

Assembling Genomes

BIO337 Systems Biology / Bioinformatics – Spring 2014

Edward Marcotte, Univ of Texas at Austin

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014



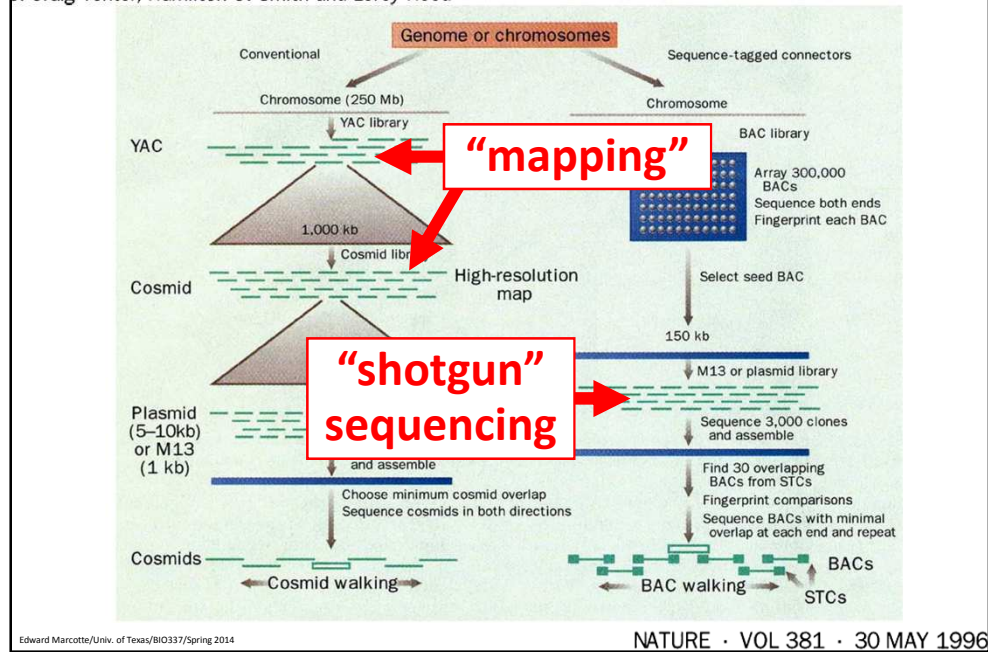
<http://www.triazzle.com>; The image from http://www.dangilbert.com/port_fun.html

Reference: Jones NC, Pevzner PA, Introduction to Bioinformatics Algorithms, MIT press

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

A new strategy for genome sequencing

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



(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)

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

NATURE · VOL 381 · 30 MAY 1996

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?

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Thinking about the basic shotgun concept

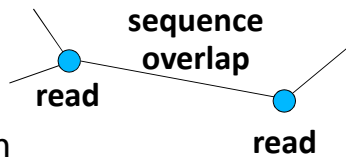
- At a high level, the first genomes were sequenced by comparing pairs of reads to find overlapping reads
- Then, building a graph (*i.e.*, a network) to represent those relationships
- The genome sequence is a “walk” across that graph

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

The “Overlap-Layout-Consensus” method

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

Layout: Construct a graph describing the overlaps

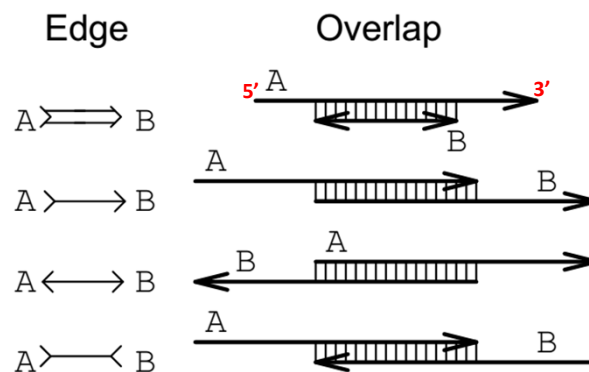


Simplify the graph
Find the simplest path through the graph

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Building an overlap graph

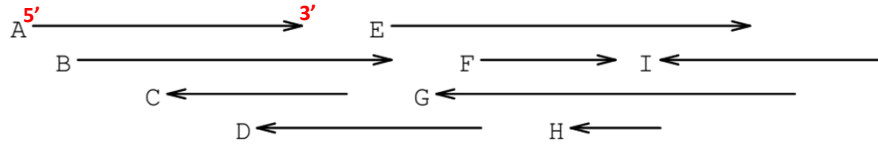


EUGENE W. MYERS. *Journal of Computational Biology*. Summer 1995, 2(2): 275-290

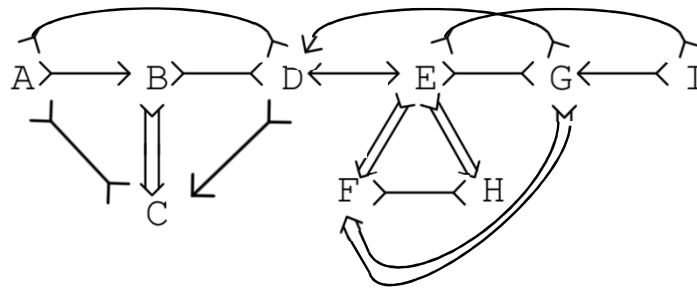
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Building an overlap graph

