**Principles and problems of de novo genome assembly**

Karin Lagesen

Norwegian Veterinary Institute

Some materials adapted from
slides provided by Lex Nederbragt

---

What is this thing called 'genome assembly'?
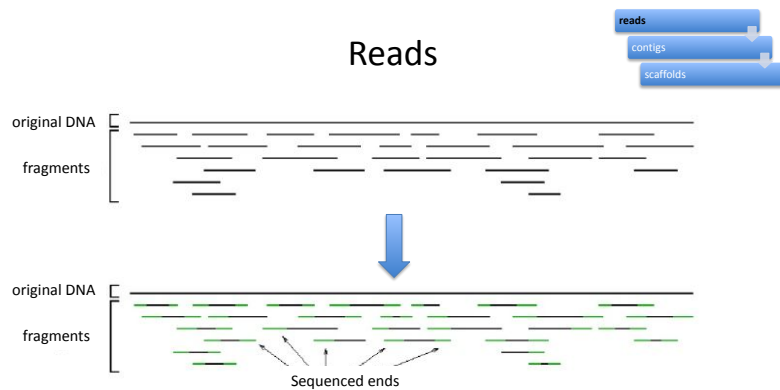
## What is a genome assembly?

A hierarchical data structure

that maps the sequence data

to a putative reconstruction of the target

Miller et al 2010, Genomics 95 (6): 315-327
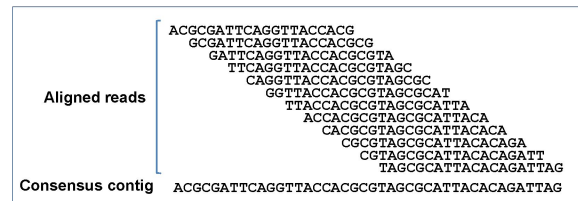
## Hierarchical structure

reads

contigs

scaffolds

# Sequence data

## Reads

original DNA

fragments

original DNA

fragments

Sequenced ends

http://www.cbcb.umd.edu/research/assembly_primer.shtml

# Contigs

## Building contigs

Aligned reads

```
ACGCGATTCAGGTTACCACG
 GCGATTCAGGTTACCACGCG
  GATTCAGGTTACCACGCGTA
    TTCAGGTTACCACGCGTAGC
     CAGGTTACCACGCGTAGCGC
      GGTTACCACGCGTAGCGCAT
       TTACCACGCGTAGCGCATTA
        ACCACGCGTAGCGCATTACA
         CACGCGTAGCGCATTACACA
          CGCGTAGCGCATTACACAGA
           CGTAGCGCATTACACAGATT
            TAGCGCATTACACAGATTAG
```

Consensus contig    ACGCGATTCAGGTTACCACGCGTAGCGCATTACACAGATTAG

# Contigs

Building contigs

**Repeat copy 1**  **Repeat copy 2**

Contig orientation?
Contig order?

**Collapsed repeat
consensus**

http://www.cbcb.umd.edu/research/assembly_primer.shtml

# Contigs

Repeats: major problem

**Repeat copy 1**  **Repeat copy 2**

http://www.cbcb.umd.edu/research/assembly_primer.shtml

# Long reads

Other read type

**Repeat copy 1**    **Repeat copy 2**

long reads

# Scaffolds

Ordered, oriented contigs

mate pairs

contigs

gap size estimate

Hierarchical structure

reads

contigs

scaffolds

ACCACGCGTAGCGCAT
ACCACGCGTAGCGCATTA
ACCACGCGTAGCGCATTACA
CACGCGTAGCGCATTACACA
CGCGTAGCGCATTACACAGA
CGTAGCGCATTACACAGA
TAGCGCATTACACAGA
ACGCGTAGCGCATTACACAGA

Assembly

How to do this?

