

Assembling Genomes

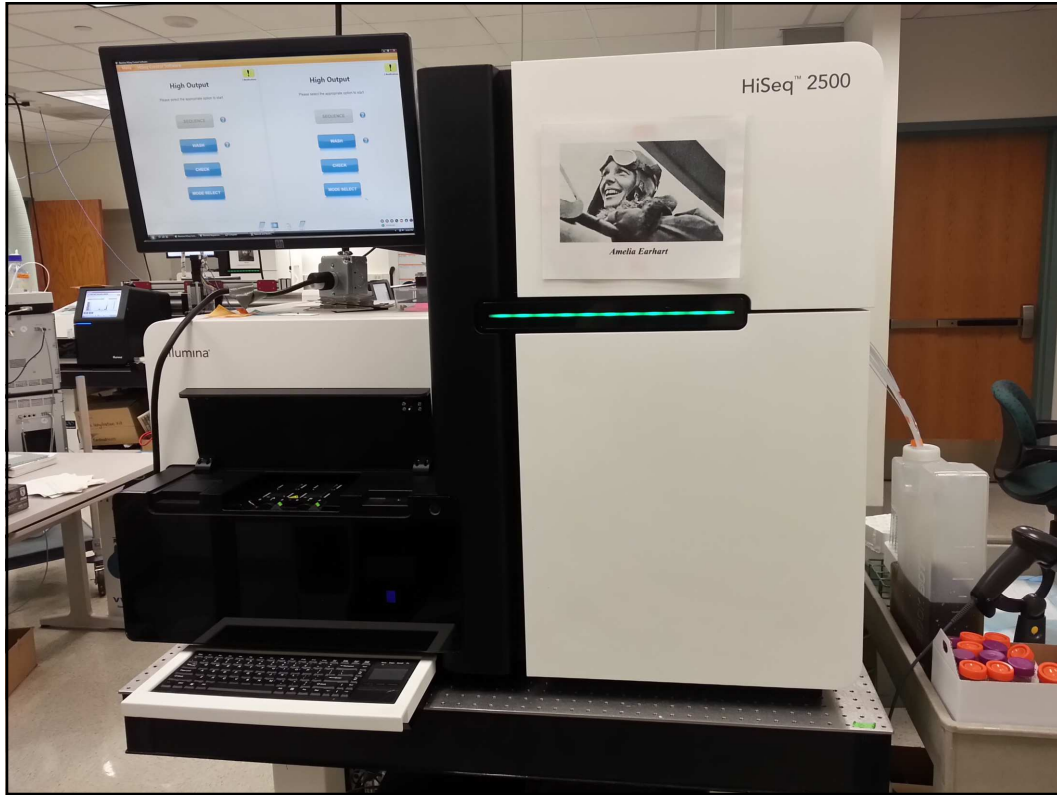
BCH394P/364C Systems Biology / Bioinformatics

Edward Marcotte, Univ of Texas at Austin

1



2



3



4

China's BGI says it can sequence a genome for just \$100

Super-cheap DNA sequencing could boost cancer screening, prenatal tests, and research into population genetics.

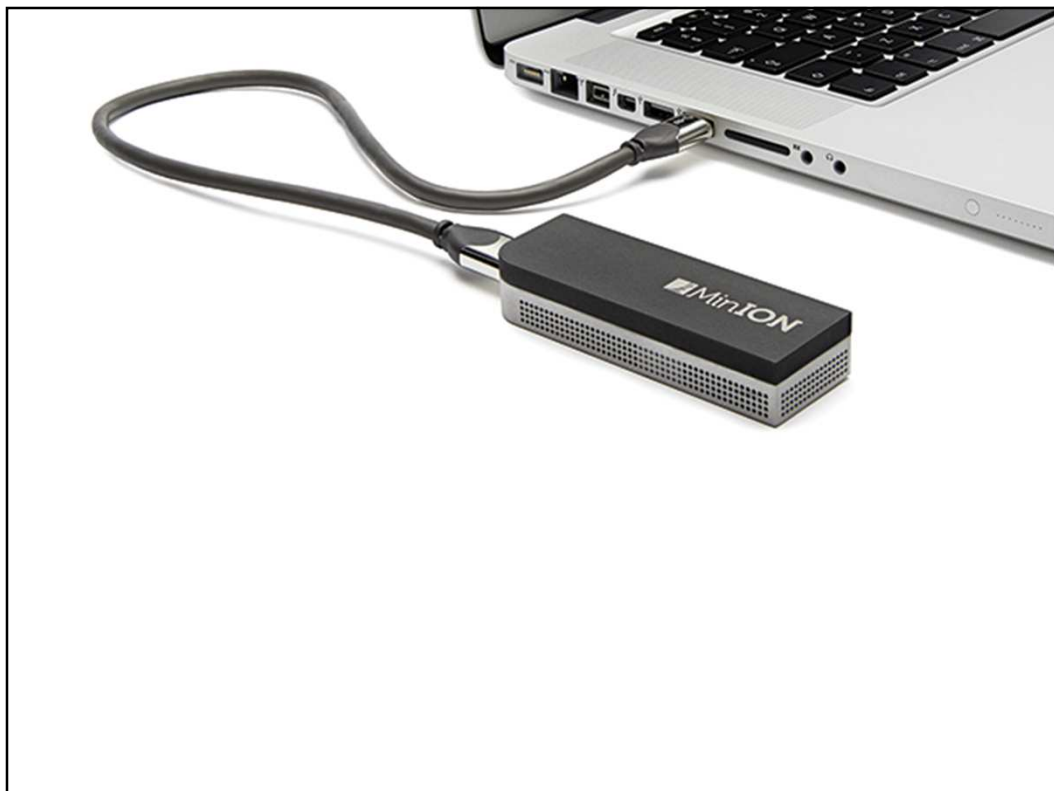
by **Antonio Regalado**

Feb 26, 2020

Using technology originally acquired in the US, the Chinese gene-giant BGI Group says it will make genome sequencing cheaper than ever, breaking the \$100 barrier for the first time.

<https://www.technologyreview.com/s/615289/china-bgi-100-dollar-genome/>

5



6

Oxford Nanopore Sequencers Have Left UK for China, to Support Rapid, Near-sample Coronavirus Sequencing for Outbreak Surveillance



English ▾

NEWS PROVIDED BY
Oxford Nanopore Technologies →
Jan 31, 2020, 10:50 ET

SHARE THIS ARTICLE



<https://www.prnewswire.com/news-releases/oxford-nanopore-sequencers-have-left-uk-for-china-to-support-rapid-near-sample-coronavirus-sequencing-for-outbreak-surveillance-300996908.html>