Reads



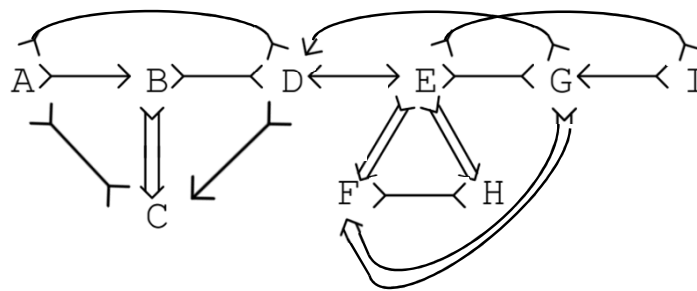
Overlap graph



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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Simplifying an overlap graph

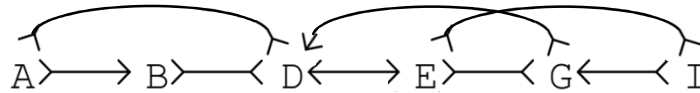


1. Remove all contained nodes & edges going to them

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Simplifying an overlap graph



2. Transitive edge removal:

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Simplifying an overlap graph



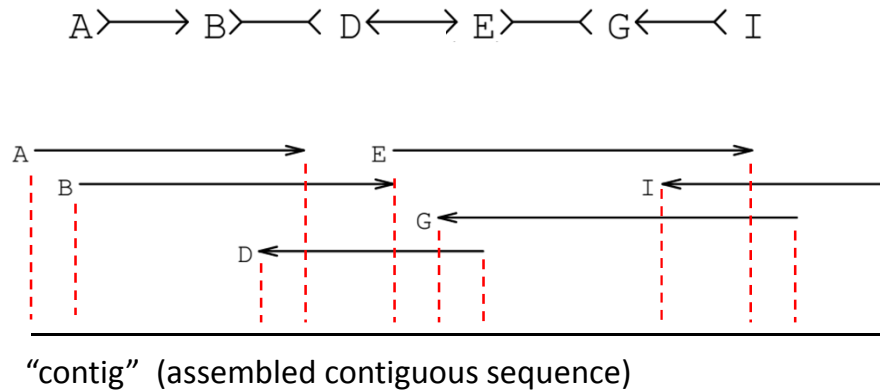
3. If un-branched, calculate consensus sequence

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

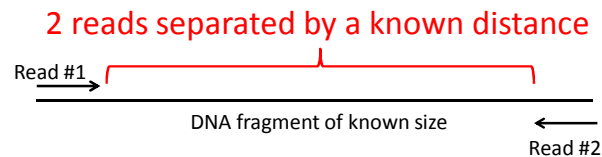
Simplifying an overlap graph



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

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



Contigs can be ordered using these paired reads



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

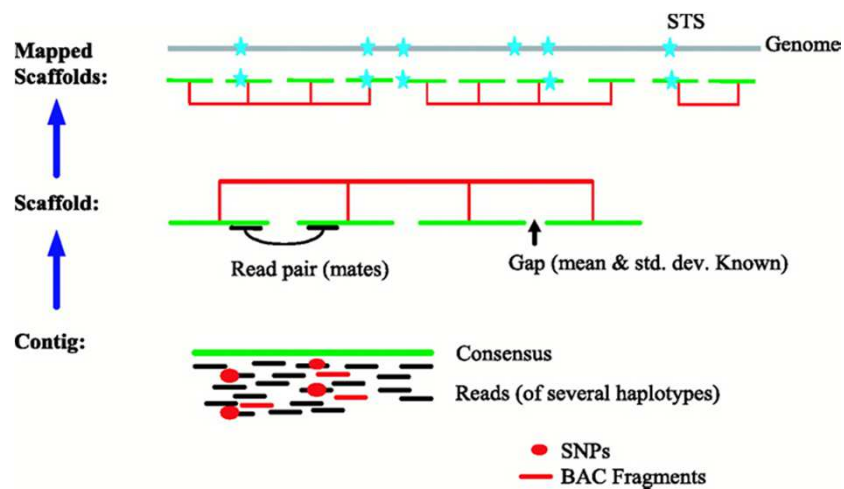


Jim Kent

David Haussler

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Whole genome Assembly: big picture



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

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
```

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

RepBase + RepeatMasker

```
taejoon@fourierseq:~/RepBase/RepBase15.05_fastas$ ls -la
total 12
drwxr-xr-x 2 taejoon taejoon 4096 Jun 10 14:54 .
drwxr-xr-x 1 taejoon taejoon 4096 Jun 10 14:54 ..
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 diarep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 mamsub.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 rodsub.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 simple.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 dtorep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 msousub.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 spurep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 fngrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 nemrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 synrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 appendix
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 fngrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 ooryep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 tmlanrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 btctrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 plnrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 tmpnemrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 cbrrrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 prirep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 tmpxenrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 celrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 prisub.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 version
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 chlrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 pseuod.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 vrtrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 cinrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 ratsub.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 zebrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 cinunc.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 mamrep.ref
-rw-r--r-- 1 taejoon taejoon  100 Jun 10 14:54 rodrep.ref
```

