Next-generation sequencing

Lecture 5

Assembly Algorithms

Main algorithm used:

- Greedy algorithms
- Overlap Layout Consensus
- De brujin graphs

Greedy Assembly



Greedy Assembly

- Advantages:
 - Simple and easy to implement
 - effective
- Disadvantages
 - Since local information is considered at each step, the assembler can be easily confused by complex repeats, leading to mis-assemblies.
 - Local approach. Easy to be trapped into a local optimal solution (local minimum).
 - Early mistakes create bad assemblies.



Overlap-layout-consensus

Try to find the Hamiltonian path:

- A path in the graph contains each node exactly once.
- Following the Hamiltonian path, combine the overlapping sequences in the nodes into the sequence of the genome
- Computationally expensive (NP-hard problem)



Overlap-layout-consensus

- Better than Greedy algorithm. It can generate correct order of contigs that the Greedy algorithms may have errors.
- No efficient algorithm to find the Hamiltonian path
- Short fragment length = very small overlap therefore many false overlaps.
- Overlap discovery is sensitive to minimum overlap length and minimum percent identity required for an overlap.
- Overlap discovery is also time consuming.
- Large number of reads + short overlap + higher error are challenging for the overlap - layout - consensus approach
- Can't assemble repeat longer than read length
- It is mostly used with Sanger or 454 data.

Assembly Algorithms

Main algorithm used:

- Greedy algorithms
- Overlap Layout Consensus
- De bruijn graphs



- Get k-bp (k-mer) subsequences for reads.
- 2. k-mers in the reads are collected into nodes and the coverage at each node is recorded. Link two kmer nodes if they have overlap.
- 3. the graph is simplified to combine nodes that are associated with the continuous linear stretches into single, larger nodes of various kmer sizes.
- error correction removes the tips and bubbles that result from sequencing errors and creates a final graph structure that accurately and completely describes in the original genome sequence.

Flicek, Nature Methods, 2009

Differences between an overlap graph and a de Bruijn graph for assembly.



Copyright © 2010 by Cold Spring Harbor Laboratory Press

De Bruijn Graphs example

- For the purposes of illustration, we can use human readable text to explore how assemblies work.
- This is an example taken from Leipzig et al 2001.
- In it he uses the opening paragraph from Dickens' "A tale of two cities".
- It is an appropriate example because like genomes, it contains strings that are repeat over and over.

"It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity,.... "

Dickens, Charles. A Tale of Two Cities. 1859. London: Chapman Hall

De Bruijn Graphs example

itwasthebestoftimesitwastheworstoftimesitwastheageofwisdomitwastheageoffoolishness...

Generate random 'reads'



How do we assemble?

fincreduli geoffoolis Itwasthebe Itwasthebe geofwisdom itwastheep epochofinc timesitwas stheepocho nessitwast wastheageo theepochof stheepocho hofincredu estoftimes eoffoolish lishnessit hofbeliefi pochofincr itwasthewo twastheage toftimesit domitwasth ochofbelie eepochofbe eepochofbe astheworst chofincred theageofwi iefitwasth ssitwasthe astheepoch efitwasthe wisdomitwa ageoffooli twasthewor ochofbelie sdomitwast sitwasthea eepochofbe ffoolishne eofwisdomi hebestofti stheageoff twastheepo eworstofti stoftimesi theepochof esitwasthe heepochof i theepochof sdomitwast astheworst rstoftimes worstoftim stheepocho geoffoolis ffoolishne timesitwas lishnessit stheageoff eworstofti orstoftime fwisdomitw wastheageo heageofwis incredulit ishnessitw twastheepo wasthewors astheepoch heworstoft ofbeliefit wastheage heepochofi pochofincr heageofwis stheageoff timereduli astheageof wisdomitwa wastheageo astheepoch olishnessi astheepoch itwastheep twastheage wisdomitwa fbeliefitw bestoftime epochofbel theepochof sthebestof lishnessit hofbeliefi Itwasthebe ishnessitwasthewa ageoffwisdo twastheage shnessitwa thebestoft itwastheag theepochof itwasthewo ofbeliefit bestoftime mitwasthea imesitwast timesitwas orstoftime estoftimes twastheage foolishnes ftimesitwa thebestoft itwastheag theepochof itwasthewo ofbeliefit bestoftime mitwasthea imesitwast timesitwas orstoftime estoftimes twastheage toolishnes is stoftimes is domitwa theworstof astheworst sitwasthe theageoffo eepochofbe theageofwi foolishnes increduli ofbeliefit chofincred beliefitwa beliefitwa wisdomitwa theworstof astheworst sitwasthewo ofbeliefit bestoftime mitwasthea imesitwast timesitwas orstoftime estoftimes twastheage mesitwasth eisdomitwa theworstof astheworst sitwasthewo ofbeliefit eepochofbe theageofwi foolishnes increduli ofbeliefit chofincred beliefitwa beliefitwa wisdomitwa eageoffool eoffoolish itwastheag mesitwasth epochofinc sitwasthee astheageoffoo estoftime mesitwasthee thebestoft oolishness heepochofb ochofbelie wastheepoc bestoftime mesitwasth ebestofti

...etc. to 10's of millions of reads

Traditional all-vs-all comparisons of datasets this size require immense computational resources.

