# GENOME ASSEMBLY
## 19 OCT 15

1

# ANNOUNCEMENTS

# RECAP

- Form Hypothesis — other class
- Collect Seq Data — What type, how much
- QC — FastQC - SolexaQA
- Error Correct — yes
- Trimming — not too harsh
- Normalize — its cool!
- QC — Once again to eval changes
- Assemble — Starting today
    - Genome v. transcriptome
- QC
- mapping
- Post-assembly
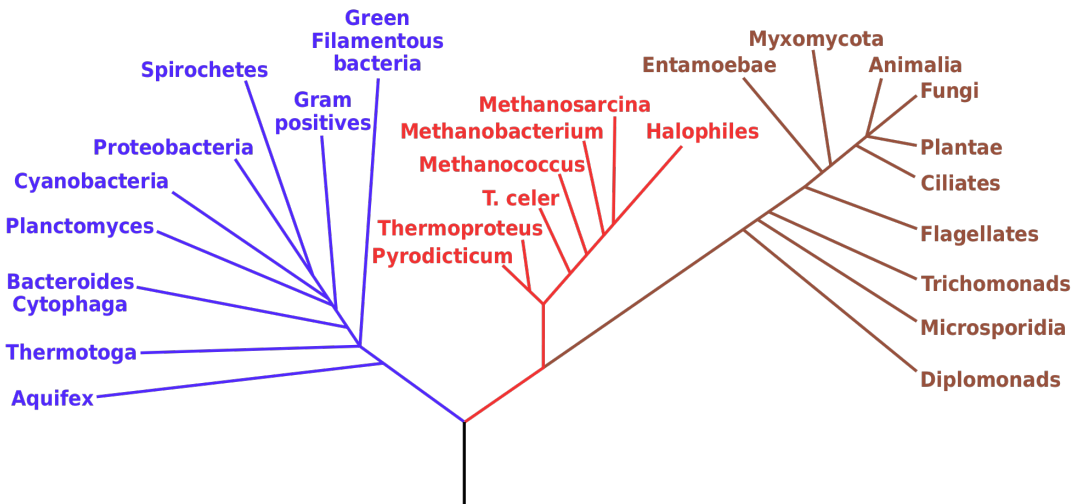    - BLAST/HMM
- Biology — other class

# WHY DO YOU WANT TO ASSEMBLE A GENOME?

# WHAT DO YOU NEED TO ASSEMBLE A GENOME?
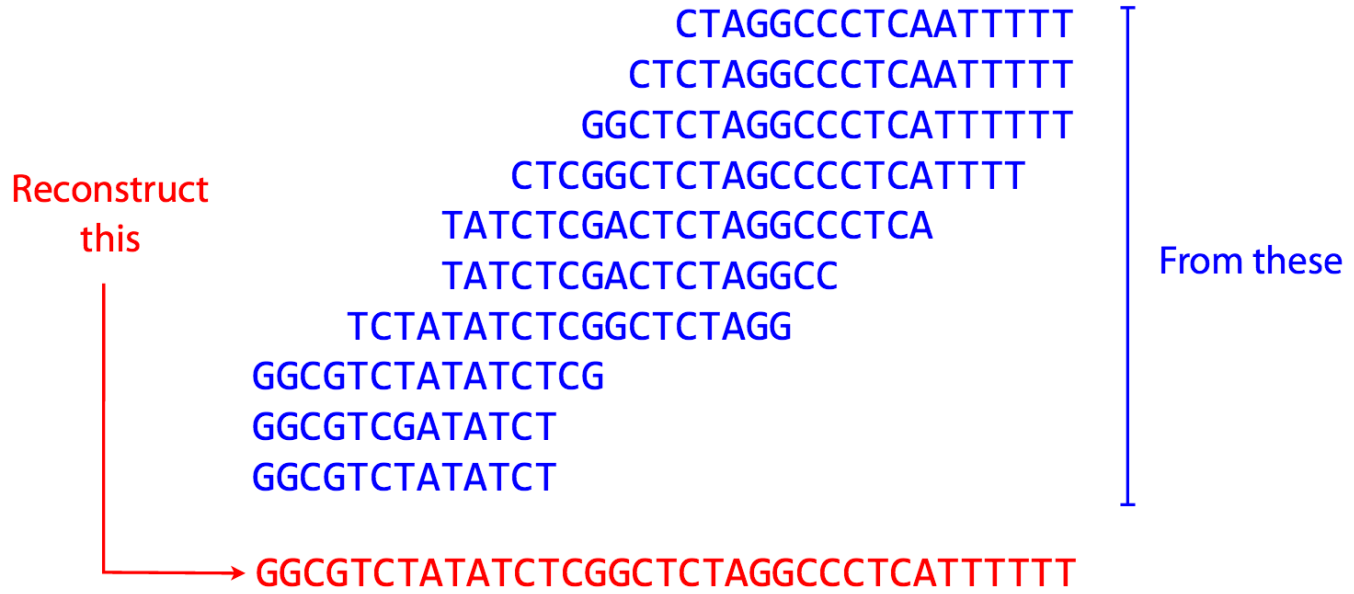
# ASSEMBLE A GENOME? GENERAL STRATEGIES

