

De novo assembly

| | |
|--|--|
| DNA assemblers | RNA assemblers |
| <ul style="list-style-type: none"> • Velvet • ABySS • SOAPdenovo • ALLPATHS-LG | <ul style="list-style-type: none"> • Trinity • Velvet/oases • TransABySS • SOAPdenovo-Trans • Rnnotator |

7

The Overlap-layout-consensus (OLC) approach

1. Pairwise alignments and overlap graph
2. Graph Layout: search of a single path in the graph (i.e. the Hamiltonian path)
3. Multiple sequence alignments and consensus

Examples: Newbler, Celera, Arachne, Edena, YASRA

Tristan Lefebvre, Cornell University

ACTGATTG

9

ACTGATTG

k-mer = 5

10

ACTGA

ACTGATTG

k-mer = 5

11

ACTGA
CTGAT

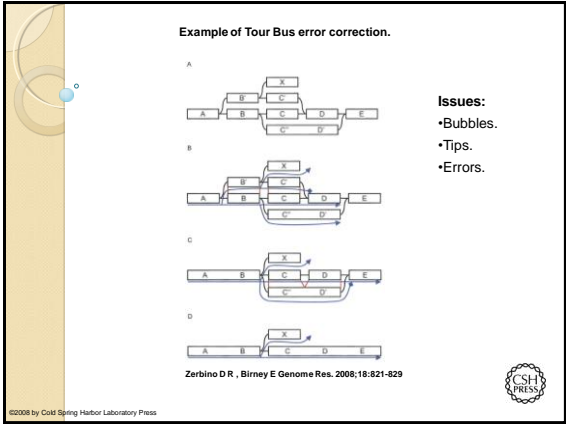
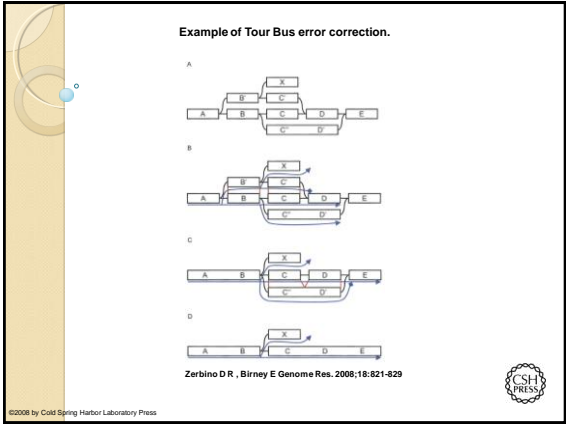
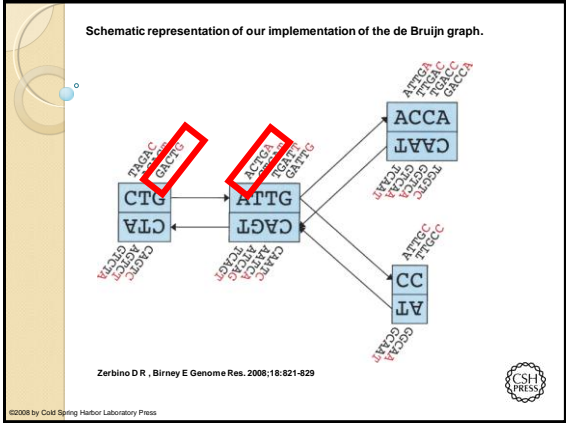
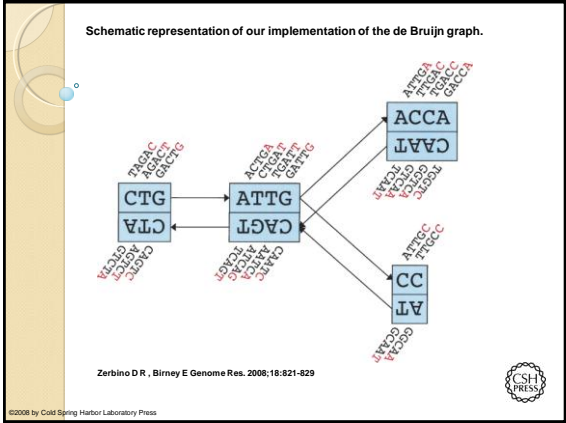
ACTGATTG

