

Genome Assembly

Michael Schatz

Sept 7, 2022

Lecture 3: Applied Comparative Genomics



Assignment 1: Chromosome Structures

Due Wednesday Sept 7 by 11:59pm

The screenshot shows a GitHub repository page for 'assignment1'. The repository contains two files: 'wheat.chrom.sizes' and 'yeast.chrom.sizes', both added as draft assignments 20 minutes ago. The main content of the page is the 'README.md' file, which includes:

Assignment 1: Chromosome Structures

Assignment Date: Wednesday, August 31, 2022
Due Date: Wednesday, Sept. 7, 2022 @ 11:59pm

Assignment Overview

In this assignment you will profile the overall structure of the genomes of several important species and then study human chromosome 22 in more detail. As a reminder, any questions about the assignment should be posted to [Piazza](#).

Question 1: Chromosome structures [10 pts]

Download the chromosome size files for the following genomes (Note these have been preprocessed to only include main chromosomes):

- Arabidopsis thaliana (TAIR10) - An important plant model species [\[info\]](#)
- Tomato (Solanum lycopersicum v4.00) - One of the most important food crops [\[info\]](#)
- E. coli (Escherichia coli K12) - One of the most commonly studied bacteria [\[info\]](#)
- Fruit Fly (Drosophila melanogaster, dm6) - One of the most important model species for genetics [\[info\]](#)
- Human (hg38) - us :) [\[info\]](#)
- Wheat (Triticum aestivum, IWGSC) - The food crop which takes up the largest land area [\[info\]](#)
- Worm (Caenorhabditis elegans, ce10) - One of the most important animal model species [\[info\]](#)
- Yeast (Saccharomyces cerevisiae, sacCer3) - an important eukaryotic model species, also good for bread and beer [\[info\]](#)

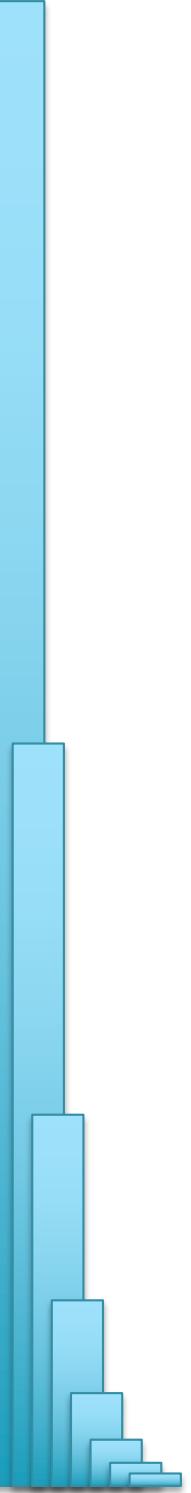
Using these files, make a table with the following information per species:

- Question 1.1. Total genome size
- Question 1.2. Number of chromosomes
- Question 1.3. Largest chromosome size and name
- Question 1.4. Smallest chromosome size and name
- Question 1.5. Mean chromosome length

Question 2. Coverage simulator [20 pts]

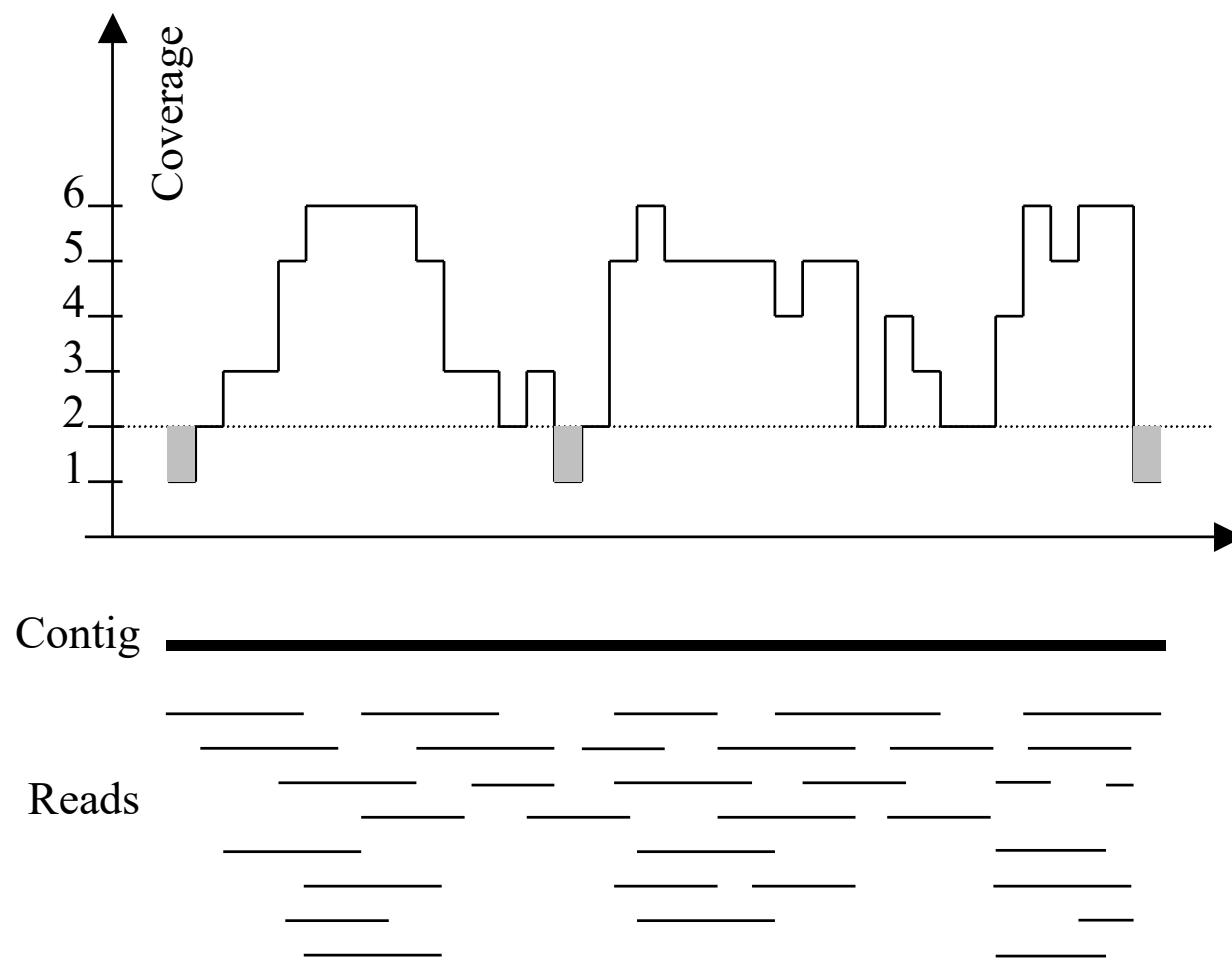
- Question 2.1. How many 100bp reads are needed to sequence a 1Mbp genome to 5x coverage?
- Question 2.2. In the language of your choice, simulate sequencing 5x coverage of a 1Mbp genome with 100bp reads and plot the histogram of coverage. Note you do not need to actually output the sequences of the reads, you can just randomly sample positions in the genome and

<https://github.com/schatzlab/appliedgenomics2022/tree/main/assignments/assignment1>
Check Piazza for questions!



