

MASSEY  
UNIVERSITY  
TE KUNINGA KI PĀRĀKOHOA  
UNIVERSITY OF NEW ZEALAND

246.201 Systems & Models: Bioinformatics I

Sebastian Schmeier  
s.schmeier@massey.ac.nz

22<sup>nd</sup> September 2014

CAN YOU DO ANYTHING RIGHT?

Copyright © Randy Glasberg - www.glasberg.com

THE ENGINE OF THE NEW  
NEW ZEALAND

## Module overview

- Week 1: Genome assembly

- Learning outcomes:

- Be able to describe the concepts regarding genome assembly
- Be able to operate comfortably the Linux command-line
- Be able to compute, investigate and evaluate the sequence quality
- Be able to compute, interpret and evaluate a whole genome assembly

2

## Module overview

- Week 2/3: Genome annotation

- Learning outcomes:

- Be able to characterize a gene and an ORF in the genomic context
- Be able to describe the concept of gene finding algorithms
- Be able to describe the concept of a sequence BLAST
- Be able to operate DNA Master
- Be able to identify genes using DNA Master
- Be able to annotate a phage genome

3

## Week 1: Overview

- DNA sequencing technologies
- Quality assessment of a sequencing run
- Genome assemblies and *de Bruijn* graphs

4

## Genomics and bioinformatics

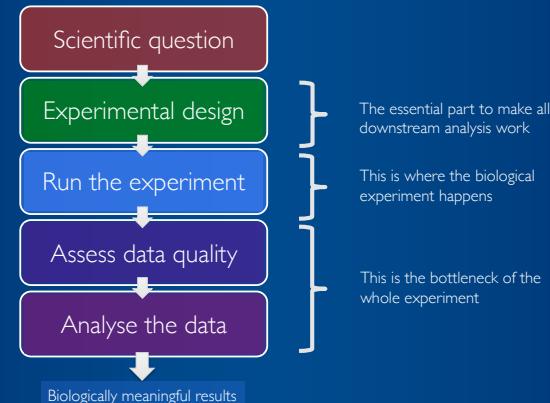
- Genomics
  - *Genomics is a discipline in genetics that applies recombinant DNA, DNA sequencing methods, and bioinformatics to sequence, assemble, and analyse the function and structure of genomes* (*The book of Wikipedia*)
- Bioinformatics
  - *The mathematical, statistical and computing methods that aim to solve biological problems using DNA and amino acid sequences and related information* (*Fredj Tekoia, Institut Pasteur*)

Genomics and bioinformatics are related disciplines

**Note!** If you want to do some genomics work, you will have to do some bioinformatics work as well

5

## Typical workflow of a genomics experiment



6

## Examples of different genomics work

- Gene expression studies
  - e.g. between two different stimuli, etc.
- Finding mutations in the DNA
  - e.g. Single nucleotide polymorphisms beneficial for fitness, etc.
- Metabarcoding studies
  - e.g. what kind of species are in my sample (e.g. natural/wild fermentation)
- Gene regulatory studies
  - e.g. transcription factor binding studies / what is regulating my genes?
  - e.g. histone modifications, DNA methylation studies
- De novo genome sequencing
  - Define the genomic context of a species

7

## De novo genome sequencing

1. DNA sequencing
2. Quality assessment
3. Genome assembly

22 Sep 2014

8

## Genome versus transcriptome

- Genome
  - The entirety of an organism's ancestral information. It is encoded either in DNA or, for many types of viruses, in RNA.
- Transcriptome
  - The set of all RNA molecules, including messenger RNA, ribosomal RNA, transfer RNA, and other non-coding RNA produced in one or a population of cells

Name	Base Pairs
HIV	9,749
E.Coli	4,600,000
Yeast	12,100,000
Drosophila	130,000,000
<b>Homo sapiens</b>	<b>3,200,000,000</b>
marbled lungfish	130,000,000,000
"Amoeba" dubia	670,000,000,000

disputed

9

## DNA sequencing

- **DNA sequencing** is the process of determining the nucleotide order of a given DNA fragment.
  - **First-generation sequencing:**
    - 1977 Sanger sequencing method development (chain-termination method)
    - 2001, Sanger method produced a draft sequence of the human genome
  - **Next-generation sequencing (NGS)**
    - Demand for low-cost sequencing has driven the development of high-throughput sequencing (or NGS) technologies that parallelize the sequencing process, producing thousands or millions of sequences concurrently
    - 2004 454 Life Sciences marketed a parallelized version of pyrosequencing

10

## Result of a sequencing run

- **Short read sequences**
  - The result of NGS technology are a collection of short nucleotide sequences (reads) of varying length (~40-400nt) depending on the technology used to generate the reads
  - Usually a reads quality is good at the beginning of the read and errors accumulate the longer the read gets → **IMPORTANT**

11

## Illumina sequencing

- **MiSeq:**
  - Bench-top sequencer
  - Produces around 30 million reads/run
  - Reads are up to 250nt
- **HiSeq:**
  - Large-scale sequencer
  - 4 billion reads/run
  - Reads up to 150nt
- The Illumina systems accumulate errors towards the end of the read sequence.