```

>SMER512 ERV1 Homo sapiens
tgaggcaggagaaatacagcaggaggaattggaagttggataaaggggagaatgagtaaaagcangagacga
gaagcaggctaaagagcgggtcgtagcaagacagagataaagaacagagaattgagcagccaaaaaaaag
taagatanaaaagaaagtgtagtaagagcccaatcgctggcctagcagaccaaaacagtaagaaggcgac
ctctcagagatgggctatgcatactagagagaaaagatctcttaaaatggcccgctagatgataatcaga
ctaataagctcatgatatggactcatgatctgatcattgtaattctaaaatttggtgtggaggtgagcgcga
agagtccaacagcacaaggggctcatgatttaagtaactgaacacccactcatcaatcaaaagacgga
ctctcgctagagataggacgctctgggaagagaagaaaaaaaacacataaaaagacccaagaatgcac
caactcgctgactgatactctttcgccagctgacccactctctctctcttgagagtgtaattctgct
taataaaacttttgcgtcttctctatctgtgtgtgtgtgtgtgtgtgtgtgtgtgtgtgtgtgtgtgtgt
ggaaactgcagcgcacaktcgtgaaca

>MIRB SINE2/tRNA Mammalia
cagagggcagctgggtgagctggaagaagacgcgggcttggagtcaggcagacctgggttgcgaatcctg
gtctctgaactctactagctgctgactctggggacagctcaactcaatctctgagctcagtttctctcat
ttgtaaaatgggataataatactgactcgcagggttggtagagataaagatagataatgcataatgcataaa
gcgctgtagcacgctctggccacagatgaagcctctaataatggtagtctattatt

>LTR45 ERV1 Homo sapiens
tgtaacgcgggagaccgccaactgggctctctgtgtgatacaaaagtctcaagttctgttggtta
ttacacagcgggcaacagtcagatgtatgaccgggctatgctgataagaaaagctttgactcttaacaa
caccagacaacaaatgattctctctctcggaacacagaagaagcggggactgaccggaaactgaaatgcgga
actcttcagaagacaaggggctcggttgcgggaagatctgggttaaaacttgcttcaactcaatctata
ccgtaattgtgtcaaaattgaaagctctcaatcagacacctgcagacacattctctaaactctctctct
gctctgtatcctctaaaactctgccacagccaaaactcgggagacagatgtgagccacactctgtctg
ctctgtgcgggttttgcaataaagctctttctctctaaaagctgtgtgcatagttgattgctgtgtgtgt
gtgtgtagcgcagcaagccatttgcctgataaca