7

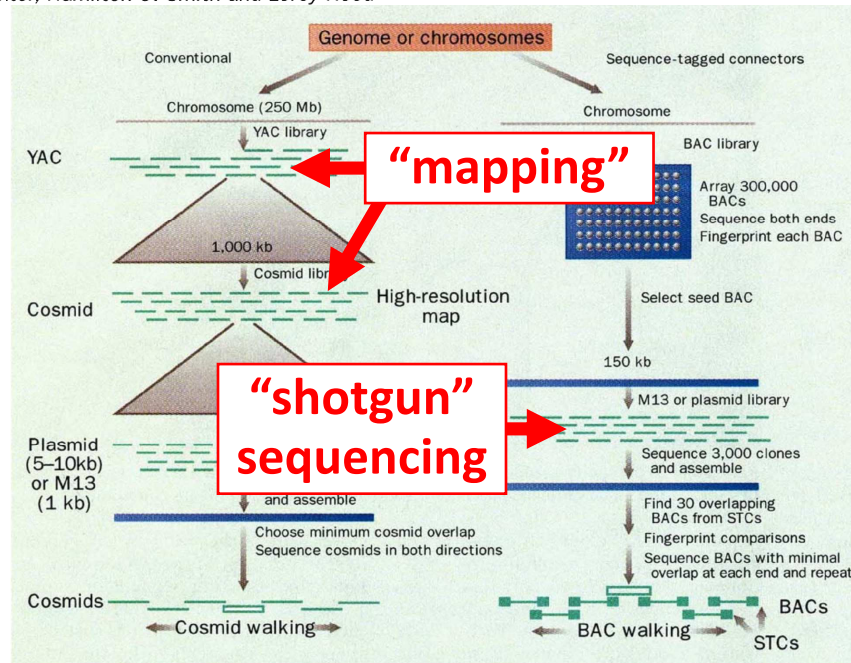


<http://www.triazle.com>; The image from http://www.danielbert.com/port_fun.html
Reference: Jones NC, Pevzner PA, Introduction to Bioinformatics Algorithms, MIT press

8

A new strategy for genome sequencing

J. Craig Venter, Hamilton O. Smith and Leroy Hood



NATURE · VOL 381 · 30 MAY 1996

9

(Translating the cloning jargon)

CLONE LIBRARIES USED FOR GENOME MAPPING AND SEQUENCING		
Vector	Human-DNA insert size range	Number of clones required to cover the human genome
Yeast artificial chromosome (YAC)	100–2,000 kb	3,000 (1,000 kb)
Bacterial artificial chromosome (BAC)	80–350 kb	20,000 (150 kb)
Cosmid	30–45 kb	75,000 (40 kb)
Plasmid	3–10 kb	600,000 (5 kb)
M13 phage	1 kb	3,000,000 (1 kb)

NATURE · VOL 381 · 30 MAY 1996

10

Thinking about the basic shotgun concept

- Start with a very large set of random sequencing reads
- How might we match up the overlapping sequences?
- How can we assemble the overlapping reads together in order to derive the genome?

11

Thinking about the basic shotgun concept

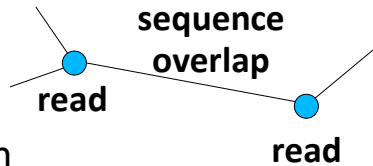
- At a high level, the first genomes were sequenced by comparing pairs of reads to find overlapping reads
- Then, building a graph (*i.e.*, a network) to represent those relationships
- The genome sequence is a “walk” across that graph

12

The “Overlap-Layout-Consensus” method

Overlap: Compare all pairs of reads
(allow some low level of mismatches)

Layout: Construct a graph describing the overlaps



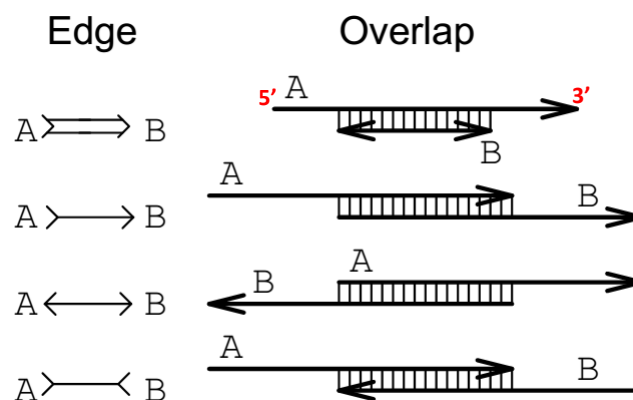
Simplify the graph

Find the simplest path through the graph

Consensus: Reconcile errors among reads along that path to find the consensus sequence

13

Building an overlap graph

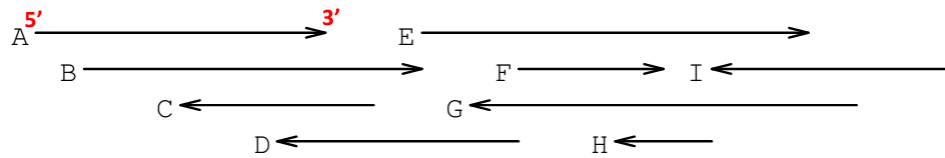


EUGENE W. MYERS. *Journal of Computational Biology*. Summer 1995, 2(2): 275-290

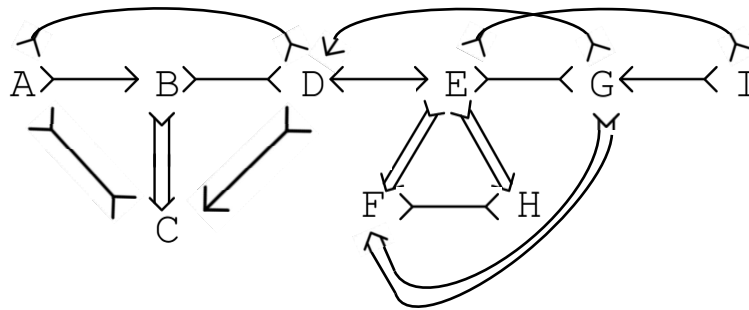
14

Building an overlap graph

Reads



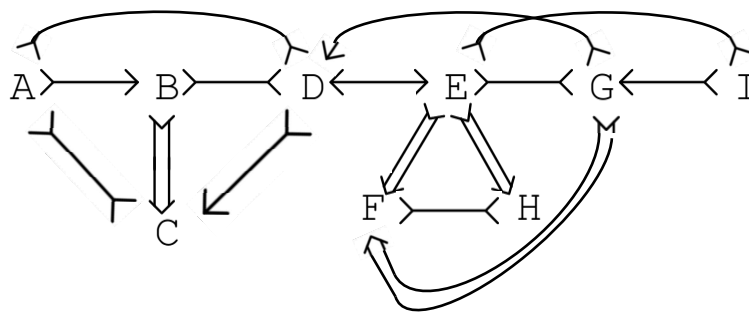
Overlap graph



EUGENE W. MYERS. *Journal of Computational Biology*. Summer 1995, 2(2): 275-290 (more or less)

15

Simplifying an overlap graph

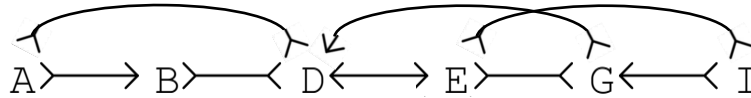


1. Remove all contained nodes & edges going to them

EUGENE W. MYERS. *Journal of Computational Biology*. Summer 1995, 2(2): 275-290 (more or less)

16

Simplifying an overlap graph



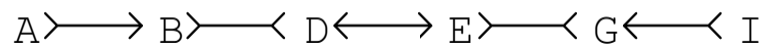
2. Transitive edge removal:

Given $A - B - D$ and $A - D$, remove $A - D$

EUGENE W. MYERS. *Journal of Computational Biology*. Summer 1995, 2(2): 275-290 (more or less)