12

### Illumina sequencing

- An Illumina flowcell is a surface to which seq. adaptors are covalently attached.
- DNA with complementary adaptors is attached, clonally amplified, and then sequenced by synthesis
- Each flowcell is subdivided into hundreds of tiles

**A**

flowcell ID  
flow in  
barcode  
flow out

**B**

flow in  
polyacrylamide-coated interior surface of flowcell  
flow out

13

### How does the output look like?

- The file-format that you will encounter soon is called FastQ

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
GATTGGGGTCAAGCAGTATCGATAATGAAATCCATTGTTCT
+
"*(((((***+))%%%++)(%%%%,1***-+*))**55CCF>>>>>
```

Sequence id  
Sequence  
Phred quality of the corresponding nucleotide (ASCII code)

14

### How does the output look like?

- FastQ: Identifier

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
GATTGGGGTCAAGCAGTATCGATAATGAAATCCATTGTTCT
+
"*(((((***+))%%%++)(%%%%,1***-+*))**55CCF>>>>>>
```

Sequence id

Casava 1.8 the format

EAS139	the unique instrument name
136	the run id
FC706VJ	the flowcell id
2	flowcell lane
2104	tile number within the flowcell lane
15343	X-coordinate of the cluster within the tile
197393	Y-coordinate of the cluster within the tile
1	the member of a pair, 1 or 2 (paired-end or mate-pair reads only)
Y	Y if the read is filtered, N otherwise
18	0 when none of the control bits are on, otherwise it is an even number
ATCACG	index sequence