Part I: Recap and Illumina Sequencing

Typical sequencing coverage



Imagine raindrops on a sidewalk

We want to cover the entire sidewalk but each drop costs \$1

If the genome is 10 Mbp, should we sequence 100k 100bp reads?

Poisson Distribution

The probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event.

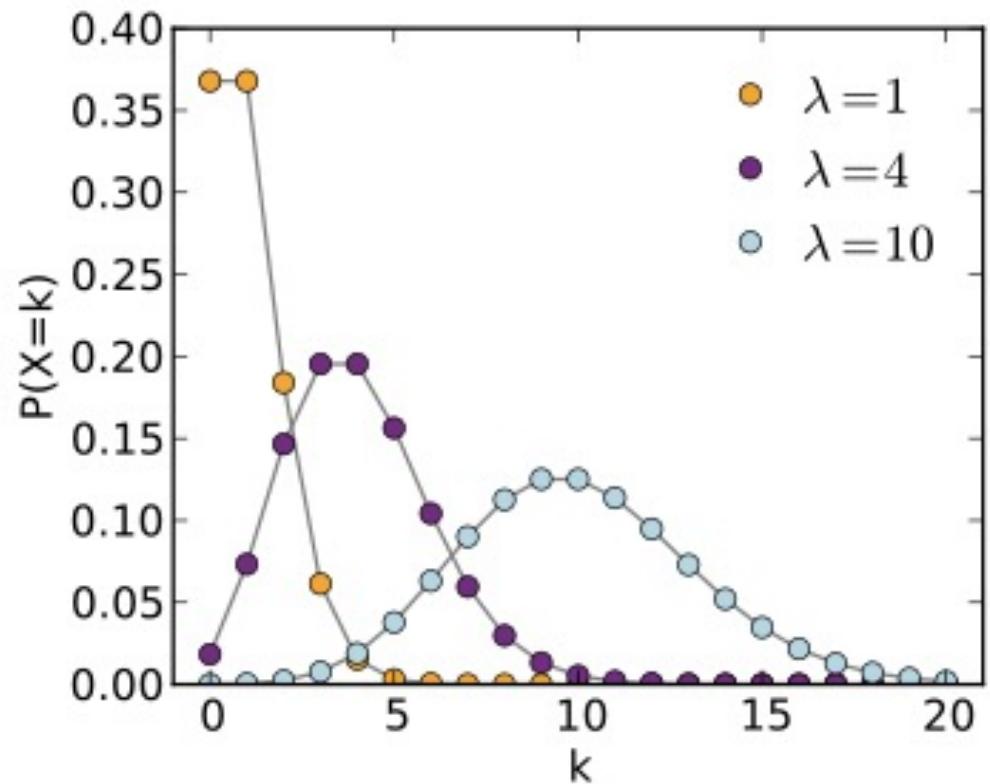
Formulation comes from the limit of the binomial equation

Resembles a normal distribution, but over the positive values, and with only a single parameter.

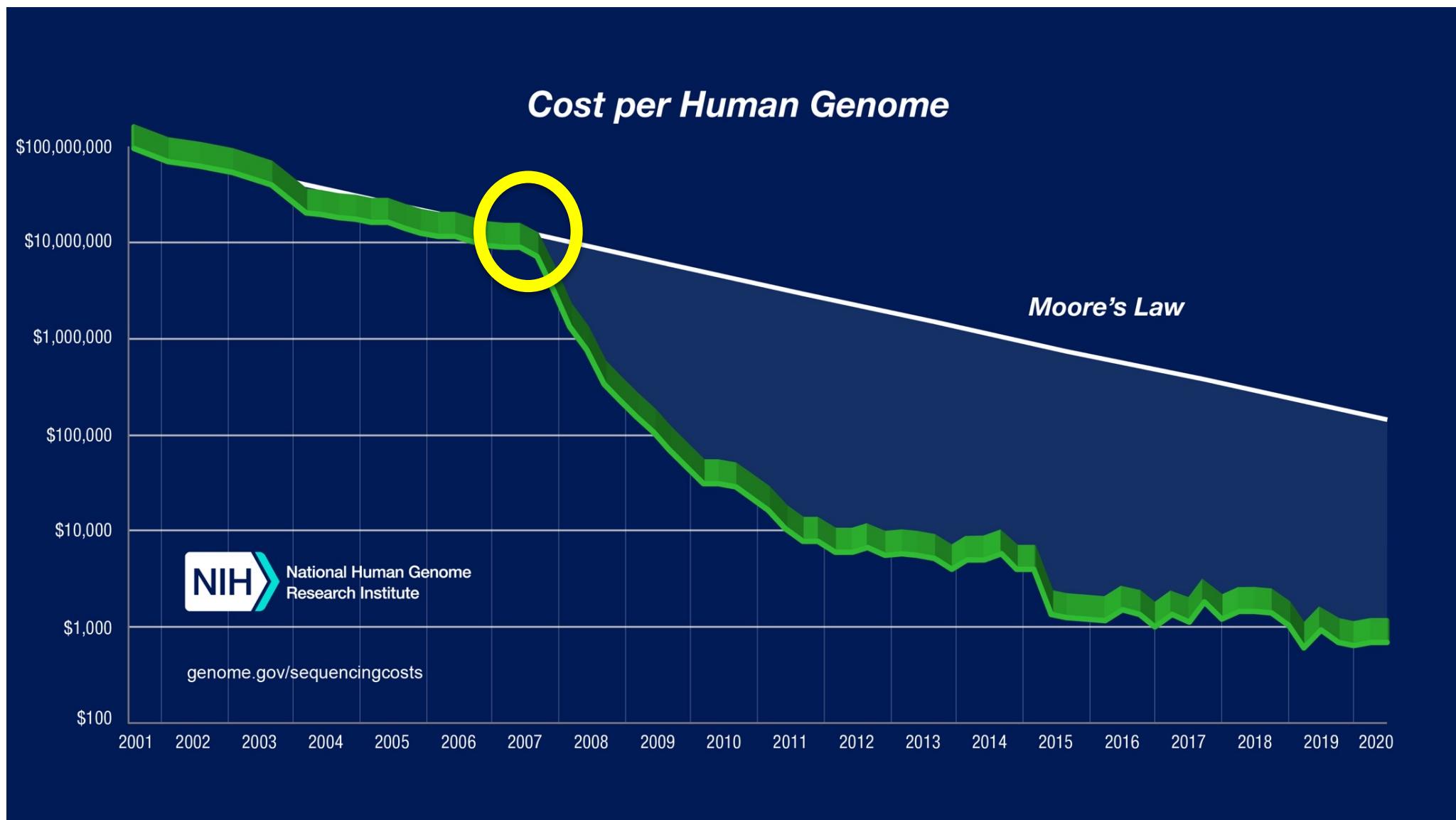
Key properties:

- ***The standard deviation is the square root of the mean.***
- ***For mean > 5, well approximated by a normal distribution***

$$P(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$



Cost per Genome

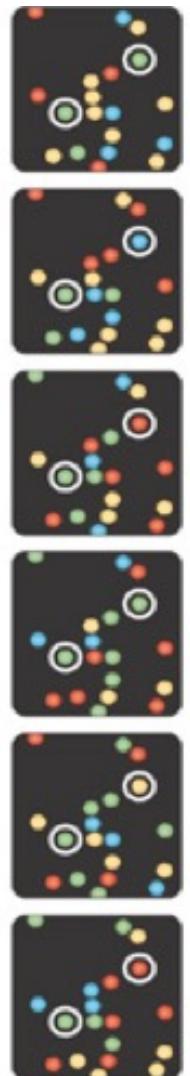
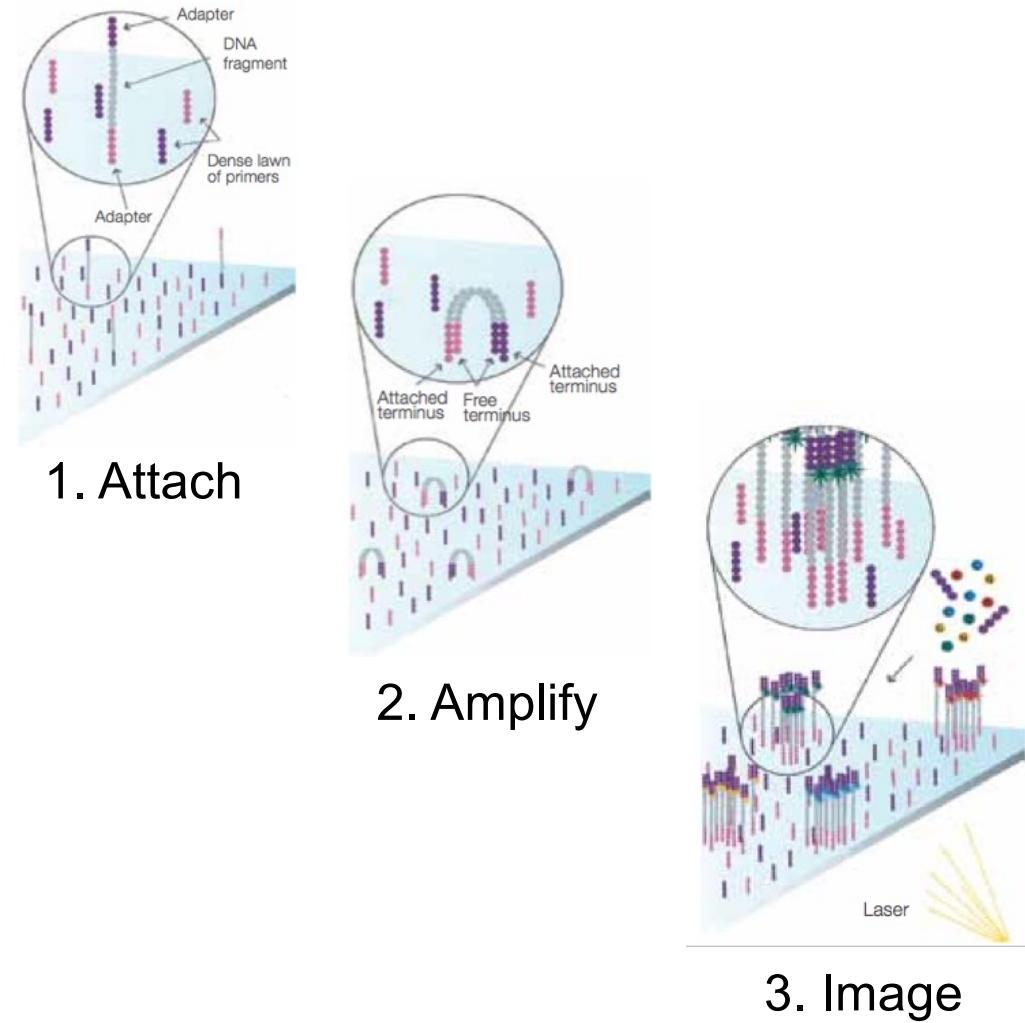


Second Generation Sequencing



Illumina HiSeq 2000
Sequencing by Synthesis

>60Gbp / day

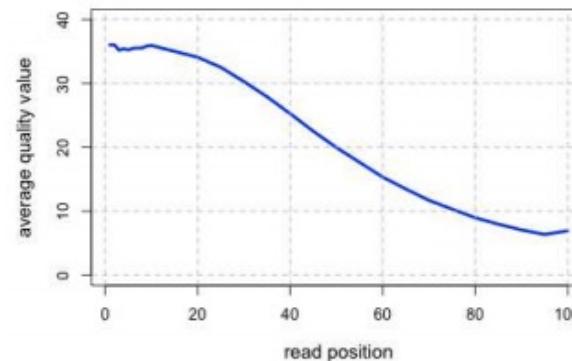


Metzker (2010) Nature Reviews Genetics 11:31-46
<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Illumina Quality