17

Simplifying an overlap graph



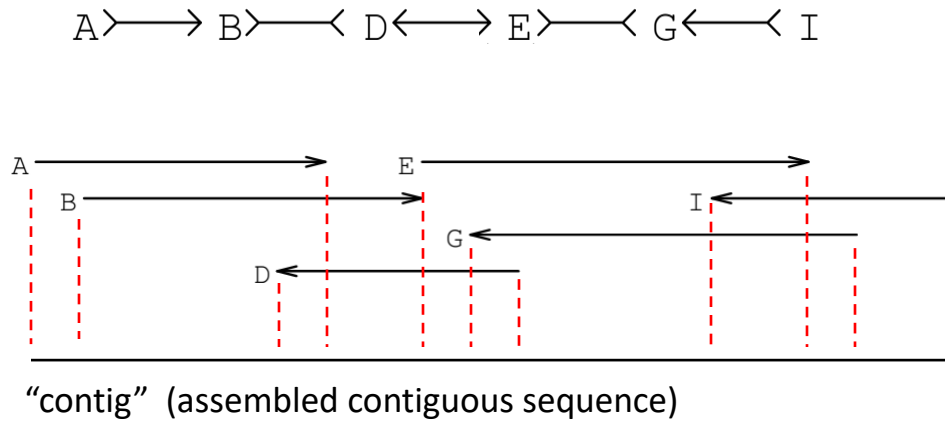
3. If un-branched, calculate consensus sequence

If branched, assemble un-branched bits and then decide how they fit together

EUGENE W. MYERS. *Journal of Computational Biology*. Summer 1995, 2(2): 275-290 (more or less)

18

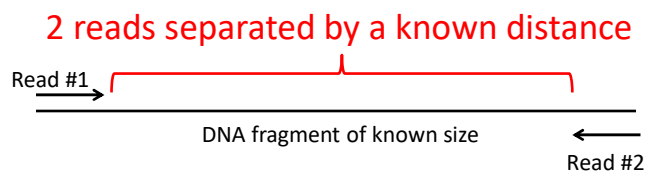
Simplifying an overlap graph



EUGENE W. MYERS. *Journal of Computational Biology*. Summer 1995, 2(2): 275-290 (more or less)

19

This basic strategy was used for most of the early genomes.
Also useful: "mate pairs"



Contigs can be ordered using these paired reads

The diagram shows two horizontal lines representing contigs, labeled "Contig #1" and "Contig #2". A red curved arrow points from Contig #1 to Contig #2, indicating the relationship between them. Below the contigs, the text "to produce 'scaffolds'" is written in red.

20

GigAssembler (used to assemble the public human genome project sequence)



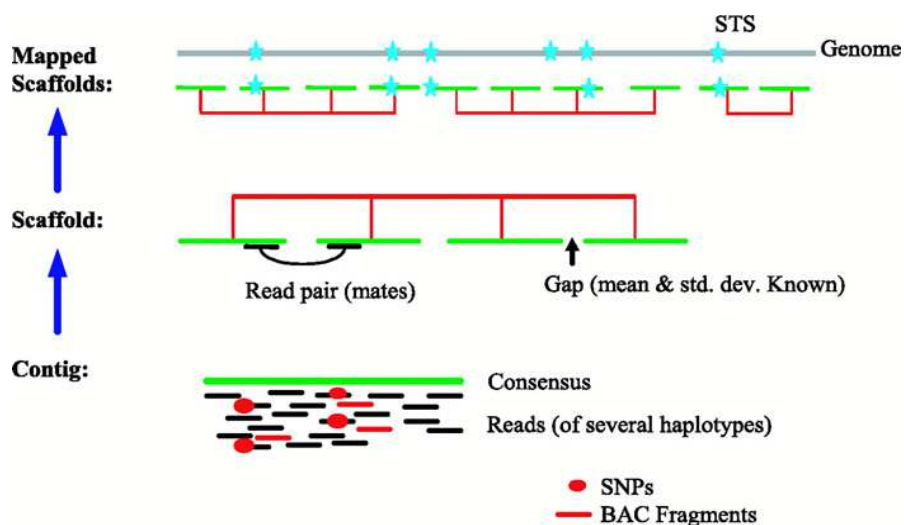
Jim Kent

David Haussler

Let's take a little walk through history to see what they did...

21

Whole genome Assembly: big picture



<http://www.nature.com/scitable/content/anatomy-of-whole-genome-assembly-20429>

22

GigAssembler – Preprocessing

1. Decontaminating & Repeat Masking.
2. Aligning of mRNAs, ESTs, BAC ends & paired reads against initial sequence contigs.
 - psLayout → BLAT
3. Creating an input directory (folder) structure.

```
chr1/  
chr1/contig1.e  
chr1/contig1.a  
chr1/contig1.c  
chr1/contig1.b  
chr1/contig1.d  
chr3/  
chr2/  
chr2/contig2.d  
chr2/contig2.b  
chr2/contig2.a  
chr2/contig2.c
```

23

RepBase + RepeatMasker

```
tajeon@fourierseq:~/RepBase/RepBase15.05.fasta$ ls -a  
.  
..  
angrep.ref  
bctrep.ref  
cbrrep.ref  
celrep.ref  
chlrep.ref  
cinrep.ref  
cinunc.ref  
dcotrep.ref  
diarep.ref  
dngrep.ref  
fngrep.ref  
fugrep.ref  
grasrep.ref  
humrep.ref  
humsub.ref  
invrep.ref  
invsref.ref  
mamrep.ref  
mamsub.ref  
mcoutrep.ref  
mousub.ref  
nemrep.ref  
nlnrep.ref  
prirep.ref  
prisub.ref  
pseudo.ref  
ratsub.ref  
rodsref.ref  
rodsref.ref  
simple.ref  
spurep.ref  
synrep.ref  
tmpnemrep.ref  
tmpnemrep.ref  
tmpnemrep.ref  
version  
vrtrep.ref  
zebrep.ref
```

```
>MER1D ERV1 Homo sapiens  
tgaggcaggagaaaatagcagaggggaattggaagtggataaaggagagaatgagtaaaagcangagagca  
gaagcaaggtaaaagaggcgggtgagcaagaagcaagataagaagcagaagttagcagcgaacacaaaag  
taagatnanaaagaagtgagtaaggagccacatggctggctagatccagacacacagtaaggagcag  
ctctcagagatgggcatgtacattagagagaaaaagtatcttaaaatgacccgtatgataatcagct  
cattaaagctcatgcatatggactgcatatcatgcatgacttaaaattatgggagggagtgacgcga  
agawgtcacacagcacaggggcatagkattaagtaactaagcaacccacatcaatcaaaaggcaga  
tgcctgctagagattaggcagccttgggaagagaagaaaaaacacataaaagacccaagtcac  
caactgacgctgattctcatttcagagagtgagccactctccctctctgagagtgtaatactgtgct  
taataaaattttgctgcttctgctgctgctgctgctgctgctgctgctgctgctgctgctgctgctgct  
ggaactgcacrgcaccakctggtaaca  
>MIRb SINE2/tRNA Mammalia  
cagaggggagcagctggtgagtggaagagcagggccttggagtcaggcagacctgggttcgaatcctg  
gctctgcaacttactagctgtgtgaccttggcgaagtcacttaacctctctgagcctcagtttctctc  
tgtaaaatggggataataatctacctcagcaggggtgtgtgtgagattaaatgagataatgcatgtaa  
gctctagcagagtgctgacacagtaagcgcctcaataaatggtagctctattatt  
>LTR45 ERV1 Homo sapiens  
tgtaacccgggaccagcgaacactgggctactctgttgatacaaaatgtcaagttacctttagtgta  
taacagagcccaaaactgcaagtcagtgacccgggcatgtgcaatagaaaaagctttagccttaacaa  
caccagaacaaatgattctccctcggaaaccaagaagacgggacatgaccggaacctgaatgccga  
actctttcagaagcaaaaggggtcgttggccgggaagatctggggctaaatctgctcaacatacctta  
ccgtaaatgttcaaatggaagccctcgaacagacccgcaagcgaacattcctaaatcctttccctt  
gcccctgtagcctttaaacttgcgccagaccccaaatcggggagagagatttagccacccctctgct  
cctgtgctggcgttttgcaataaagcctttctttctcaaaagcgtgtgcatagttattggtctctgt  
gtgcatcaggcagcaagccattgtctcgataaca  
>MER80B hAT Homo sapiens  
cagggctcttaaacagaggttccatggatggccttcaggaggtctgtgaacccctggaattatatacaa  
aaatgtgtgtatgtgcatatgtatgttttctggggagaggggttcagctttcatagattctcaa  
aggggtctatgatctmaaaaggttaagaagccctg
```

24

GigAssembler: Build merged sequence contigs (“rafts”)




Figure 1 Two sequences overlapping end to end. The sequences are represented as dashes. The aligning regions are joined by vertical bars. End-to-end overlap is an extremely strong indication that two sequences should be joined into a contig.

