

Discovery of Structural Variants: Introduction

Alexander Schönhuth



Bielefeld University
April 28, 2021

Organization

ORGANIZATION

- ▶ Organization and introduction: *today*
- ▶ Full literature list available: *by April 30*
- ▶ How to present (brief): *May 5*
- ▶ How to write (brief): *May 12*
- ▶ **Presentations:** *from May 19:*
 - ▶ Each presentation 30-45 minutes
 - ▶ We can do two presentations per week, if this suits best
- ▶ **Technical Report:** *after presentation:*
 - ▶ Each report 8-15 pages
 - ▶ Optimally, report profits from feedback provided after presentation
 - ▶ Drafts can be submitted for discussion
 - ▶ Improving drafts based on feedback

Genetic Variants

GENETIC VARIANTS

Until 2006

Single nucleotide polymorphisms (SNPs)

CCCAGCACTTTGGGAGG	C	CAAGGTGGGGGGAGGAAAT	T	GCTTAAGCCCAGGAGT	Reference
CCCAGCACTTTGGGAGG	T	CAAGGTGGGGGGAGGAAAT	A	GCTTAAGCCCAGGAGT	New Genome

GENETIC VARIANTS

Until 2006

Single nucleotide polymorphisms (SNPs)

CCCAGCACTTTGGGAGG**C**CAAGGTGGGGGGAGGAAAT**T**GCTTAAGCCCAGGAGT Reference
CCCAGCACTTTGGGAGG**T**CAAGGTGGGGGGAGGAAAT**A**GCTTAAGCCCAGGAGT New Genome

From 2006

Structural Variants

Deletion

CCCAGCACTTTGGGAGGCCAAGGTG**GGGGGAG**GAAATTGCTTAAGCCCAGGAGT Reference
CCCAGCACTTTGGGAGGCCAAGGTG**G**GAAATTGCTTAAGCCCAGGAGT New Genome

Insertion

CCCAGCACTTTGGGAGGCCAAGGTGGGGGGAGGAAATTGCTTAAGCCCAGGAGT Reference
CCCAGCACTTTGGGAG**AGTT**ATGCCAAGGTGGGGGGAGGAAATTGCTTAAGCCCAGGAGT New Genome

Translocation

CCCAGCACTTTGGGAG**GCCA**AGGTGGGGGGAGGAAAT**TG**CTTAAGCCCAGGAGT Reference
CCCAGCACTTTGGGAG**AGGTGGGGGGAGGAAATGCCA**TGCTTAAGCCCAGGAGT New Genome

Further variations: inversions, duplications, ...

Somatic Variants

SOMATIC VARIANTS

CANCER ↔ CONTROL

Single Nucleotide Polymorphisms (SNPs)

CAGCATTGAAATA **A** GGCACAT **C** CGAA Cancer Genome

CAGCATTGAAATA T GGCACAT **C** CGAA Control Genome

CAGCATTGAAATA T GGCACAT G CGAA Reference

Somatic SNP **Germline SNP**

Analysis of two samples necessary

How to Discover Genetic Variants?

GENETIC VARIANTS: MODES OF DISCOVERY

RE-SEQUENCING

- ▶ Sequence DNA of genome of interest
- ▶ Align resulting reads against reference genome
- ▶ Note down all differences

DE NOVO ASSEMBLY

- ▶ Sequence DNA of genome of interest
- ▶ Connect resulting reads to form full-length genome
- ▶ Note down differences as per full-length comparison with reference genome

SOMATIC VARIANTS

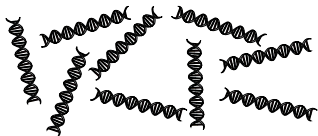
- ▶ Note down differences between cancer and control as well

NEXT GENERATION SEQUENCING

1. Extract Donor Genome DNA



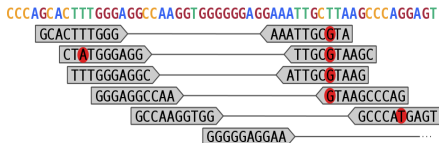
2. Break into fragments



3. Sequence fragments



4. Map against reference genome



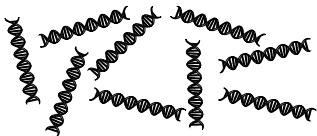
- ▶ For **reference guided variant discovery**, start from 4.
- ▶ For **de novo assembly**, start from 3.