Algorithms

# Algorithms

All are graph-based

Node                    Node

Edge, often
directed

**Graph-theory!**

# Algorithms

All are graph-based

Read 1                  Read 2

Overlap

Aligned reads
```
ACGCGATTCAGGTTACCACG
  GCGATTCAGGTTACCACGCG
    GATTCAGGTTACCACGCGTA
      TTCAGGTTACCACGCGTAGC
        CAGGTTACCACGCGTAGCGC
          GGTTACCACGCGTAGCGCAT
            TTACCACGCGTAGCGCATTA
              ACCACGCGTAGCGCATTACA
                CACGCGTAGCGCATTACACA
                  CGCGTAGCGCATTACACAGA
                    CGTAGCGCATTACACAGATT
                      TAGCGCATTACACAGATTAG
```
Consensus contig  ACGCGATTCAGGTTACCACGCGTAGCGCATTACACAGATTAG

**Graph-theory!**

# Algorithms

Hamiltonian path
  — a path that contains all the nodes

# Algorithms

Overlap calculation (alignment)
  — computationally intensive

# Algorithms

Path through the graph

☐contig



Read 1       Read 2       Read 3       Read 4

**Overlap**     **Overlap**     **Overlap**

```
ACGCGATTCAGGTTACCACG
 GCGATTCAGGTTACCACGCG
  GATTCAGGTTACCACGCGTA
   TTCAGGTTACCACGCGTAGC
    CAGGTTACCACGCGTAGCGC
     GGTTACCACGCGTAGCGCAT
      TTACCACGCGTAGCGCATTA
       ACCACGCGTAGCGCATTACA
        CACGCGTAGCGCATTACACA
         CGCGTAGCGCATTACACAGA
          CGTAGCGCATTACACAGATT
           TAGCGCATTACACAGATTAG
```
**Aligned reads**

**Consensus contig**   ACGCGATTCAGGTTACCACGCGTAGCGCATTACACAGATTAG

---

# Algorithms

Many flavors



Two most used
- Overlap Layout Consensus
- de Bruijn graph

http://www.waialuasodaworks.com/images/flavors2009.jpg

# Overlap-Layout-Consensus

Developed for Sanger-type reads (longer reads)



# Overlap-Layout-Consensus

Steps
- Overlap computation
- Layout: graph simplification
- Consensus: sequence

# Overlap-Layout-Consensus

Overlap phase: find "similar enough" reads

Comparing all against all: expensive

Trick for finding "similar enough" reads:

- Split reads into k-mers

  K-mer: substring of length *k* from a longer string

  `ACGCGATTCAGGTTACCACG`

- Make list over which read has which k-mers

- If two reads share k-mers, test for similarity

---

# Overlap-Layout-Consensus

**A** Read Layout

```
R₁:  GACCTACA
R₂:    ACCTACAA
R₃:      CCTACAAG
R₄:       CTACAAGT
A:         TACAAGTT
B:          ACAAGTTA
C:           CAAGTTAG
X:         TACAAGTC
Y:          ACAAGTCC
Z:           CAAGTCCG
```
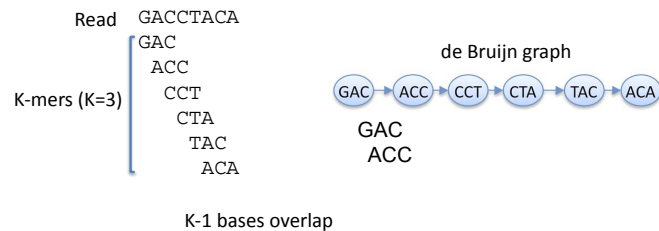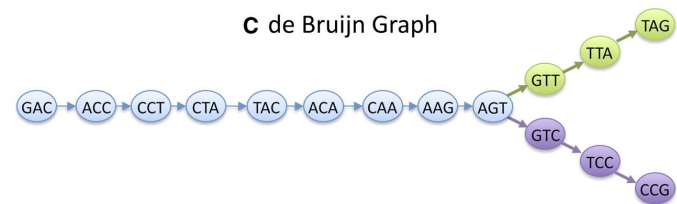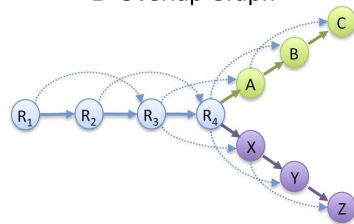
**B** Overlap Graph