25

[illegible]

Sequencing quality (Phred Score)

$$Q = -10 \log_{10} P \quad \text{Base-calling Error Probability}$$

or

$$P = 10^{\frac{-Q}{10}}$$

Phred quality scores are logarithmically linked to error probabilities

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90 %
20	1 in 100	99 %
30	1 in 1000	99.9 %
40	1 in 10000	99.99 %
50	1 in 100000	99.999 %

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

27

GigAssembler: Build merged sequence contigs (“rafts”)

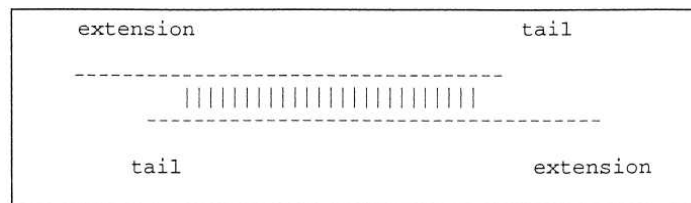


Figure 2 Two sequences with tails. The nonaligning regions on either side can be classified into ‘extensions’ and ‘tails.’ Short tails are fairly common even when two sequences should be joined into a contig because of poor quality sequence near the ends and occasional chimeric reads. Long tails, however, are generally a sign that the alignment is merely due to the sequences sharing a repeating element.

28

GigAssembler: Build merged sequence contigs (“rafts”)

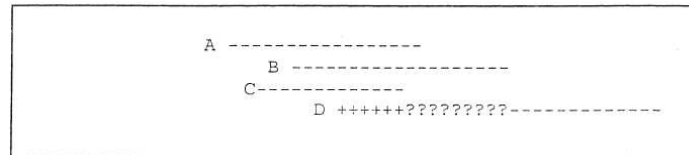


Figure 3 Merging into a raft. A contig (‘raft’) of three sequences: A, B, and C has already been constructed by GigAssembler. The program now examines an alignment between sequence C and a new sequence, D, to see whether D should also be added to the raft. The parts of D marked with +s are compatible with the raft because of the C/D alignment. The program must also check that the parts of D marked with ?s are compatible with the raft by examining other alignments.

29

GigAssembler: Build sequenced clone contigs (“barges”)

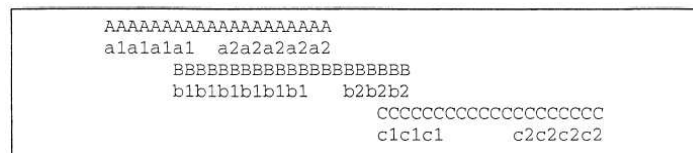


Figure 4 Three overlapping draft clones: A, B, and C. Each clone has two initial sequence contigs. Note that initial sequence contigs a1, b1, and a2 overlap as do b2 and c1.

30

GigAssembler: Build a “raft-ordering” graph

```

AAAAAAAAAAAAAAAAAAAA
a1a1a1a1 a2a2a2a2a2
BBBBBBBBBBBBBBBBBBBB
b1b1b1b1b1b1 b2b2b2
CCCCCCCCCCCCCCCCCCCC
c1c1c1 c2c2c2c2

```

Figure 4 Three overlapping draft clones: A, B, and C. Each clone has two initial sequence contigs. Note that initial sequence contigs a1, b1, and a2 overlap as do b2 and c1.

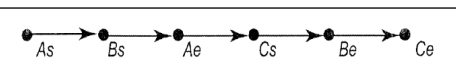


Figure 5 Ordering graph of clone starts and ends. This represents the same clones as in Fig. 4. (As) The start of clone A; (Ae) the end of clone A. Similarly Bs, Be, Cs, and Ce represent the starts and ends of clones B and C.

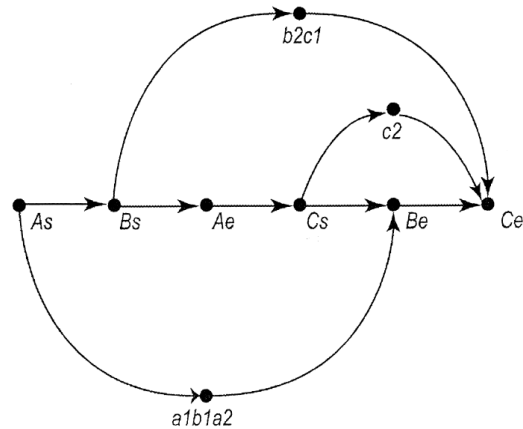


Figure 6 Ordering graph after adding in rafts. The initial sequence contigs shown in Fig. 4 are merged into rafts where they overlap. This forms three rafts: a1b1a2, b2c1, and c2. These rafts are constrained to lie between the relevant clone ends by the addition of additional ordering edges to the graph shown in Fig. 5.

31

GigAssembler: Build a “raft-ordering” graph

- Add information from mRNAs, ESTs, paired plasmid reads, BAC end pairs: building a “bridge”
 - Different weight to different data type: (mRNA ~ highest)
 - Conflicts with the graph as constructed so far are rejected.
- Build a sequence path through each raft.
- Fill the gap with N's.
 - 100: between rafts
 - 50,000: between bridged barges

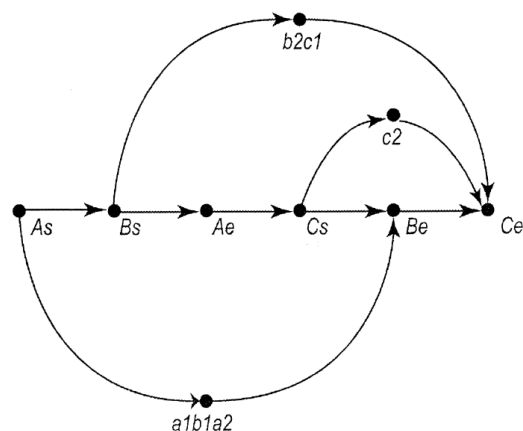


Figure 6 Ordering graph after adding in rafts. The initial sequence contigs shown in Fig. 4 are merged into rafts where they overlap. This forms three rafts: a1b1a2, b2c1, and c2. These rafts are constrained to lie between the relevant clone ends by the addition of additional ordering edges to the graph shown in Fig. 5.