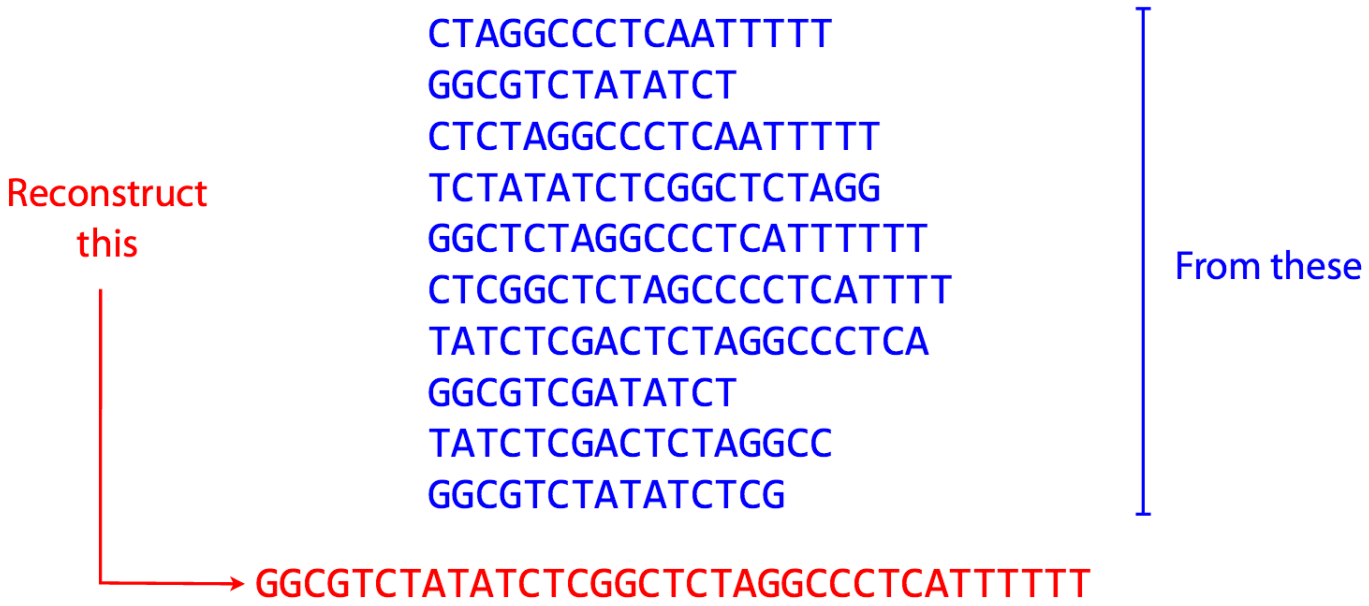| Genome size | Unlimited $$ | Typical |
|---|---|---|
| >10Mb | | |
| 10Mb - 100Mb | | |
| > 100 Mb | | |

# GENOME SIZES

# ASSEMBLY

Assume sequencing produces such a large # fragments that almost all genome positions are *covered* by many fragments...