De Bruijn solution: Construct a graph efficiently

Step 1: "Kmerize" the data

Reads:	theageofwi	sthebestof	astheageof	worstoftim	imesitwast	
Kmers :	the	sth	ast	wor	ime	
(k=3)	hea	the	sth	ors	mes	
	eag	heb	the	rst	esi	
	age	ebe	hea	sto	sit itw twa	
	geo	bes	eag	tof		
	eof	est	age	oft		
	ofw	sto	geo	fti	was	
	fwi	tof	eof	tim	ast	

.....etc for all reads in the dataset

Step2 Build the graph

Look for k-1 overlaps: given by the reads



.....etc for all 'kmers' in the dataset

step3: simplify the graph



- The final step is to remove redundancy, result in the final De Bruijn Graph representation of our genome.
- the overlaps between reads are implicit in the graph, so all the millions v.s millions of comparisons are not required.
- On the downside, information is lost as repetitive sequences are "collapsed" into a single representation.

De Bruijn Graphs step4: Create contigs



Find the Hamiltonian path or cycle in the De Bruijn graph. Each path in the simplified De Bruijn graph is a Hamiltonian path. A Hamiltonian path is a contig.

Common Problems

- Spurs: dead-end sequences
- Bubbles: divergent paths that then converge
- Frayed rope: convergent then divergent paths
- Cycles: paths convergent upon themselves



Resolve graph complexity





Strengths and problems of De Bruijn approach

Strengths:

- No need to calculate the overlaps
- Size of the final graph is proportional to the genome size
- successfully for very short reads (<50bp)

Problems:

- The main drawback to the de Bruijn approach is the loss of information caused by decomposing a read into a path of k-mers.
- require an enormous amount of computer space
- Can only resolve k long repeat
- Loose connectivity when create the contigs

Strengths and problems of De Bruijn approach



Schlebusch, 2012



Pros: correctly links two sequences without having to compute overlap score. (above case) Cons: two sequences are linked without any real overlap. (left case)

De Bruijn Assemblers

- Euler: http://nbcr.sdsc.edu/euler/, Sanger, 454, 2001-2006
- Velvet: <u>http://www.ebi.ac.uk/~zerbino/velvet/</u>, small genomes, Sanger, 454, Solexa, SOLiD, 2007-2009 (very good for small genome)
- ABySS: <u>http://www.bcgsc.ca/platform/bioinfo/software/abyss</u>, large genome, Solexa, SOLiD, 2008-2011 (for very large genome)
- SOAP-denovo: <u>http://soap.genomics.org.cn/soapdenovo.html</u>, Solexa, 2009
- ALLPATH-LG: <u>http://www.broadinstitute.org/software/allpaths-lg/blog/</u>, large genome, Solexa, SOLiD, 2011 (very good performance bu require 2 lib of different insert sizes)
- IDBA-UD: <u>http://i.cs.hku.hk/~alse/hkubrg/projects/idba_ud/</u>, Sanger, 454, Solexa, 2010 (metagenomic, doesn't rely on coverage to remove error)

Comparison of Assembly tools

Algorithm Feature	Greedy Assemblers	OLC Assemblers	DBG Assemblers
Approaches to graph construction			
Implicit Reads as graph nodes K-mers as graph nodes Simple paths as graph nodes Multiple values of K Multiple overlap stringencies	SSAKE, SHARCGS, VCAKE SHARCGS	CABOG, Newbler, Edena	Euler, Velvet, ABySS, SOAP AllPaths Euler
Approaches to graph reduction			
Filter overlaps Greedy contig extension Collapse simple paths Erosion of spurs Transitive overlap reduction Bubble smoothing Bubble detection Reads separate tangled paths Break at low coverage Break at high coverage High coverage indicates repeat Special use of long reads	SSAKE, SHARCGS, VCAKE	CABOG CABOG, Newbler CABOG, Edena Edena Edena CABOG CABOG Shorty	Euler, Velvet, SOAP Euler, Velvet, AllPaths, SOAP Euler, Velvet, SOAP AllPaths Euler, SOAP Velvet, SOAP Euler Velvet Velvet Velvet
Graph partitions			
Partition by K-mers Partition by scaffolds			ABySS AllPaths