Schatz M C et al. Genome Res. 2010;20:1165-1173

## de Bruijn graphs

Developed outside of DNA-related work
- Best solution for short(er) reads

Read    GACCTACA

K-mers (K=3)
GAC
ACC
CCT
CTA
TAC
ACA

de Bruijn graph

GAC → ACC → CCT → CTA → TAC → ACA

GAC
ACC

K-1 bases overlap

## Graphs

**C** de Bruijn Graph

GAC → ACC → CCT → CTA → TAC → ACA → CAA → AAG → AGT → GTT → TTA → TAG
AGT → GTC → TCC → CCG

Schatz M C et al. Genome Res. 2010;20:1165-1173

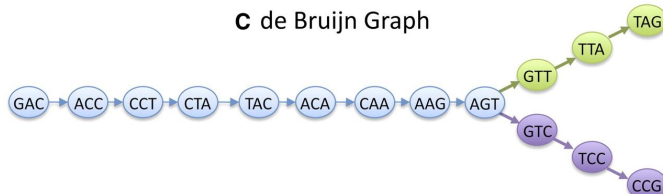# Graphs

**A** Read Layout

$R_1$: GACCTACA
$R_2$: ACCTACAA
$R_3$: CCTACAAG
$R_4$: CTACAAGT
A: TACAAGTT
B: ACAAGTTA
C: CAAGTTAG
X: TACAAGTC
Y: ACAAGTCC
Z: CAAGTCCG

**B** Overlap Graph

**C** de Bruijn Graph

Schatz M C et al. Genome Res. 2010;20:1165-1173

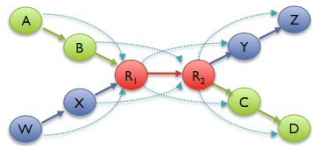# Graphs

Simplify the graph

Add scaffolding information
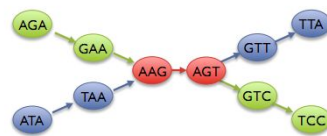
## de Bruijn Graphs

### Overlap Graph



Long read assemblers

- Repeats depends on read length
- Read coherency, placements kept
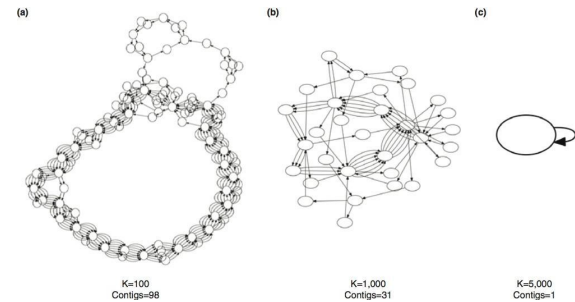- Tangled by high coverage

### de Bruijn Graph

Short read assemblers

- Repeats depends on word length
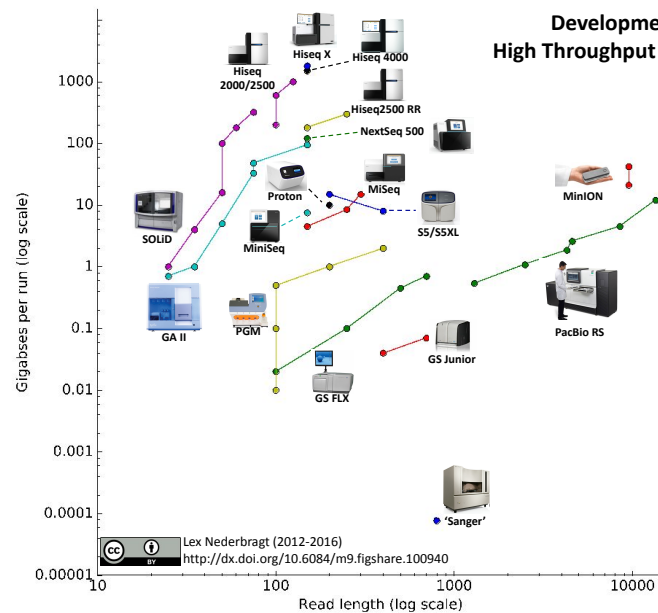- Read coherency, placements lost
- Robust to high coverage