Reconstruct this

From these

```
                              CTAGGCCCTCAATTTTT
                             CTCTAGGCCCTCAATTTTT
                            GGCTCTAGGCCCTCATTTTTT
                           CTCGGCTCTAGCCCCTCATTTT
                          TATCTCGACTCTAGGCCCTCA
                          TATCTCGACTCTAGGCC
                       TCTATATCTCGGCTCTAGG
                    GGCGTCTATATCTCG
                    GGCGTCGATATCT
                    GGCGTCTATATCT
```

GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT

# ASSEMBLY

...but we don't know what came from where

Reconstruct
this

CTAGGCCCTCAATTTTT
GGCGTCTATATCT
CTCTAGGCCCTCAATTTTT
TCTATATCTCGGCTCTAGG
GGCTCTAGGCCCTCATTTTTT
CTCGGCTCTAGCCCCTCATTTT
TATCTCGACTCTAGGCCCTCA
GGCGTCGATATCT
TATCTCGACTCTAGGCC
GGCGTCTATATCTCG

From these

GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT

# ASSEMBLY

Key term: *coverage*.  Usually it's short for *average coverage*: the average number of reads covering a position in the genome.

```
                    CTAGGCCCTCAATTTTT
                   CTCTAGGCCCTCAATTTTT
                  GGCTCTAGGCCCTCATTTTTT
                 CTCGGCTCTAGCCCCTCATTTT
                TATCTCGACTCTAGGCCCTCA           177 nucleotides
                TATCTCGACTCTAGGCC
              TCTATATCTCGGCTCTAGG
          GGCGTCTATATCTCG
          GGCGTCGATATCT
          GGCGTCTATATCT
          GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT      35 nucleotides
```

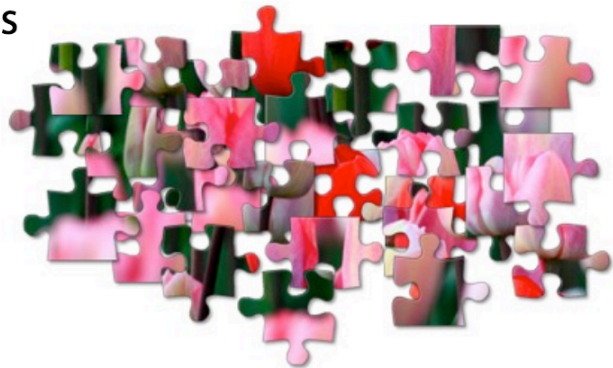Average coverage = 177 / 35 ≈ 7x

# OTHER ASSEMBLY TERMS

Unitig

Contig

scaffold

# ASSEMBLY

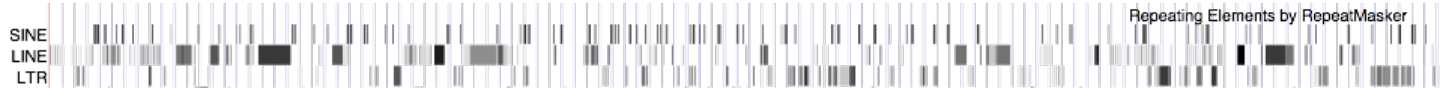- Complicated by:

Reads



+

Reference genome



How to assemble puzzle without the benefit of knowing what the finished product looks like?

Input DNA

# ASSEMBLY

- Complicated by:

# ASSEMBLY

- Workflow:

# ASSEMBLY

- 3 assembly strategies:

# ASSEMBLY

- OLC Assembly

**Overlap** — **Build overlap graph**

Layout — Bundle stretches of the overlap graph into *contigs*

Consensus — Pick most likely nucleotide sequence for each contig

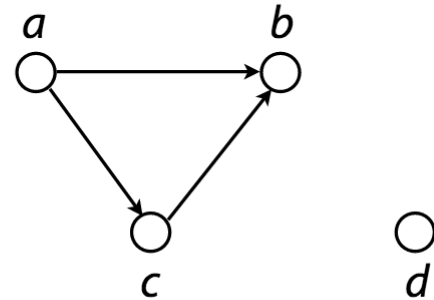# ASSEMBLY

- OLC Assembly: Characteristics

# Assembly

Directed graph $G(V, E)$ consists of set of *vertices, V* and set of *directed edges, E*

Directed edge is an *ordered pair* of vertices.
First is the *source*, second is the *sink*.