Miller, genomics, 2010, 95(6):315-27

Comparison of Assembly tools

Algorithm Feature	Greedy Assemblers	OLC Assemblers	DBG Assemblers
Modeled features of reads			
Base substitutions			Euler, AllPaths, SOAP
Homopolymer miscount		CABOG	
Concentrated error in 3' end			Euler
Flow space		Newbler	
Color space		Shorty	Velvet
Removal of erroneous reads			
Based on K-mer frequencies			Euler, Velvet, AllPaths
Based on K-mer freq and QV			AllPaths
For multiple values of K			AllPaths
By alignment to other reads		CABOG	
By alignment and QV	SHARCGS		
Correction of erroneous base calls			
Based on K-mer frequencies			Euler, SOAP
Based on Kmer freq and QV			AllPaths
Based on alignments		CABOG	

Comparison of the effect of various coverage depths on average contig length



Read length =35

Read length =75

Lin Y et al. Bioinformatics 2011;27:2031-2037

Couple issues

- Try different assemblers and compare their results.
- Need a big fat memory computer (from 16GB to 1TB).
- Running time is long: from several hours to several days.

			U			
		Rur	ntime (s)			
	Bench.Seq	E.coli	C.ele	H.sap-2	H.sap-3	
	(Length: bp)	(4.6M)	(20.9M)	(50.3M)	(100.5M)	
	SSAKE	2,776				
	VCAKE	1,672	16,742			
	Euler-sr	1,689	11,961	29,622		
SE	Edena	895	8,450	17,043		
	Velvet	205	1,003	2,786	6,098	
	ABySS	265	1,300	3,307	6,608	
	SOAPdenovo	62	253	560	1,029	
	SSAKE	9,163				
PE	Euler-sr	1,455	15,068			
	Velvet	229	1,351	55,581		
	ABySS	458	3,081	9,199	21,683	
	SOAPdenovo	78	374	889	2,257	

_ . _ _ /_ __ .

Running time

Lin Y et al. Bioinformatics 2011;27:2031-2037

'-' denotes runtime is too long (>10 days) to get assembly results

RAM used

			(112)		
	Bench.Seq	E.coli	C.ele	H.sap-2	H.sap-3
	(Length: bp)	(4.6M)	(20.9M)	(50.3M)	(100.5M)
	SSAKE	9,933			
	VCAKE	4,099	17,408		
	Euler-sr	1,536	7,065	13,312	
SE	Edena	1,741	7,557	30,720	
	Velvet	1,229	4,045	9,830	22528
	ABySS	1,126	3,993	8,909	18432
	SOAPdenovo	935	2,867	8,089	18227
	SSAKE	16,384			
	Euler-sr	1,638	7,578		
PE	Velvet	1,331	5,324	30,720	
	ABySS	950	4,505	9,830	18,432
	SOAPdenovo	1,638	5,939	10,342	19,456

RAM (MB)

'-' denotes the RAM of computer is not enough

Lin Y et al. Bioinformatics 2011;27:2031-2037

Whole genome sequencing

- *De Novo* whole genome sequencing
- Mapping assembly (Reference-guided assembly) (Resequencing)

Paired-end sequencing

- Paired-End sequencing (for Mate-pairs)
 - Sequence two ends of a fragment of known size.



- Currently fragment length (insert size) can range from 200 bps 10,000 bps
- Paired-end sequencing is helpful for assembly and locating repeat. It also can detect rearrangements, including insertions and deletions (indels) and inversions.
- As paired end reads are more likely to align to a reference, the quality of the entire data set improves



Paired-end sequencing by Illumina

Solid-phase amplification and Cyclic reversible termination

A simple modification to the standard single-read DNA library preparation.

Both the forward and reverse template strands of each cluster can be sequenced.

Mate-pair libraries



De Novo sequencing

- New species/strains
- Challenge of assembly with short reads
 - 8x coverage of 3 GB genome = 750 million fragments
 - Exponential problem for all-vs-all algorithm (overlap)
- Big problem with repeats
- Assemble contigs, fill gaps
- Paired-end reads are essential

Shotgun Sequencing



- Breaking the genome into a collection of small DNA fragments
- Sequencing.
- Reconstitute the genome.