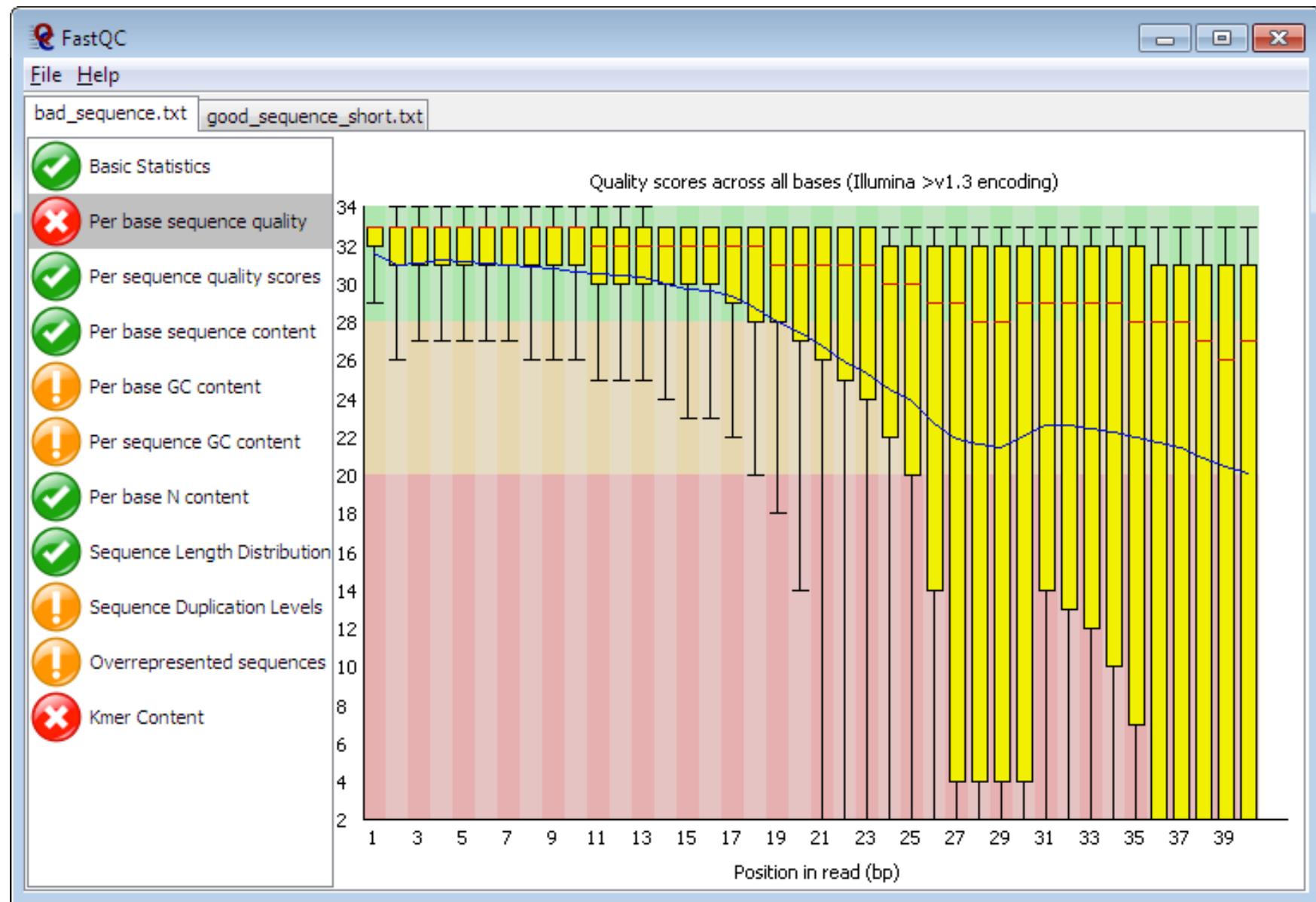
QV	p _{error}
40	1/10000
30	1/1000
20	1/100
10	1/10

$$Q_{\text{sanger}} = -10 \log_{10} p$$



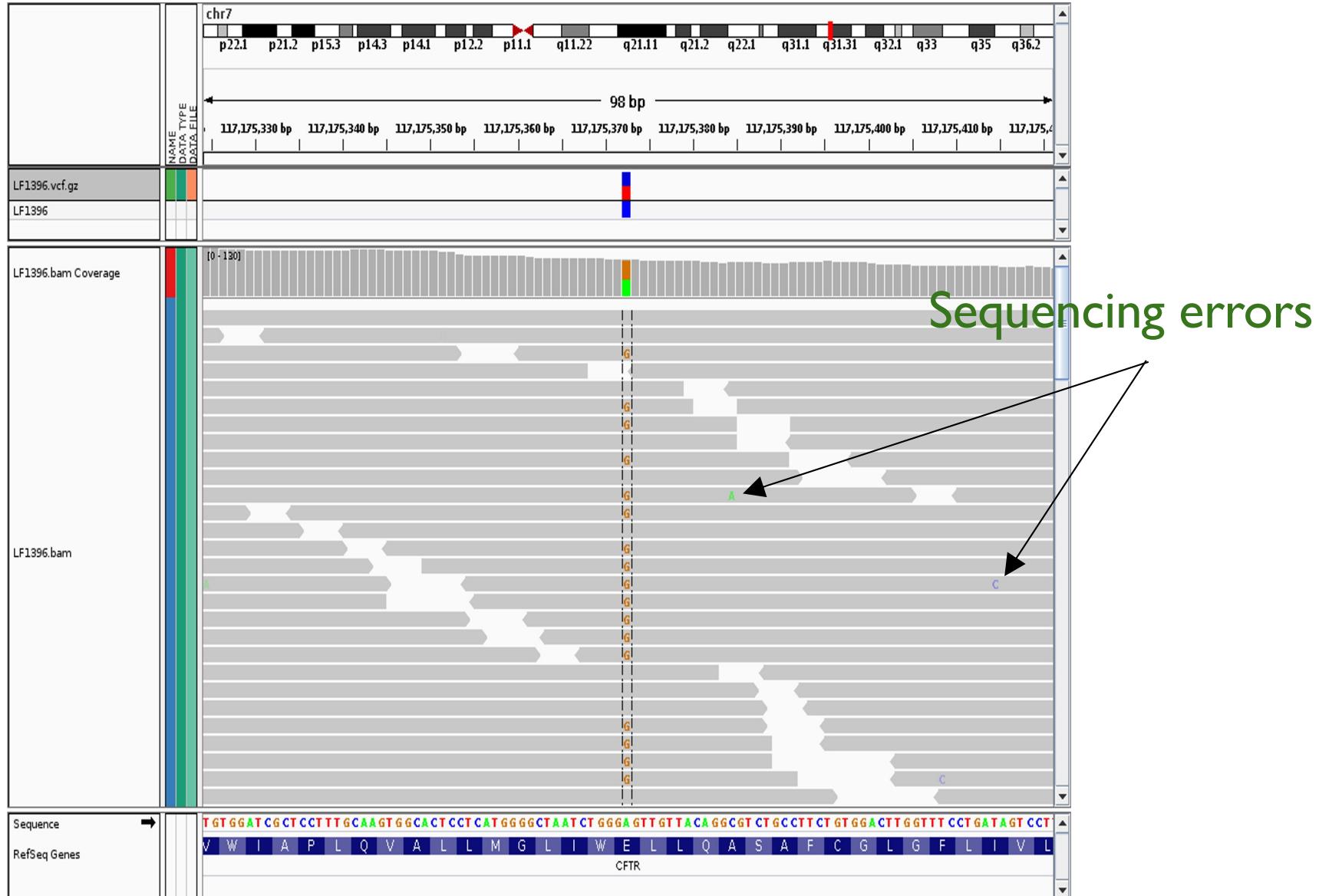
S - Sanger Phred+33, raw reads typically (0, 40)
X - Solexa Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)
with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (**bold**)
(Note: See discussion above).
L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)

FASTQC: Are my data any good?

