

Genome Assembly Continued

---

---

---

---

---

---

---

---

Genome Assembly Exercise

```
graph TD; A[Download fastq files from server] --> B[Quality assessment using FastQC]; B --> C[Create hash table in Velvet]; C --> D[Assembly genome using Velvet]; D --> E[Visualize data in Tablet];
```

---

---

---

---

---

---

---

---

Sequences and Annotations

**Some common genomics data file formats**

- FASTA: gene and genome sequence format (fa).
- FASTQ: high-throughput sequencing (HTS) data format (fastq).
- GFF: feature format typically general feature format (gff).
- SAM/BAM: HTS read alignment data (sam or bam).

---

---

---

---

---

---

---

---

FASTA Format

**FASTA:** DNA sequence alignment software. The software gave rise to the fasta format, now ubiquitous sequence file format.

```
>sequence1_description
AGCTAGCATGCACTAGCAGTACGATGCGAGTACAGGTAGGAGTAGGGGCTTACGATGCTA
CCCCGGACTACGGGAGTCCCGATTCACGGGGATCGAGGAGGAGCCATGAGGAGATTCATCTTA
TCGAGGAGACTACTACTCTCTCTACTACTACTACTACTTACCCCTCTTAGGGTTCATTAATTGCTGGTAG
GATCGAGGATTCAGAGGATTCGAGGACTGAGGAGACTTACTACTATGAGGAGACTACTT
>sequence2_description
CTCTAGCATGCACTAGCAGTACGATGCGAGTACAGGTAGGAGTAGGGGCTTACGATGCTA
CCCCGGACTACGGGAGTCCCGATTCACGGGGATCGAGGAGGAGCCATGAGGAGATTCATCTTA
TCGAGGAGACTACTACTCTCTCTACTACTACTACTACTTACCCCTCTTAGGGTTCATTAATTGCTGGTAG
GATCTTCTTAGGAGGATTCAGAGGATTCAGGAGGAGTACGAGGAGTACTACTACTATGAGGAGACTACT
TC
* DNA, RNA, or amino acid sequence
```

---

---

---

---

---

---

---

---

---

---

FASTQ Format

```
Index sequence
Read 1 ① ② ③ ④ ⑤ ⑥ ⑦ ⑧ ⑨ ⑩ ⑪ ⑫ ⑬ ⑭ ⑮ ⑯ ⑰ ⑱ ⑲ ⑳ ㉑ ㉒
1 @H04T0FP1.1248.CS50RMACX.5:1101:1241:2095:1:R10 @ATCACG
2 CA0CC0C0C0T0C0TATCC0G0GACTCG0BATTCT0G0G0TC0C0AG0A0CTC0A
3 +
4 CC0FFFFFHHBBLJ3J0BLJJ3J1JJ3JJGG0FFFFFAB0B0HFFHFF#FD0
Read 2 ① ② ③ ④ ⑤ ⑥ ⑦ ⑧ ⑨ ⑩ ⑪ ⑫ ⑬ ⑭ ⑮ ⑯ ⑰ ⑱ ⑲ ⑳ ㉑ ㉒
1 @H04T0FP1.1248.CS50RMACX.5:1101:1241:2154:1:R10 @ANCCG
2 TCAATATTTCATAGGGTATCTCGAATTCTCGGGTCCAGAACTCCAGT
3 +
4 CC0FFFFFHHBBLJ3J0PFIJ3J2JJ3JJJJ3JJJFPBH1JJBHBJPBLJ1
Read 3 ① ② ③ ④ ⑤ ⑥ ⑦ ⑧ ⑨ ⑩ ⑪ ⑫ ⑬ ⑭ ⑮ ⑯ ⑰ ⑱ ⑲ ⑳ ㉑ ㉒
1 @H04T0FP1.1248.CS50RMACX.5:1101:1461:2205:1:R10 @ATCACG
2 CAAGAGACTCTCTCTAGATATGGAATCTCGGGTCCAGAACTCCAGT
3 +
4 CC0FFFFFHHBBLJ3J2JJ3JJ3JJ3JJ3JJ3JJ3JJ3JJ3JJ3JJ1FPJL3JJ
```

- Line 1: sequence ID, description, and index; begins with @
- Line 2: sequence; contains only A, C, T, G, and N
- Line 3: optional sequence ID; begins with +
- Line 4: signal quality of each base, cryptic code, ASCII\_base 33 or 64

---

---

---

---

---

---

---

---

---

---

Compressing and Decompressing Large Files

```
gzip...compress or decompress a file.
usage (compress file) $ gzip file.txt
usage (decompress a gz file) $ gzip -d file.gz

tar...combine files into a single archive -- a tarball.
usage $ tar cf archive_name file1 file2 file3
usage (extract files from a tarball) $ tar xf archive_name.tar

gzip+tar... compress files and combine into a single archive
usage $ tar cvzf newname.tar file1 file2 fileN
usage (extract and decompress) $ tar xvzf file.tar.gz
```

**See unix cheat sheet for additional options.**  
 Datasets can be stored in a variety of locations: portable external hard drives, storage arrays, etc.

---

---

---

---

---

---

---

---

---

---