

De novo assembly of NGS data

Sandra Smit
WUR bioinformatics



De novo assembly. Why?



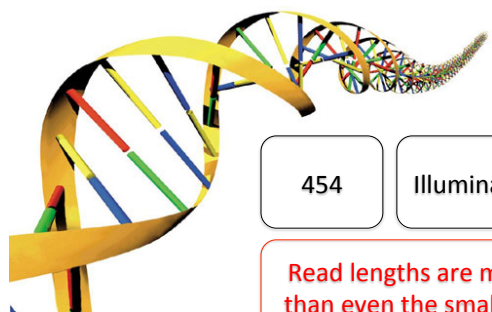
PacBio

454

Illumina

SOLiD

De novo **assembly**. Why?



PacBio

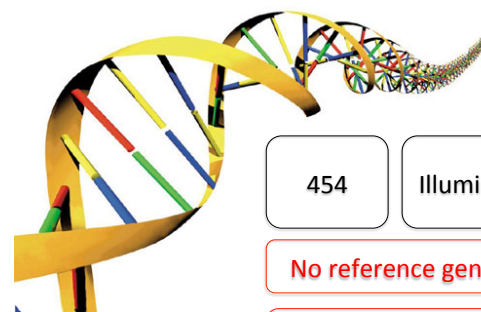
454

Illumina

SOLiD

Read lengths are much shorter than even the smallest genome

De novo **assembly**. Why?



PacBio

454

Illumina

SOLiD

No reference genome available

Not all reads map to the reference

Assemble this!

quickb			umpso
ckbrow			thequ
		thelaz	
sove	the	uickb	bro
helaz	ydo	azy	nfox
ownfo	umps	ckbro	kbr
heq	azyd	lazydo	oxjum
quic	fox	thel	erthel
	nfo		ver

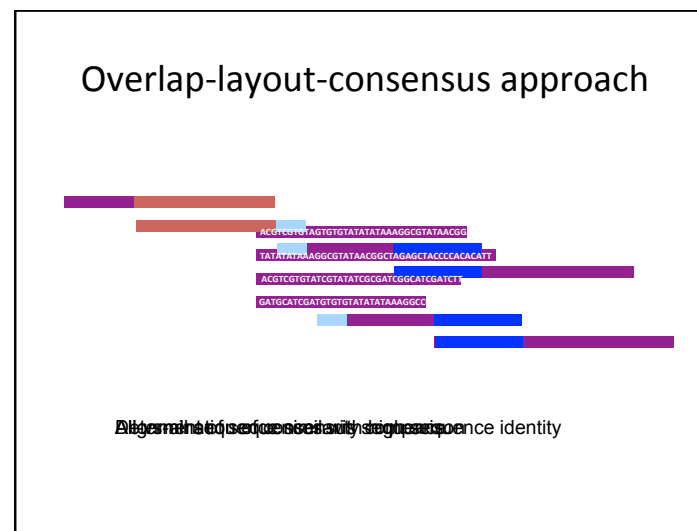
Find overlap

quickb			umpso
ckbrow			thequ
		thelaz	
sove	the	uickb	bro
helaz	ydo	azy	nfox
ownfo	umps	ckbro	kbr
heq	azyd	lazydo	oxjum
quic	fox	thel	erthel
	nfo		ver

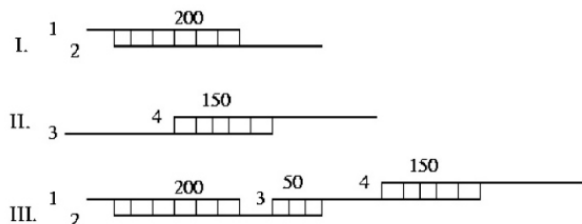
Layout of the reads

the		oxjum
thequ		umps
heq		umpso
qui		sove
quic		ver
quickb		erthel
uickb		rthela
ckbro		thel
ckbrow		thelaz
kbr		helaz
bro		lazydo
ownfo		azy
nfo		azyd
nfox		ydo
fox		ydog