>MER80B hAT Homo sapiens
cagggctctttcaacagaggtctatgtgggcttcaggaggtctgtgaacctctgaaattatatacaaa
aaatgtgtgtatgtgcataatgattttcttgggagaggggttcatagttttcacagattctcaaa
aagggctatgatctmaaaagaaataaagacccgtg

```

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

GigAssembler: Build merged sequence contigs (“rafts”)

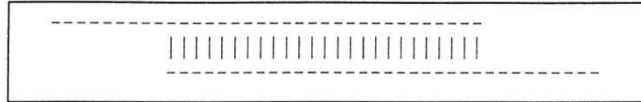
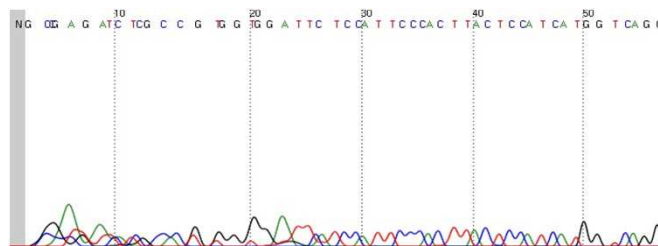


Figure 1 Two sequences overlapping end to end. The sequences are represented as dashes. The aligning regions are joined by vertical bars. End-to-end overlap is an extremely strong indication that two sequences should be joined into a contig.

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Sequencing quality (Phred Score)



>gnl|tij2299297598 name:fw01a01.x1 NCBI Accession: [AC243936](#) Mate pair: [2299297599](#)

Quality score:	not available	>0 - <20	>=20 - <40	>=40 - <60	>=60 - <80	>=80 - <100
0	3	6	6	6	6	8
1	6	6	6	6	6	8
2	6	6	6	6	6	8
3	6	6	6	6	6	8
4	6	6	6	6	6	8
5	6	6	6	6	6	8
6	6	6	6	6	6	8
7	6	6	6	6	6	8
8	6	6	6	6	6	8
9	6	6	6	6	6	8
10	6	6	6	6	6	8
11	6	6	6	6	6	8
12	6	6	6	6	6	8
13	6	6	6	6	6	8
14	6	6	6	6	6	8
15	6	6	6	6	6	8
16	6	6	6	6	6	8
17	6	6	6	6	6	8
18	6	6	6	6	6	8
19	6	6	6	6	6	8
20	6	6	6	6	6	8
21	6	6	6	6	6	8
22	6	6	6	6	6	8
23	6	6	6	6	6	8
24	6	6	6	6	6	8
25	6	6	6	6	6	8
26	6	6	6	6	6	8
27	6	6	6	6	6	8
28	6	6	6	6	6	8
29	6	6	6	6	6	8
30	6	6	6	6	6	8
31	6	6	6	6	6	8
32	6	6	6	6	6	8
33	6	6	6	6	6	8
34	6	6	6	6	6	8
35	6	6	6	6	6	8
36	6	6	6	6	6	8
37	6	6	6	6	6	8
38	6	6	6	6	6	8
39	6	6	6	6	6	8
40	6	6	6	6	6	8
41	6	6	6	6	6	8
42	6	6	6	6	6	8
43	6	6	6	6	6	8
44	6	6	6	6	6	8
45	6	6	6	6	6	8
46	6	6	6	6	6	8
47	6	6	6	6	6	8
48	6	6	6	6	6	8
49	6	6	6	6	6	8
50	6	6	6	6	6	8
51	6	6	6	6	6	8
52	6	6	6	6	6	8
53	6	6	6	6	6	8
54	6	6	6	6	6	8
55	6	6	6	6	6	8
56	6	6	6	6	6	8
57	6	6	6	6	6	8
58	6	6	6	6	6	8
59	6	6	6	6	6	8
60	6	6	6	6	6	8
61	6	6	6	6	6	8
62	6	6	6	6	6	8
63	6	6	6	6	6	8
64	6	6	6	6	6	8
65	6	6	6	6	6	8
66	6	6	6	6	6	8
67	6	6	6	6	6	8
68	6	6	6	6	6	8
69	6	6	6	6	6	8
70	6	6	6	6	6	8
71	6	6	6	6	6	8
72	6	6	6	6	6	8
73	6	6	6	6	6	8
74	6	6	6	6	6	8
75	6	6	6	6	6	8
76	6	6	6	6	6	8
77	6	6	6	6	6	8
78	6	6	6	6	6	8
79	6	6	6	6	6	8
80	6	6	6	6	6	8
81	6	6	6	6	6	8
82	6	6	6	6	6	8
83	6	6	6	6	6	8
84	6	6	6	6	6	8
85	6	6	6	6	6	8
86	6	6	6	6	6	8
87	6	6	6	6	6	8
88	6	6	6	6	6	8
89	6	6	6	6	6	8
90	6	6	6	6	6	8
91	6	6	6	6	6	8
92	6	6	6	6	6	8
93	6	6	6	6	6	8
94	6	6	6	6	6	8
95	6	6	6	6	6	8
96	6	6	6	6	6	8
97	6	6	6	6	6	8
98	6	6	6	6	6	8
99	6	6	6	6	6	8
100	6	6	6	6	6	8

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Sequencing quality (Phred Score)

$$Q = -10 \log_{10} P \leftarrow \begin{array}{l} \text{Base-calling} \\ \text{Error} \\ \text{Probability} \end{array}$$

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 %

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

GigAssembler: Build merged sequence contigs ("rafts")

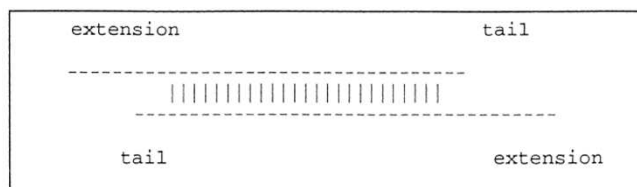


Figure 2 Two sequences with tails. The nonaligning regions on either side can be classified into 'extensions' and 'tails.' Short tails are fairly common even when two sequences should be joined into a contig because of poor quality sequence near the ends and occasional chimeric reads. Long tails, however,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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

GigAssembler: Build merged sequence contigs (“rafts”)

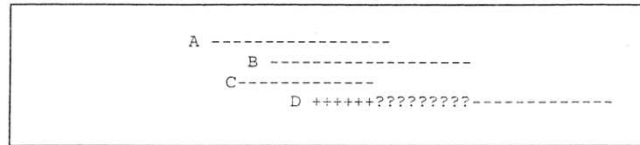


Figure 3 Merging into a raft. A contig (“raft”) of three sequences: A, B, and C has already been constructed by GigAssembler. The program now examines an alignment between sequence C and a new sequence, D, to see whether D should also be added to the raft. The parts of D marked with +s are compatible 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.

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

GigAssembler: Build sequenced clone contigs (“barges”)

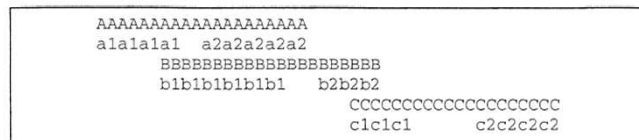


Figure 4 Three overlapping draft clones: A, B, and C. Each clone has two initial sequence contigs. Note that initial sequence contigs a1, b1, and a2 overlap as do b2 and c1.

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

GigAssembler: Build a “raft-ordering” graph

```

AAAAAAAAAAAAAAAAAAAA
a1a1a1a1  a2a2a2a2a2
BBBBBBBBBBBBBBBBBBBB
b1b1b1b1b1  b2b2b2
CCCCCCCCCCCCCCCCCCCC
c1c1c1c1  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.