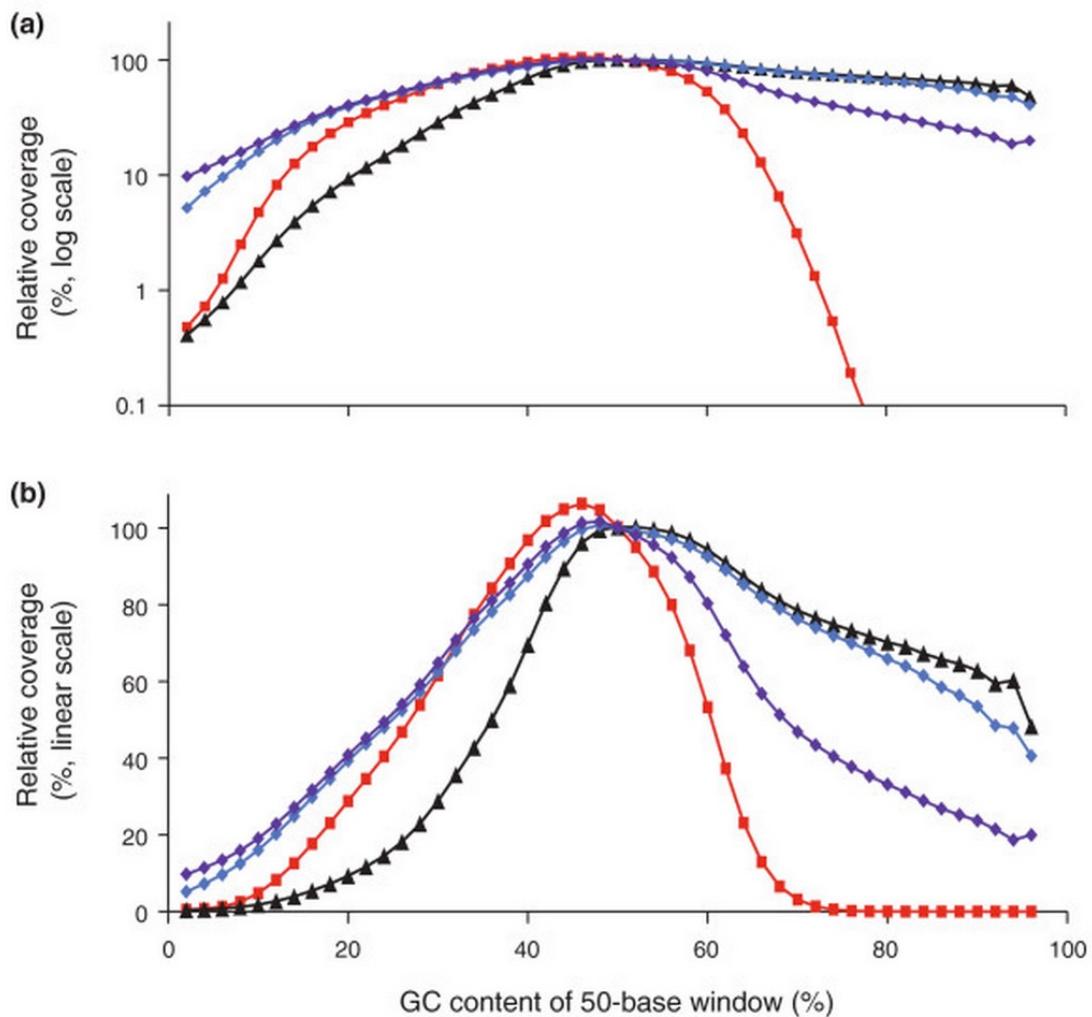


<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Sequencing errors fall out as noise (most of the time)



Beware of GC Biases



Illumina sequencing does not produce uniform coverage over the genome

- Coverage of extremely high or extremely low GC content will have reduced coverage in Illumina sequencing
- Biases primarily introduced during PCR; lower temperatures, slower heating, and fewer rounds minimize biases
- This makes it very difficult to identify variants (SNPs, CNVs, etc) in certain regions of the genome

Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries.

Aird et al. (2011) *Genome Biology*. 12:R18.

Question?

We would love to generate
longer and longer reads with this technology

What can we do?

Paired-end and Mate-pairs

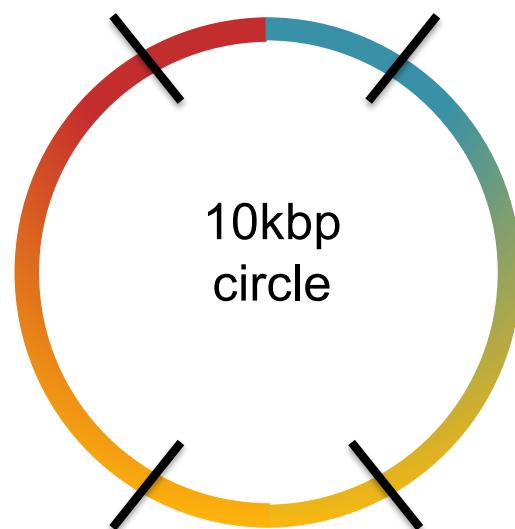
Paired-end sequencing

- Read one end of the molecule, flip, and read the other end
- Generate pair of reads separated by up to 500bp with inward orientation



Mate-pair sequencing

- Circularize long molecules (1-10kbp), shear into fragments, & sequence
- Mate failures create short paired-end reads



2x100 @ ~10kbp (outies)



2x100 @ 300bp (innies)



FASTQ Files



```
@SEQ_ID  
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTTT  
+  
! ' * ( ( ( (****+) ) %%%++ ) (%%% ) . 1***-+* ' ) ) **55CCF>>>>cccccccc65
```

@Identifier
Sequence
+Separator
Quality Values
...

Illumina Sequencing Summary

Advantages:

- Best throughput, accuracy and read length for any 2nd gen. sequencer
- Fast & robust library preparation



Disadvantages:

- Inherent limits to read length (practically, 150bp)
- Some runs are error prone
- Requires amplification, sequences a population of molecules

Illumina HiSeq

~3 billion paired 100bp reads
~600Gb, \$10K, 8 days
(or “rapid run” ~90Gb in 1-2 days)

Illumina X Ten / NovaSeq

~6 billion paired 150bp reads
1.8Tb, <3 days, ~1000 / genome(\$\$)
(or “rapid run” ~90Gb in 1-2 days)

Illumina NextSeq

One human genome in <30 hours



NASDAQ: ILMN

[All](#) [News](#) [Maps](#) [Images](#) [Videos](#) [More](#)

Tools

About 578,000 results (0.78 seconds)



Illumina, Inc.

NASDAQ: ILMN

[Overview](#)[News](#)[Compare](#)[Financials](#)

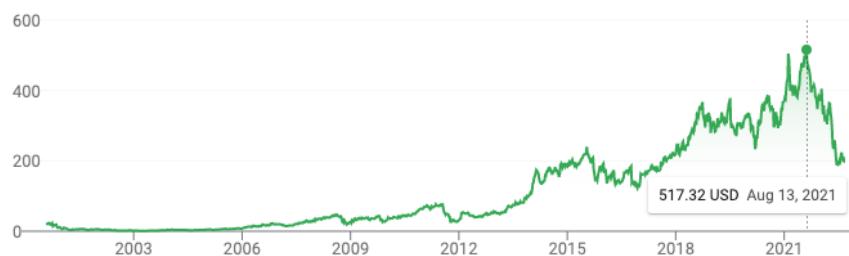
Market Summary > Illumina, Inc.

