# DNA sequencing + Genome sequencing

Goal :



File : >seq1
ACTAGGTA    ] read1
>seq
TGACTAG. ,   ] read2
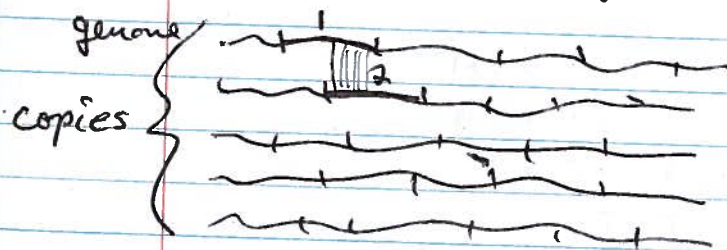⋮

Illumina Sequencing : See video.

Limitation : - Length of a read is limited ( Illumina ≤ 300bp )

---

## Genome Sequencing + Assembly

Goal: Get entire DNA seq. of a genome (long)
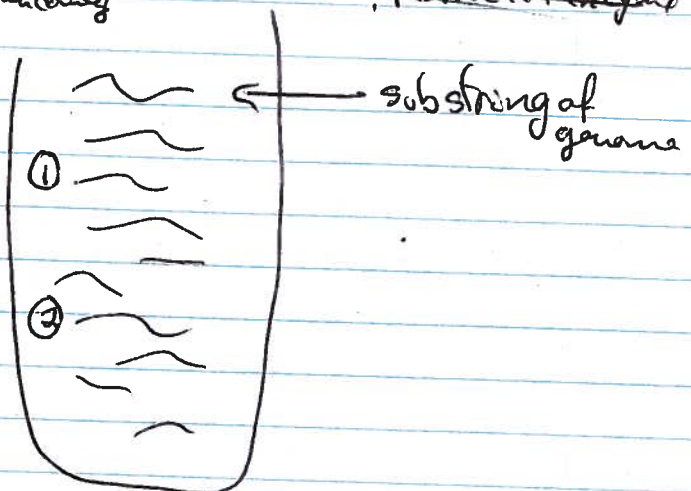
Problem: Seq. machines produce reads that short

### Shotgun sequencing

genome
copies {



① Generate many copies of genome
② Cut seq, into small pieces
   randomly ↳ sonication, restriction enzyme



substring of genome

③ Sequence DNA fragments
   ACTAGG
   TAGCC...
   ⋮

④ Assemble reads into genome

True genome : A C T A G C T T C T T A G C C T T
(unknown)

| | |
|---|---|
| Read 1 | C T T C T T |
| Read 2 | A C T A G C |
| Read 3 | A G C C T T |
| Read 4 | C T T A G C |
| Read 5 | T A G C T T |
| Read 6 | T T C T T A |

Problem: Shortest Superstring Problem

Given: Set of reads $R_1 \ldots R_n$, of length $L = six$

Find: Sequence $G$ such that
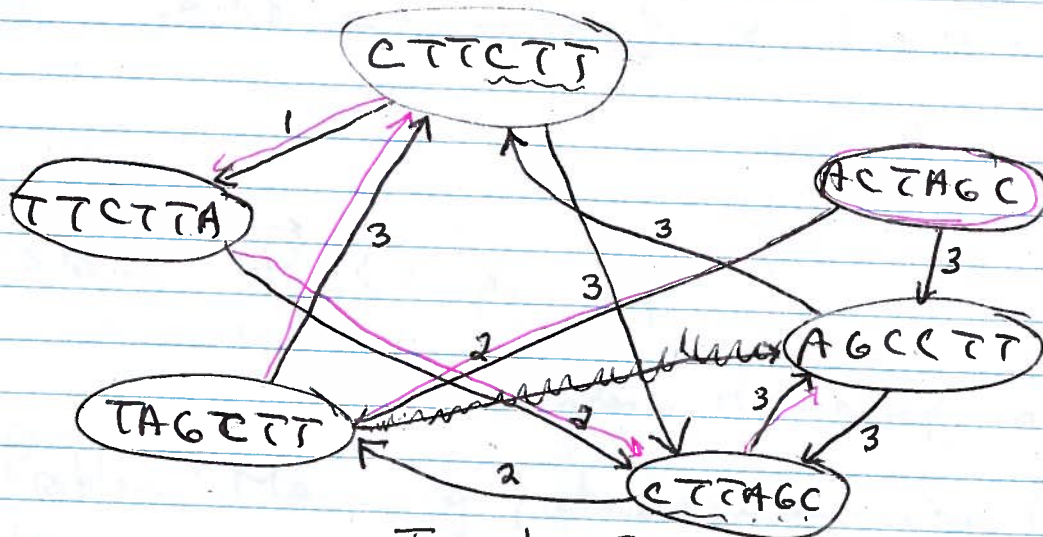  1) Each $R_i$ is a substring of $G$
  2) $G$ is a short as possible

Assumption: No sequencing error.

# Overlap-Layout-Consensus Approach

Build Graph: V: set of reads

E: overlaps btw reads of at least $k$ bases $= 3$



Goal: Find shortest Hamiltonian Path in G

Traveling Salesperson Problem

→ NP-Complete Problems

↓ Smallest total weight

→ Path that visits each vertex exactly once

ACTAGC
  TAGCTT
    CTTCTT
     → TTCTTA
      CTTAGC
        AGCCTT
_____

weight: $2+3+1+2+3$
    $= 11$

Predicted Genome ACTAGCTTCTTAGCCTT