

# Assembling Genomes

**BCH364C/391L Systems Biology / Bioinformatics – Spring 2015**

**Edward Marcotte, Univ of Texas at Austin**

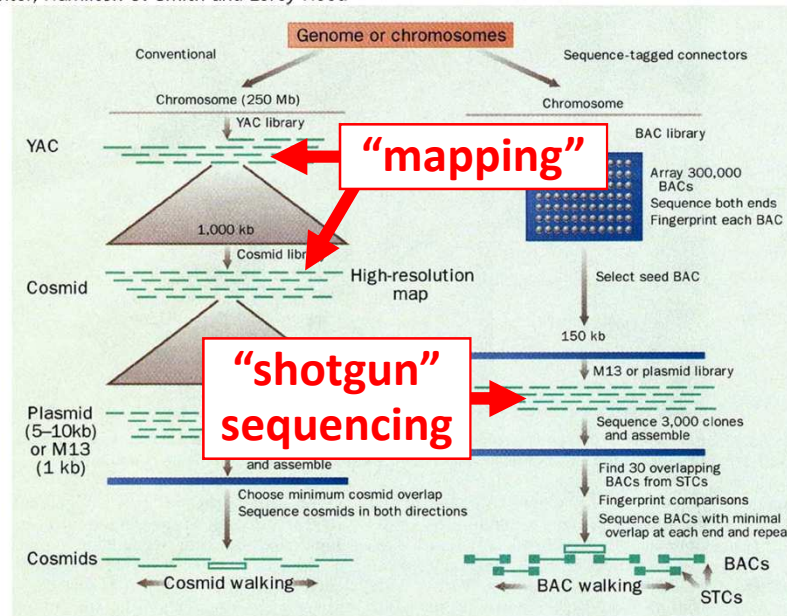
Edward Marcotte/Univ. of Texas/BCH364C-391L/Spring 2015



<http://www.triazzle.com>; The image from [http://www.dangilbert.com/port\\_fun.html](http://www.dangilbert.com/port_fun.html)  
Reference: Jones NC, Pevzner PA, Introduction to Bioinformatics Algorithms, MIT press

# A new strategy for genome sequencing

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



NATURE · VOL 381 · 30 MAY 1996

## (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

### **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?

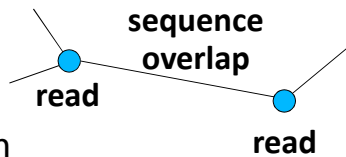
### **Thinking about the basic shotgun concept**

- 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

## The “Overlap-Layout-Consensus” method

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

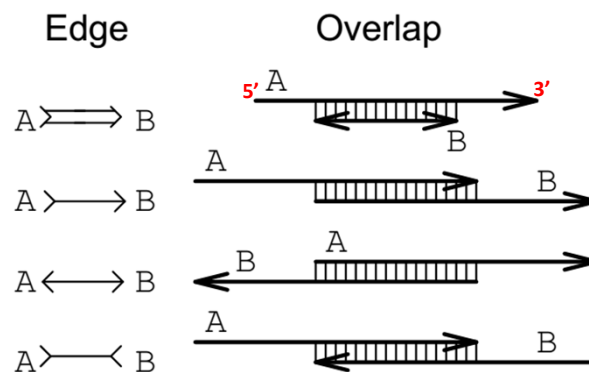
**Layout:** Construct a graph describing the overlaps



Simplify the graph  
Find the simplest path through the graph

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

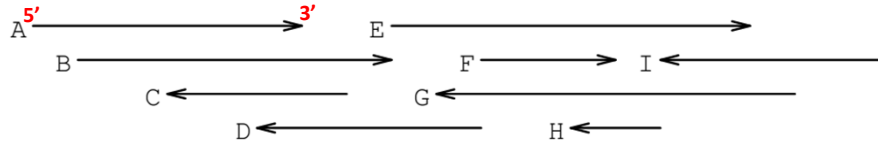
## Building an overlap graph



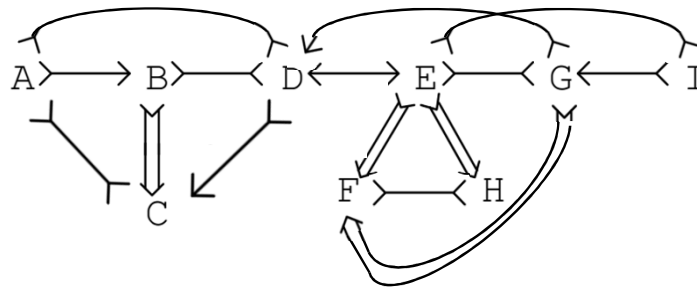
EUGENE W. MYERS. *Journal of Computational Biology*. Summer 1995, 2(2): 275-290