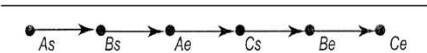


Figure 5 Ordering graph of clone starts and ends. This represents the same clones as in Fig. 4. (As) The start of clone A; (Ae) the end of clone A. Similarly Bs, Be, Cs, and Ce represent the starts and ends of clones B and C.

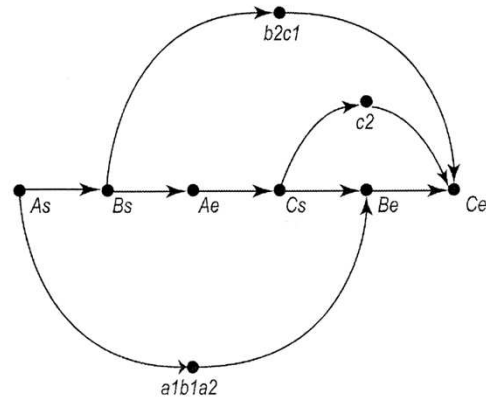


Figure 6 Ordering graph after adding in rafts. The initial sequence contigs shown in Fig. 4 are merged into rafts where they overlap. This forms three rafts: a1b1a2, b2c1, and c2. These rafts are constrained to lie between the relevant clone ends by the addition of additional ordering edges to the graph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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

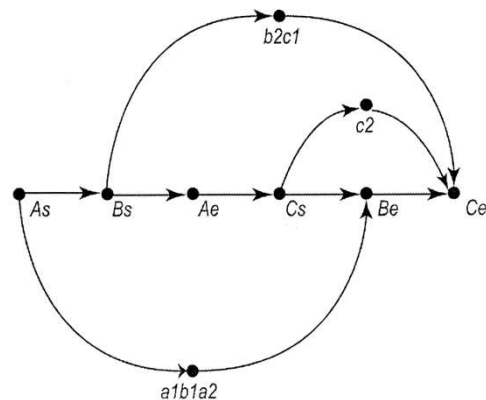


Figure 6 Ordering graph after adding in rafts. The initial sequence contigs shown in Fig. 4 are merged into rafts where they overlap. This forms three rafts: a1b1a2, b2c1, and c2. These rafts are constrained to lie between the relevant clone ends by the addition of additional ordering edges to the graph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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

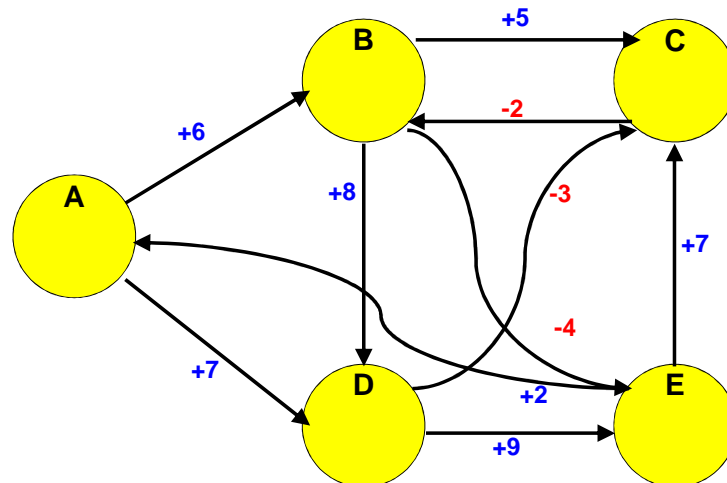
Finding the shortest path across the ordering graph using the Bellman-Ford algorithm

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Find the shortest path to all nodes.

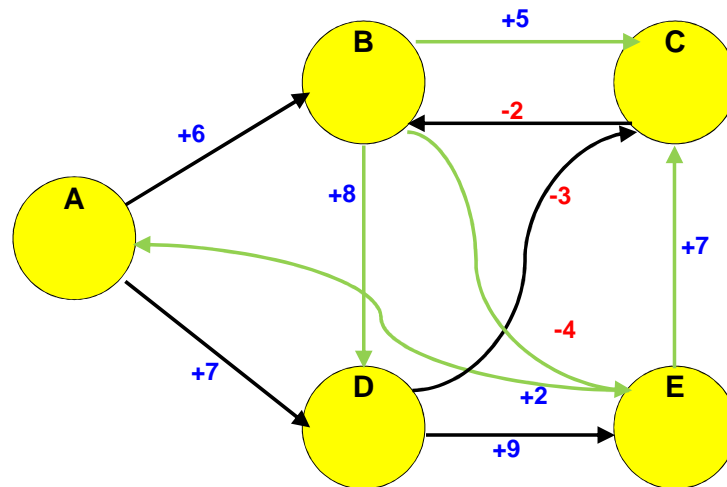
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Find the shortest path to all nodes.

