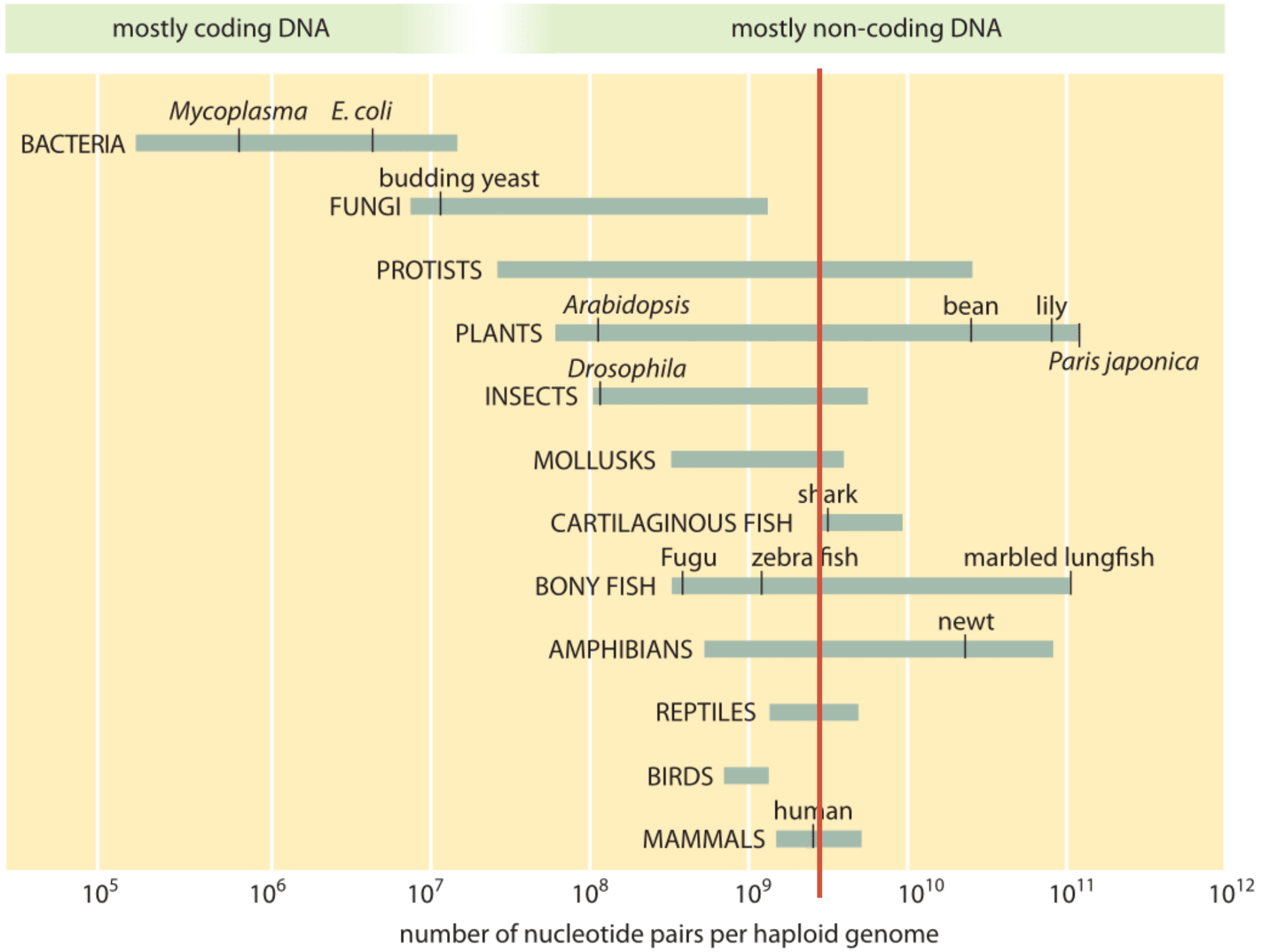# Bioinformatics

## CS300
## Genome Sequencing and Assembly

**Fall 2019**
**Oliver BONHAM-CARTER**

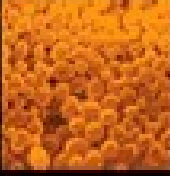# What is a Genome?

