

ShadeLab Bioinformatics Workshop

Week 1

Sequencing platforms, amplicon (16S rRNA, ITS) processing
with USEARCH

<https://github.com/ShadeLab/workshop>

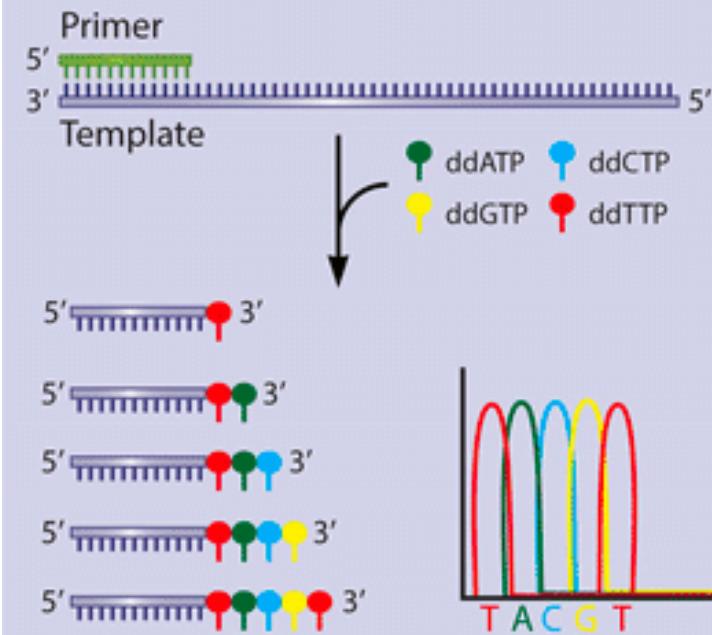
Checklist

- **GitHub**
 - git commands *clone, add, commit, push, pull*
- **HPCC**
 - directories
 - commands *ssh, cd, ls, mkdir, rm, cp, mv, scp*
- Install R and RStudio before next workshop

Sequencing platforms

First Generation

Shotgun Sequencing

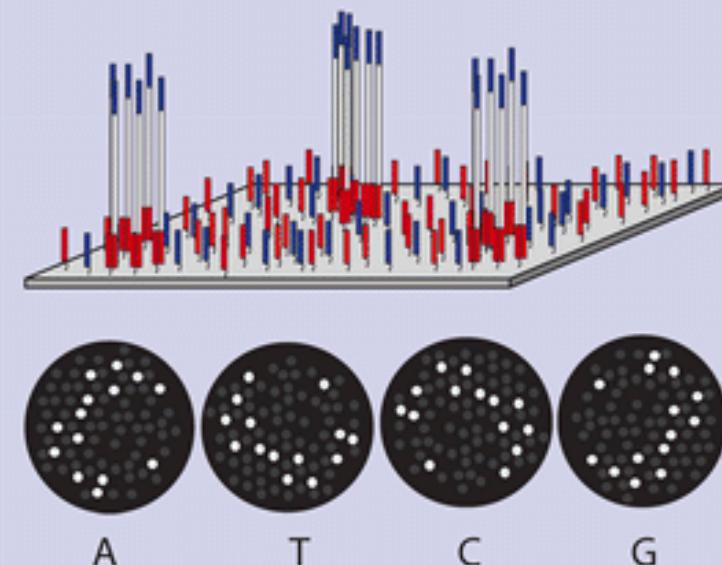


- Sequencing by synthesis
- High accuracy
- Long read lengths
- Relatively small amount of data generated

e.g., ABI capillary sequencer (ABI)