NEXT GENERATION SEQUENCING

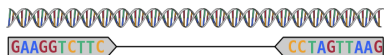
1. Extract Donor Genome DNA



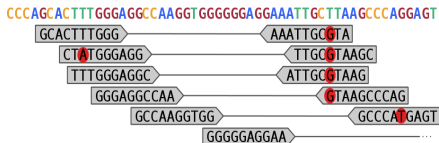
2. Break into fragments



3. Sequence fragments



4. Map against reference genome



- ▶ For **reference guided variant discovery**, start from 4.
- ▶ For **de novo assembly**, start from 3.

Re-Sequencing

RE-SEQUENCING: VARIANT DISCOVERY

Evaluate signals emerging from aligned reads

SNP'S AND SMALL INDELS

- ▶ Look at alignments of single reads with reference

RE-SEQUENCING: VARIANT DISCOVERY

Evaluate signals emerging from aligned reads

SNP'S AND SMALL INDELS

- ▶ Look at alignments of single reads with reference

STRUCTURAL VARIANTS

- ▶ Variants may still yield signals in alignments directly
- ▶ Variants give rise to signals in paired-end alignments

RE-SEQUENCING: VARIANT DISCOVERY

Evaluate signals emerging from aligned reads

SNP'S AND SMALL INDELS

- ▶ Look at alignments of single reads with reference

STRUCTURAL VARIANTS

- ▶ Variants may still yield signals in alignments directly
- ▶ Variants give rise to signals in paired-end alignments

Example: Read Pair Signals

DISCOVERING INDELS

Reference genome

CCCAGCACTTTGGGAGGCCAA**GGTGGGGGAGG**AAATTGCTTAAGCCCAGGAGT

DISCOVERING INDELS

Reference genome

CCCAGCACTTTGGGAGGCCAA**GGTGGGGGAGG**AAATTGCTTAAGCCCAGGAGT

Sequenced genome

CCCAGCACTTTGGGAGGCCAAAATTGCTTAAGCCCAGGAGT



DISCOVERING INDELS

Reference genome

CCCAGCACTTTGGGAGGCCAA**GGTGGGGGAGG**AAATTGCTTAAGCCCAGGAGT

Sequenced genome

CCCAG**CACTTTGGGAGGCCAA**AAATTGCTTAAGCCCAGGAGT

Fragment

DISCOVERING INDELS

Reference genome

CCCAGCACTTTGGGAGGCCAA**GGTGGGGGGAGG**AAATTGCTTAAGCCCAGGAGT

Sequenced genome

CCCAG**GCAC**TTTGGGAGGCCAA**AAATTGCTTAAGCCCAG**GGAGT

GGACTTTGGG ————— **TTAAGCCCAG**

DISCOVERING INDELS

Reference genome

CCCAGCACTTTGGGAGGCCAA**GGTGGGGGGAGG**AAATTGCTTAAGCCCAGGAGT

GGACTTTGGG

TTAAGCCCAG

Sequenced genome

CCCAGCACTTTGGGAGGCCAA**AAATTGCTTAAGCCCAG**GAGT

GGACTTTGGG

TTAAGCCCAG

too long!

- **Insertions:** alignment length **too short**

STATISTICAL QUANTIFICATION

FRAGMENT LENGTH DISTRIBUTION

Reference genome

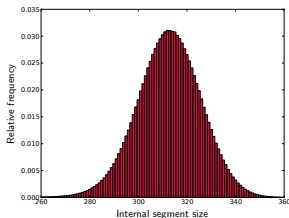
CCCAGCACTTTGGGAGGCCAA**GGTGGGGGAGG**AAATTGCTTAAGCCCAGGAGT
GGACTTTGGG TTAAGCCCAG

Sequenced genome

CCCAGCACTTTGGGAGGCCAA**AAATTGCTTAAGCCCAG**GAGT
GGACTTTGGG TTAAGCCCAG

too long!

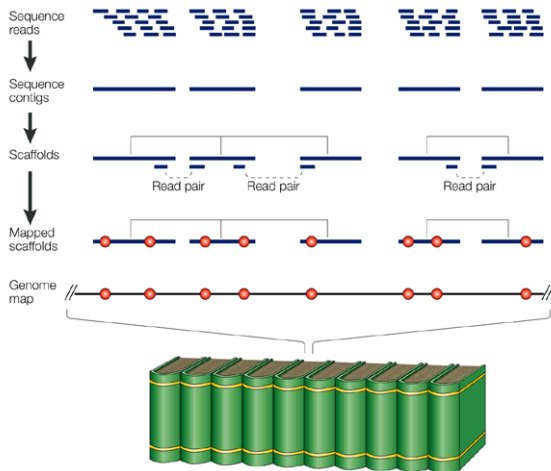
- ▶ **fragment length** $\sim \mathcal{N}_{\mu, \sigma}$
- ▶ alignment length L : the greater $|L - \mu|$, the more significant



