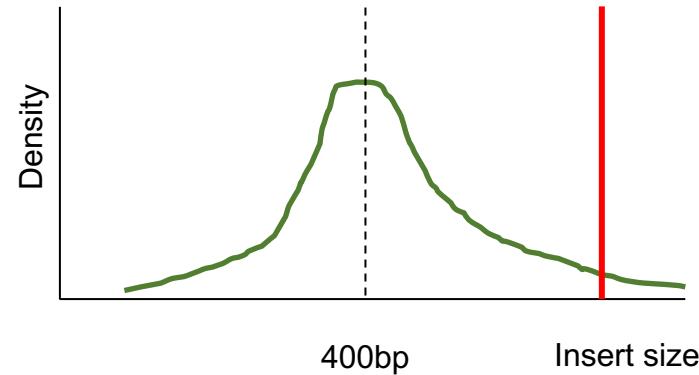
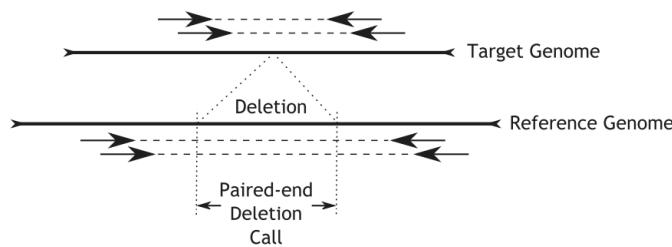


Abnormal insert size

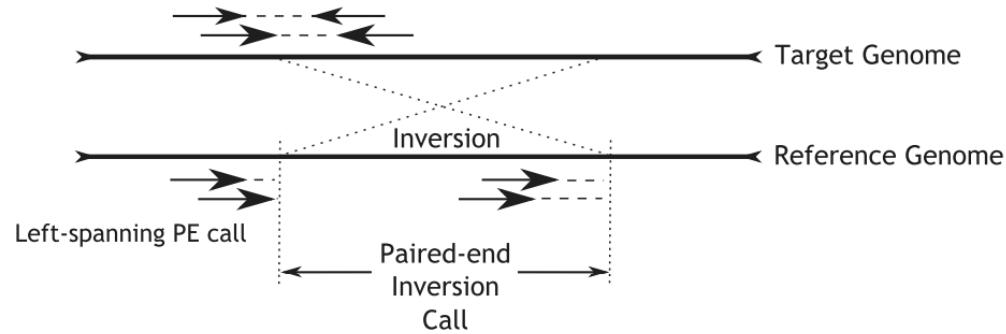


- Discordantly mapped read-pairs that either have an abnormal orientation or an insert size greater than the expected range.
- Default insert size cutoff is three standard deviations from the median insert size.

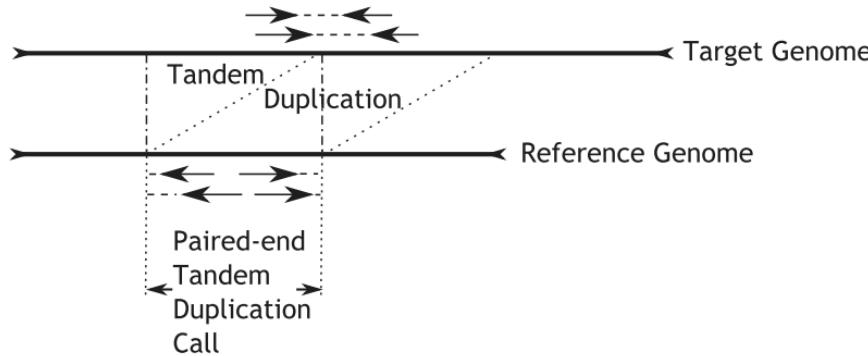
# Deletion: default orientation & extremely big insert size