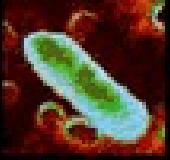- An organism's complete set of DNA, including all of its genes, regulatory regions, non-coding regions, etc.

- An organism's complete set of genetic instructions

mostly coding DNA       mostly non-coding DNA

BACTERIA — *Mycoplasma*   *E. coli*

FUNGI — budding yeast

PROTISTS

PLANTS — *Arabidopsis*   bean   lily   *Paris japonica*

INSECTS — *Drosophila*

MOLLUSKS

CARTILAGINOUS FISH — shark

BONY FISH — Fugu   zebra fish   marbled lungfish

AMPHIBIANS — newt

REPTILES

BIRDS

MAMMALS — human

$10^5$   $10^6$   $10^7$   $10^8$   $10^9$   $10^{10}$   $10^{11}$   $10^{12}$

number of nucleotide pairs per haploid genome

# What Is In a Genome?

| Organism | Number of genes in the genome |
|---|---|
| Myscoplasma genitalium | 517 |
| Saccharomyces cerevisiae | 6,275 |
| Arabidopsis thaliana | ~ 20,000 |
| Caenorhabditis elegans | 19,099 |
| Haemophilus influenzae | 1,743 |
| Drosophila melanogaster | 13,601 |
| Neisseria meningitdis | 2,158 |
| Homo sapiens | 20,000–25,000 |

# Genome Projects

- Goals:
  - Determine complete genome sequence of an organism
  - Annotate protein-coding genes and other important genome-encoded features

# Genome Projects

- Projects:
  - Over 15,000 <u>genome projects</u> in progress or completed



https://www.ncbi.nlm.nih.gov/genome/browse/

# Genome Projects

- Contrast genetic material of populations to determine ancestry



https://en.wikipedia.org/wiki/1000_Genomes_Project#Human_genetic_variation

# Genome Projects: Data

- Locate genes for proteins in sequences.



https://www.ncbi.nlm.nih.gov/genome/browse/

# Genome Projects: Data

- Protein metadata

# Human Genetic Variation

- Having diverse human genetic information helps to spot genetic conditions

- Genetic drift: a random fluctuation in the population frequency of a trait
  - Occurring in subsequent generations and would result in the loss of all variation in the absence of external influence

# Detection By Comparison

- Early detection of genetic problems by being able to compare genomes to some "wild-type" genome.



Hapsburg jaw



Ellis-Van Creveld
syndrome, a sixth finger

# Genome Sequencing

- Bases are recorded as little peaks

- Reads = Small segments of DNA from sequencer machine

- Contigs = Segments of partially combined reads

# Genome Sequencing

- The technology works by "exploding" DNA into smaller, manageable pieces

- Then it recombines pieces (*Reads*) into bigger pieces (*Contigs*)

- And then it combines contigs bigger chunks like a jigsaw puzzle



Cloned genomes

Multiple genomes are sheared into variable sized segments

Unordered sequenced segments

# Shredded Book Reconstruction

- Dickens accidentally shreds first printing of Tale of Two Cities

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

the best of times, it

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

the best of times, it

best of times, it was

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

the best of times, it

best of times, it was

Assembly Parameter:
100% identify across 4 words

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

was the best of times,

the best of times, it

best of times, it was

Assembly Parameter:
100% identify across 4 words

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

It was the best of

was the best of times,

the best of times, it

best of times, it was

Assembly Parameter:
100% identify across 4 words

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

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

# Tale of Two Cities
# Charles Dickens

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 winter of despair, we had everything before us, we had nothing before us, we were all going direct to Heaven, we were all going direct the other way - in short, the period was so far like the present period, that some of its noisiest authorities insisted on its being received, for good or for evil, in the superlative degree of comparison only.