Second Generation

Massively Parallel Sequencing

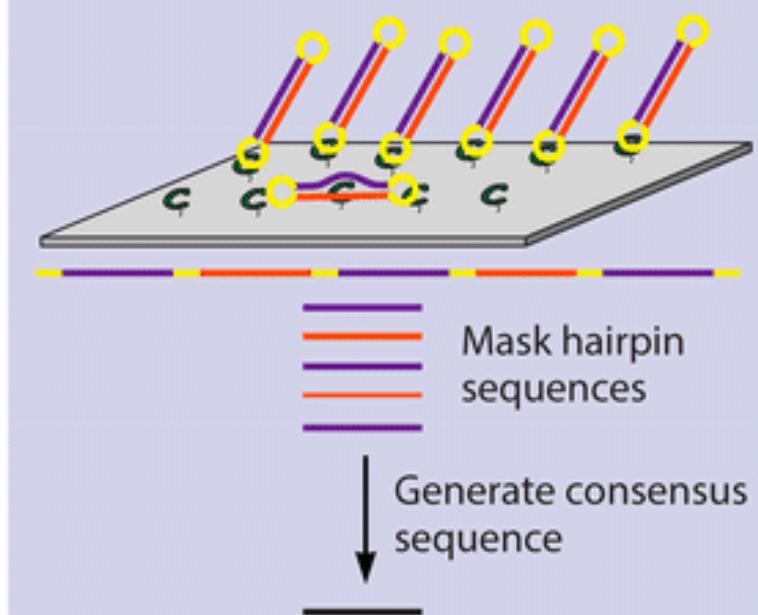


- Sequencing by synthesis
- Amplified templates are generated during sequencing, reducing the requirements for starting material
- High accuracy
- Short read lengths

e.g., MiSeq (Illumina), Ion Torrent (Thermo Fisher Scientific)

Third Generation

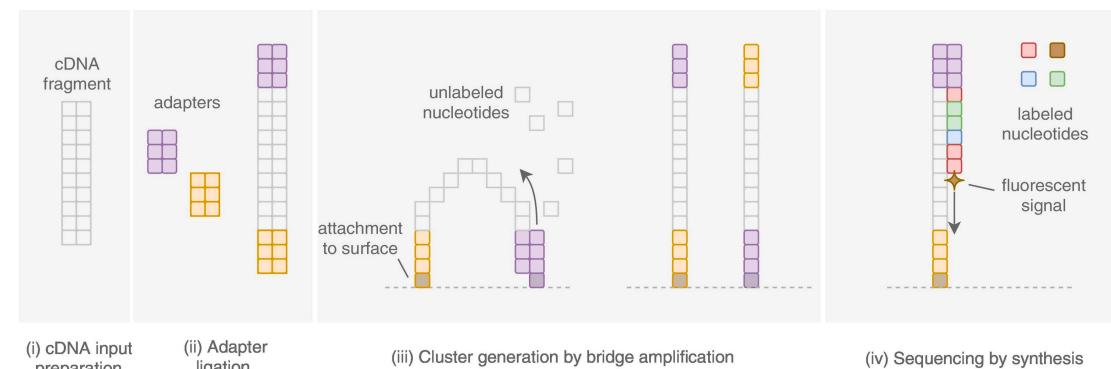
Single-molecule Sequencing



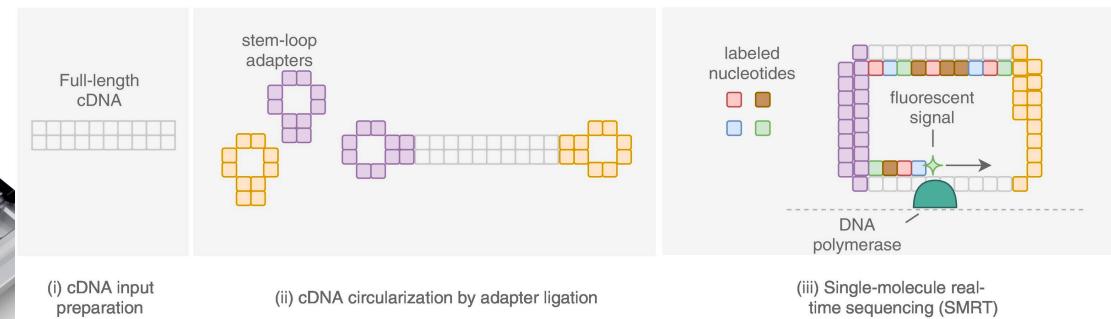
- Single-molecule templates
- Low accuracy
- Long read lengths

e.g., Single-Molecule Real-Time (SMRT) — Sequencing (Pacific Biosciences), MinION (Oxford Nanopore Technologies)