thequickbrownfoxjumpsoverthelazydog

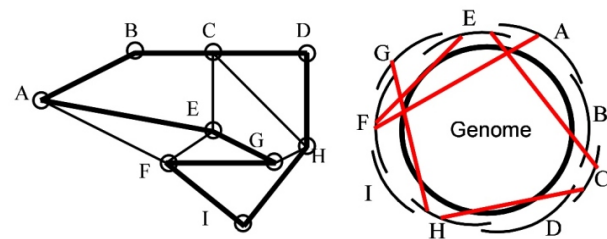


Greedy extension



Greedy join reads that are most similar to each other

Graph-based OLC



Find Hamiltonian path

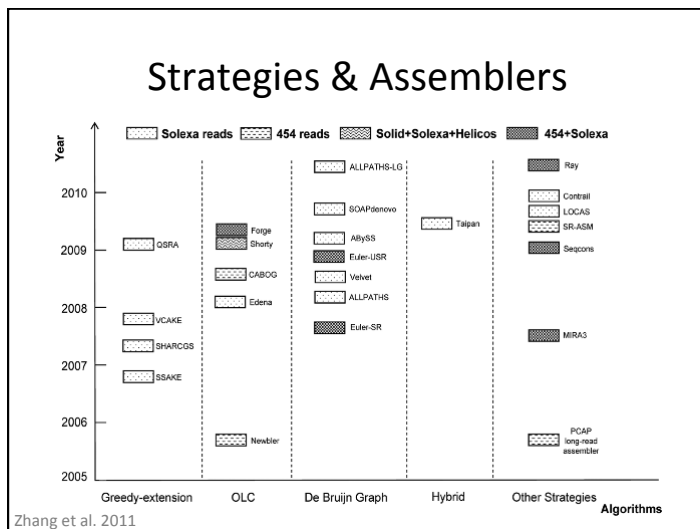
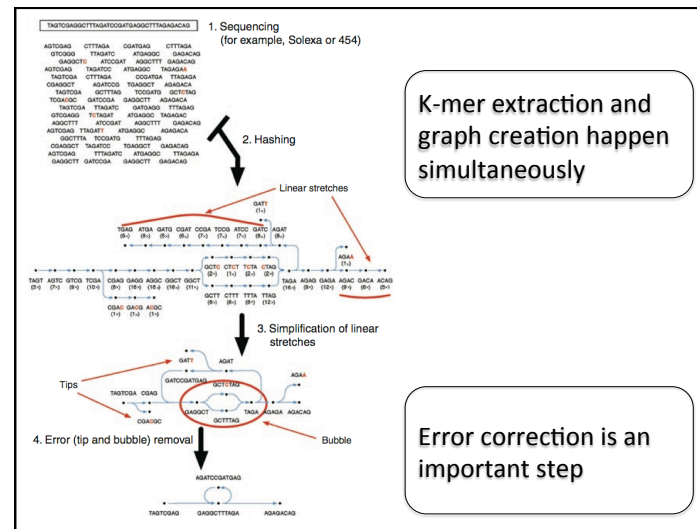
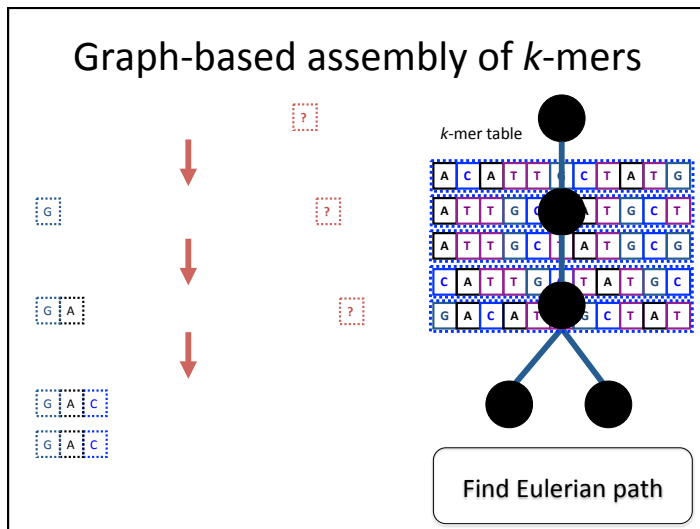
The de Bruijn graph approach