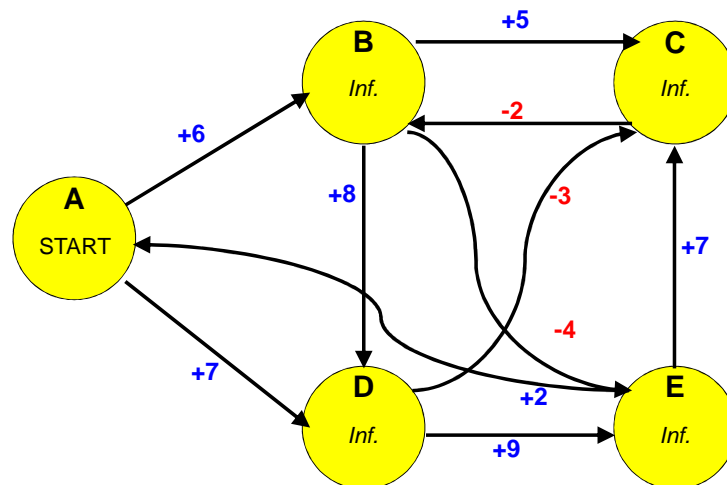
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Find the shortest path to all nodes.

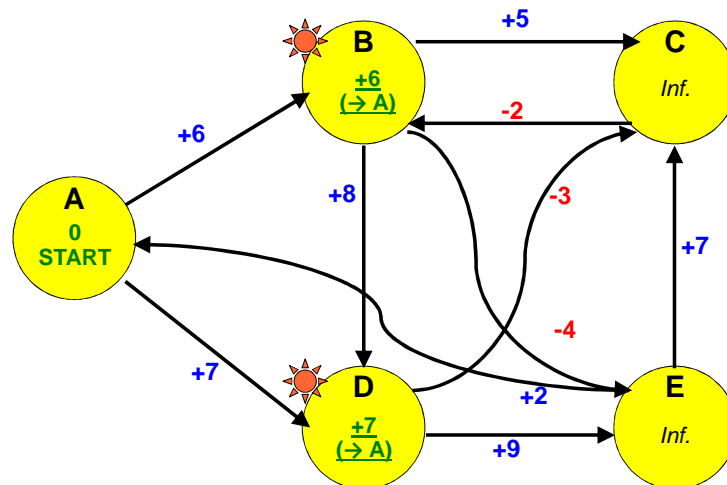
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Find the shortest path to all nodes.

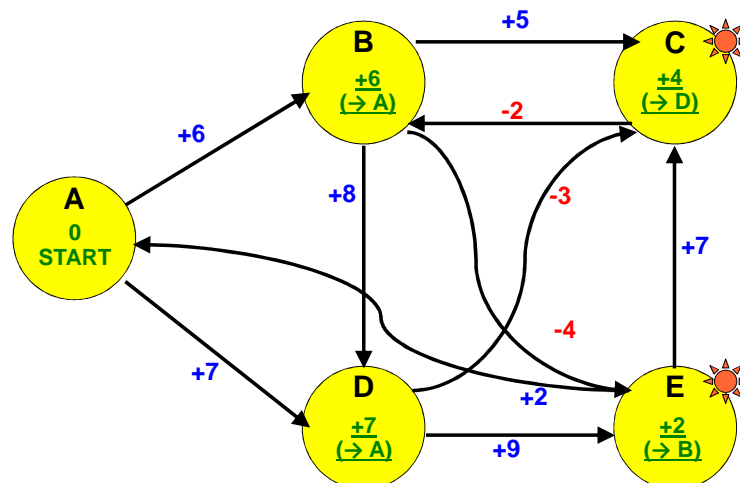
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Find the shortest path to all nodes.

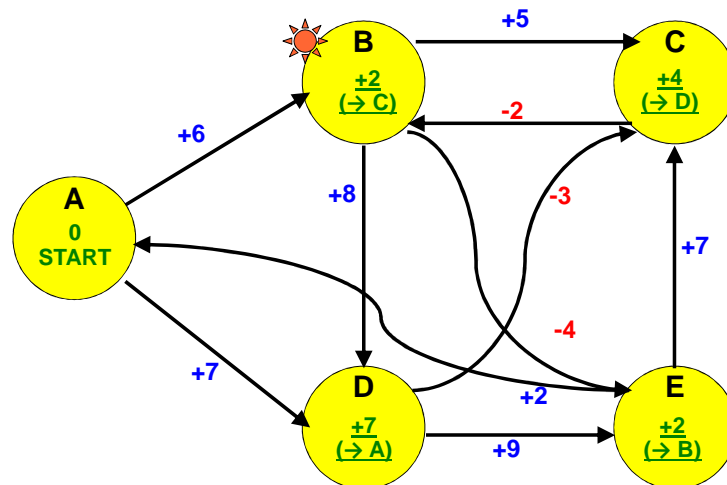
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Find the shortest path to all nodes.

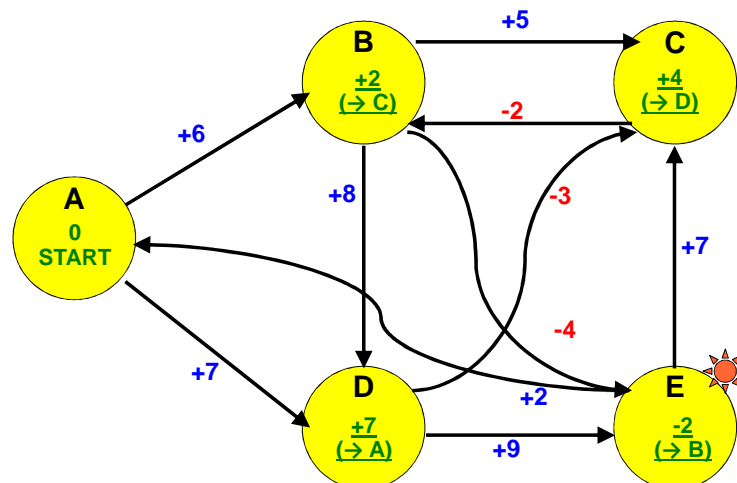
Take every edge and try to relax it ($N - 1$ times where N is the count of nodes)