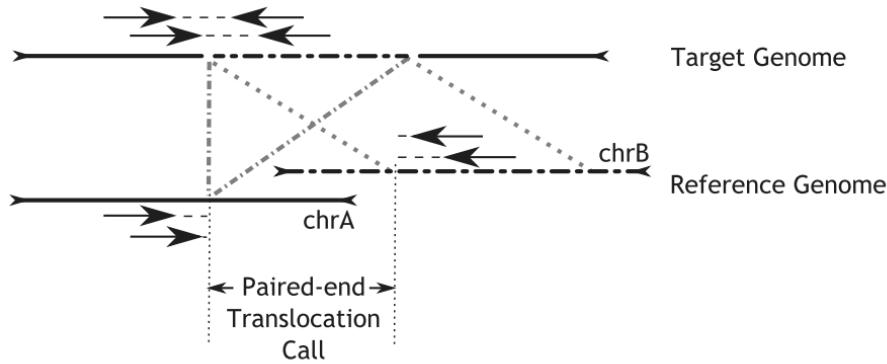
# Inversion: abnormal orientation



## Tandem duplications: first and second read changed their relative order

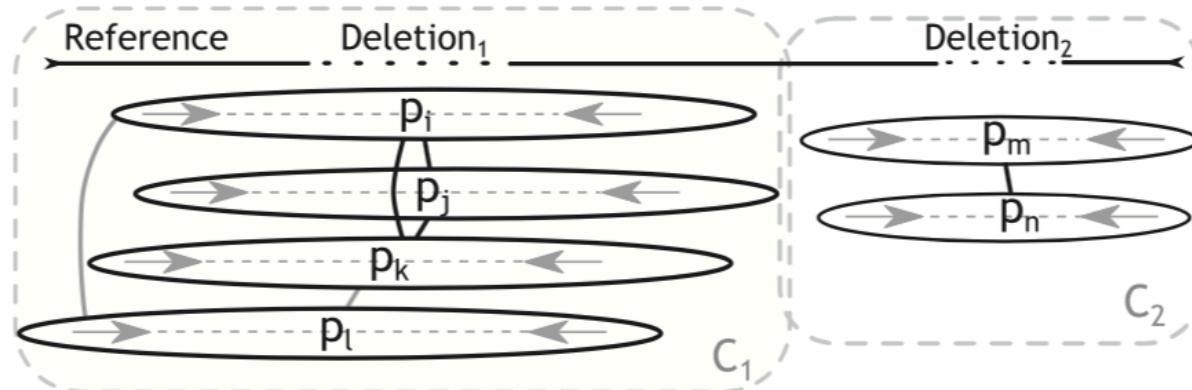


# Translocation: paired-ends mapping to different chromosomes



## Collect discordant pairs and make clusters

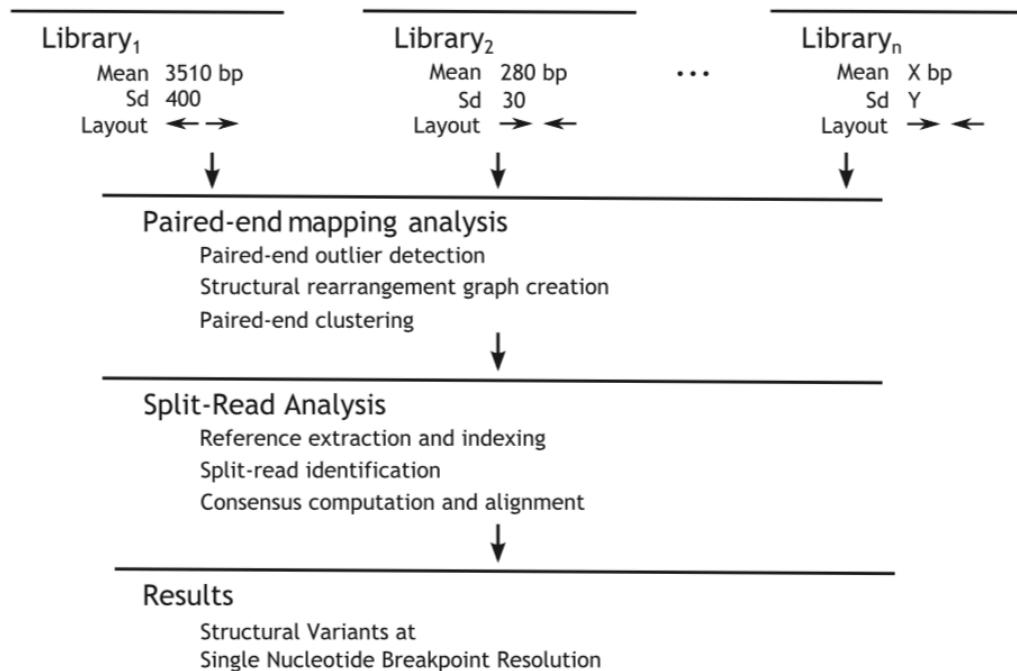
It gives breakpoint-containing genomic intervals!