32

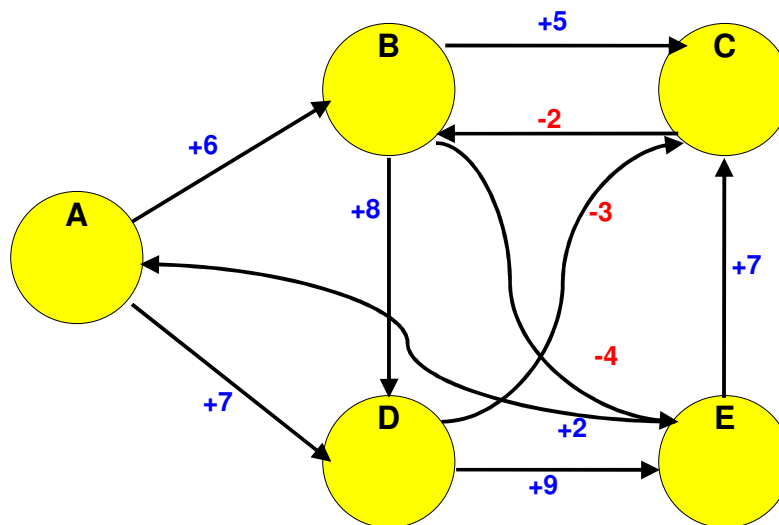
Finding the shortest path across the ordering graph using the Bellman-Ford algorithm

<http://compprog.wordpress.com/2007/11/29/one-source-shortest-path-the-bellman-ford-algorithm/>

33

Find the shortest path to all nodes.

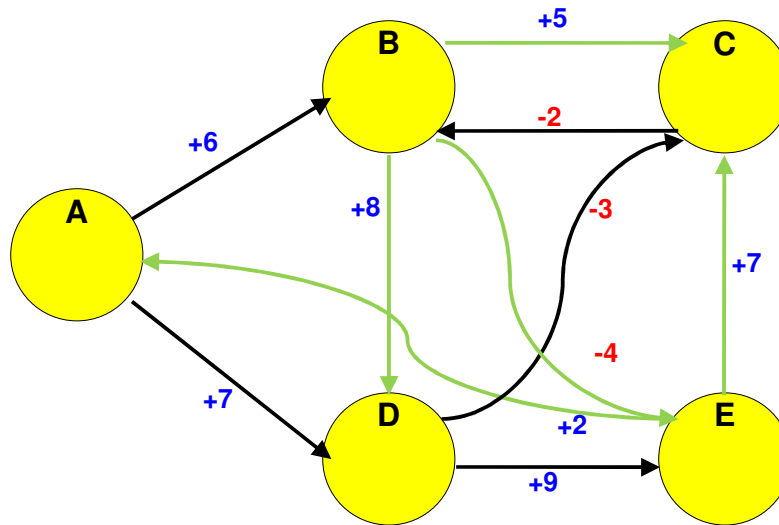
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



34

Find the shortest path to all nodes.

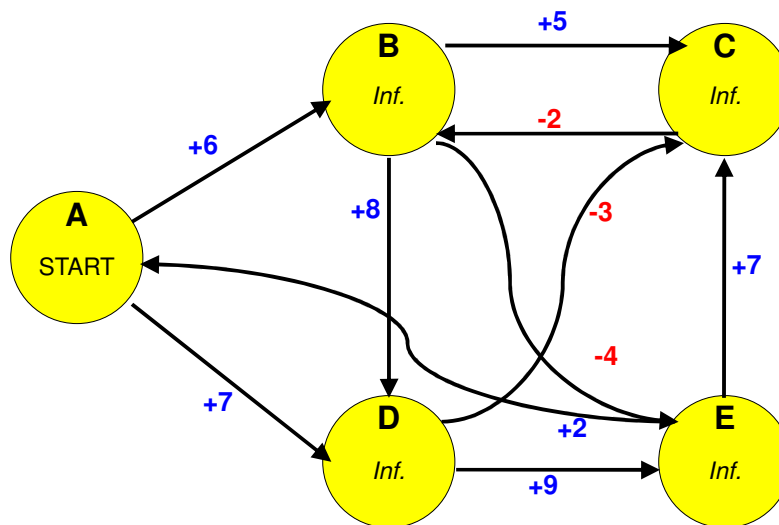
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



35

Find the shortest path to all nodes.

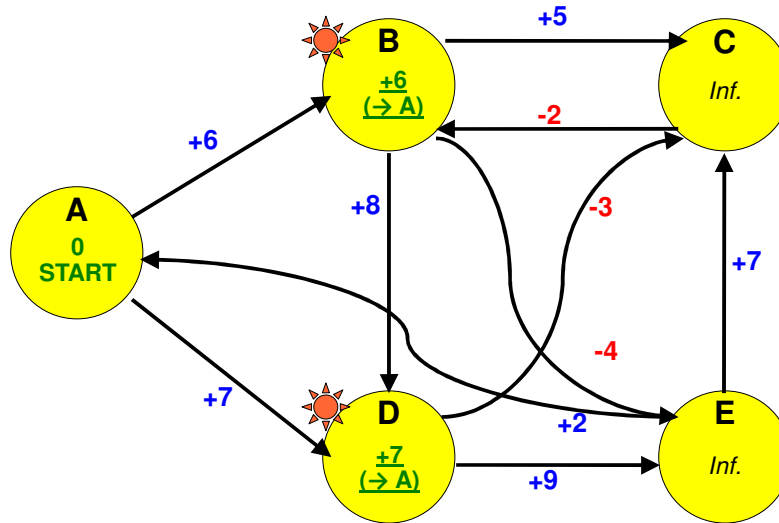
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



36

Find the shortest path to all nodes.

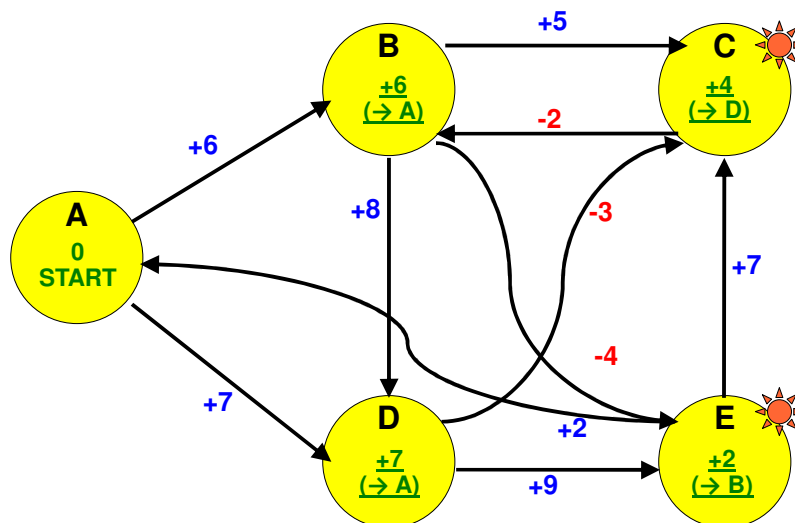
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



37

Find the shortest path to all nodes.

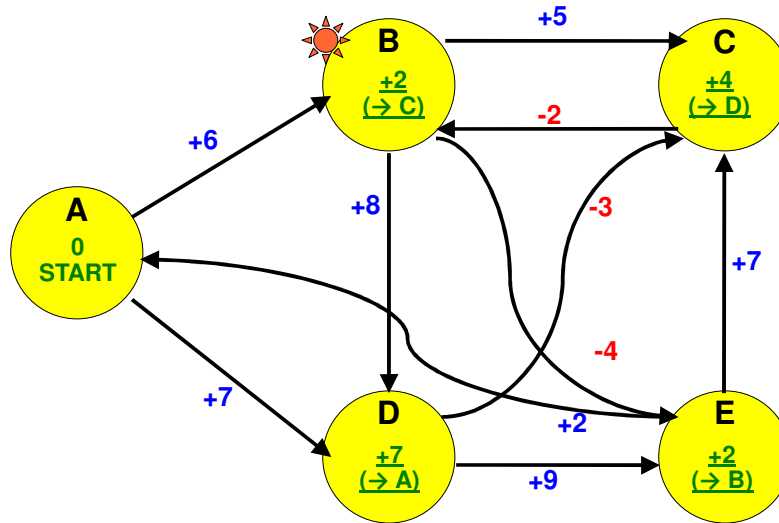
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



38

Find the shortest path to all nodes.

