# Structural Variants and Copy-Number Variants

**Tobias Rausch**
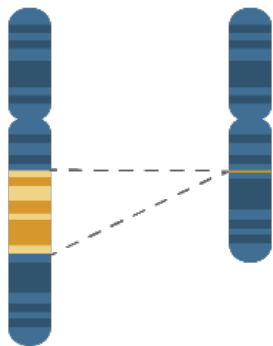
**European Molecular Biology Laboratory (EMBL)**

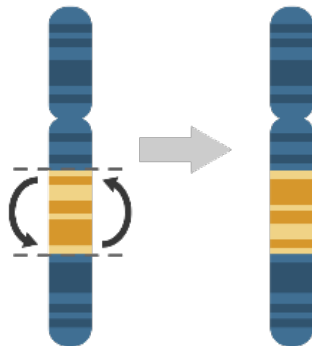25 June 2024

EMBL

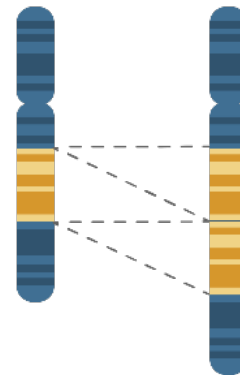# Somatic and Germline Structural Variants (SVs)

# Cancers harbor a wide Range of Chromosome Abberrations
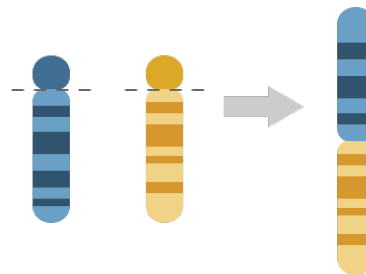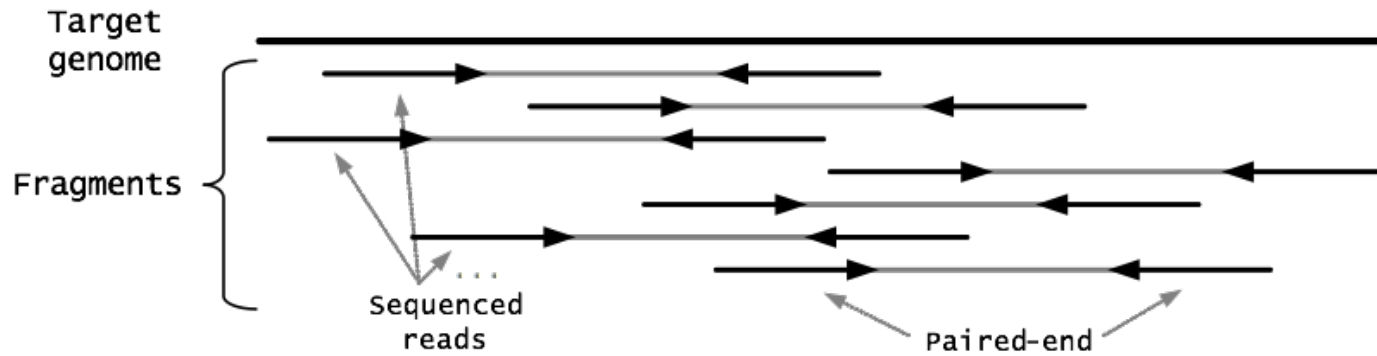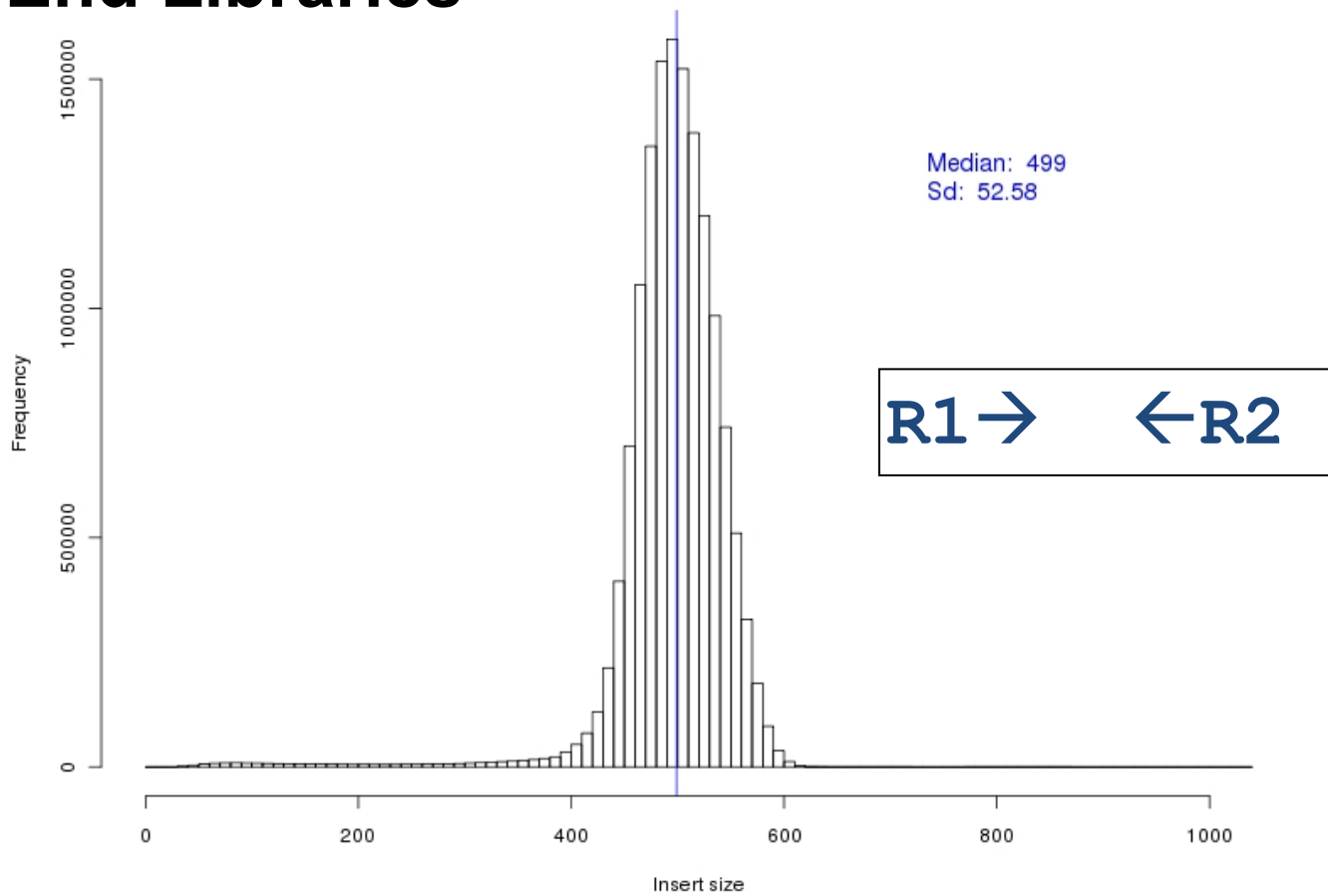
# Structural Variants (SVs)