Genome Assembly

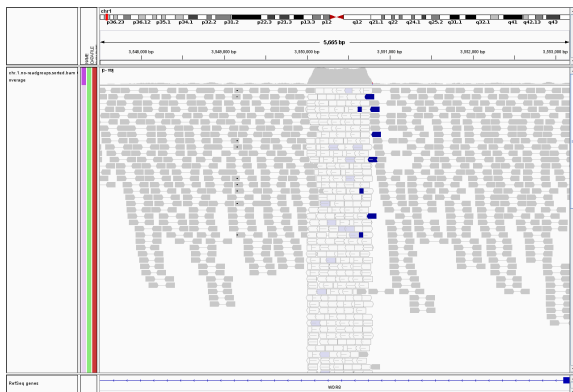
GENOME ASSEMBLY

STEP BY STEP



From Contigs To Scaffolds: Why?

REPETITIVE SEQUENCE



- ▶ Repetitive areas disturb the problem decisively
- ▶ Make **contigs** from reads exhibiting unique sequence
- ▶ Make **scaffold** of contigs, to bridge repetitive sequence

GENERATING CONTIGS

- ▶ De novo assembly programs usually deliver contigs
- ▶ Scaffolding requires data from additional, unconventional sources
- ▶ We will focus on contig generation in the following

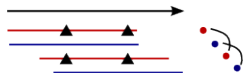
GENERATING CONTIGS

- ▶ De novo assembly programs usually deliver contigs
- ▶ Scaffolding requires data from additional, unconventional sources
- ▶ We will focus on contig generation in the following

GENERATING CONTIGS: ASSEMBLY PARADIGMS

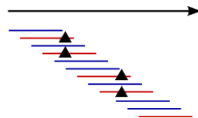
Overlap graphs

nodes: sequencing reads
edges: approximate
suffix-prefix overlaps

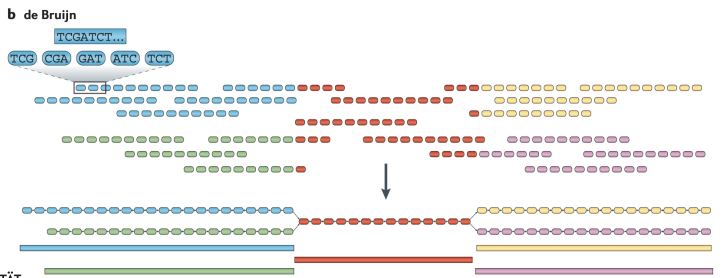
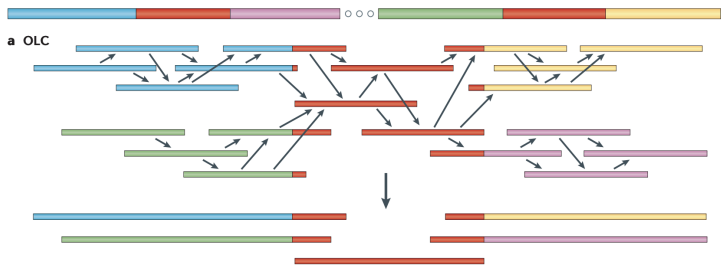


De Bruijn graphs

length k substrings
exact overlaps of
length $k-1$



ASSEMBLY PARADIGMS



GENOME ASSEMBLY PARADIGMS

Overlap-Layout-Consensus Paradigm

- ▶ Natural approach: identify overlapping fragments
- ▶ Seminal paper: [Kececioglu, Myers, 1995]
- ▶ In use at Celera for human genome assembly

GENOME ASSEMBLY PARADIGMS

Overlap-Layout-Consensus Paradigm

- ▶ Natural approach: identify overlapping fragments
- ▶ Seminal paper: [Kececioglu, Myers, 1995]
- ▶ In use at Celera for human genome assembly

De Bruijn Graph Paradigm

- ▶ Construct *de Bruijn graph* from sequence fragments
- ▶ “*Cut pieces into even smaller pieces to solve the puzzle*”
- ▶ Seminal paper: [Pevzner, Tang, Waterman, PNAS 2001]
- ▶ Predominant paradigm for NGS based genome assembly

De Bruijn Graphs

DE BRUIJN GRAPHS

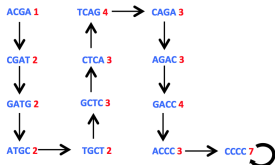
Genome: ACGATGCTCAGACCCCCCCC
 Short reads: ACGATGCTCAGA CTCAGACCC AGACCCC CCCCCC

k-mers:

```

    ACGAT      CTCAG      AGACC      CCCCC
     CGATG    TCAGA    GACCC    CCCCC
      GATGC    CAGAC    GACCC    CCCCC
        ATGCT   AGACC
         TGCTC   GACCC
          GCTCA
           CTCAG
            TCAGA
    
```