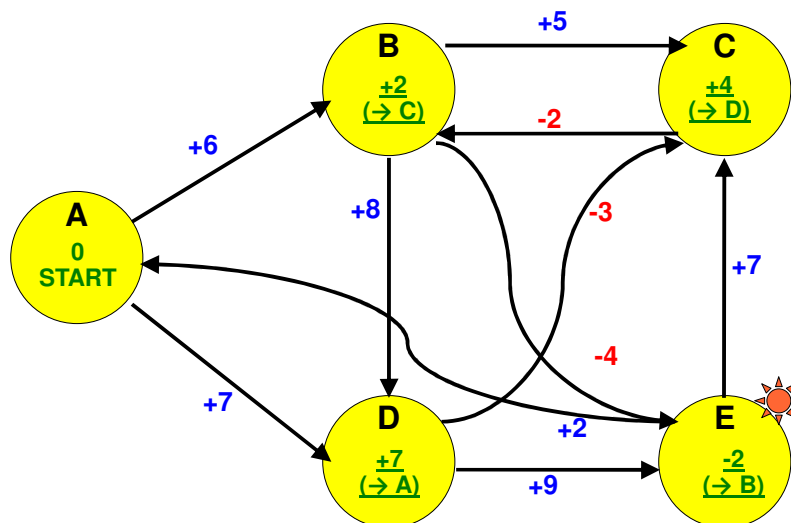
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



39

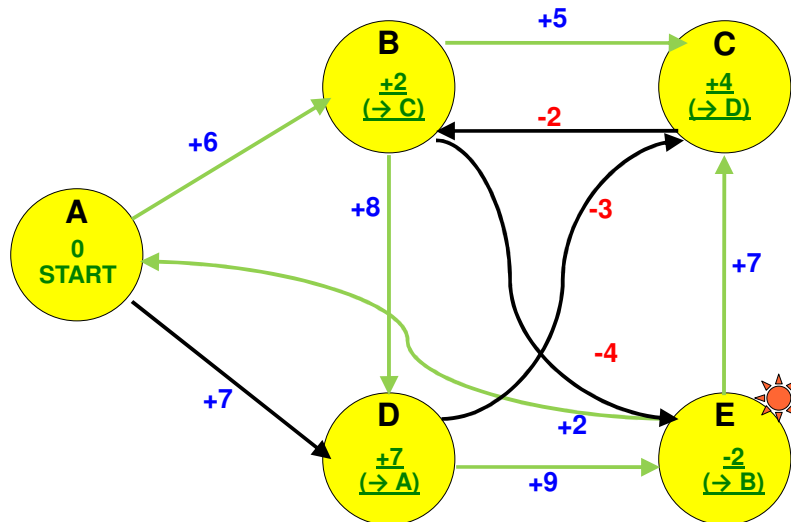
Find the shortest path to all nodes.

Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



40

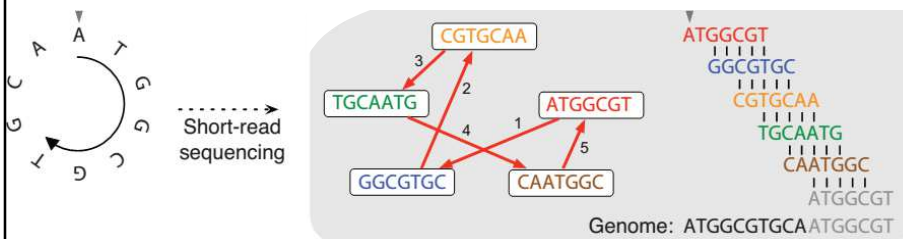
Answer: A-D-C-B-E



41

Modern assemblers now work a bit differently, using so-called **DeBruijn graphs**:

Here's what we saw before:



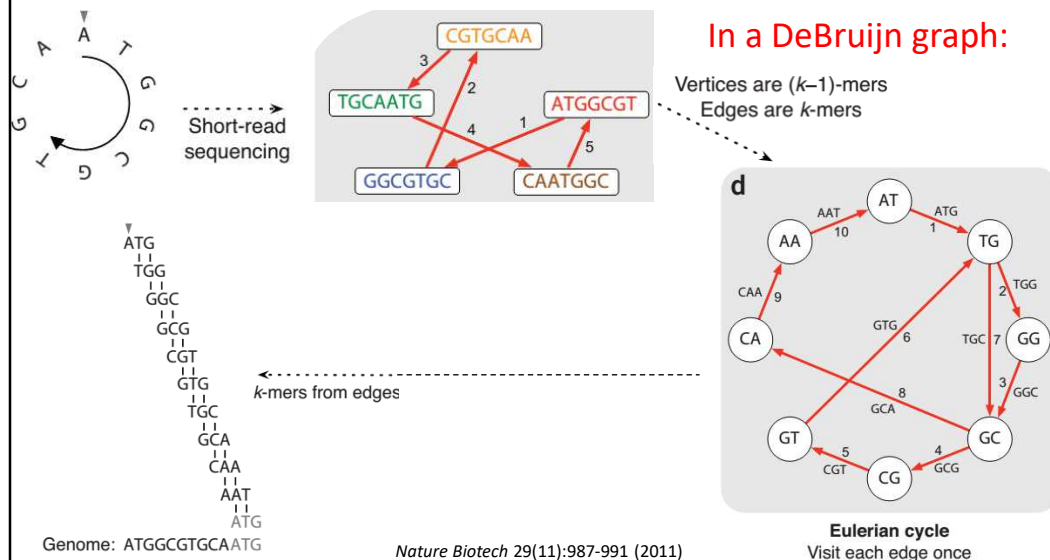
In Overlap-Layout-Consensus:

Nodes are reads
Edges are overlaps

Nature Biotech 29(11):987-991 (2011)

42

Modern assemblers now work a bit differently, using so-called **DeBruijn graphs**:

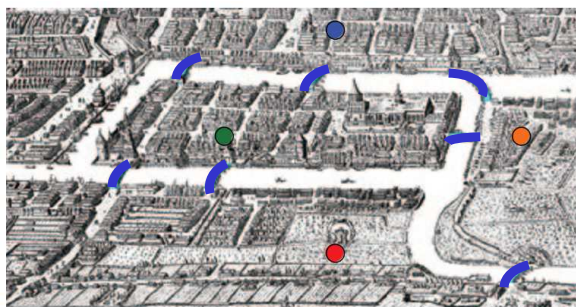


43

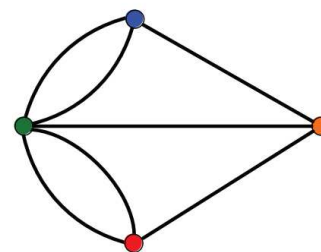
Why Eulerian?