Mike Schatz

## Read length matters

5.2 Mb circular genome, infinite error-free reads



(a) K=100 Contigs=98

(b) K=1,000 Contigs=31

(c) K=5,000 Contigs=1

Developments in High Throughput Sequencing. Lex Nederbragt (2012-2016) http://dx.doi.org/10.6084/m9.figshare.100940

Why is genome assembly such a difficult problem?

## 1) Repeats

**Repeat copy 1**

**Repeat copy 2**

Repeats break up assembly

**Collapsed repeat consensus**

## 2) Diploidy

Differences between sister chromosomes

'heterozygosity'

## 2) Diploidy

Region 1 → Polymorphic region 2 → Region 4

Region 1 → Polymorphic region 3 → Region 4

Homozygous    Heterozygous    Homozygous

## 3) Polyploidy

Haploid (N)    Diploid (2N)

Triploid (3N)    Tetraploid (4N)

http://en.wikipedia.org/wiki/Polyploidy

# 4) Lots of tools to choose from

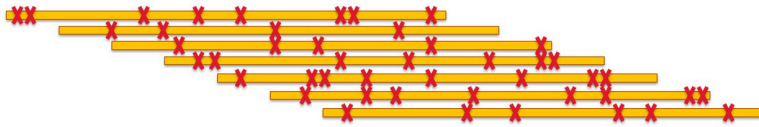| First generation 1 kb High accuracy Sanger | ARACHNE [70], Atlas [71], CAP3 [72], Celera [73], Euler [74], JAZZ [75], Minimus [76], MIRA [77], phrap [78], Phusion [79], SUTTA [80], TIGR [81] |
|---|---|
| Second generation 25-300 bp High accuracy 454, IonTorrent, Solexa, SOLiD | ABySS [82, 83], ALLPATHS [84], BASE [85], CABOG [86], Edena [87], EPGA [88], Euler-SR [89], Gossamer [90], IDBA [91], ISEA [92], JR-Assembler [93], LightAssembler [94], Meraculous [95], MIRA [77], Newbler [96], PCAP [97], PERGA [98], Platanus [99], PE-Assembler [100], QSRA [101], Ray [102], Readjoiner [103], SGA [104], SHARGCS [105], SOAPdenovo [106], SOAPdenovo2 [107], SPAdes [108], SparseAssembler [109], SSAKE [110], SUTTA [80], Taipan [111], VCAKE [112], Velvet [113] |
| Third generation 10-100,000+ kb PacBio CLR, Nanopore | Canu [114], FALCON [115], Flye [116], HINGE [117], MECAT [118], MECAT2 [118], miniasm [119], NECAT [120], NextDenovo [121], Ra [122], Raven [123], Shasta [124], SMARTdenovo [125], wtdbg [126], wtdbg2 [127] |
| 15-25 kb High accuracy PacBio HiFi, | Flye [116], HiCanu [128], hifiasm [129], IPA [130], LJA [131], mdBG [132], MBG [133], NextDenovo [121], Peregrine [134], Raven [123], wtdbg2 [127] |

Guiglielmoni, *et. al.* Peer Comm. J. 2022

Assembly with noisy single molecule sequencing data
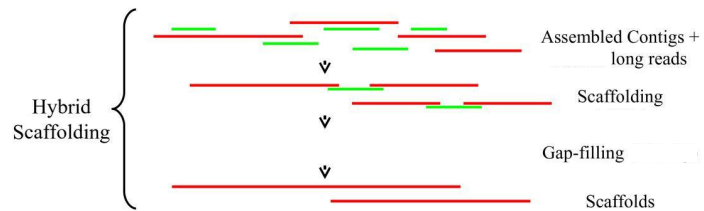
## Usage of long reads