- Solving the assembly problem with De Bruijn graphs
 - Nodes are k -mers
 - Edges are overlaps between k -mers on $k-1$ positions
- A k -mer is a (short) substring of a read
 - A read of n bp consists of $(n - k + 1)$ k -mers
 - Example: a read of 75 bp consists of 45 (overlapping) 31-mers

```
TGACCAAGT
TGAC
GACC
ACCA
CCAG
....
```

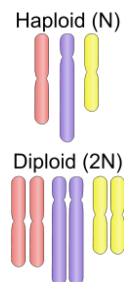
Generating k -mers from reads



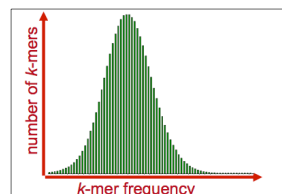
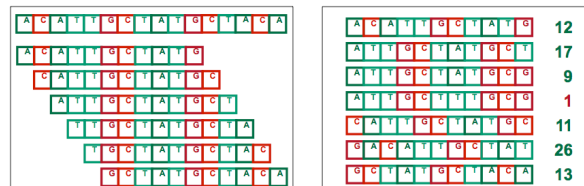


Complicating factor: ploidy

- Ploidy is the number of sets of chromosomes in a biological cell
- Diploid: two sets of chromosomes
 - Human cells
- Polyploidy means more than two sets of chromosomes per nucleus
- Tetraploid: four sets of chromosomes
 - Common in plants

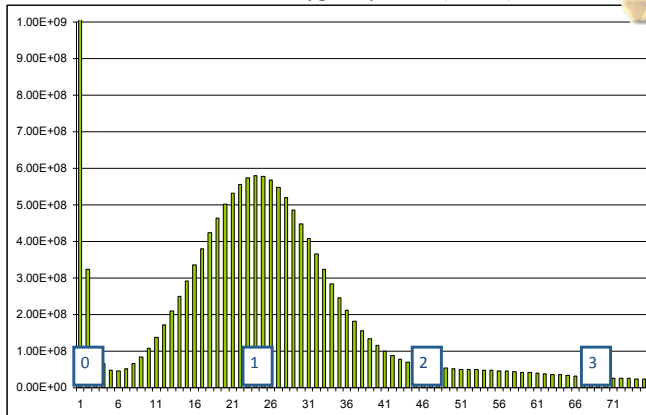


Complicating factor: heterozygosity

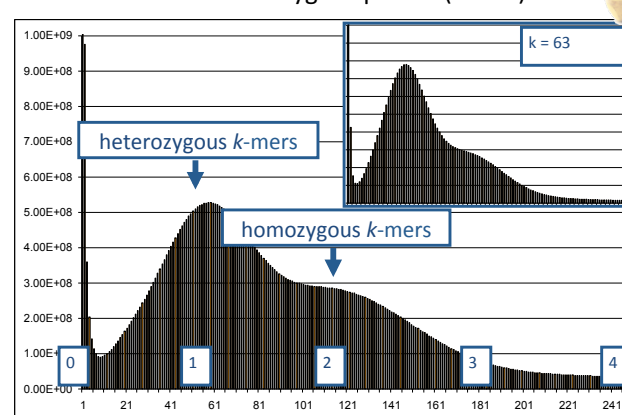


- Assumptions
- | the reads are randomly sampled from the genome
 - | the k-mer table represents the complete genome

k-mers – homozygous potato (k = 31)

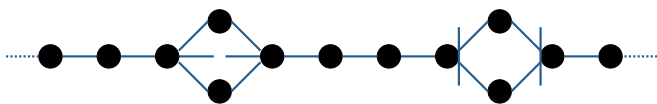


k-mers – heterozygous potato (k = 31)