Assembly Pipeline



Shotgun sequencing statistics

Typical contig coverage



Lander-Waterman statistics

L = read length G = genome size N = number of reads c = coverage = (NL / G)T = minimum detectable overlap $\sigma = 1 - T/L$

 $E(\# \text{ of islands}) = \text{Ne}^{-c\sigma}$ E(island size) = L((e^{c\sigma} - 1) / c + 1 - σ) contig = island with 2 or more reads





Example

Genome size: 1 Mbp Read Length: 600

c	Ν	#islands	#contigs	bases not in	bases not in
				any read	contigs
1	1,667	655	614	698	367,806
3	5,000	304	250	121	49,787
5	8,334	78	57	20	6,735
8	13,334	7	5	1	335

Experimental data

X coverage	# ctgs	% > 2X	avg ctg size (L-W)	max ctg size	# ORFs
1	284	54	1,234 (1,138)	3,337	526
3	597	67	1,794 (4,429)	9,589	1,092
5	548	79	2,495 (21,791)	17,977	1,398
8	495	85	3,294 (302,545)	64,307	1,762
complete	1	100	1.26 M	1.26 M	1,329

Numbers based on artificially chopping up the genome of *Wolbachia pipientis* dMel

Errors in Lander-Waterman Estimate

Lander-Waterman has errors:

- repeats
- GC/AT rich regions
- other low complexity regions
- cloning biases in shotgun libraries

Expected average contig length for a range of different read lengths and coverage values.





Copyright © 2010 by Cold Spring Harbor Laboratory Press

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

		Input sequence									
0					Nie of	Deed	Dein	Contigs			
organism/genome size	Assembler/status ^a	Туре	Pair size	Avgerage read (bp)	NO. OF reads	coverage ^b	coverage ^c	No.	N50	Max	Total
Human (<i>H.</i> sapiens)/3.0 Gb	ABYSS published 2009	GA	210 bp	35–46	3.5 B	45×	120×	2.76 M	1.5 kb	18.8 kb	2.18 Gb
Grapevine (V. vinifera)/500 Mb	Myriad published 2007	Sanger Sanger Sanger 454	2–10 kb 40 kb 120 kb None	579 460 369 169	5.95 M 144 k 68 k 12.5 M	6.9× 0.13× 0.02× 4.2×	21× 4.4× 4.2×	58,611	18.2 kb	238 kb	531 Mb ^d
Cucumber (C. sativus)/367 Mb	RePS2 published 2009	Sanger Sanger Sanger	2–6 kb 40 kb 140 kb	439 496 551	2.08 M 339 K 33.2 k	3.35× 0.46× 0.04×	9.9× 16.7× 5.6×	62,412	19,807	NR	226 Mb
		GĂ	200 bp	42	282 M	32.5×	76.8×	NR	2.6 kb	NR	204 Mb
		GA GA	400 bp 2 kb	44 53	173 M 105 M	20.6× 15.3×	94.4× 286×	NR	12.5 kb	NR	190 Mb
Panda (A. <i>melanoleura</i>)/2.4 Gb	SOAPdenovo published 2010	GA GA GA GA	150 500 2 kb 5 kb 10 kb	45 67 71 38 35	1.31 B 917 M 397 M 505 M 254 M	24.5× 25.5× 11.8× 8.0× 3.7×	43.3× 90.2× 192× 533× 571×	200,604	36,728	434,635	2.25 Gb
Strawberry (F. vesca)/220 Mb	CABOG and Velvet announced	454 454 454 454 GA SOLiD	None None 2.5 kb 20 kb None 2 kb	209 368 193 236 76 25	7.73 M 787 M 2.39 M 1.58 M 36 M 1.30 M	7.3× 13.2× 2.1× 1.7× 12.4× 0.14×		16,487	28,072	215,349	202 Mb
Turkey (<i>M.</i> gallopavo)/1.1 Gb	CABOG announced	454 454 GA GA	3 kb 20 kb None 180 bp None	180 195 366 74 74	6 M 2 M 13 M 200 M 200 M	1× 0.3× 4× 13× 13> Sch	8× 18× 	128,271 t al. Gen	12,594 ome Re	90 kb s. 2010;2	931 Mb 20:1165-11

One more example

For yeast 12Mbp

- read length: 200-400 bp
- coverage: 50X (how many reads do we need?)
- paired-end read insert size: 8kb (better to make multiple libraries with different insert sizes.)

Assembly Pipeline



- Scaffolding groups contigs into subsets with known order and orientation.
- Nodes are contigs
- Directed edge is between two nodes if they are adjacent in the genome.

• Mate pairs , if in different contigs, have a chance of being neighbors.





Scaffolding Algorithm

- Find all connected components
- Find a consistent orientation for all nodes in the graph (all contigs).
 - Nodes (contigs) have two types of edges
 - Same orientation
 - Different orientation
 - Make sure linked contigs have consistent orientation.
 - Optimization problem find the smallest number of edges to be removed so that all contigs have consistent orientation.
- Find the Hamiltonian path again.