203.70 USD

[+ Follow](#)

+182.55 (932.33%) ↑ all time

Sep 7, 1:11 PM EDT • Disclaimer

[1D](#) [5D](#) [1M](#) [6M](#) [YTD](#) [1Y](#) [5Y](#) [Max](#)

Open	201.00	Mkt cap	31.75B	CDP score	B
High	202.98	P/E ratio	-	52-wk high	469.87
Low	197.60	Div yield	-	52-wk low	173.45

[More about Illumina, Inc. →](#)

About

[illumina.com](#)

Illumina, Inc. is an American company. Incorporated on April 1, 1998, Illumina develops, manufactures, and markets integrated systems for the analysis of genetic variation and biological function. [Wikipedia](#)

CEO: Francis deSouza (2016–)

Founded: 1998

Headquarters: San Diego, CA

Number of employees: 9,825 (2022)

Revenue: 3.239 billion USD (2020)

Founders: Larry Bock, Anthony Czarnik

Subsidiaries: GRAIL, Solexa, Emedgene Inc, Bluebee Holding B.V., MORE

Outstanding shares

Earnings history

Dividend yield

About 578,000 results (0.78 seconds)



Illumina, Inc.

NASDAQ: ILMN

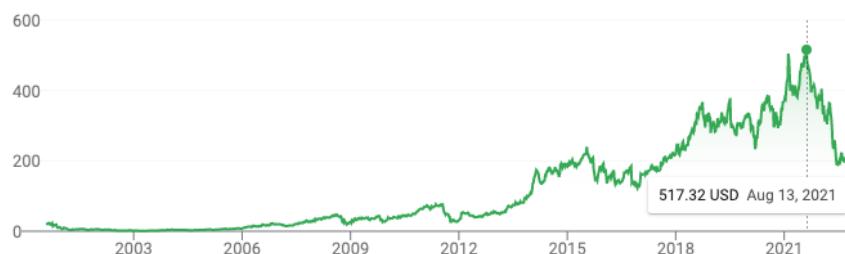
[Overview](#)[News](#)[Compare](#)[Financials](#)

Market Summary > Illumina, Inc.

203.70 USD

+182.55 (932.33%) ↑ all time

Sep 7, 1:11 PM EDT • Disclaimer

[1D](#) [5D](#) [1M](#) [6M](#) [YTD](#) [1Y](#) [5Y](#) [Max](#)

Open	201.00	Mkt cap	31.75B	CDP score	B
High	202.98	P/E ratio	-	52-wk high	469.87
Low	197.60	Div yield	-	52-wk low	173.45

[More about Illumina, Inc. →](#)[Feedback](#)

Inbox - michael.schatz@... | appliedgenomics2022@... | (2019) Illumina Sequencing | NASDAQ: ILMN - Google | Scientists Finish the Human Genome | Other Bookmarks

The New York Times

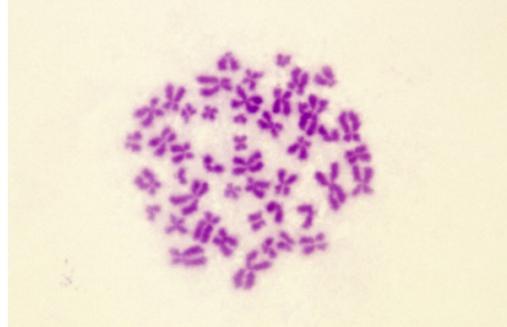
SCIENCE

MATTER

Scientists Finish the Human Genome at Last

The complete genome uncovered more than 100 new genes that are probably functional, and many new variants that may be linked to diseases.

Give this article



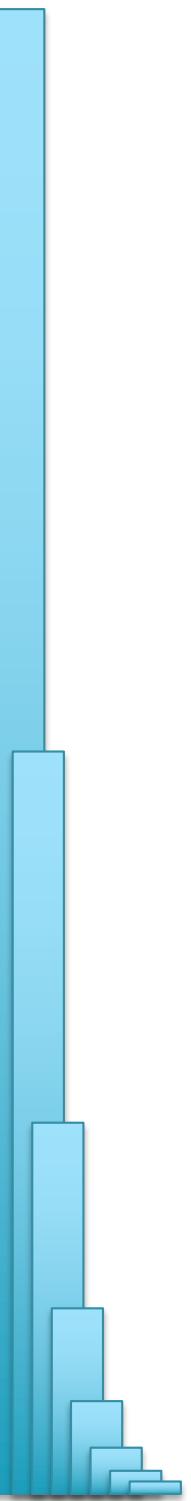
A century ago, scientists knew that genes were spread across 23 pairs of chromosomes. But pinpointing any single gene and deciphering its sequence was a struggle that could have consumed a career. Michael Abbey/Science Source

By Carl Zimmer

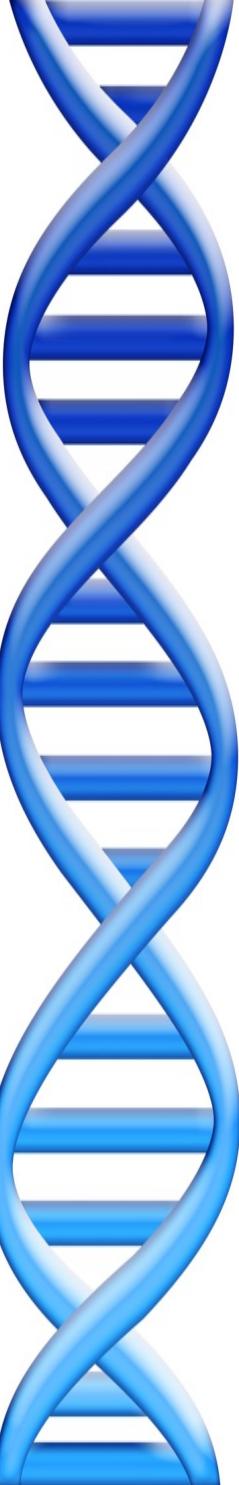
Published July 23, 2021 Updated July 26, 2021

Two decades after the draft sequence of the human genome was unveiled to great fanfare, a team of 99 scientists has finally deciphered the entire thing. They have filled in vast gaps and corrected a long list of errors in previous versions, giving us a new view of our DNA.

The consortium has posted six papers online in recent weeks in which they describe the full genome. These hard-sought data, now under review by scientific journals, will give scientists a deeper understanding of how DNA influences risks of disease, the scientists say, and how cells keep it neatly organized.