From Leonhard Euler's solution in 1735 to the 'Bridges of Königsberg' problem:

Königsberg (now Kaliningrad, Russia) had 7 bridges connecting 4 parts of the city. **Could you visit each part of the city, walking across each bridge only once, & finish back where you started?**



(Visiting every edge once = an *Eulerian* path)

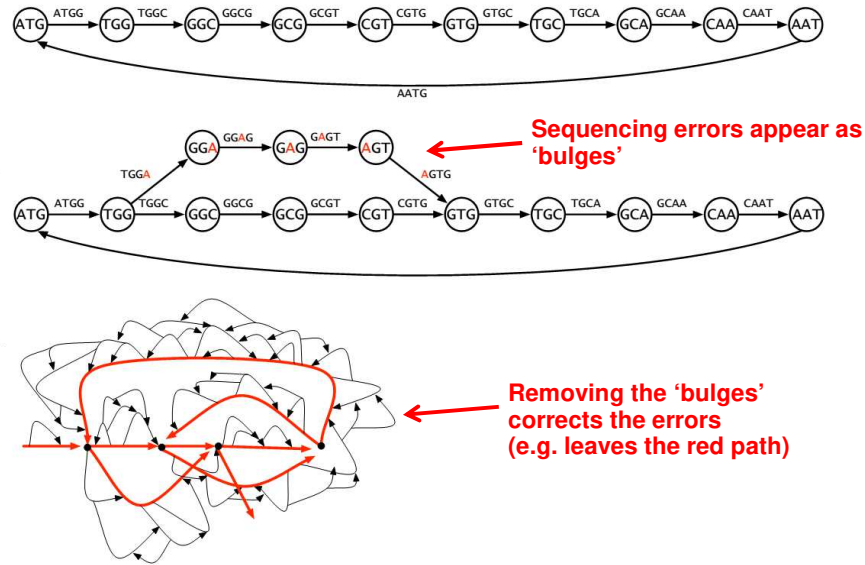


Euler conceptualized it as a graph:
Nodes = parts of city
Edges = bridges

Nature Biotech 29(11):987-991 (2011)

44

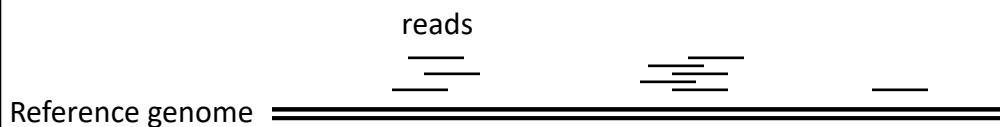
DeBruijn graph assemblers tend to have nice properties, e.g. correcting sequencing errors & handling repeats better



Nature Biotech 29(11):987-991 (2011)

45

Once a reference genome is assembled, new sequencing data can 'simply' be mapped to the reference.



46

Mapping reads to assembled genomes

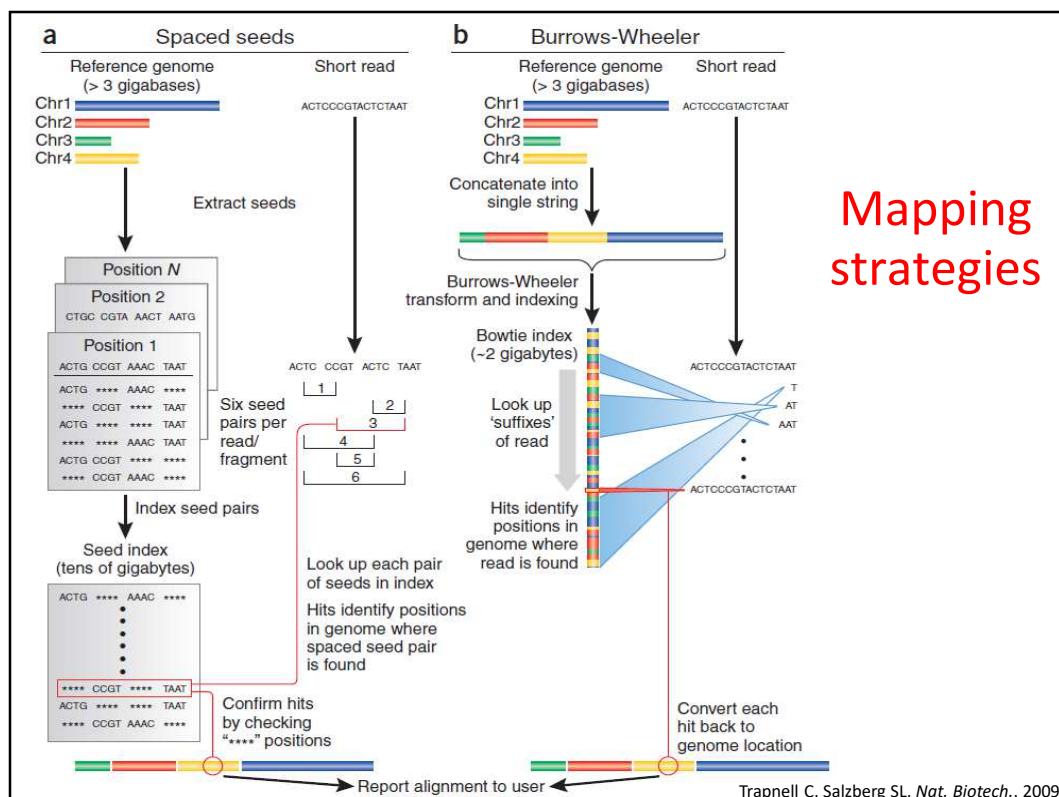
Table 1 A selection of short-read analysis software

Program	Website	Open source?	Handles ABI color space?	Maximum read length
Bowtie	http://bowtie.cbcb.umd.edu	Yes	No	None
BWA	http://maq.sourceforge.net/bwa-man.shtml	Yes	Yes	None
Maq	http://maq.sourceforge.net	Yes	Yes	127
Mosaik	http://bioinformatics.bc.edu/marthlab/Mosaik	No	Yes	None
Novoalign	http://www.novocraft.com	No	No	None
SOAP2	http://soap.genomics.org.cn	No	No	60
ZOOM	http://www.bioinfor.com	No	Yes	240

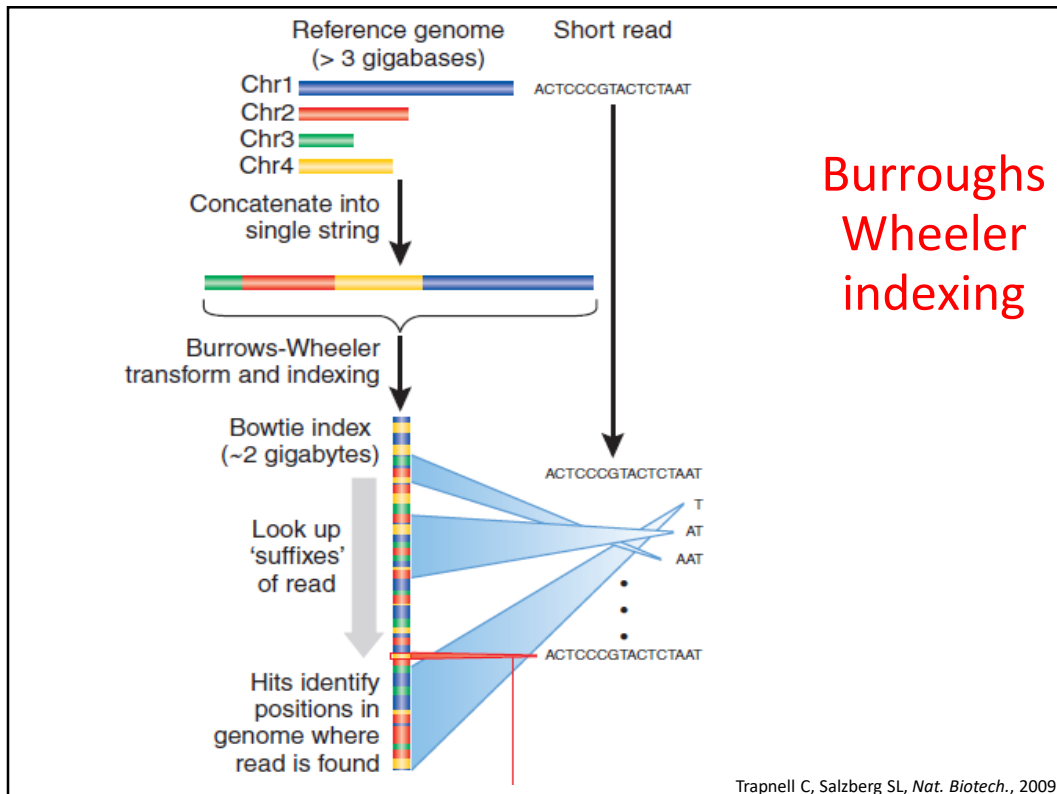
The list is a little longer now! e.g. see https://en.wikipedia.org/wiki/List_of_sequence_alignment_software#Short-Read_Sequence_Alignment