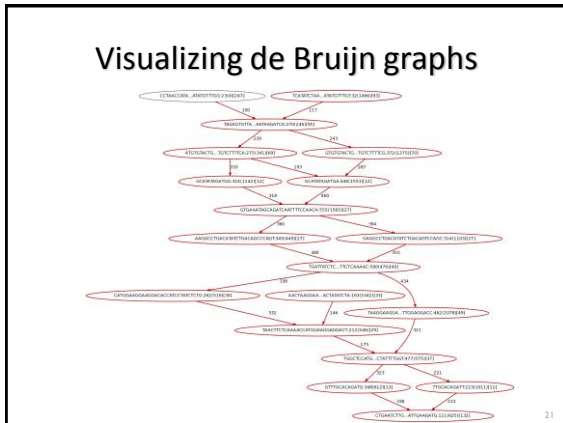
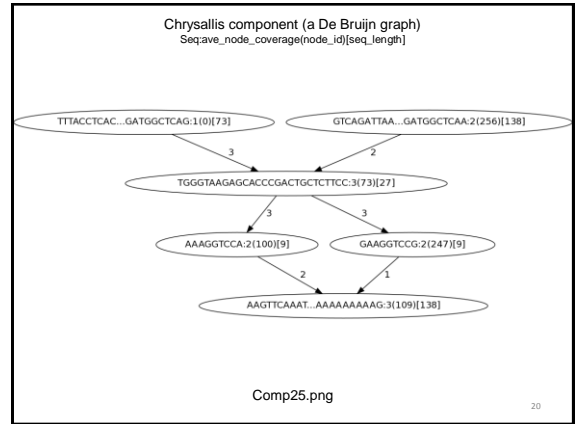
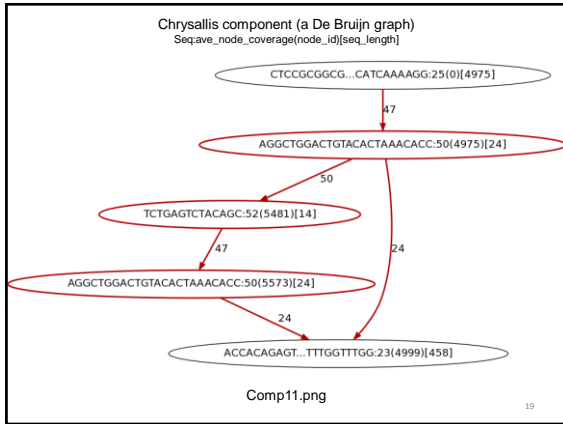
k-mer = 5

12

Brian Knaus, US Forest Service

Botany 2012 Intro to Next Generation Sequencing Workshop





Read Mapping

- Bowtie/bowtie2
- BWA
- MAQ
- Stampy
- RUM (RNA)

22

MAQ: Mapping and Assembly with Qualities

- Hash based assembler.
- Compiled code (C).
- Seeded alignment.
 - Uses first 28 bp (highest quality region).
 - Allows for up to 2 mismatches in this region (~7% sequence divergence).
 - After finding these regions it extends through the rest of the read length.
- Currently limited to 63 bp reads.
- Ungapped alignment.

23

MAQ vs. Bowtie

- Maq = 0.49 million reads mapped per hour.
- Bowtie = 29.5 million reads mapped per hour.

24

Bowtie

- Reference:
 - Burrows-Wheeler transform.
 - FM index (Feragina and Manzini).
- Reads from 4 to 1,024 bp long.
- Greedily finds matches.

25

Fig. 1. Prefix trie of string "GGOGGG". Symbol Δ marks the start of the string. The two numbers in a node give the SA interval of the string represented by the node (see Section 2.3). The dashed line shows the route of the brute-force search for a query string "LGL", allowing at most one mismatch. Edge labels in squares mark the mismatches to the query in searching. The only hit is the bold node [1, 1] which represents string "GGG".

Li and Durbin. 2009. *Bioinformatics* 25: 1754-1760.

26

Genomic references.
 Splice junctions?

27

Strand specificity

Parkhomchuck et al. 2009. Transcriptome analysis by strand-specific sequencing of complementary DNA. *Nucleic Acids Research* 37(18):e123

28

Graphical options

- Galaxy
- Easy terminal alternative (ETA)
- CLC Genomics Workbench (commercial)
- Geneious (commercial)