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

Repeats pile up – actual placement of each individual fragment unknown

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

Assembly Parameter:
100% identify across 4 words

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

Repeats pile up – actual placement of each individual fragment unknown

Repeats can cause ambiguity and prevent proper assembly

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 age

times, it was the worst

Assembly Parameter:
100% identify across 4 words

times, it was the age

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

worst of times, it was

wisdom, it was the age

was the worst of times,

was the best of times,

was the age of foolishness,

was the age of wisdom

times, it was the worst

Repeats pile up – actual placement of each individual fragment unknown

Repeats can cause ambiguity and prevent proper assembly

It was the best of

was the best of times,

the best of times, it

best of times, it was

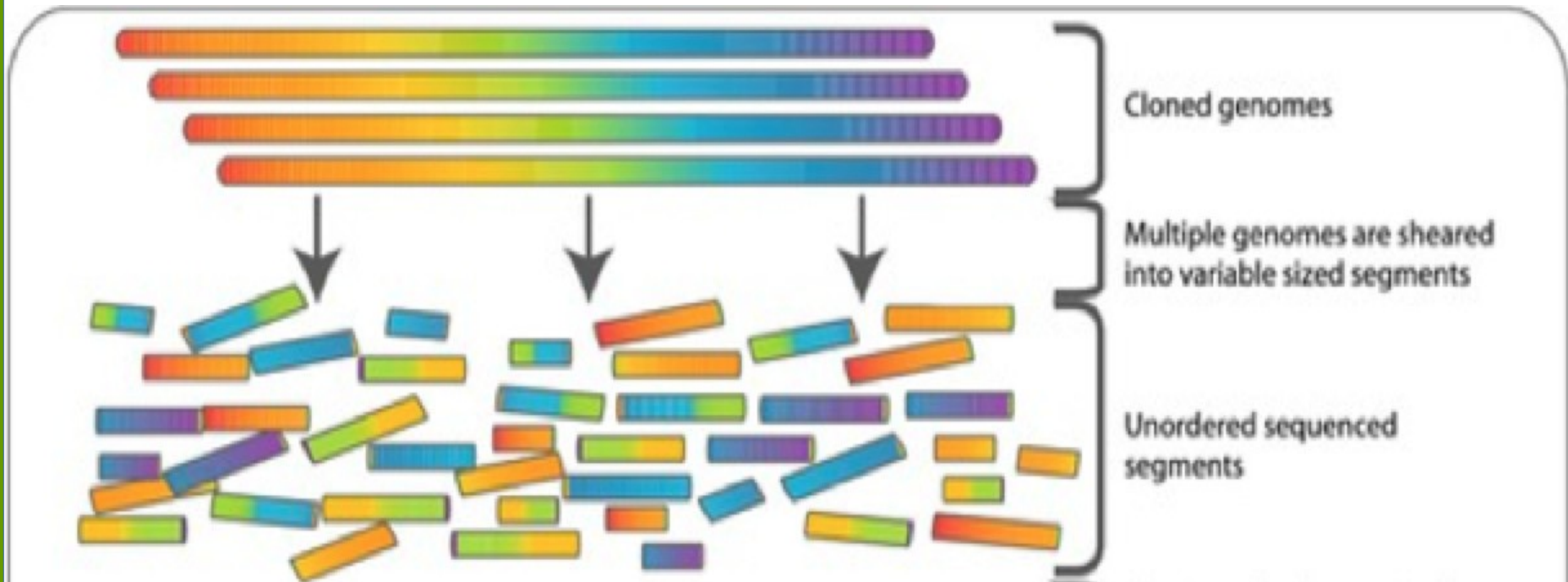of times, it was the

of times, it was the

times, it was the age

times, it was the worst

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

Assembly Parameter:
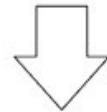100% identify across 4 words