## Building an overlap graph

### Reads

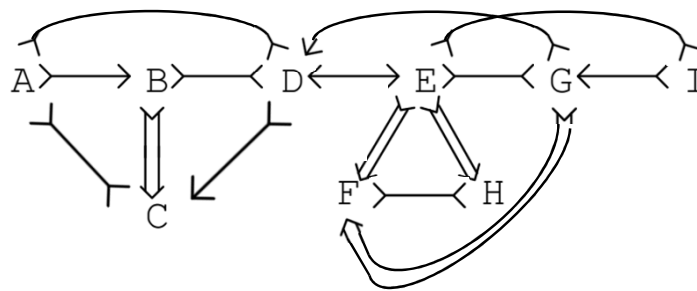


### Overlap graph



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

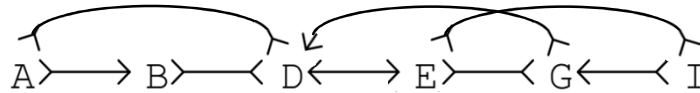
## Simplifying an overlap graph



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)

## Simplifying an overlap graph



### 2. Transitive edge removal:

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

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

## Simplifying an overlap graph

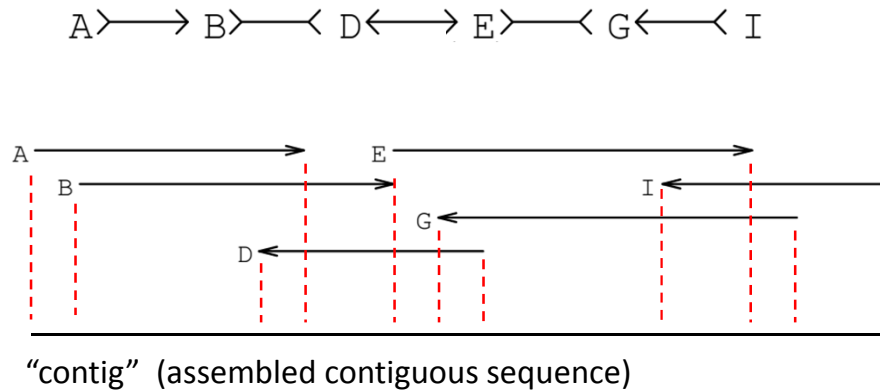


### 3. If un-branched, calculate consensus sequence

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

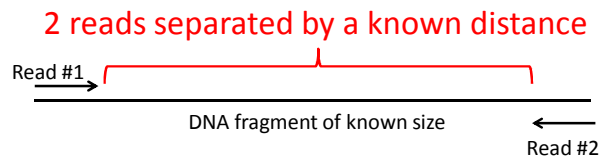
EUGENE W. MYERS. *Journal of Computational Biology*. Summer 1995, 2(2): 275-290 (more or less)

## Simplifying an overlap graph



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



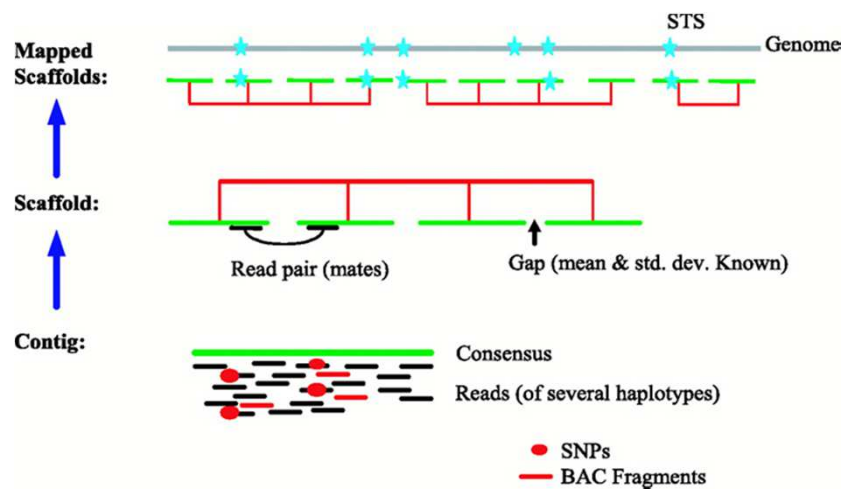
## GigAssembler (used to assemble the public human genome project sequence)



Jim Kent

David Haussler

## Whole genome Assembly: big picture



<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.