**Now we need Split-read analysis to determine breakpoint!**

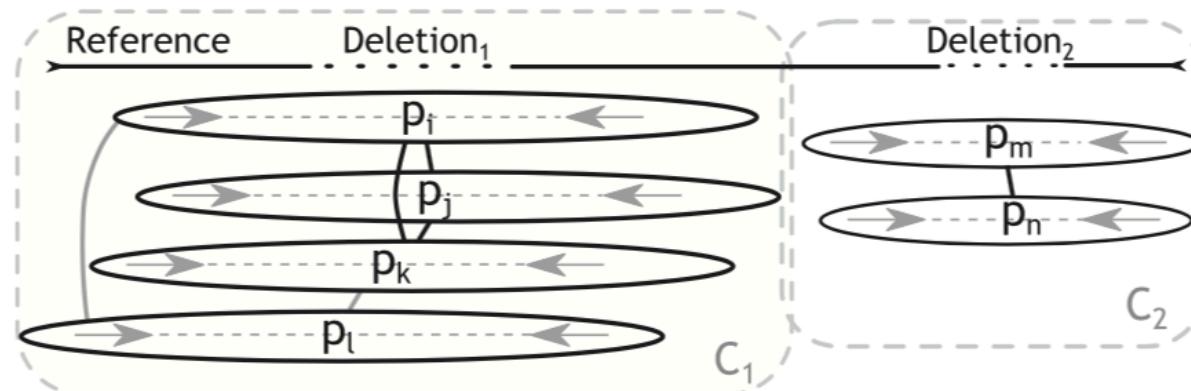
# Two components of DELLY analysis

- Paired-end mapping analysis
- Split-read analysis



# Searching for the single-anchored paired ends

A single-anchored paired-end is a read pair where one read maps to the reference and the other read is unmapped

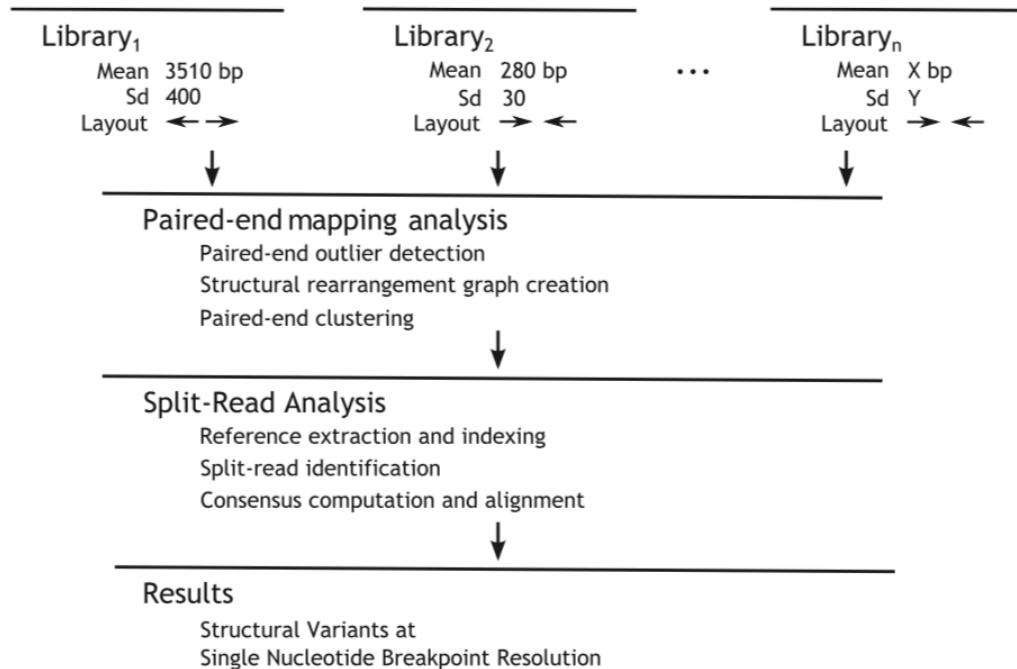


Breakpoint



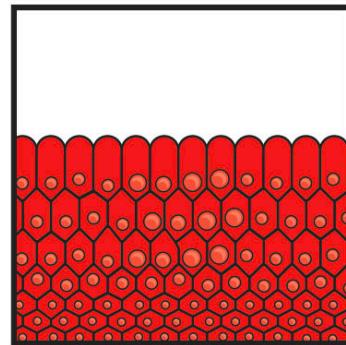
# Two components of DELLY analysis gives SV call

