#### Next-generation sequencing

Lecture 4



## NGS

- Introduction to the background
- NGS workflow and accuracy
- Data format, quality control, data management
- Assembly: sequence assembly refers to aligning and merging fragments of a much longer DNA sequence in order to reconstruct the original sequence.
- RNA-seq
  - Aligner
  - Analysis tools
  - Applications, such as MiRNA
- Chip-seq
  - Applications

# Goals of Assembly

- Reconstruction of unknown genomes
  - viruses, bacteria
  - Eucaryotes (individual genomes)
  - metagenomics (environmental samples, microbial communities)
- RNA-Seq => transcriptomes
  - for organisms without reference genomes
  - novel transcripts
  - fusion genes
  - viral integration
- Local assembly for detection of insertions and genomic rearrangements
  - unmapped reads from WGS
  - transposable elements
  - viral integration

# Next-generation sequencing

- Much higher throughput (1-4gbps / day)
- Lower cost / base pair
- Inherent ability to do paired-end (mate-pair) sequencing

# Assembly

- sequence assembly refers to aligning and merging fragments of a much longer DNA sequence in order to reconstruct the original sequence.
- Taking many copies of a book, passing each of them through a shredder with a different cutter, and piecing the text of the book back together just by looking at the shredded pieces.
- Combines short sequencing reads into contigs based on sequence similarity and overlap between reads.
- Find the shortest common sequence of a set of reads

# Assembly

- **Spanner**: single read that spans a repeat instance with sufficient unique sequence on either side of the repeat.
- Contig: contiguous sequence formed by several overlapping reads with no gaps.
- Supercontig (scaffold): ordered and oriented set of contigs, usually by mate pairs. Relative distance known => fill gaps between contigs with "NNNNNNNN..."
- Consensus sequence: sequence derived from the multiple alignment of reads in a contig.



# **Challenges for Assembly**

Repeats or similar parts in the genome. The reads originating from different copies of a repeat appear identical to the assembler and cause assembly errors.



- Lose DNA fragments during library preparation
- Bias during amplification
- Very short fragment lengths (25-200bps)
- High error rate

# Assembly

- Assembly algorithms
- De novo whole genome assembling strategies
- Mapping assembling strategies

# Assembly algorithms

- Assembly, finding an optimal path that connect all short reads (or part of short reads) once time, is NP-hard problem.
- No efficient solution.

– We need to use some approximation algorithm.

### **Assembly Algorithms**

Main algorithm used:

- Greedy algorithms
- Overlap Layout Consensus
- De brujin graphs

# Greedy Assembly



# Alternative Greedy Algorithm

- Instead of calculating overlaps of all reads:
  - 1. Start with a random read as an initial contig (seed)
  - 2. Go over all unassembled reads, pick the one that best fits the 3' end of the contig and elongate the contig by this read.
  - 3. Repeat step 2 until no elongation is possible anymore.
  - 4. Repeat step 2 for the 5' end of the reverse complement of the contig.
  - 5. Stop in case of a conflict (fork)
    - if two reads that do not overlap with each other could elongate the contig equally well.

# 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.



# Greedy Assemblers

- TIGR Assembler: <u>ftp://ftp.jcvi.org/pub/software/assembler/</u>, Sanger, 2003
- SSAKE: <u>http://www.bcgsc.ca/platform/bioinfo/software/ssake</u>, small genomes, Solexa/Illumina, 2007
- SHARCGS: <u>http://sharcgs.molgen.mpg.de/</u>, small genomes, Solexa/Illumina, SOLiD, Sanger, 454, 2007
- VCAKE: <u>http://sourceforge.net/projects/vcake</u>, small genomes, Solexa/Illumina, 2007
- Phrap: <u>http://www.phrap.org/</u>, Sanger, 454, Solexa, 1995-2008

### **Assembly Algorithms**

Main algorithm used:

- Greedy algorithms
- Overlap Layout Consensus
- De brujin graphs

# Graph model for Assembly

(a)



Graphs are made up by vertices (nodes) and edges (links). G=(V,E) V: a finite set of vertices. E: edges of the graph

aaccqq

ccqqtt

Graph for Assembly:

- Graph model. A node is a read or a k-mer subsequence
- Edge: if there are over lap between two k-mer subsequence



Miller, 2010

### Overlap-layout-consensus

Main entity: read Relationship between reads: overlap



#### Graph model:

- A node is a
- Edge: if there are overlap between two reads





#### Overlap-layout-consensus Step 1: Find Overlapping Reads

- What reads are intersecting ?
- Issues:
  - Need efficient alignment algorithm (parallelization and index based strategies)
  - 2. Doesn't scale well when number of read is high
  - 3. Use seed based alignment with extension



Overlap-layout-consensus Step 2: Construct overlap graph

- A graph is constructed:
  - Nodes are reads
  - Edges represent overlapping reads
- Need to simplify graph, such as remove redundant nodes and edges.



#### Overlap-layout-consensus Step 3: Consensus stage, Find Contigs

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)



# Travelling salesman problem

 The travelling salesman problem (TSP) asks the following question: Given a list of cities and the distances between each pair of cities, what is the shortest possible route that visits each city exactly once and returns to the origin city?

# Circular genome

• Hamiltonian circuit: visit each node (city) exactly once, returning to the start B C D





# 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.

#### Overlap-layout-consensus

- Celera (CABOG): <u>http://www.jcvi.org/cms/research/projects/celera-assembler/overview/</u>, large genome, Sanger, 454,Solexa, 2004/2010
- Newbler: <a href="http://www.454.com/">http://www.454.com/</a>, 454, Sanger, 2009
- Mira: <u>http://sourceforge.net/apps/mediawiki/mira-assembler/</u> Sanger, 454, Solexa, 1998/2011
- Edena: <u>http://www.genomic.ch/edena.php</u> Snger, 454, 2008/2013

### **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
	geo	bes	eag	tof	itw
	eof	est	age	oft	twa
	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)</li>

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: <a href="http://nbcr.sdsc.edu/euler/">http://nbcr.sdsc.edu/euler/</a>, 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		
Graph partitions					
Partition by K-mers Partition by scaffolds			ABySS AllPaths		

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

# 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		
Runtime (s)					
	Bench.Seq	E.coli	C.ele	H.sap-2	H.sap-3
	(Length: bp)	(4.6M)	(20.9M)	(50.3M)	(100.5M)
SE	SSAKE	2,776			
	VCAKE	1,672	16,742		
	Euler-sr	1,689	11,961	29,622	
	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
PE	SSAKE	9,163			
	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

	Bench.Seq	E.coli	C.ele	H.sap-2	H.sap-3
	(Length: bp)	(4.6M)	(20.9M)	(50.3M)	(100.5M)
SE	SSAKE	9,933			
	VCAKE	4,099	17,408		
	Euler-sr	1,536	7,065	13,312	
	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
PE	SSAKE	16,384			
	Euler-sr	1,638	7,578		
	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)