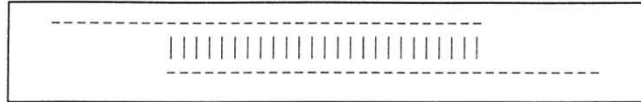
```
chr1/  
chr1/contig1.e  
chr1/contig1.a  
chr1/contig1.c  
chr1/contig1.b  
chr1/contig1.d  
chr3/  
chr2/  
chr2/contig2.d  
chr2/contig2.b  
chr2/contig2.a  
chr2/contig2.c
```

# RepBase + RepeatMasker

```
tajeon@fourierseq:~/RepBase/RepBase15.05.fasta$ ls -la  
.  
..  
angrep.ref  
appendix  
athrep.ref  
bctrep.ref  
cbrrep.ref  
celrep.ref  
chlrep.ref  
cinunc.ref  
dcotrep.ref  
drorrep.ref  
fngrep.ref  
grasrep.ref  
humrep.ref  
humsub.ref  
invrep.ref  
mamrep.ref  
mamsub.ref  
mousub.ref  
nemrep.ref  
nryrep.ref  
plnrep.ref  
prirep.ref  
prsub.ref  
pseudoref  
ratsub.ref  
rodsub.ref  
rodrrep.ref  
simple.ref  
spurep.ref  
synrep.ref  
tmplanrep.ref  
tmpnemrep.ref  
tmpxenrep.ref  
version  
vrtrep.ref  
zebrep.ref
```

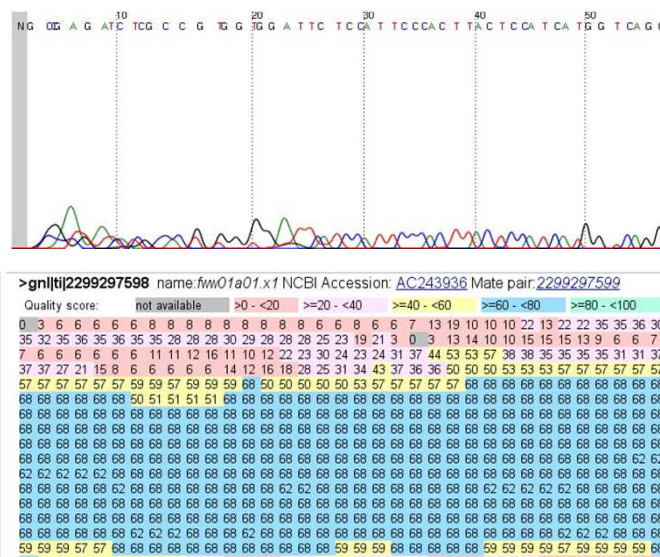
```
>MER51D ERV1 Homo sapiens  
tgaggcaggagaaaatagcagagggaattggaattggataaaggagaaatgaagtaaaagcangagagca  
gaagcaaggtaaaagagcgggtgagcaagaagcaagataagaagcagaagttagcagcagcaaaacaaaag  
taagatnanaaagaagtgagtaaggagccacatggctggctagatccagaccaaccagtaaggggcag  
ctcctcagagatggcagtgatgacattagagagaaaaatcttcttaaaatgacccgtatgataatcagct  
cattaaagctcatgcatatggactgcataatcatgcatgtacttaaaattatgggatggaggtgacgcgca  
agawgtcacagcacacagggccatagkattaaagtaactaagcaaccacccatcaatcaaaaagcaga  
tgctggctagagattaggcagccttgggaagagaagaaaaaacacataaaagacccaaagacacac  
caaaactgcgctgactctcttgcagaggtcagccactctccctctctgagaggtgaataactgtgt  
taataaaacttttgcgtcttgcctatctgtgtgtcttgccttccaattcttcttgggacaccaagagcct  
ggaaactgcacrgcaccactgtgtaaca  
>MIRb SINE2/trna Mammalia  
cagaggggacgctgggtgagtgaaagagacagggccttggagtcaggcagacgtgggttcgaatcctg  
gctctgccacttactagctgtgtgaccttgggcaagtcacttaacctctctgagcctcagtttctctatc  
tgtaaaatggggataataatcactcctgcaggggtgtgtgaggattaaatgagataatgcatgttaa  
gcgcttagcacagtcctggcacacagtaagcgtcaataaatggtagctctattatt  
>LTR45 ERV1 Homo sapiens  
tgtaaccgaggacagcccaactggcctactctgttgataacaaaatgtcaagttaccttgtagtta  
taacagagcccaaaactgcagtcagtgtagccggcagtgcaatagaaaaagcttgcaccttaacaa  
caccagaaaccaatgattctccctcggaaccaagaagacgggacatgacgggaacctgaatgccgga  
actctttcagagacaaaggggtcgttggccgggaagatctggggctaaaaatctgcctcaacatacctta  
ccgtaaatgttcaaatgtgaagcctccaaatcagacctgccaagccaactcctaaatcctttccctt  
gccctctgaccccttaaaacttgcacagcccaaatcggggagacagattgagccacacctctgtct  
cctgtggtggcgggtttgcaataaagcctttctttctcaaaagctgggtgacatagttatgtgctctgt  
gtgcatcaggcagcaagccatttgcgtgataaca  
>MER80B HAT Homo sapiens  
cagggcttcttaaccagaggttccatggagggcttcaggaggtctgtgaacctctgaaattatatacaa  
aaatgttgtatagtgcataatgtattttctggggagaggggttcagctttcatcagattctcaaa  
aggggtctatgactmaaaagggttaagaagccctg
```