- Paired-end mapping analysis
- Split-read analysis

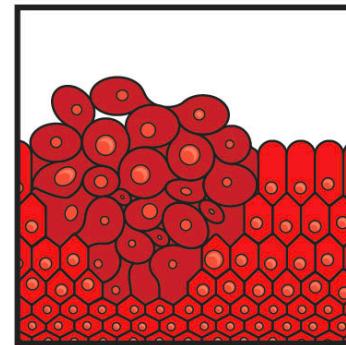


# Let's try Delly analysis to identify somatic SVs!

**Real world**



Normal cells



Cells forming a tumour

**Model system**

*RPE1-WT*



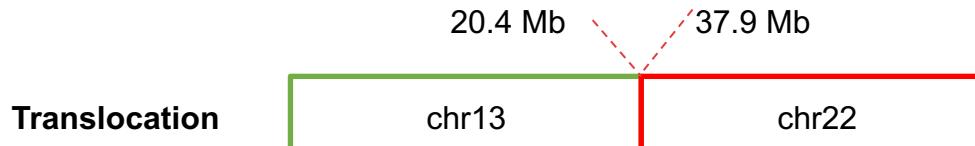
**vs**

*RPE1-BM510*



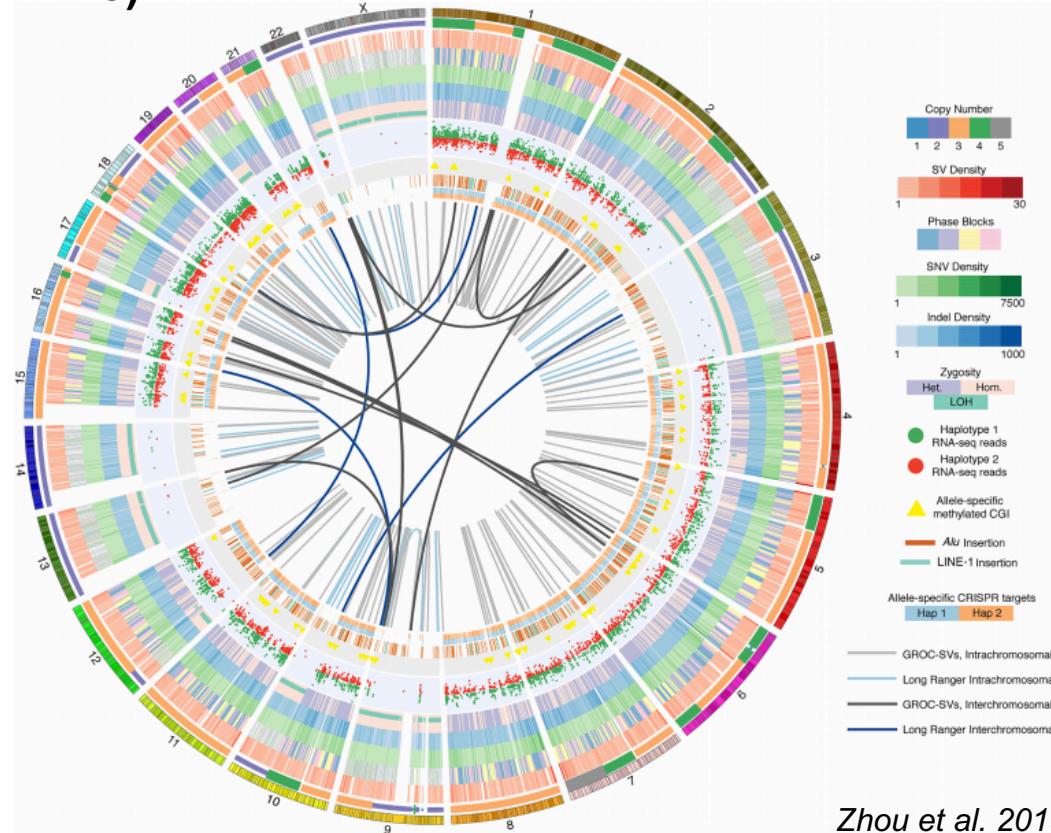
# Delly example clonal somatic translocation (Pure BM510 VAF100% vs WT)