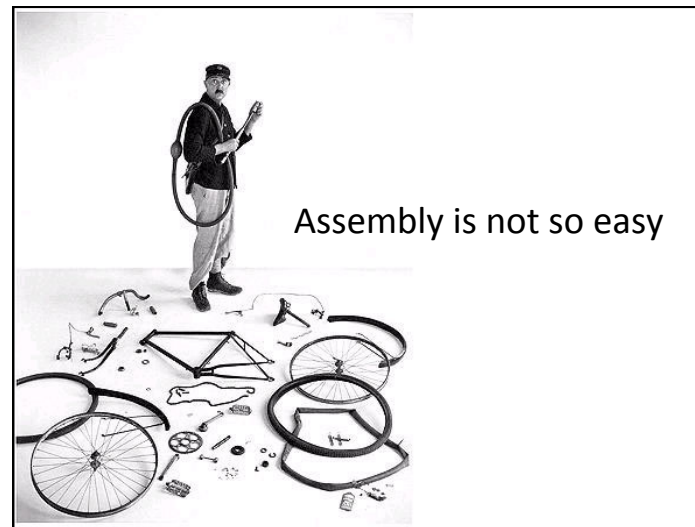


Complicating factor: heterozygosity

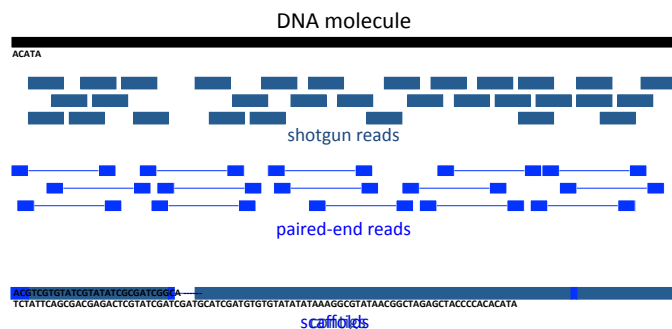
- Heterozygous k -mers create “bubbles” in the graph
 - Every SNP and in/del creates a bubble
- Solutions
 - Pinching bubbles: loss of diversity
 - Breaking graph: loss of contiguity




Assembly is not so easy



Assembly result



Assembly quality

- Number and length of contigs/scaffolds
 - N50 
 - Number of gaps/Ns
 - GC content
 - Coverage and completeness
 - Structural consistency
 - Paired-end data
 - Physical and genetic maps
 - Error rate
 - Known sequences (BAC clones, ESTs, etc.)
- 50% of assembly
N50 length: 18 Kb
N50 index: 4

Basic stats: exercise

Contigs (length)

250
150
120
100
90
80
80
70
40
20

- Calculate
 - N50 contig size and index
 - N90 contig size and index

Basic stats: exercise

Contigs (length)

250
150
120
100
90
80
80
70
40
20

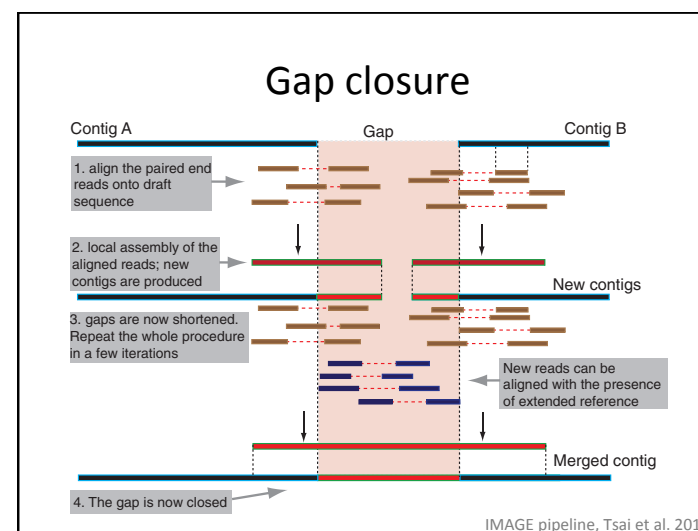
- Total assembly size = 1000
- Sort from large to small
- N50 lookup value = 500
- N50 contig size = 120
- N50 contig index = 3
- N90 contig size = 70
- N90 contig index = 8

Genome finishing

From initial draft to complete genome

- Closing gaps
 - IMAGE pipeline, Tsai et al. 2010
 - GapFiller, Boetzer and Pirovano 2012
- Establish order and orientation of contigs/scaffolds
 - E.g. FISH experiments
- Base error correction
 - k*-mer correction
- Contamination removal

Chain et al., Science 2009
Genome Project Standards in a New Era of Sequencing



Sequencing strategy

Sequencing strategy: which sequencing data is necessary for my project?
How can we balance the benefits and the costs?