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

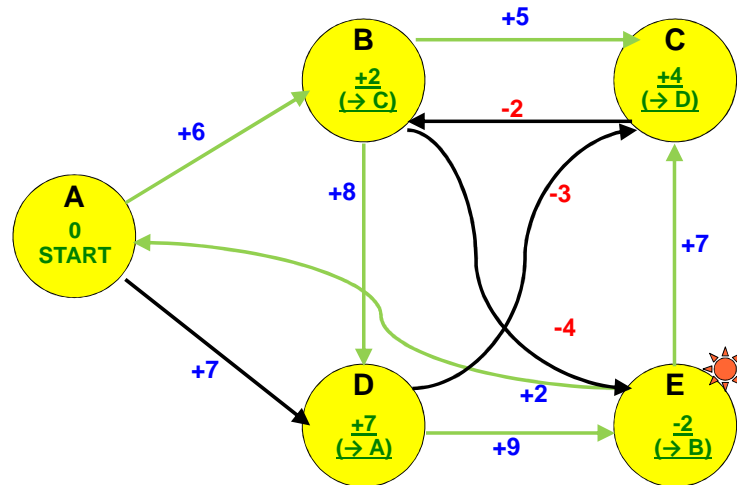
Find the shortest path to all nodes.

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



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

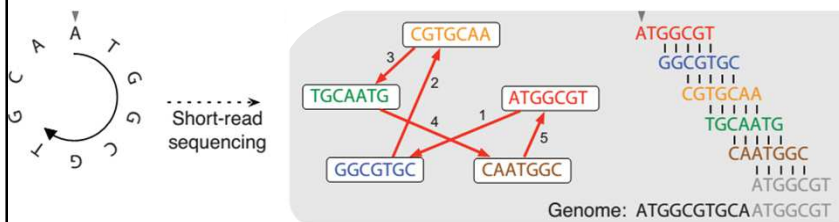
Answer: A-D-C-B-E



Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Here's what we saw before:



In Overlap-Layout-Consensus:

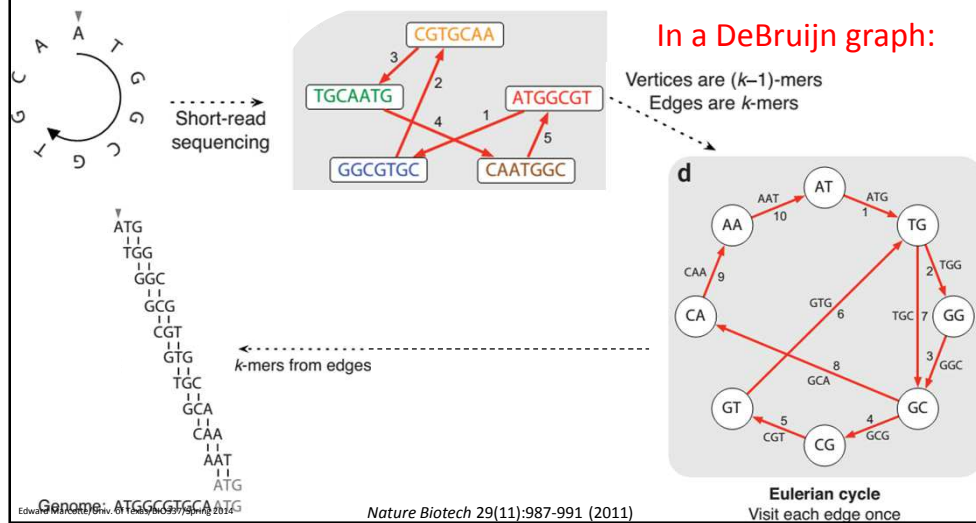
Nodes are reads

Edges are overlaps

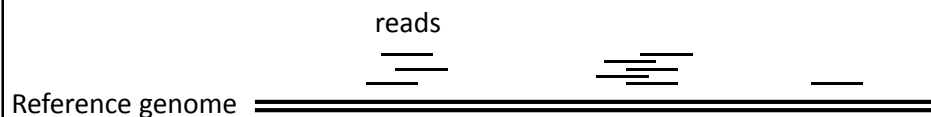
Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