Part 2: De novo genome assembly



Outline

1. ***Assembly theory***

- Assembly by analogy

2. ***Practical Issues***

- Coverage, read length, errors, and repeats

3. ***Next-next-gen Assembly***

- Canu: recommended for PacBio/ONT project

4. ***Whole Genome Alignment***

- MUMmer recommended

Shredded Book Reconstruction

- Dickens accidentally shreds the first printing of A Tale of Two Cities
 - Text printed on 5 long spools

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 best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

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 best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

- How can he reconstruct the text?

- 5 copies \times 138,656 words / 5 words per fragment = 138k fragments
- The short fragments from every copy are mixed together
- Some fragments are identical

It was the best of
age of wisdom, it was
best of times, it was
it was the age of
it was the age of
it was the worst of
of times, it was the
of times, it was the
of wisdom, it was the
the age of wisdom, it
the best of times, it
the worst of times, it
times, it was the age
times, it was the worst
was the age of wisdom,
was the age of foolishness,
was the best of times,
was the worst of times,
wisdom, it was the age
worst of times, it was

Greedy Reconstruction

It was the best of
was the best of times,
the best of times, it
best of times, it was
of times, it was the
of times, it was the
times, it was the worst
times, it was the age

The repeated sequence make the correct reconstruction ambiguous

- It was the best of times, it was the [worst/age]

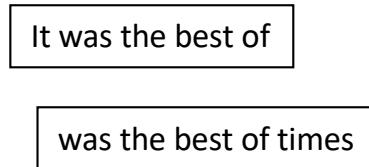
Model the assembly problem as a graph problem

How long will it take to compute the overlaps?

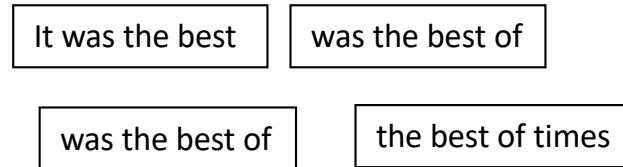
de Bruijn Graph Construction

- $G_k = (V, E)$
 - V = Length- k sub-fragments
 - E = Directed edges between consecutive sub-fragments
 - Sub-fragments overlap by $k-1$ words

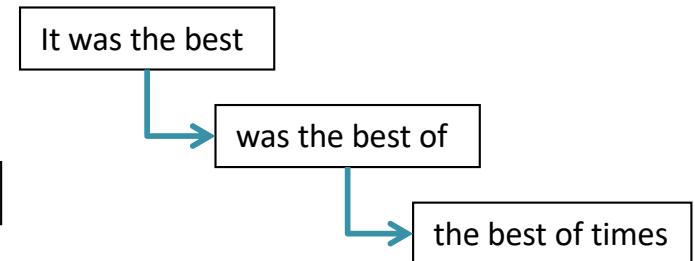
Fragments $|f|=5$



Sub-fragment $k=4$



Directed edges (overlap by $k-1$)



– Overlaps between fragments are implicitly computed

How to pronounce:

https://forvo.com/word/de_briuin/

de Bruijn, 1946
Idury et al., 1995
Pevzner et al., 2001

de Bruijn Graph Assembly

It was the best

was the best of

the best of times,

best of times, it

of times, it was

times, it was the

it was the worst

was the worst of

the worst of times,

worst of times, it

it was the age

was the age of

the age of foolishness

the age of wisdom,

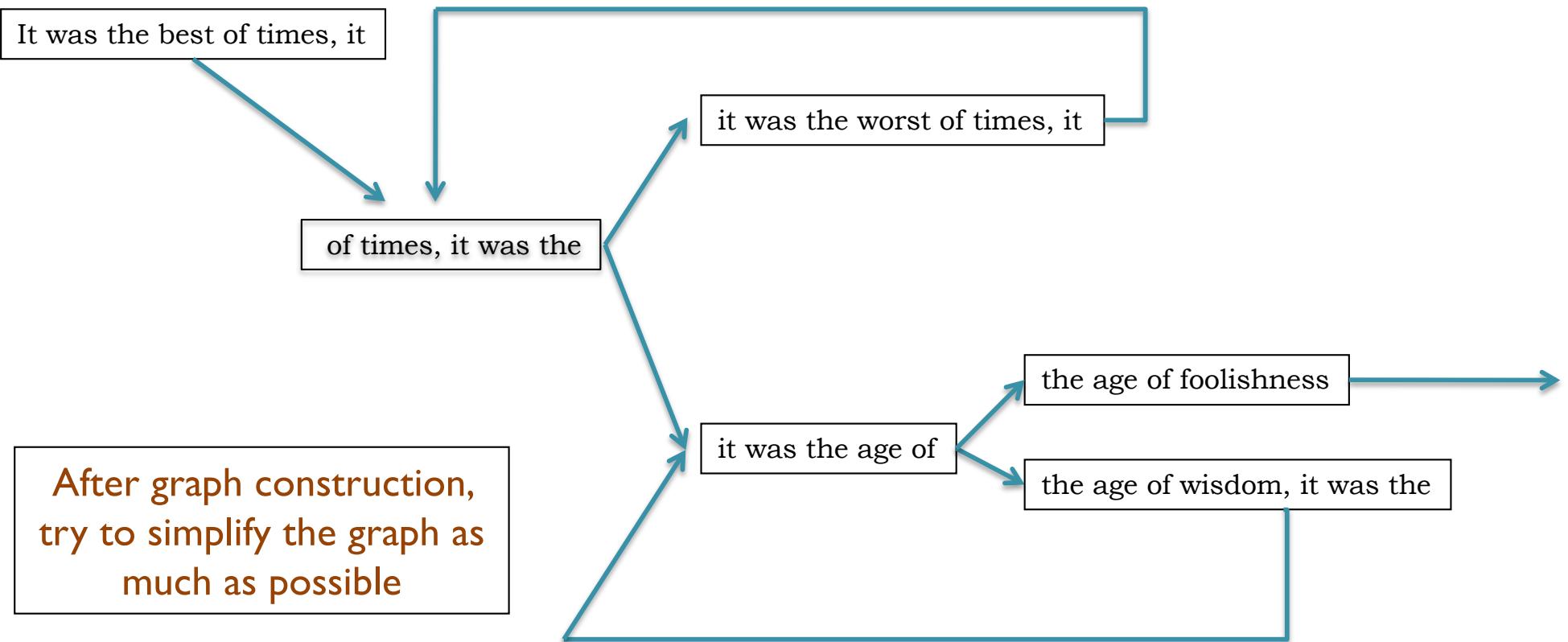
age of wisdom, it

of wisdom, it was

wisdom, it was the

After graph construction,
try to simplify the graph as
much as possible

de Bruijn Graph Assembly



The full tale

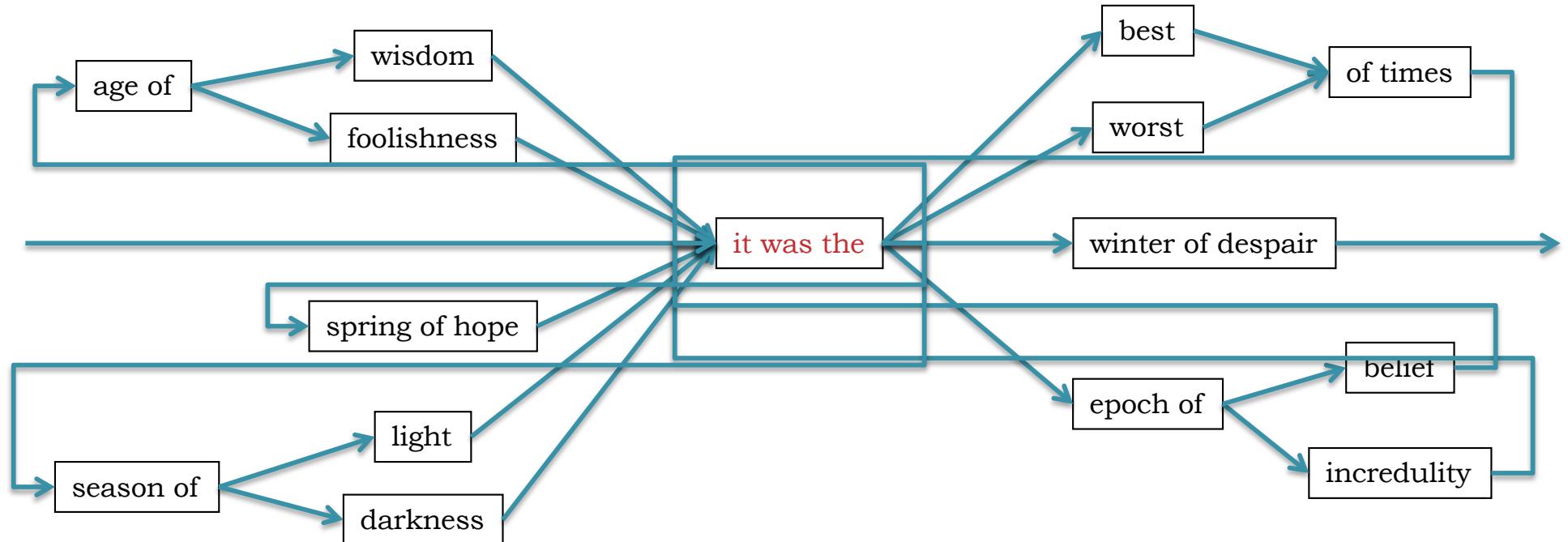
... it was the best of times it was the worst of times ...

... it was the age of wisdom it was the age of foolishness ...

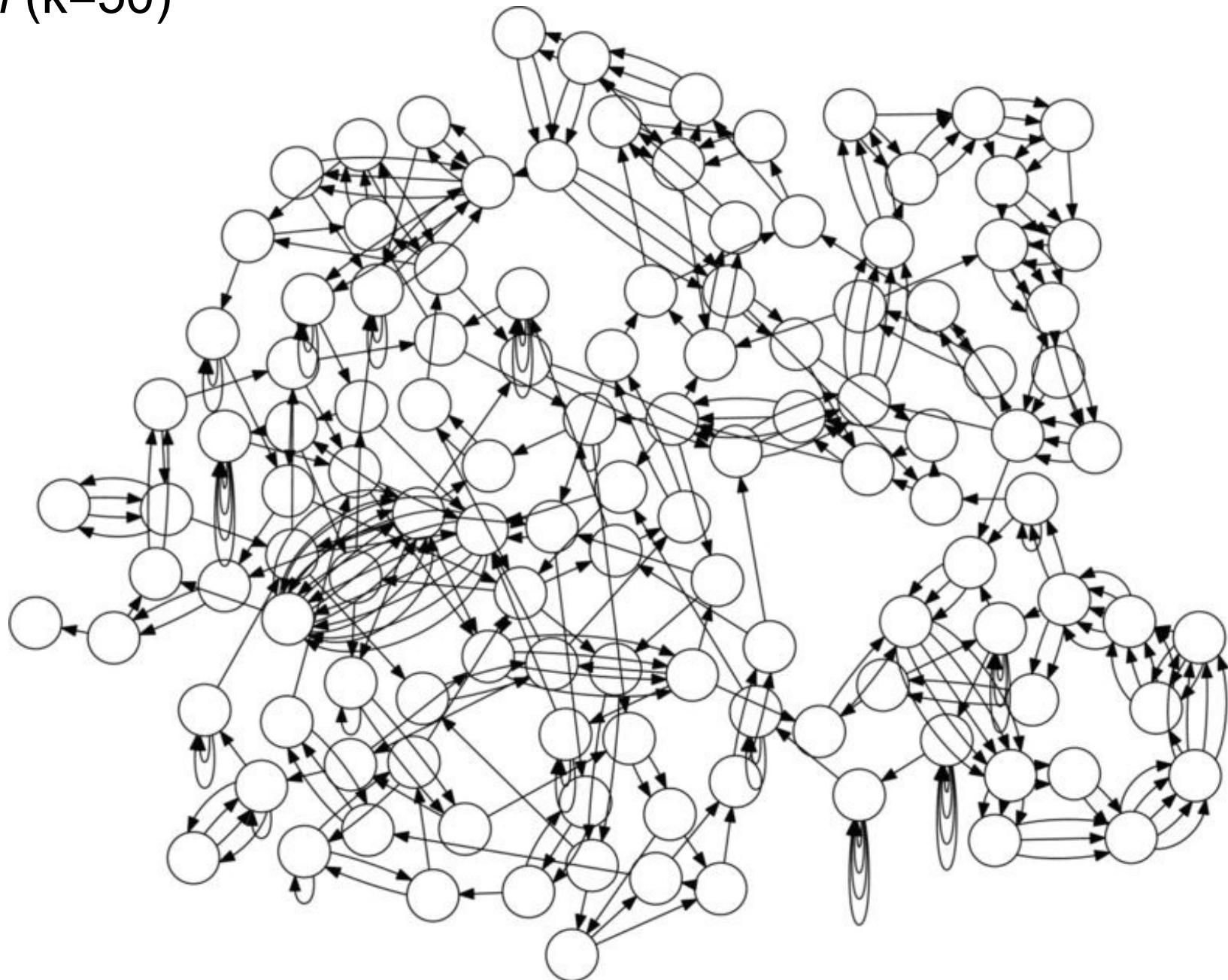
... it was the epoch of belief it was the epoch of incredulity ...

... it was the season of light it was the season of darkness ...

... it was the spring of hope it was the winder of despair ...



E. coli ($k=50$)



Reducing assembly complexity of microbial genomes with single-molecule sequencing
Koren et al (2013) Genome Biology. 14:R101 <https://doi.org/10.1186/gb-2013-14-9-r101>