De Bruijn graph:



Assembled Contigs: ACGATGCTCAGACCCC

Advantages:

- ▶ Removes redundancy among reads
- ▶ Easy to construct and store

[Lu, Shen, Warren and Walter, 2016]

DE BRUIJN GRAPHS

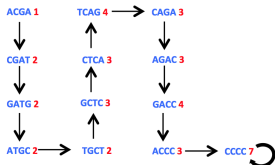
Genome: ACGATGCTCAGACCCCCCCC
 Short reads: ACGATGCTCAGA CTCAGACCC AGACCCC CCCCCC

k-mers:

```

    ACGAT      CTCAG      AGACC      CCCCC
     CGATG    TCAGA      GACCC      CCCCC
      GATGC    CAGAC      GACCC      CCCCC
        ATGCT  AGACC
          TGCTC AGACC
            GCTCA GACCC
              CTCAG
                TCAGA
    
```

De Bruijn graph:



Assembled Contigs: ACGATGCTCAGACCCC

Advantages:

- ▶ Removes redundancy among reads
- ▶ Easy to construct and store

Issues:

- ▶ loses information about linked mutations
- ▶ works only well with error-corrected reads

[Lu, Shen, Warren and Walter, 2016]

DE BRUIJN GRAPHS

SUMMARY

- ▶ Superior for NGS based consensus genome assembly
- ▶ For polyploid genome assembly: *colored de Bruijn graphs* (not treated here, but seem to have certain limitations)
- ▶ (Likely severe) limitations for polyploid genome assembly:
 - ▶ Error correction removes lowly frequent true mutations
 - ▶ Loose information about linked mutations

Overlap-Layout-Consensus


OVERLAP-LAYOUT-CONSENSUS (OLC)

- ▶ **Overlap:** Construct overlap graph from sequencing reads
- ▶ **Layout:** Compute *contigs*, longer stretches of overlapping reads using the overlap graph
- ▶ **Consensus:** Compute consensus sequence for contigs

OLC AND POLYPLOID GENOMES

SUMMARY

Advantages

- ▶ Edges  overlapping reads from the same haplotype
- ▶ Error correction *after graph construction*

Disadvantages

- ▶ Overlap graph construction time and space consuming
- ▶ Error correction: new ideas required

Overlap Graph Construction

OVERLAP GRAPH CONSTRUCTION

Without reference genome

For each pair of reads, determine

- ▶ how they optimally overlap
- ▶ whether the resulting overlap indicates that the two reads are from the same haplotype (statistically significantly likely!)

Naive approaches infeasible!

Literature

LITERATURE REFERENCES

REVIEWS

- ▶ *Structural Variant Discovery: “Genome structural variation discovery and genotyping”*, by Alkan et al., Nature Reviews Genetics, 2011
- ▶ *De Novo Assembly: “Genetic variation and the de novo assembly of human genomes”*, Chaisson et al., Nature Reviews Genetics, 2015
- ▶ *Long-range sequencing: “Piercing the dark matter: bioinformatics of long-range sequencing and mapping”*, Sedlazeck, et al., Nature Reviews Genetics, 2018
- ▶ *Somatic Structural Variant Discovery: “Structural variant detection in cancer genomes: computational challenges and perspectives for precision oncology”*, van Belzen et al., Nature Precision Oncology, 2021

LITERATURE REFERENCES

PAPERS

- ▶ *Varlociraptor*: Koester et al., Genome Biology, 2020 (somatic variants)
- ▶ *Delly*: Rausch et al., Bioinformatics, 2012 (structural variants, also somatic)
- ▶ *FreeBayes*: Garrison and Marth, arXiv:1207.3907, 2012 (haplotype-aware small variants)
- ▶ *Hifiasm*: Cheng et al., Nature Methods, 2021 (de novo assembly of “HiFi” reads)
- ▶ *MiniMap2*: Heng Li, Bioinformatics, 2018 (overlap graph construction)
- ▶ *Lancet*: Narzisi et al., Communications Biology, 2018 (somatic variants)
- ▶ *GRIDSS*: Cameron et al., Genome Research, 2017 (rearrangement type variants)
- ▶ *Sniffles*: Sedlazeck et al., Nature Methods, 2018 (structural variants from long reads)
- ▶ *NanoVar*: Tham et al., Genome Biology, 2020 (structural variants from long reads)