# Paired-End Sequencing

# Paired-End Libraries



Median: 499
Sd: 52.58

R1→    ←R2

# SV Discovery Approaches

# Paired-end mapping

# Copy-Number Variants (CNVs)

# Human karyotype

Normal Karyotype

Cancer Karyotype (NSCLC cell line D117)



EMBL

# CNV detection technologies



| Tech: | FISH | Array CGH | Genotype arrays | WGS |
|-------|------|-----------|-----------------|-----|
| #: | <10 | 30-100K | 100K-2M | 3G! |

Resolution

# Tumor / Normal Read-Depth Ratio

- Read counting in windows for tumor and normal data



- Log2 ratio for each window
- Chromosome-wide plot

$$\log_2 \frac{\#\text{Reads}_{Disease}}{\#\text{Reads}_{Normal}}$$



EMBL

# Copy-Number Variants

- Can vary the gene dosage of a tumor suppressor or oncogene

- Aneuploidy or non-reciprocal translocations are one form of CNV

- Rare pathogenic germline CNVs can affect known cancer predisposition genes

- Recurrent deletions or duplications indicate a selective advantage

**Aneuploid**

BRCA1 Deletion chr17:41,262,855-41,272,521

p13.2   p12   p11.2 p11.1   q11.2   q12   q21.2   q21.32   q22   q23.1   q24.1   q24.3   q25.2

41,260 kb          41,270 kb

Germline

Tumor

BRCA1

# Copy-Number Variation

# Somatic Structural Variants

# Childhood Brain Tumor Medulloblastoma

- Li-Fraumeni syndrome
  - Germline TP53 mutation

• WHO Grade I Ependymoma

• SHH-medulloblastoma
• Myelodysplastic syndrome (MDS)

• Choroid plexus carcinoma (CPC)



Deletion-type
In-tandem-type
Head-to-head-type
Tail-to-tail-type

log2-ratio (read-depth tumor vs. control)

LFS-MB1, chr3

Chromosomal coordinate [Mb]

# Somatic DNA alterations



Tumor-specific SNVs, non-synonymous (in coding regions)

Deletion (simple)

Tandem dup. (simple)

Inversion (simple)

$Log_2$ read-depth ratio

Complex alterations

Deletion-type

In-tandem-type

Head-to-head-type

Tail-to-tail-type

Inter-chromosomal

EMBL

# Complex DNA alterations forming double-minute chromosome

# Validation of double-minute and co-localization of distant segments



Inter-chromosomal connections validated by PCR

Co-localization of distal segments on chr3 by FISH

# Chromothripsis

EMBL

# Predisposing Germline Structural Variants

# Cancer Predisposition



- Heterozygous germline deletion
- Loss of wildtype copy in the tumor

# Germline SV detection using short-reads is largely incomplete!

# Long-reads and T2T references for SV discovery

Short-reads: 100bp-300bp

Long-reads: 1,000bp-20,000Kbp, few >>20Kbp



Nanopore sequencing

Linear reference genome (GRCh38)