Trapnell C, Salzberg SL, *Nat. Biotech.*, 2009

47



48



49

Burroughs-Wheeler transform indexing

BWT is often used for file compression (like bzip2), here used to make a fast 'lookup' index in a genome

BWT = 'reversible block-sorting'

Input	SIX.MIXED.PIXIES.SIFT.SIXTY.PIXIE.DUST.BOXES	
	↓ Forward BWT	This sequence is more compressible
Output	TEXYDST.E.IXIXIXSSMPPS.B..E.S.EUSFXDIIOIIT	
	↓ Reverse BWT	
Recovered input	SIX.MIXED.PIXIES.SIFT.SIXTY.PIXIE.DUST.BOXES	

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

50

Burroughs-Wheeler transform indexing

Input

^BANANA |

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

51

Burroughs-Wheeler transform indexing

**All
Rotations**

^BANANA |
| ^BANANA
A | ^BANAN
NA | ^BANA
ANA | ^BAN
NANA | ^BA
ANANA | ^B
BANANA | ^

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

52

Burroughs-Wheeler transform indexing

Sorting All Rows in Alphabetical Order

ANANA | ^B
ANA | ^BAN
A | ^BANAN
BANANA | ^
NANA | ^BA
NA | ^BANA
^BANANA |
| ^BANANA

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

53

Burroughs-Wheeler transform indexing

Taking Last Column

ANANA | ^**B**
ANA | ^BA**N**
A | ^BANAN**N**
BANANA | ^
NANA | ^BA**A**
NA | ^BANA**A**
^BANANA |
| ^BANANA**A**

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

54

Burroughs-Wheeler transform indexing

Output Last Column
BNN^AA A

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

55

Burroughs-Wheeler transform indexing

Transformation				
Input	All Rotations	Sorting All Rows in Alphabetical Order	Taking Last Column	Output Last Column
<div> ^BANANA </div>	<div> ^BANANA ^BANANA A ^BANAN NA ^BANA ANA ^BAN NANA ^BA ANANA ^B BANANA ^ </div>	<div> ANANA ^B ANA ^BAN A ^BANAN BANANA ^ NANA ^BA NA ^BANA ^BANANA ^BANANA </div>	<div> ANANA ^B ANA ^BAN A ^BANAN BANANA ^ NANA ^BA NA ^BANA ^BANANA ^BANANAA </div>	<div> BNN^AA A </div>

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

56

BWT is remarkable because it is
reversible.

Any ideas as how you might reverse it?

57

Burroughs-Wheeler transform indexing

Input
BNN^AA A

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

58

Burroughs-Wheeler transform indexing

Add 1	Sort 1	Add 2	Sort 2
B N N ^ A A A	A A A B N N ^ 	BA NA NA ^B AN AN ^ A	AN AN A BA NA NA ^B ^
Write the sequence as the last column	Sort it...	Add the columns...	Sort those...

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

59

Burroughs-Wheeler transform indexing

Add 3	Sort 3	Add 4	Sort 4
BAN NAN NA ^BA ANA ANA ^B A ^	ANA ANA A ^ BAN NAN NA ^BA ^B	BANA NANA NA ^ ^BAN ANAN ANA ^BA A ^B	ANAN ANA A ^B BANA NANA NA ^ ^BAN ^BA
Add the columns...	Sort those...	Add the columns...	Sort those...

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

60

Burroughs-Wheeler transform indexing

Add 5	Sort 5	Add 6	Sort 6
BANAN NANA NA ^B ^BANAN ANANA ANA ^ ^BAN A ^BA	ANANA ANA ^ A ^BA BANAN NANA NA ^B ^BANAN ^BAN	BANANA NANA ^ NA ^BA ^BANAN ANANA ANA ^B ^BANAN A ^BAN	ANANA ANA ^B A ^BAN BANANA NANA ^ NA ^BA ^BANAN ^BANAN
Add the columns...	Sort those...	Add the columns...	Sort those...

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

61

Burroughs-Wheeler transform indexing

Add 7	Sort 7	Add 8
BANANA NANA ^B NA ^BAN ^BANANA ANANA ^ ANA ^BA ^BANAN A ^BANA	ANANA ^ ANA ^BA A ^BANA BANANA NANA ^B NA ^BAN ^BANANA ^BANAN	BANANA ^ NANA ^BA NA ^BANA ^BANANA ANANA ^B ANA ^BAN ^BANANA A ^BANAN
Add the columns...	Sort those...	Add the columns...

The row with the "end of file" character at the end is the original text

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

62

Burroughs-Wheeler transform indexing

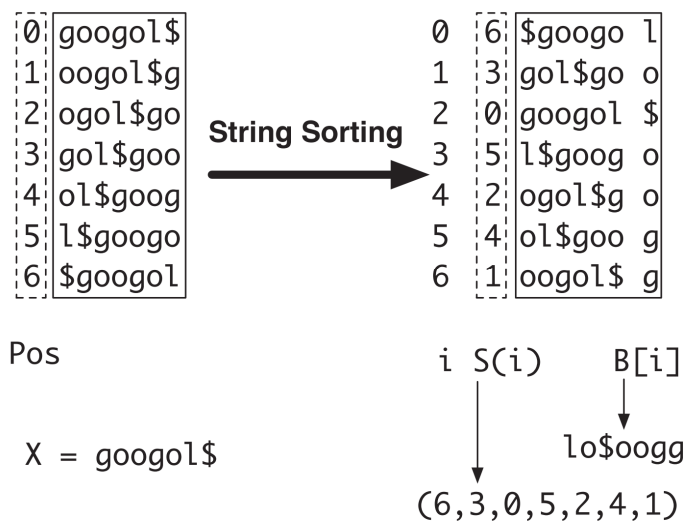
Output
^BANANA

The row with the "end of file" character at the end is the original text

http://en.wikipedia.org/wiki/Burrows-Wheeler_transform

63

The Burroughs-Wheeler transform leads naturally to a suffix array



Li & Durbin, doi:10.1093bioinformatics/btp324/

64

Li & Durbin, doi:10.1093/bioinformatics/btp324/



Burroughs Wheeler indexing

Trapnell C, Salzberg SL, *Nat. Biotech.*, 2009