# 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.

## Sequencing quality (Phred Score)



## Sequencing quality (Phred Score)

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

or

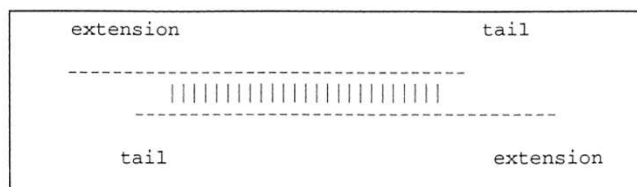
$$P = 10^{\frac{-Q}{10}}$$

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 %

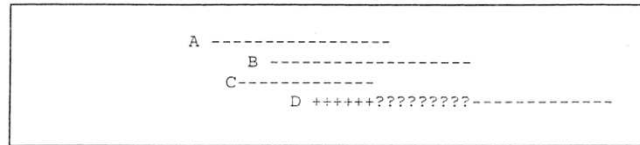
[http://en.wikipedia.org/wiki/Phred\\_quality\\_score](http://en.wikipedia.org/wiki/Phred_quality_score)

## GigAssembler: Build merged sequence contigs (“rafts”)



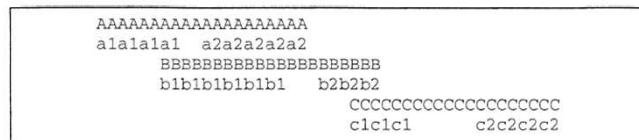
**Figure 2** Two sequences with tails. The nonaligning regions on either side can be classified into ‘extensions’ and ‘tails.’ Short tails are fairly common even when two sequences should be joined into a contig because of poor quality sequence near the ends and occasional chimeric reads. Long tails,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

## 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.

## 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.

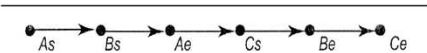
# GigAssembler: Build a “raft-ordering” graph

```

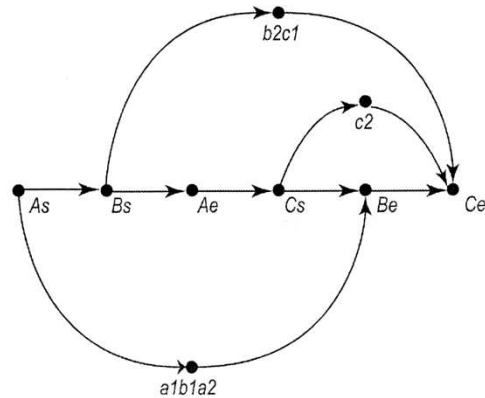
AAAAAAAAAAAAAAAAAAAA
a1a1a1a1 a2a2a2a2a2
BBBBBBBBBBBBBBBBBBBB
b1b1b1b1b1 b2b2b2
CCCCCCCCCCCCCCCCCCCC
c1c1c1 c2c2c2c2

```

**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.



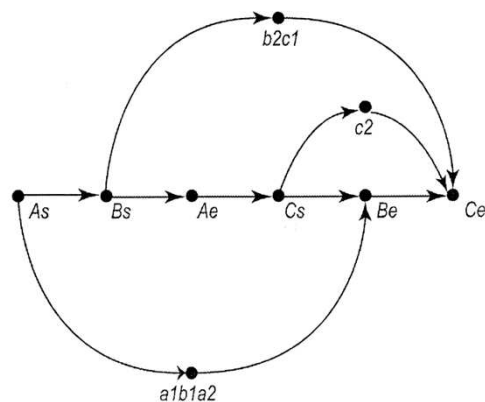
**Figure 5** Ordering graph of clone starts and ends. This represents the same clones as in Fig. 4. (As) The start of clone A; (Ae) the end of clone A. Similarly Bs, Be, Cs, and Ce represent the starts and ends of clones B and C.



**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.

# 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



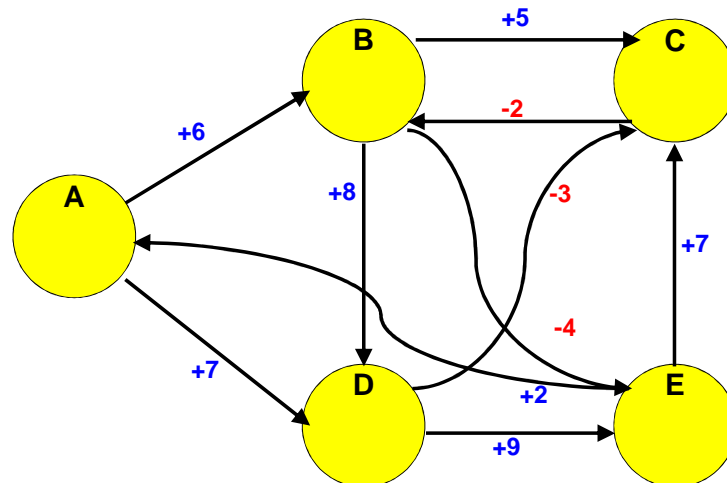
**Figure 6** Ordering graph after adding in rafts. The initial sequence contigs shown in Fig. 4 are merged into rafts where they overlap. This forms three rafts: a1b1a2, b2c1, and c2. These rafts are constrained to lie between the relevant clone ends by the addition of additional ordering edges to the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