- Problem: higher error rates
- Overlaps more difficult/expensive to find
- OLC more commonly used than for 2nd generation data

## Long read assembly strategies

- Alt 0: Scaffolding short read asms

- Alt 1: Correct reads, then assemble

- Alt 2: Assemble reads, then correct

## Scaffolding and gap closing (hybrid)



Hybrid Scaffolding

Assembled Contigs + long reads
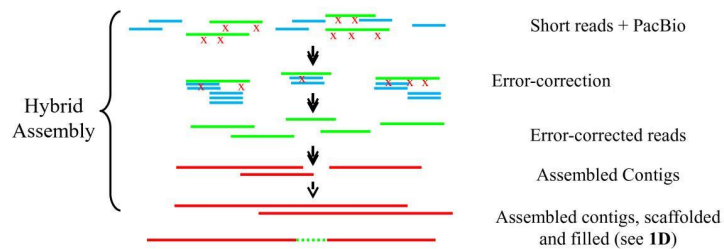
Scaffolding

Gap-filling

Scaffolds

Powers *et.al.*, BMC genomics 2013
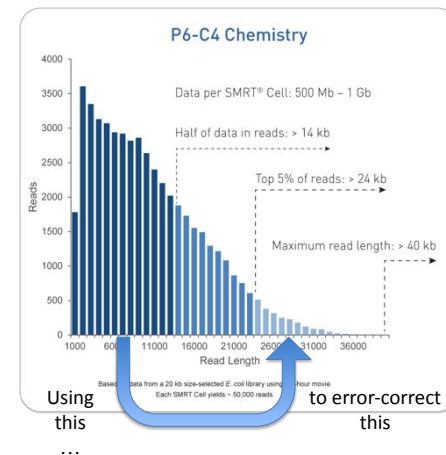
## Correct, assemble

- Do pairwise comparison, find shorter reads that support the longer
- Align supporting reads, correct longer reads
- Overlap-Layout-Consensus on corrected reads
- Polish assembly
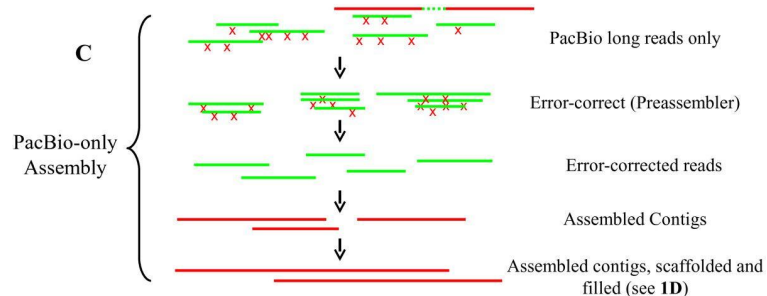
# Mapping and error correcting (hybrid)



Short reads + PacBio

Error-correction

Error-corrected reads

Assembled Contigs

Assembled contigs, scaffolded and filled (see **1D**)

Hybrid Assembly

Powers *et.al.*, BMC genomics 2013

# Hierarchical approach (self-correcting)



**P6-C4 Chemistry**

Data per SMRT® Cell: 500 Mb – 1 Gb

Half of data in reads: > 14 kb

Top 5% of reads: > 24 kb

Maximum read length: > 40 kb

Reads

Read Length

Based on data from a 20 kb size-selected *E. coli* library using 4 hour movie.
Each SMRT Cell yields ~ 50,000 reads.

Using this … to error-correct this

https://genome.duke.edu/cores-and-services/sequencing-and-genomic-technologies/pacbio

## Short read error correction



C

PacBio-only Assembly

PacBio long reads only

Error-correct (Preassembler)

Error-corrected reads

Assembled Contigs

Assembled contigs, scaffolded and filled (see **1D**)

Powers *et.al.*, BMC genomics 2013

## Assemble, correct

- Compare reads, find overlaps
- Assemble reads, knowing things will be wrong
- Align reads to assembly
- Correct assembly

Questions?