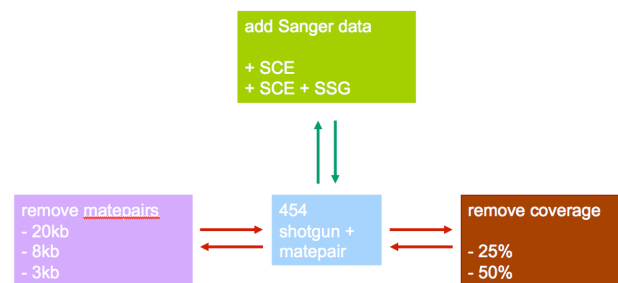
Technologies have different properties and error profiles

What is the contribution of different input data to an assembly?

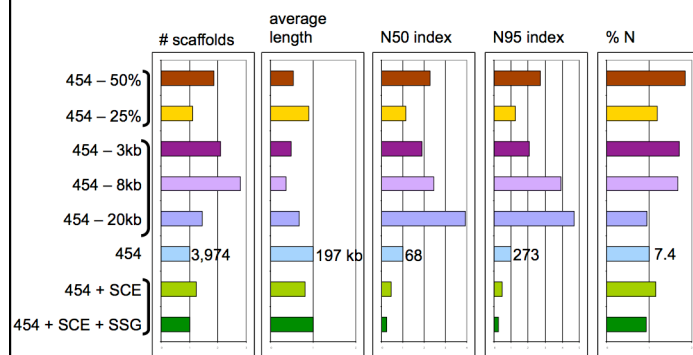
	Throughput	Length	Quality	Costs
Sanger	8 Mb/day	800 nt	10^{-4} - 10^{-5}	~500\$/Mb
454/Roche	750 Mb/day	400 nt	10^{-3} - 10^{-4}	~20\$/Mb
Illumina	5,000 Mb/day	100 nt	10^{-2} - 10^{-3}	~0.50\$/Mb
SOLID	5,000 Mb/day	50 nt	10^{-2} - 10^{-3}	~0.50\$/Mb
Helicos	5,000 Mb/day	32 nt	10^{-1}	<0.50\$/Mb

Kircher & Kelso, 2010

Sequencing strategy

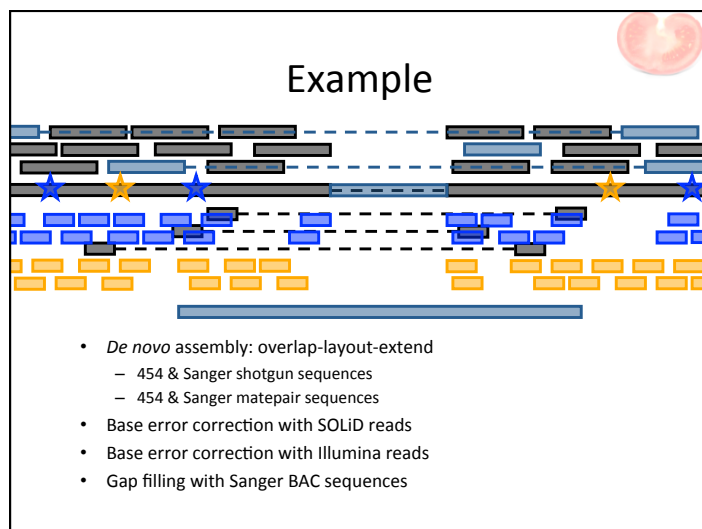


Sequencing strategy



Recommendations

- Longer reads (Sanger, 454) improve assembly quality
- Deep coverage by reads with lengths longer than common repeats
- Mixture of short and long insert sizes for paired-end data



Summary: assembly workflow

NGS sequencing/data production

Data cleaning/filtering/error correction

De novo assembly into contigs and (super)scaffolds

Assembly finishing/gap filling

Assembly validation

Literature

Assembly of large genomes using second-generation sequencing
Schatz et al.
Genome Res. 2010


Assembly algorithms for next-generation sequencing data
Miller et al.
Genomics 2010


High-throughput DNA sequencing -- concepts and limitations
Kircher and Kelso
Bioessays 2010

Acknowledgements

- Jack Leunissen
- Roeland van Ham
- Erwin Datema
- Jan van Haarst



 The International Tomato Genome Sequencing Consortium

 Potato Genome Sequencing Consortium