## 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/>

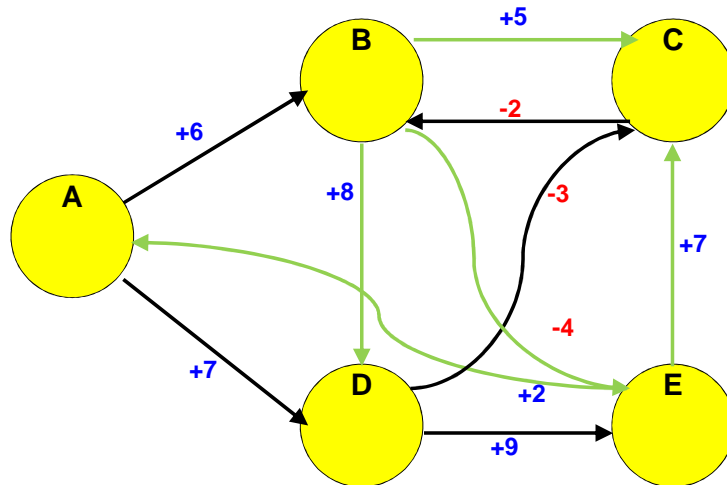
### Find the shortest path to all nodes.

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



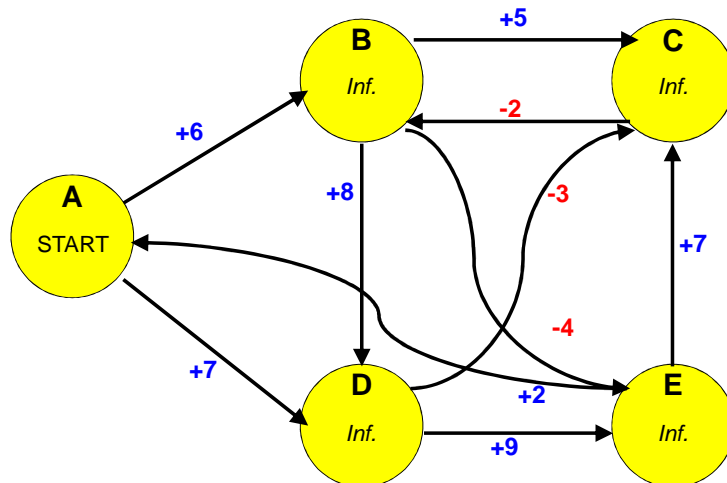
## Find the shortest path to all nodes.

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



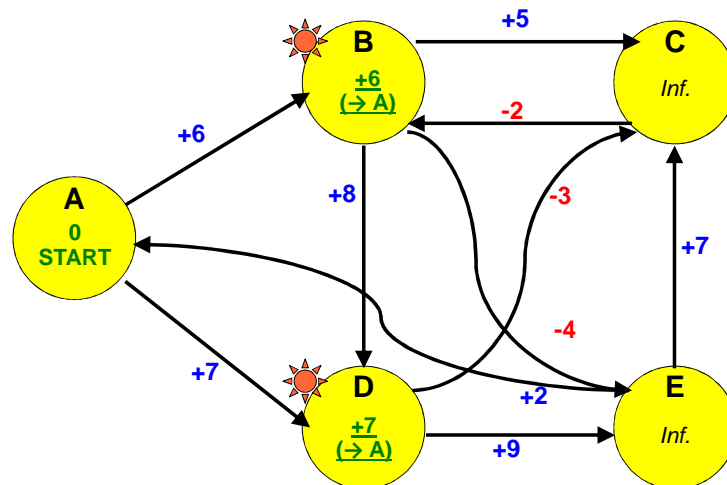
## Find the shortest path to all nodes.