- CHROM, POS : [chr22, 37934425], [chr13, 20412250]
- MAPQ : mapping quality [60]
- GT : genotype [0/1] [0/0]
- GL : log10-scaled genotype likelihoods for RR, RA, AA genotypes [-23.6655,0,-90.8635]  
[0,-7.5218,-88.7961]
- GQ : Genotype Quality [10000] [75]
- FT : Per-sample genotype filter [PASS] [PASS]
- RC : Raw high-quality read counts for the SV [0] [0]
- RCL : Raw high-quality read counts for the left control region [0] [0]
- RCR : Raw high-quality read counts for the right control region [0] [0]
- CN : Read-depth based copy-number estimate for autosomal sites [-1] [-1]
- DR : # high-quality reference pairs [25] [23]
- DV : # high-quality variant pairs [5] [0]
- RR : # high-quality reference junction reads [27] [25]
- RV : # high-quality variant junction reads [9] [0]
- PRECISE/IMPRECISE :SVs refined using split-reads [PRECISE]
- PASS/LowQual : genotype quality [PASS]



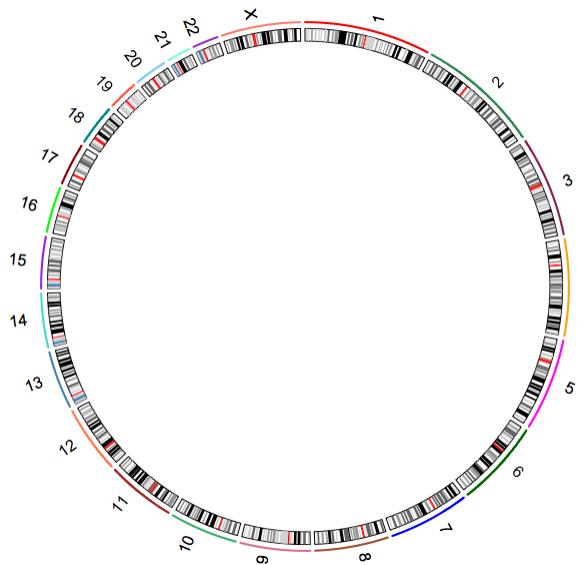
*Can we visualize it with more fancy way?*