15

### How does the output look like?

- FastQ: Phred base quality

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
GATTGGGGTCAAGCAGTATCGATAATGAAATCCATTGTTCT
+
"*(((((***+))%%%++)(%%%%,1***-+*))**55CCF>>>>>>
```

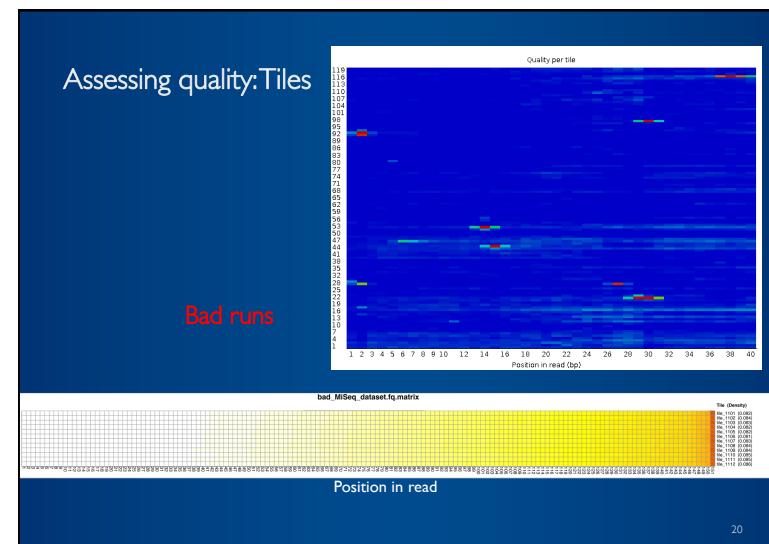
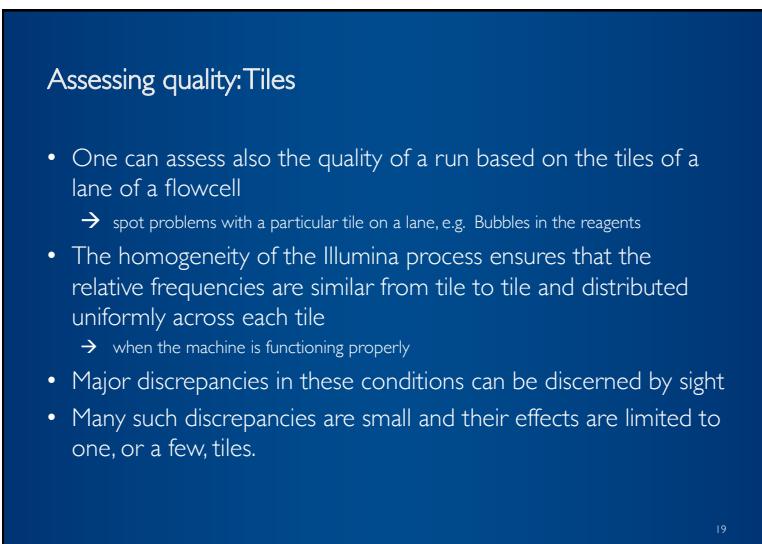
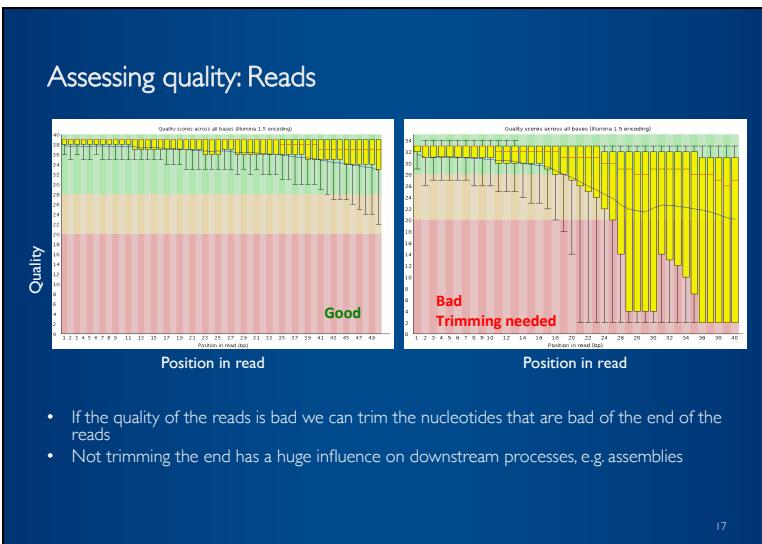
Phred quality of the corresponding nucleotide (ASCII code)

- One ASCII character per nucleotide.
- Encodes for a quality  $Q = -10 \log_{10}(P)$ , where P is the error probability

The Relationship Between Quality Score and Base Call Accuracy		
Quality Score	Probability of Incorrect Base Call	Inferred Base Call Accuracy
10 (Q10)	1 in 10	90%
20 (Q20)	1 in 100	99%
30 (Q30)	1 in 1000	99.9%

Q	ASCII	P
1	#	0.79433
2	\$	0.63096
3	%	0.50119
4	&	0.39011
5	,	0.31623
6	(	0.25119
7	)	0.19953
8	:	0.15097
9	*	0.12589
10	+	0.10000
11	,	0.07943

16



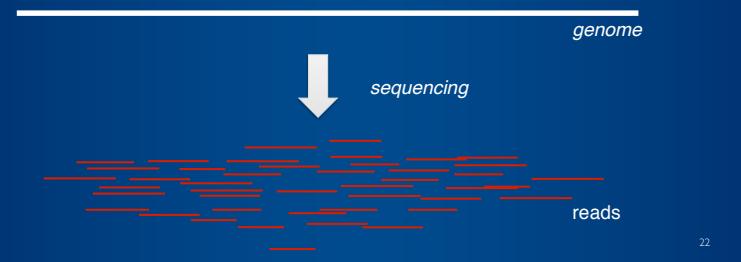
### Assessing quality: Final

- After assessing the quality we would try to remove all bp from the ends that do not fulfil a certain quality
- Thus, we work with a adjusted set of sequencing reads for which we are more certain that they represent correct nt sequences from the genome

21

### *De novo* genome assembly

- The process of generating a new genome sequence from NGS genome sequence reads based on assembly algorithms
- Assembly involves joining short sequence fragments together into long pieces – contigs



22

ATCG

GTGGC

GCGTG

TGGCA

23

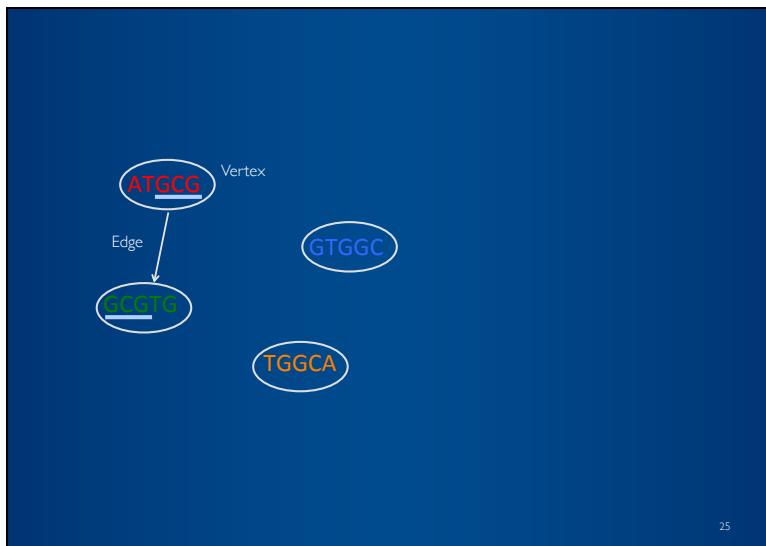
ATCG

GTGGC

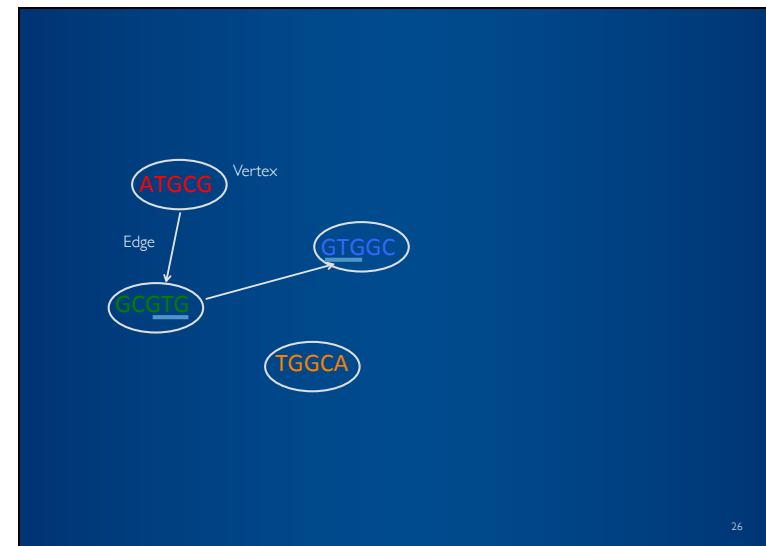
GCGTG

TGGCA

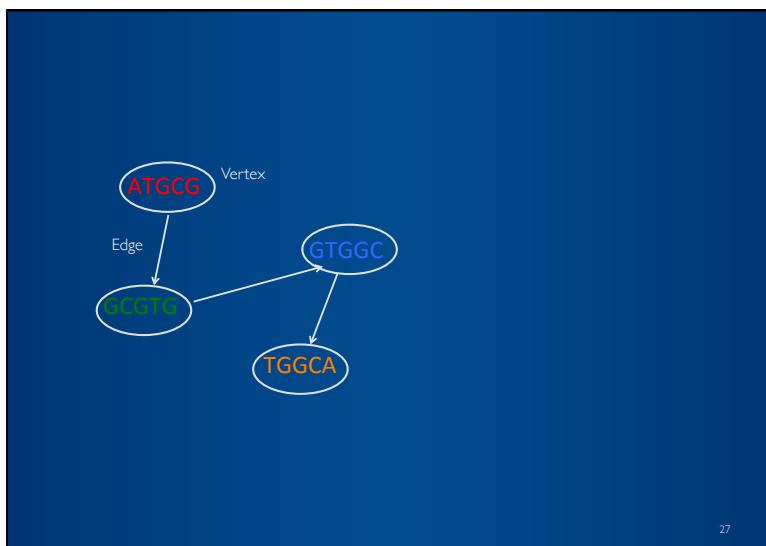
24



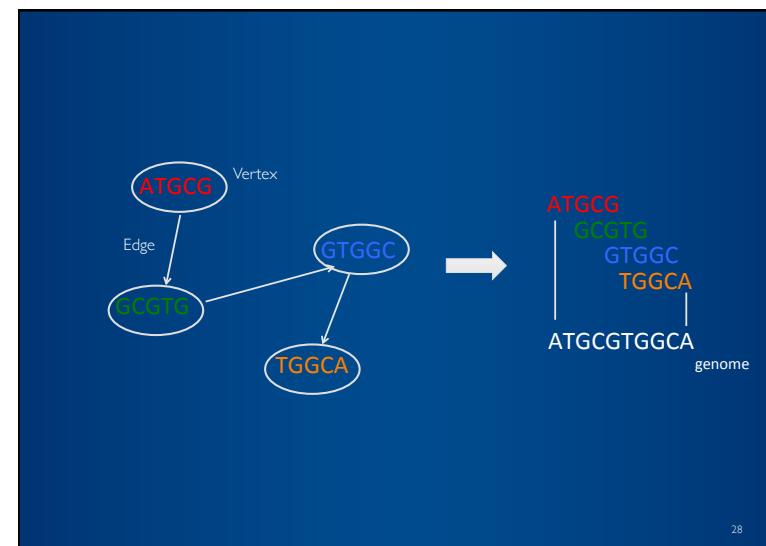
25



26



27



28

## The fragment assembly problem

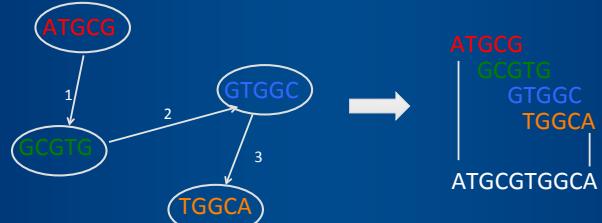