# Genome Sequencing



Cloned genomes

Multiple genomes are sheared into variable sized segments

Unordered sequenced segments

# Coverage
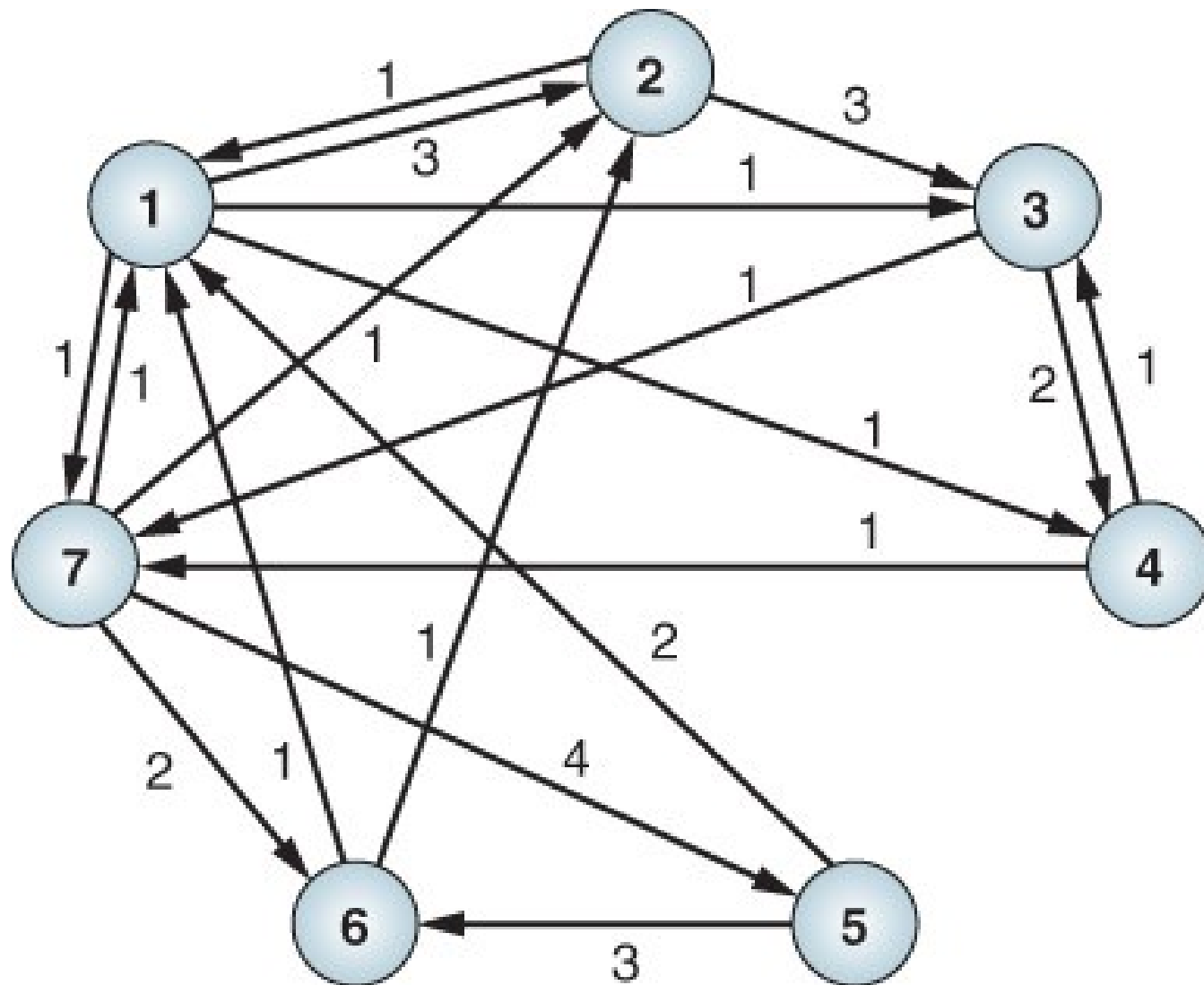
# Assembling a Contig

**Table 8.3  Overlaps for a hypothetical set of sequence reads.**

| Fragments | Overlaps (Length) |
|-----------|-------------------|
| 1. TACCTTG | 2 (3), 3 (1), 4 (1), 7 (1) |
| 2. TTGAT | 1 (1), 3 (3) |
| 3. GATATGG | 4 (2), 7 (1) |
| 4. GGAG | 3 (1), 7 (1) |
| 5. CTCTA | 1 (2), 6 (3) |
| 6. CTAGT | 1 (1), 2 (1) |
| 7. GCTCT | 1 (1), 2 (1), 5 (4), 6 (2) |

For each sequence, we name an overlap with another sequence by number and number of overlaps

# Assembling a Contig:
## graph representation

# Your Turn to Investigate!!!

- Investigate the reads included in your sandbox (file: ***SH_reads.txt***).

- Questions:
  - What is the quote?
  - How did you determine this quote?

GitHub Activity Repository:
https://classroom.github.com/a/DuBJW7yi
Due at 12:30 on 31 Oct. 2019
Create a directory: **act2**
Create work file: **act2/readWork.md**