Graph pan-genome



**?**

Human Pangenome
Reference Consortium

EMBL

# Oxford Nanopore Technologies (ONT) and Pacific Biosciences (PacBio)

## Oxford Nanopore Sequencing



1,000bp – 20,000bp reads but some >>20Kbp
~1 error in 100 bases

## Pacific Biosciences Sequencing



1,000bp – 20,000bp reads
~1 error in 10,000 bases

EMBL

# Long read applications

- *De novo* genome Assembly
- Haplotype-resolved genome analysis
- Structural variant (SV) discovery
  - Repetitive SVs
  - Complex SVs
- Resolving genome structure
  - Derivative chromosomes in cancer

# Pan-genome graphs

- A succinct representation of a set of reference genomes



Haplotype-resolved human assemblies

Pan-genome graph

Human Pangenome Reference Consortium (HPRC)
- 44 samples (88 haplotypes) + GRCh38 + CHM13

EMBL

# Pan-genome graphs

- HPRC pan-genome graph: 90 haplotypes (44 samples, GRCh38, CHM13)

- How to incorporate all types of variation?
  - Coarse-grained pangenome graph (structural variants only)
    - 751M on disk: **391,950 segments** (S); 566,204 links (L); 3,198,196,033bp
  - Fine-grained pangenome graph (including small variants)
    - 8.6G on disk: **81,415,956 segments** (S); 112,955,105 links (L); 3,287,932,785bp

Alignment to pan-genome graph

EMBL

# Genome variation discovery using long reads and graph genomes

# 1000 Genomes ONT Project



1,019 samples sequenced with ONT
- ~15x coverage
- **Structural variant calling** using pan-genome graphs

EMBL

# Variant Calling Strategy



**1** Map to linear references and HPRC graph

CHM13 & GRCh38

HPRC_mg

**2** Call structural variants

CHM13 & GRCh38

Delly — v1 v2 | v3 v4 v5
Sniffles — v1 v2 v6 | v4 v5

HPRC_mg

SVarp — v1 | v6 | v2 v3 | v4 v5

**3** Create pseudo-haplotypes

GRCh38

v1 v6
v2 | v3 v4 v5

**4** Construct augmented graph

HPRC_mg_44+966

**5** Augmented graph-based genotyping

HPRC_mg_44+966

giggles and phasing

h1 v1 v2 | v3 v4 v5
h2 v2 | v4 v5

h1 v1 | v2 v6 | v4
h2 v1 | v4

h1 v2 | v3 v4 v5
h2 v2 v6 | v4 v5

— Read    ▬ Reference (linear or graph)    v# Structural Variation    ⊔ Added sequence    h1/2 Haplotype

EMBL

# Long-reads facilitate the discovery of sequence-resolved insertions

# Insertion SV classes: VNTR variation is abundant!



| SV Class | | Resolved INS | DEL | SV Size | MAF |
|---|---|---|---|---|---|
| VNTR | | 26,904 | 16,338 | | |
| Alu | | 20,469 | 2,068 | | |
| Tandem Duplication | | 14,001 | | | |
| L1 | | 4,202 | 227 | | |
| SVA | | 2,571 | 65 | | |
| Non-Canonical MEI | | 1,101 | | | |
| Processed Pseudogene | | 143 | 12 | | |
| NUMT | | 136 | | | |
| solo-LTR | | 25 | | | |
| HERVK | | 6 | | | |

MAF≤1%   MAF>1%

EMBL

# Improved resolution of complex SVs and Inversions

Dotplot



Reference genome

Long – reads

# Repeat-mediated SVs

- A large fraction of deletions (other than VNTR) is repeat-mediated (35%)



Carsten Hain

# Cancer Predisposing SVs

Hypothetical Example: Deletion that affects geneX



BRCA1 Deletion chr17:41,262,855-41,272,521

- Pan-cancer cohort (e.g. 300 breast cancer samples)
  - 18 out of 300 samples are a carrier
  - Allele frequency: ~6%

- 1000 Genomes cohort (2504 samples)
  - 5 out of 2504 samples are a carrier
  - Allele frequency: ~0.2%

- SV may confer a higher risk for breast cancer but be aware of many possible confounders!

  – Sex, Related individuals, Population structure, …

    - All 5 carriers have European ancestry and the cohort of Europeans is much smaller than 2504 samples

  – Technical confounders: Low vs. high-coverage, different insert size, error rate

EMBL

# Structural variants affecting genes



GTEx Portal:
Genotype-Tissue Expression
https://gtexportal.org/home/

# Summary - Challenges in SV and CNV Calling

- Comprehensive detection usually requires long-reads
- Breakpoints from short-reads are often imprecise (e.g., is a fusion gene in-frame?)
- Copy-number baseline in cancer is not necessarily copy-number 2
- Incomplete copy-number and SV polymorphism map
- Very incomplete understanding of SV mechanisms and structure
  - Repeat-mediated SVs
  - Role of centromeres and telomeres?
- Short-read methods tend to have a very high false positive rate
  - Long-read methods are still being developed
- Complex rearrangements are difficult to disentangle with short reads
- Assembling a cancer genome is currently NOT possible



EMBL

# Thank you for your attention!