- Given: A set of reads (strings)  $\{s_1, s_2, \dots, s_n\}$
- Do: Determine a large string  $s$  that "best explains" the reads
- What do we mean by "best explains"?
- What assumptions might we require?

29

## Shortest superstring problem

- Objective: Find a string  $s$  such that
  - all reads  $s_1, s_2, \dots, s_n$  are substrings of  $s$
  - $s$  is as short as possible
- Assumptions:
  - Reads are 100% accurate
  - Identical reads must come from the same location on the genome
  - "best" = "simplest"

30



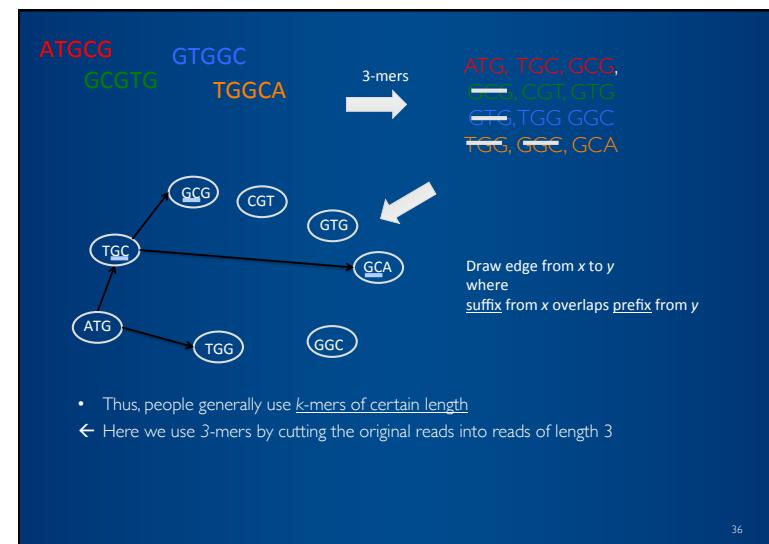
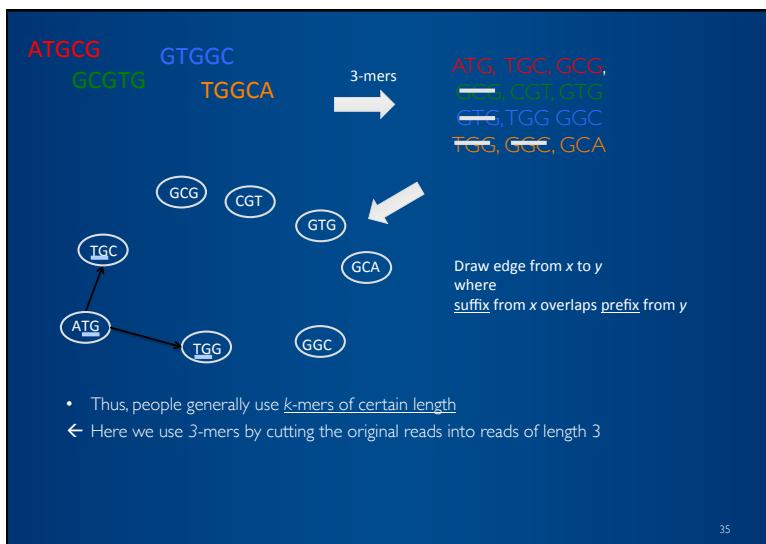
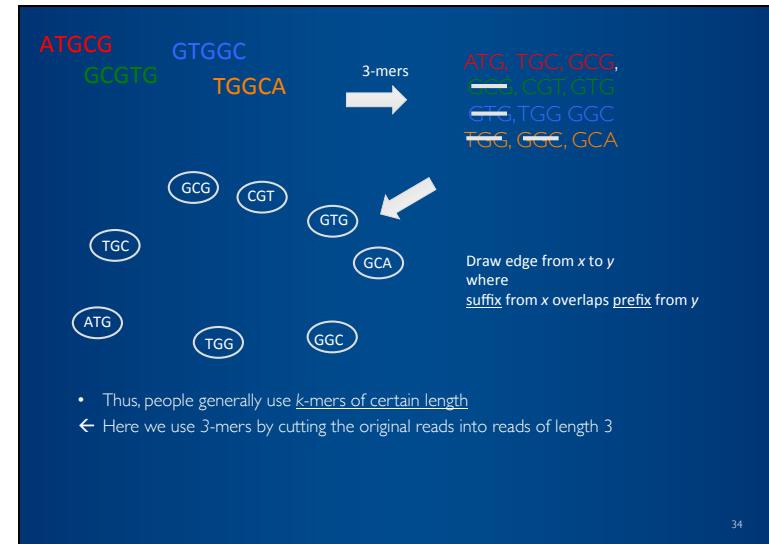
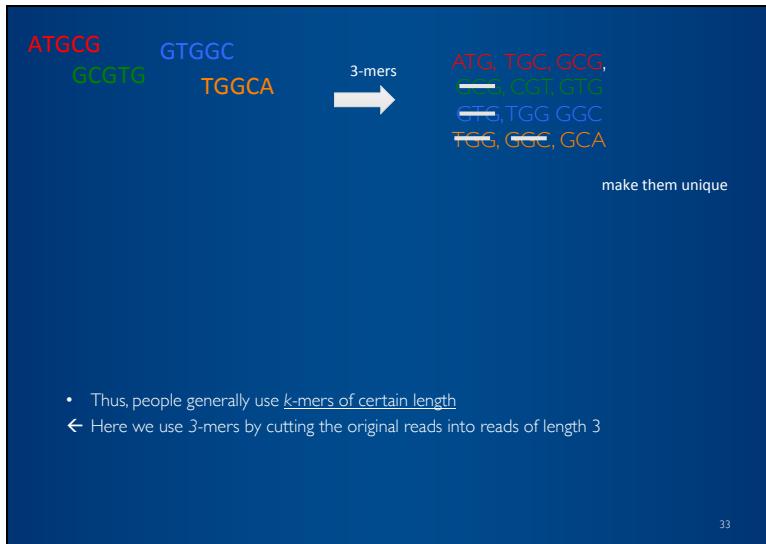
- The assumption is that all substrings are represented
- Even modern sequencers that generate 100nt reads do not cover all possible 100-mers

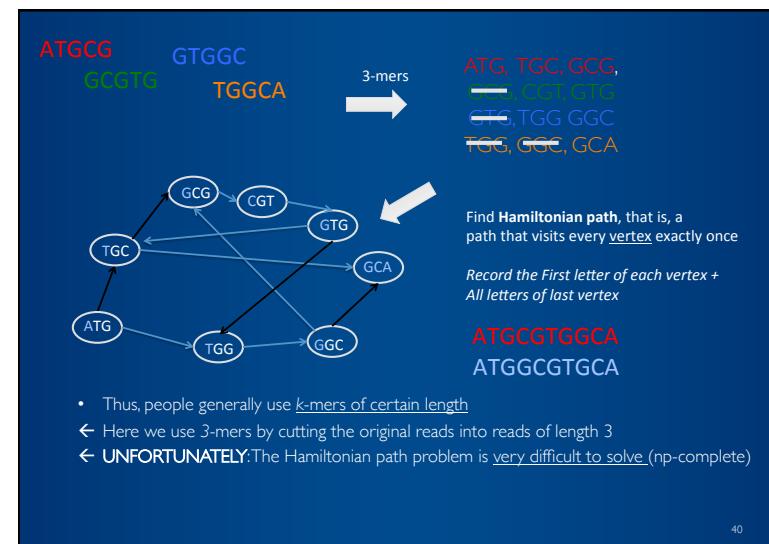
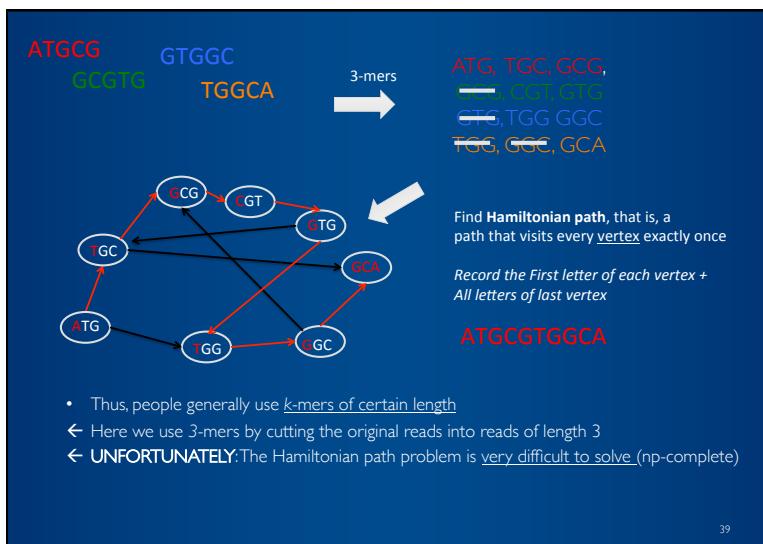
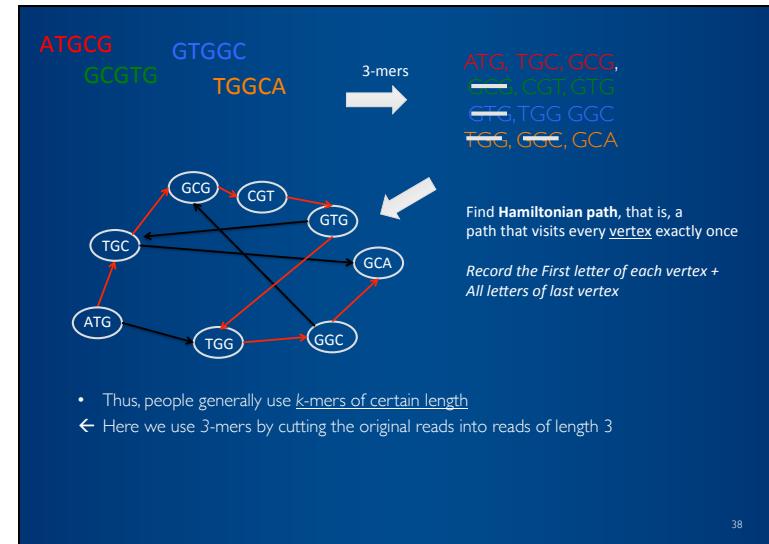
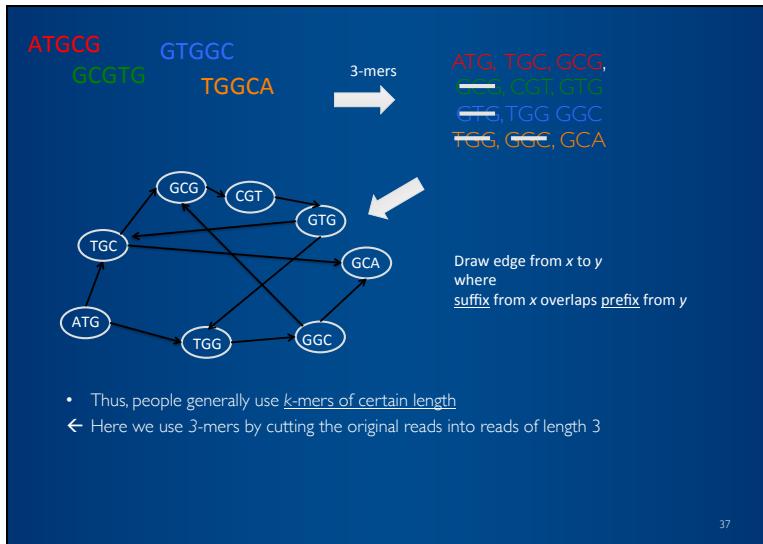
31



- Thus, people generally use  $k$ -mers of certain length
- Here we use 3-mers by cutting the original reads into reads of length 3

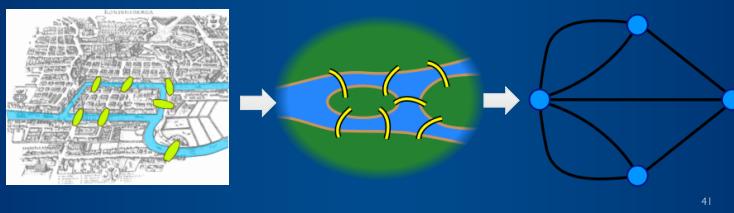
32





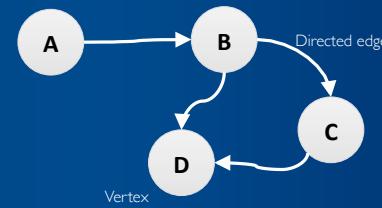
## Seven bridges of Königsberg

- In 1735 Leonhard Euler was presented with the following problem:
  - find a walk through the city that would cross each bridge once and only once
  - He proved that a connected graph with undirected edges contains an Eulerian cycle exactly when every node in the graph has an even number of edges touching it.
  - For the Königsberg Bridge graph, this is not the case because each of the four nodes has an odd number of edges touching it and so the desired stroll through the city does not exist.



## Assembly as a graph theoretical problem

- The degree of a vertex: # of edges connected to it