Scaffolding software

- Some assembly software, such velvet, can do scaffolding as well.
- **Bambus** <u>http://www.cbcb.umd.edu/software/bambus</u>
- SSPACE -

http://www.baseclear.com/landingpages/basetools-awide-range-of-bioinformatics-solutions/sspacev12/

• **GRASS** - <u>http://code.google.com/p/tud-scaffolding/</u>

Additional techniques for orientation

- Physical mapping. Using information from Bacterial Artificial Chromosome (BAC)-based physical maps. Physical maps are built by clustering together of BACs sharing portions of a DNA "fingerprint," which is a pattern of DNA fragments of various sizes.
- Using markers along a DNA strand as independent information for scaffolding software. Markers are known sequences of nucleotides and tags. Markers are searched in the contigs.
- Using large scale maps of landmarks that lie along the the chromosomal DNA.

- Additional information is also useful:
 - Sequences of closely related organisms are also used as scaffolding information.

Example: aligning scaffolds of a mouse genome to the human genome

Scaffolding: Issues

- Errors in length of inserts (affecting distances between clone mates)
- Physical mapping is error prone.
- first builds a sequence based on linking information with high confidence, then factors in linking information with lower confidence.

Assembly Pipeline



The variability in repetitiveness among species species.

The ratio == the percentage of the genome that is covered by unique sequences of length k or longer.

The figure shows how much of each genome would be covered by *k*mers (reads) that occur exactly once.



The k-mer uniqueness ratio for five well-known organisms and one single-celled human parasite.

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

Repeat Control Issues

 Assembly programs should detect repeats in the assembly process and not after.

Incorrect genome reconstruction

• Assemblers should try to resolve correctly as many repeats as possible.

-Avoid intensive human labor

Repeat Control – When? & How?

- **pre-assembly:** find fragments that belong to repeats
 - statistically (most existing assemblers)
 - repeat database (*RepeatMasker*)
- **during assembly:** detect "tangles" indicative of repeats (Pevzner, Tang, Waterman 2001)
- **post-assembly:** find repetitive regions and potential mis-assemblies.
 - *Reputer, RepeatMasker*
 - "unhappy" mate-pairs (too close, too far, misoriented)

Detecting repeats pre-assembly:

- Statistical methods
 - Assemblers assume that reads are sampled uniformly at random.
 - Significant deviations from average coverage flagged as repeats.
 - frequent k-mers are ignored
 - "arrival" rate of reads in contigs compared with theoretical value.
- (e.g., 800 bp reads & 8x coverage reads "arrive" every 100 bp)

Detecting repeats during assembly

- Example: In Euler assembly program
 - Finds repeats by complex parts of the graph constructed during the assembly process.
 - Researchers look into these complex areas to try and resolve repeats.
 - Assemblers can use clone mate information to find incorrect assemblies. This is based on finding clone-mate pairs too close or too far from one another. ("unhappy" mate-pairs)

Detecting repeats post-assembly: Mis-assembled repeats



Repeat resolution

- Assemblers deduce that areas covered by a large number of reads may show an over-collapsed repeat.
- Problems with this samples are not uniformly distributed (for example, non-random libraries and poor clonability regions). leads to false positives.
- Repeats with low copy number are missed leads to false negatives.

Repeat resolution

- Techniques for repairing sequencing errors during repeat resolution
 - find clusters of reads where the clusters share differences.
 - For example, four reads contain an A , four contain a B. it is likely that the first four reads are from one copy and the last four from a different one.
 - Drawbacks are if certain areas of the sequence have low coverage.
 - Difficult to separate from true polymorphism

Discussion: Virtual genome assembly

- Plant mitochondrion genome 500,000 bp DNA circular
- How can you get mitochondria DNA? What problems do we need to concern for this step?
- For DNA fragmenting, what sizes of DNA fragments will you use? A. 1Kbp, B. 5kbp, C. both
- Pair-ended or single ended?
- What depth do you sequence? how many lanes do you need if you use illumina hiseq 2000? or how many reads do you need to get?
- Which assembler will you use? Why?
- What computer do you used to do assemble? A. 4GB laptop B. 50GB workstation C. computer cluster in HCC
- According to your estimate, how long does it take for assemble? A. 30 minutes B.2 hours C. 12 hours D. 4 days
- What software do you used to do scaffold? how long does it take?
- What is longest gap in one scaffold? How do you fill gaps?
- How do you determine if your assembled genome is good enough?
- how do you annotate genes?