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



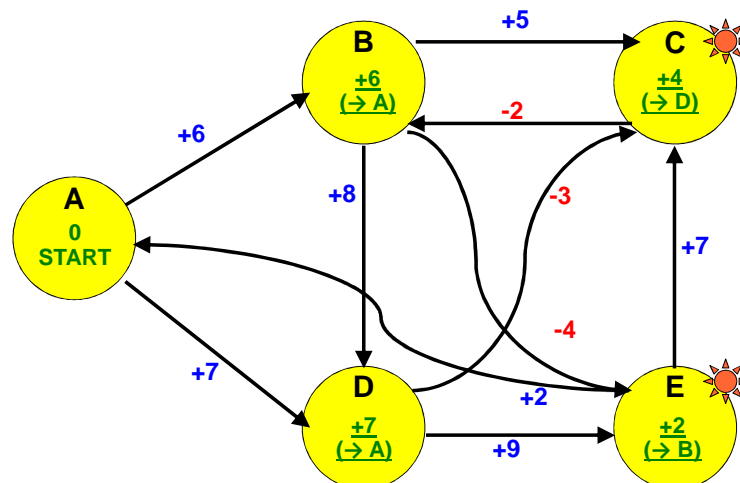
## Find the shortest path to all nodes.

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



## Find the shortest path to all nodes.

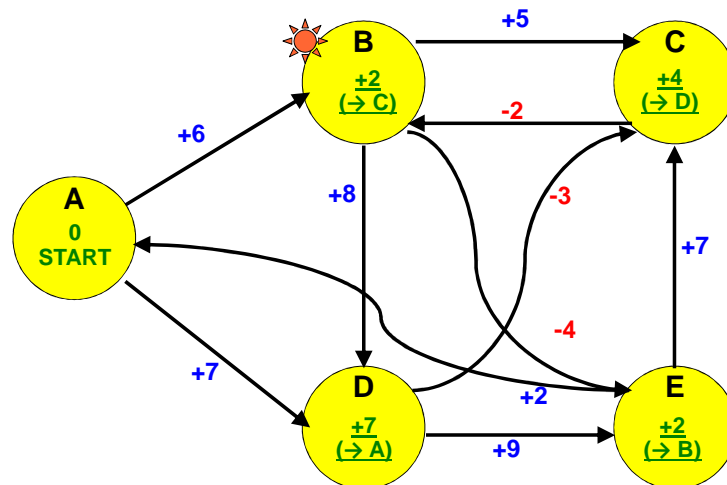
Take every edge and try to relax it ( $N - 1$  times where  $N$  is the count of nodes)





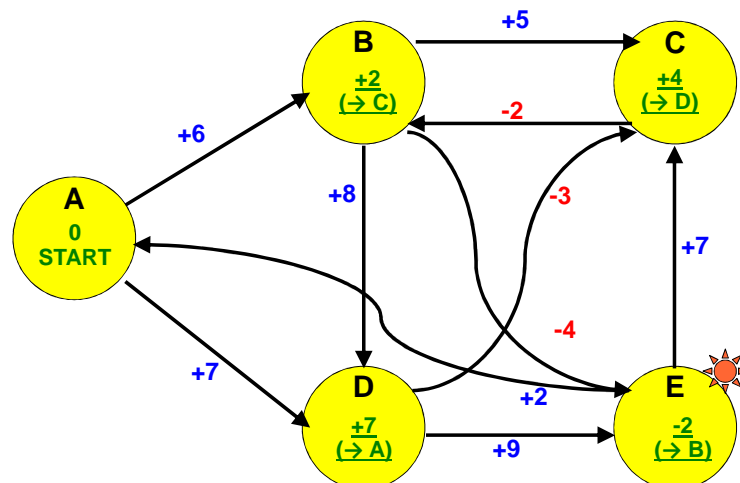
## Find the shortest path to all nodes.

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

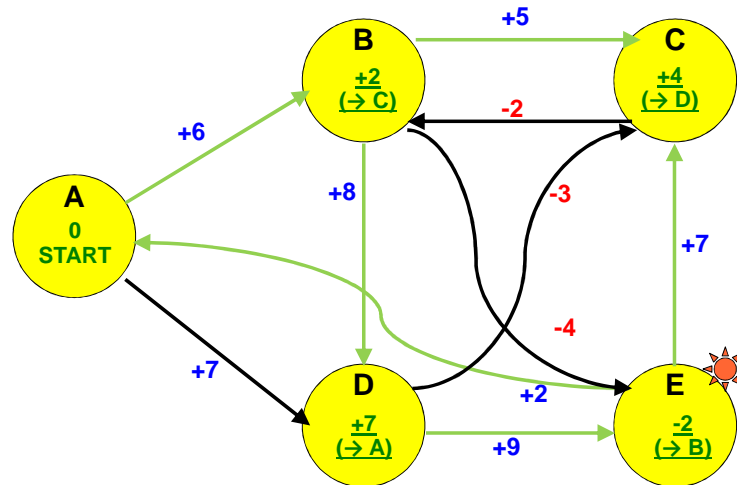


## Find the shortest path to all nodes.

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

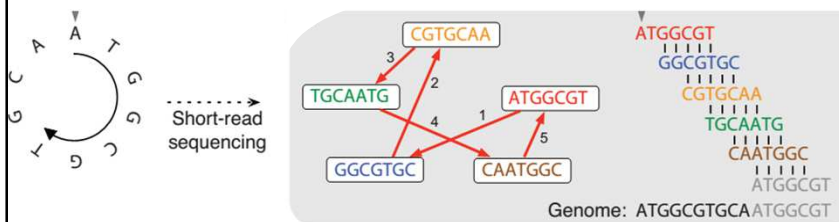


Answer: A-D-C-B-E



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

Here's what we saw before:

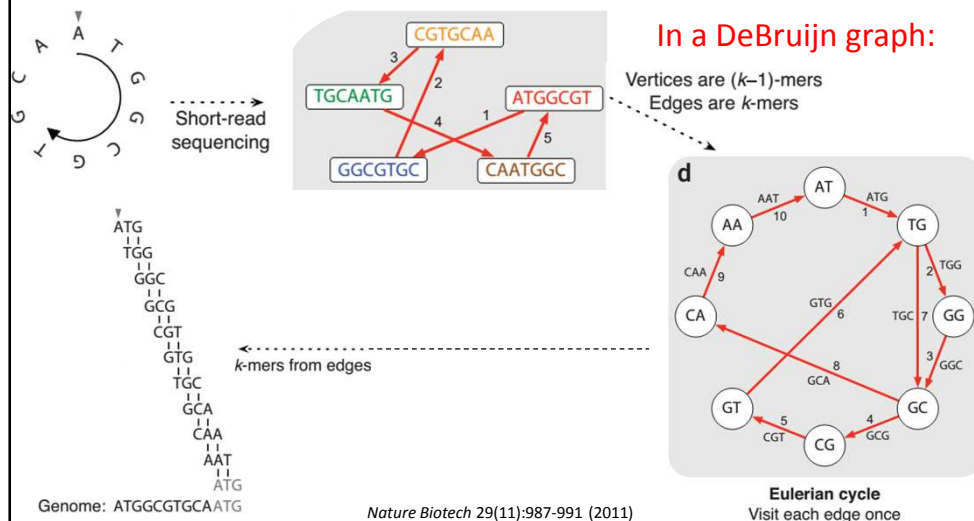


In Overlap-Layout-Consensus:

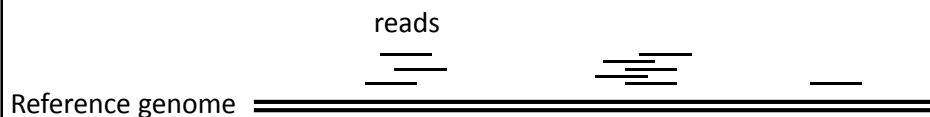
Nodes are reads

Edges are overlaps

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.