- outdegree: # of outgoing edges
- indegree: # of ingoing edges
- degree(B)?
- outdegree(B)?
- indegree(D)?



42

## Seven bridges of Königsberg II

- The case of directed graphs is similar:
  - A graph in which indegrees are equal to outdegrees for all nodes is called 'balanced'.
  - Euler's theorem states that a connected directed graph has an Eulerian cycle if and only if it is balanced.
- Mathematically/computationally finding Eulerian path is much easier than Hamiltonian  
→ we need to reformulate our assembly problem

43

We construct a de Bruijn graph:

- edges represent k-mers
- vertices correspond to (k-1)-mers

- Form a node for every distinct prefix or suffix of a k-mer
- Connect vertex x to vertex y with a directed edge if some k-mer (e.g., ATG) has prefix x (e.g., AT) and suffix y (e.g., TG), and label the edge with this k-mer.

k-mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct (k-1)-mers:

44

We construct a ***de Bruijn graph***:

- edges represent k-mers
- vertices correspond to  $(k-1)$ -mers

1. Form a node for every distinct prefix or suffix of a k-mer
2. Connect vertex x to vertex y with a directed edge if some k-mer (e.g., ATG) has prefix x (e.g., AT) and suffix y (e.g., TG), and label the edge with this k-mer.

k-mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct  $(k-1)$ -mers:



45

We construct a ***de Bruijn graph***:

- edges represent k-mers
- vertices correspond to  $(k-1)$ -mers

1. Form a node for every distinct prefix or suffix of a k-mer
2. Connect vertex x to vertex y with a directed edge if some k-mer (e.g., ATG) has prefix x (e.g., AT) and suffix y (e.g., TG), and label the edge with this k-mer.

k-mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct  $(k-1)$ -mers:



46

We construct a ***de Bruijn graph***:

- edges represent k-mers
- vertices correspond to  $(k-1)$ -mers

1. Form a node for every distinct prefix or suffix of a k-mer
2. Connect vertex x to vertex y with a directed edge if some k-mer (e.g., ATG) has prefix x (e.g., AT) and suffix y (e.g., TG), and label the edge with this k-mer.

k-mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct  $(k-1)$ -mers:



47

We construct a ***de Bruijn graph***:

- edges represent k-mers
- vertices correspond to  $(k-1)$ -mers

1. Form a node for every distinct prefix or suffix of a k-mer
2. Connect vertex x to vertex y with a directed edge if some k-mer (e.g., ATG) has prefix x (e.g., AT) and suffix y (e.g., TG), and label the edge with this k-mer.

k-mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct  $(k-1)$ -mers:



48

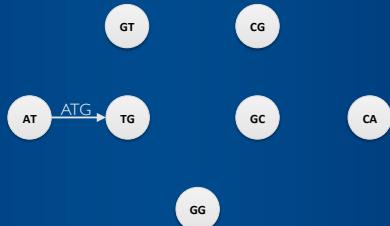
We construct a **de Bruijn graph**:

- edges represent k-mers
- vertices correspond to (k-1)-mers

- Form a node for every distinct prefix or suffix of a k-mer
- Connect vertex x to vertex y with a directed edge if some k-mer (e.g., ATG) has prefix x (e.g., AT) and suffix y (e.g., TG), and label the edge with this k-mer.

k-mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct (k-1)-mers:



49

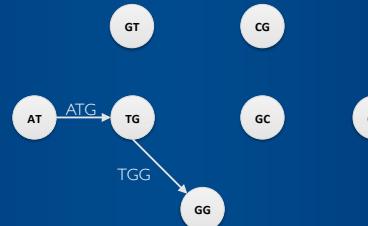
We construct a **de Bruijn graph**:

- edges represent k-mers
- vertices correspond to (k-1)-mers

- Form a node for every distinct prefix or suffix of a k-mer
- Connect vertex x to vertex y with a directed edge if some k-mer (e.g., ATG) has prefix x (e.g., AT) and suffix y (e.g., TG), and label the edge with this k-mer.

k-mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct (k-1)-mers:



50

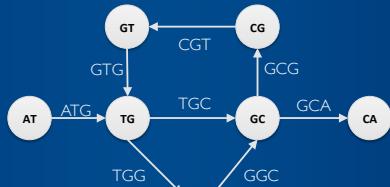
We construct a **de Bruijn graph**:

- edges represent k-mers
- vertices correspond to (k-1)-mers

- Form a node for every distinct prefix or suffix of a k-mer
- Connect vertex x to vertex y with a directed edge if some k-mer (e.g., ATG) has prefix x (e.g., AT) and suffix y (e.g., TG), and label the edge with this k-mer.

k-mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct (k-1)-mers:



51

Can we find a DNA sequence containing all k-mers?

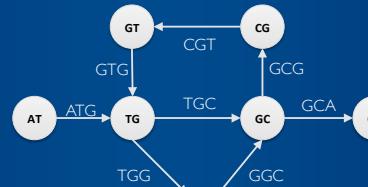
→ In a de Bruijn graph, can we find a path that visits every edge of the graph exactly once?

→ Eulerian path

- a vertex v is semibalanced if  $|\text{indegree}(v) - \text{outdegree}(v)| = 1$
- a connected graph has an Eulerian path if and only if it contains at most two semibalanced vertices

k-mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct (k-1)-mers:



52

Can we find a DNA sequence containing all  $k$ -mers?

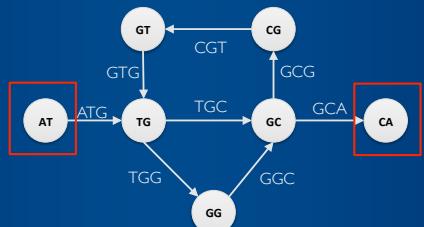
→ In a de Bruijn graph, can we find a path that visits every edge of the graph exactly once?

→ Eulerian path

- a vertex  $v$  is semibalanced if  $|indegree(v) - outdegree(v)| = 1$
- a connected graph has an Eulerian path if and only if it contains at most **two** semibalanced vertices

$k$ -mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct  $(k-1)$ -mers:



53

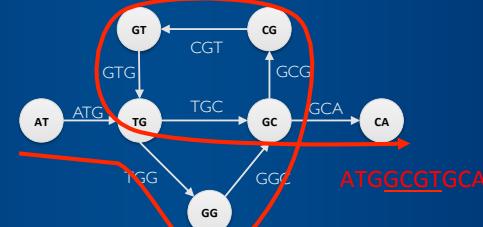
Can we find a DNA sequence containing all  $k$ -mers?

→ In a de Bruijn graph, can we find a path that visits every edge of the graph exactly once?

→ Eulerian path

$k$ -mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct  $(k-1)$ -mers:



54

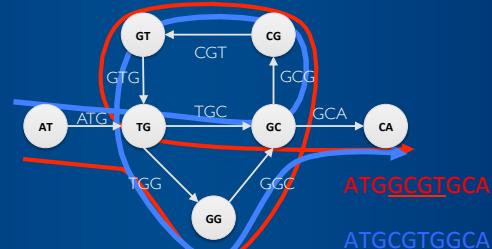
Can we find a DNA sequence containing all  $k$ -mers?

→ In a de Bruijn graph, can we find a path that visits every edge of the graph exactly once?

→ Eulerian path