Vertex is drawn as a circle

Edge is drawn as a line with an arrow
connecting two circles



Vertex also called *node* or *point*

Edge also called *arc* or *line*

Directed graph also called *digraph*

$V = \{ a, b, c, d \}$
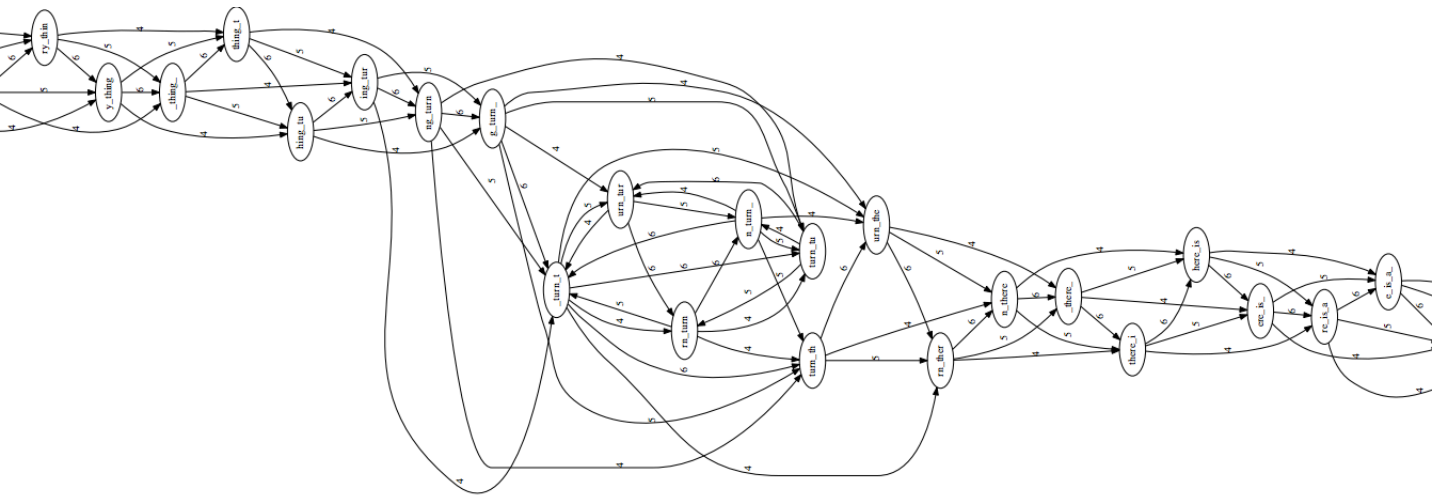
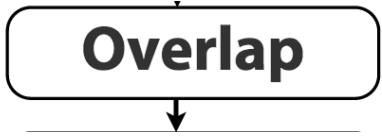$E = \{ (a, b), (a, c), (c, b) \}$

Source    Sink

# ASSEMBLY

**Overlap** Build overlap graph

to_every_thing_turn_turn_turn_there_is_a_season
L=4, k=7

# ASSEMBLY

**Overlap** **Build overlap graph**

Vertices (reads): { *a*: CTCTAGGCC, *b*: GCCCTCAAT, *c*: CAATTTTT }

Edges (overlaps): { (*a*, *b*), (*b*, *c*) }

*a*: CTCTAGGCC →₃ *b*: GCCCTCAAT →₄ *c*: CAATTTTT

```
CTCTAGGCC              GCCCTCAAT
   | | |                  | | | |
   GCCCTCAAT              CAATTTTT
```