# Variety of SV information can be included in the circos plot (K562 cell line)

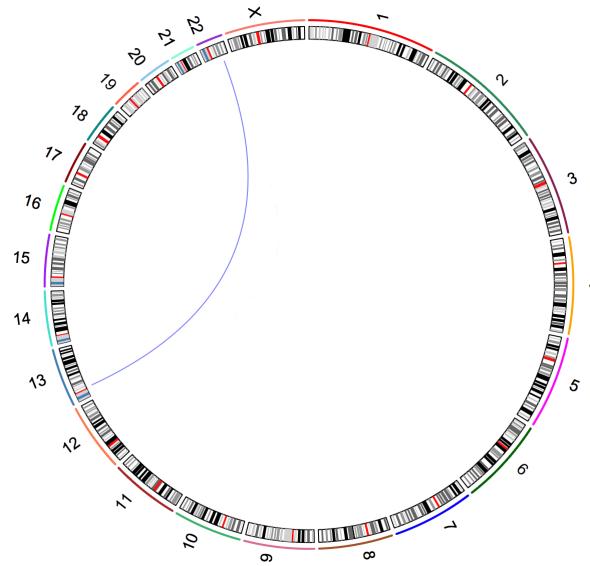


Zhou et al. 2019

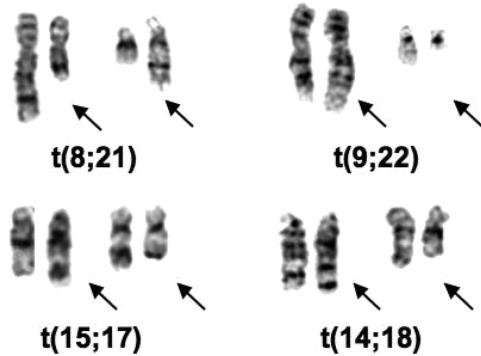
# Visualize SVs using circos plot



# Visualize SVs using circos plot

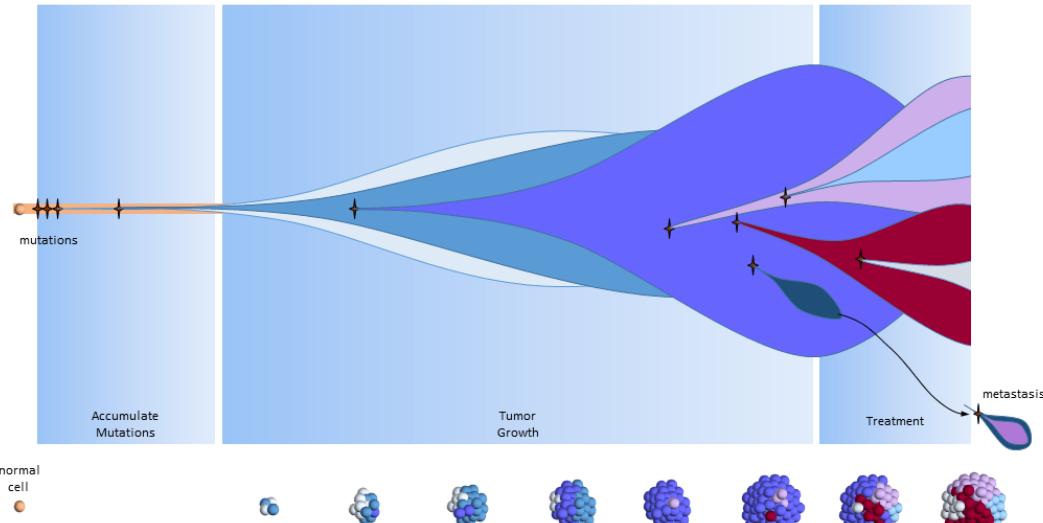


# But what about sub-clonal variation?



(Rowley, J., D., *Blood* 2008)

10-40% of certain AML subtypes  
are driven by translocations

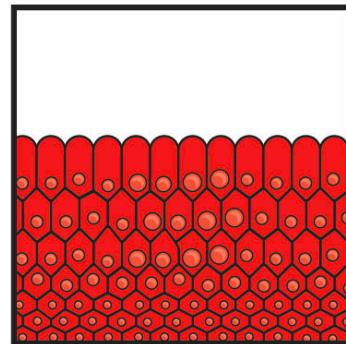


In theory, VAFs as low as 2.5% can be identified; assuming:

1. Known target locus/loci
2. Extremely high sequencing depth (20,000 - >100,000X)

Can we detect subclonal translocation using same method we practiced today?

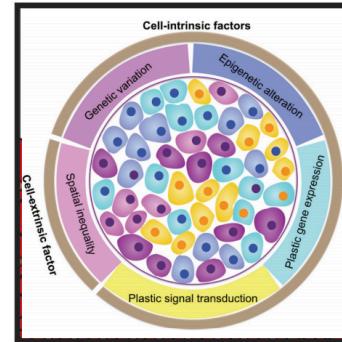
Real  
world



Normal cells

Model  
system

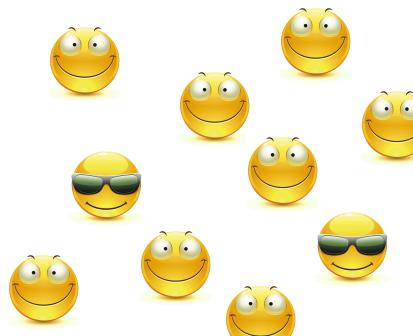
WT



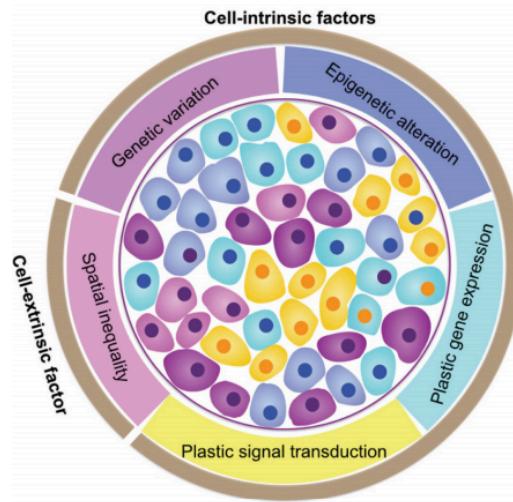
Cells forming a tumour

*VAF20% BM510*

vs



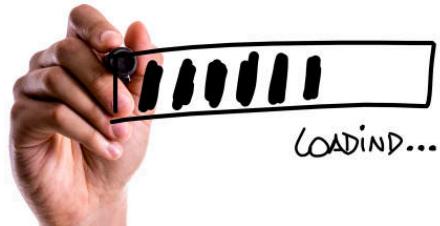
# What would be the best strategy to detect subclonal SVs?



- If certain variants have low VAF, it can be difficult to be confidently detected by bulk DNA sequencing approaches, what would be the reason?
- Question1: How can we more confidently detect subclonal SV event?
- Question2: How can we know which subset of cells carry SVs?

*Thank you for  
listening*

**COMING SOON**



**Day2  
Single-cell  
approach**