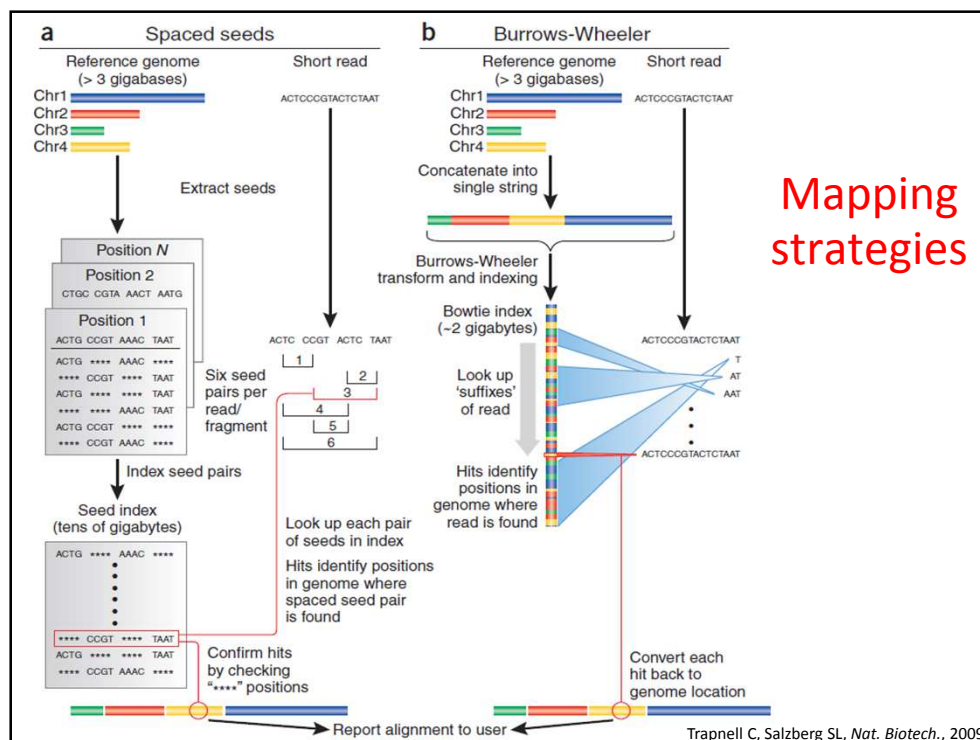


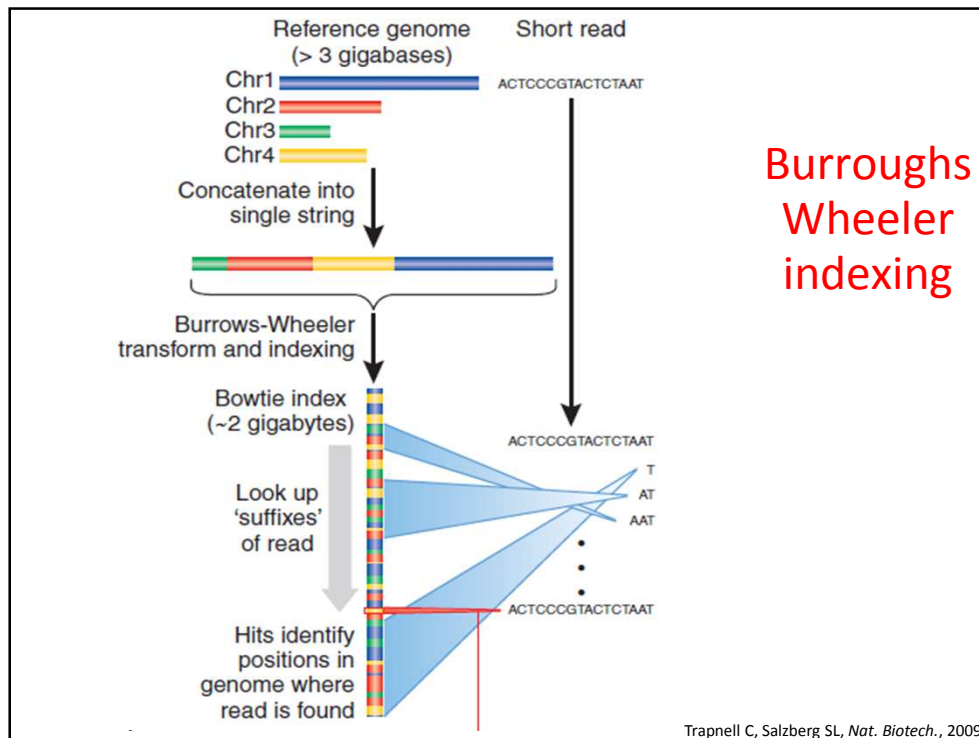
# Mapping reads to assembled genomes

**Table 1 A selection of short-read analysis software**

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

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

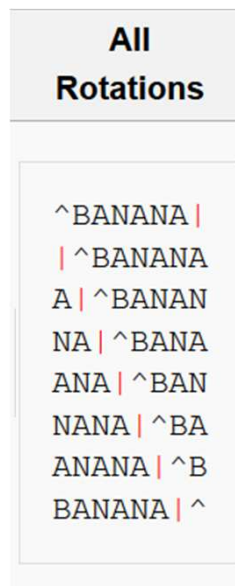
[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](http://en.wikipedia.org/wiki/Burrows-Wheeler_transform)

## Burroughs-Wheeler transform indexing



[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](http://en.wikipedia.org/wiki/Burrows-Wheeler_transform)

## Burroughs-Wheeler transform indexing



[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](http://en.wikipedia.org/wiki/Burrows-Wheeler_transform)

## Burroughs-Wheeler transform indexing

### Sorting All Rows in Alphabetical Order

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

[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](http://en.wikipedia.org/wiki/Burrows-Wheeler_transform)

## Burroughs-Wheeler transform indexing

### Taking Last Column

ANANA | ^**B**  
ANA | ^**B**AN  
A | ^**B**ANAN  
BANANA | ^  
NANA | ^**B**A  
NA | ^**B**ANA  
^**B**ANANA |  
| ^**B**ANANA

[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](http://en.wikipedia.org/wiki/Burrows-Wheeler_transform)

## Burroughs-Wheeler transform indexing

Output	Last Column
$BNN^{\wedge}AA \mid A$	

[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](http://en.wikipedia.org/wiki/Burrows-Wheeler_transform)

## Burroughs-Wheeler transform indexing

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

[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](http://en.wikipedia.org/wiki/Burrows-Wheeler_transform)



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

***Any ideas as how you might reverse it?***

### Burroughs-Wheeler transform indexing

Input
BNN^AA   A

[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](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...

[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](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...

[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](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   ^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](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 <u>^BANANA  </u> 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](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

[http://en.wikipedia.org/wiki/Burrows-Wheeler\\_transform](http://en.wikipedia.org/wiki/Burrows-Wheeler_transform)