$k$ -mers: ATG, TGG, TGC, GTG, GGC, GCA, GCG, CGT

Distinct  $(k-1)$ -mers:



55

## Underlying assumptions

- Four hidden assumptions that do **not** hold for next-generation sequencing  
We took for granted that:

1. we can generate all  $k$ -mers present in the genome
2. all  $k$ -mers are error free
3. each  $k$ -mer appears at most once in the genome
4. the genome consists of a single chromosome

56

## Underlying assumptions

- Four hidden assumptions that do not hold for next-generation sequencing  
We took for granted that:
  - we can generate all k-mers present in the genome
  - all k-mers are error free
- That is the reason that we do not choose the longest possible k-mer
- The smaller the k-mer the higher the possibility that we see all k-mers
- Errors:

ATGGC**G**TGCA

Mostly unaffected k-mers  
100% affected k-mers

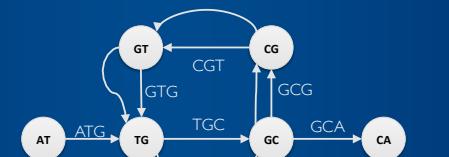
57

Each k-mer appears at most once in the genome → repeats

- This is most often not true
- This is known as k-mer multiplicity

k-mers:  
ATG, GCA,  
TGC, **TGC**, GTG, **GTG**, GCG, **GCG**, CGT, **CGT**

Distinct (k-1)-mers:



**ATGGCTGCGTGCA**

58

## Questions?

### References

How to apply de Bruijn graphs to genome assembly. Phillip E C Compeau, Pavel A Pevzner & Glenn Tesler. Nature Biotechnology 29, 987– 991 (2011) doi:10.1038/nbt.2023 Published online 08 November 2011

Sequence Assembly, Lecture by Mark Craven (craven@biostat.wisc.edu), BMI/CS 576 (www.biostat.wisc.edu/bmi576/), Fall 2011

### Assignment : Due date: September 29<sup>th</sup>

Write a one paragraph summary of how de Bruijn graph assemblers can assemble genome sequences from short read sequences such as Illumina. This should be written at a conceptual level (explanatory).

Sebastian Schmeier  
s.schmeier@massey.ac.nz