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



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



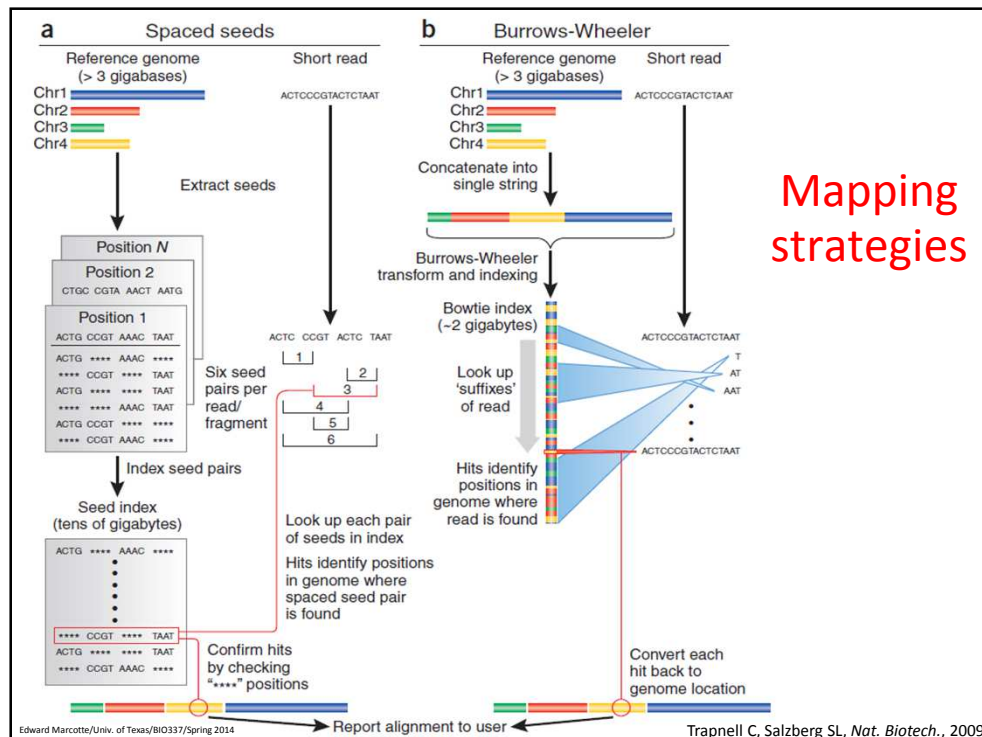
Mapping reads to assembled genomes

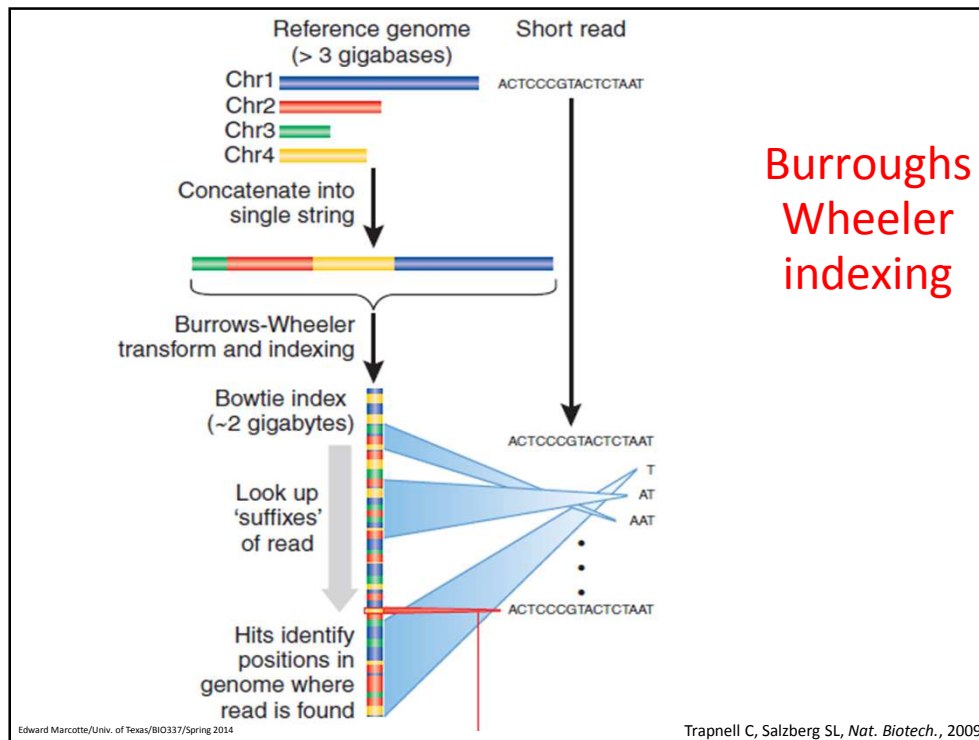
Table 1 A selection of short-read analysis software

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.bioinform.com	No	Yes	240

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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





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.IXIXXSSMPPS.B..E.S.EUSFXDIIIOIIT

Reverse BWT

Recovered input SIX.MIXED.PIXIES.SIFT.SIXTY.PIXIE.DUST.BOXES

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Burroughs-Wheeler transform indexing

Input

^BANANA |

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Burroughs-Wheeler transform indexing

**All
Rotations**

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Burroughs-Wheeler transform indexing

Sorting All Rows in Alphabetical Order

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Burroughs-Wheeler transform indexing

Taking Last Column

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

**BWT is remarkable because it is
reversible.**

Any ideas as how you might reverse it?

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

Burroughs-Wheeler transform indexing

Input
BNN^AA A

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

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...

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

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...

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

Burroughs-Wheeler transform indexing

Output
^BANANA

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

Edward Marcotte/Univ. of Texas/BIO337/Spring 2014

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