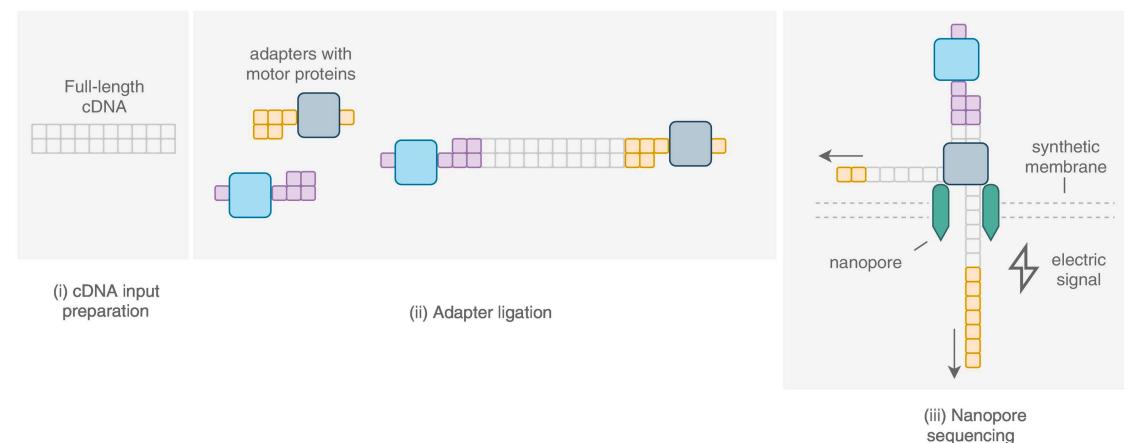
(A) Illumina RNA-Seq



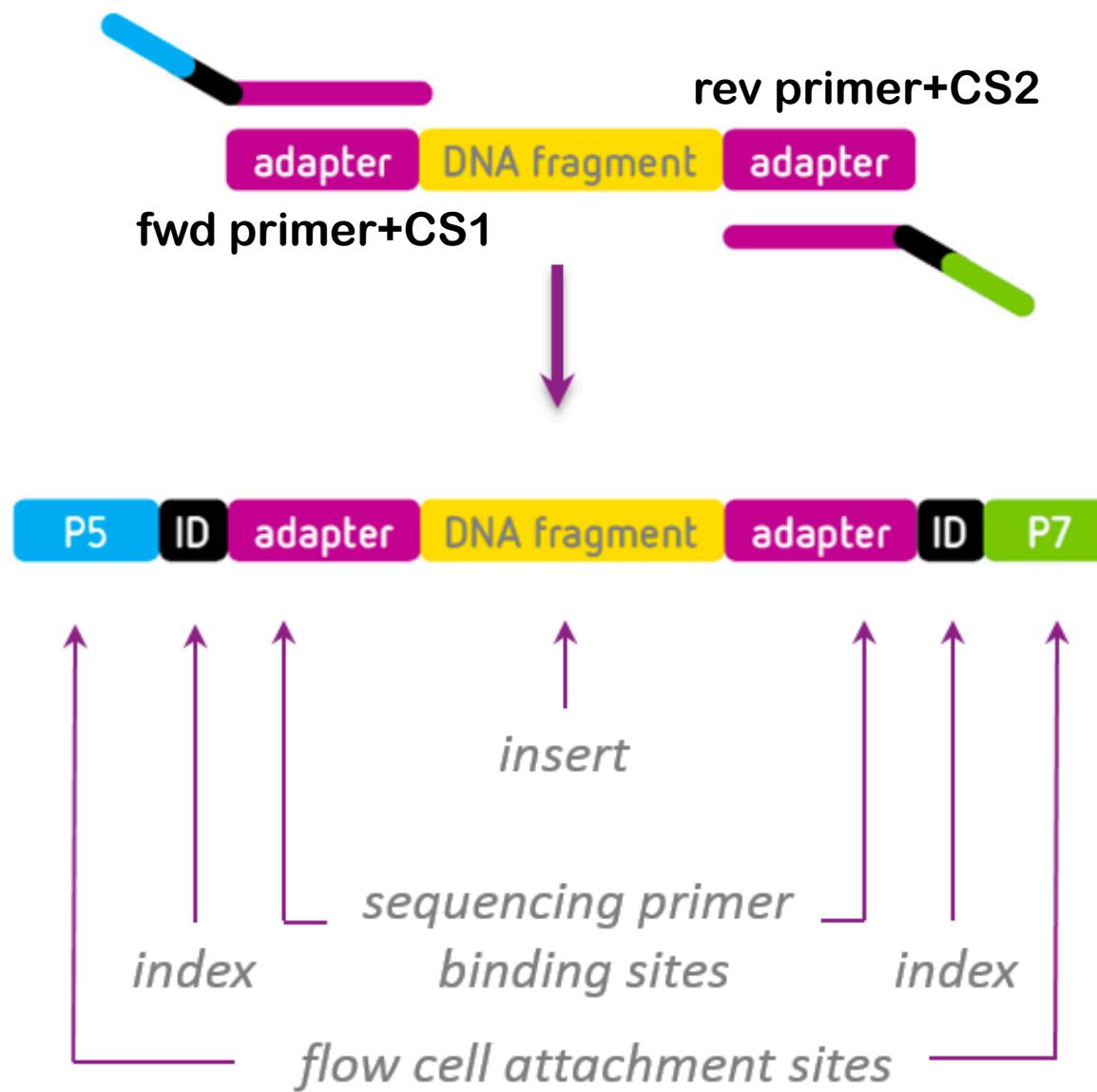
(B) PacBio Iso-Seq



(C) MinION 1D cDNA Sequencing



Illumina amplicon structure – library



**All extra information
needs to be
removed!**
(in some cases even
PhiX signal – Illumina
internal control)

RTSF (Illumina) sequencing results

... but first how to submit your sample:

<https://my.labguru.com/knowledge/experiments/7540>

<https://rtsf.natsci.msu.edu/genomics/submitting-samples/>

You will receive an email from RTSF containing the link to retrieve the files and QC report.

Sequencing is complete for samples submitted to the RTSF Genomics Core, project ID STO8258 (Aggregate_diversity). This project included a mixture of sample types. There were 16 samples of microbial metagenomic DNA submitted for 16 amplicon library preparation. The V4 hypervariable region was amplified using dual indexed, Illumina compatible primers 515f/806r following the Schloss lab protocol (Kozich, JJ, et al. 2013). PCR products were batch normalized using an Invitrogen SequlPrep DNA Normalization plate and the products recovered from the plate were pooled. The pool was QC'd and quantified using a combination of Qubit dsDNA HS, Agilent 4200 TapeStation HS DNA1000 and Kapa Illumina Library Quantification qPCR assays.

In addition, 36 primary PCR products prepared in your lab were also provided. Products from three (3) different targets were generated. These 1° PCR products were generated using primers with Fluidigm CS1/CS2 universal oligomers at their 5' ends. The Genomics Core performed 2° PCR using dual indexed, Illumina compatible primers which targeted the CS1/CS2 ends of the 1° PCR products. PCR products were batch normalized using an Invitrogen SequlPrep DNA Normalization plate and the products recovered from the plate were pooled. The pool was QC'd and quantified using a combination of Qubit dsDNA HS, Agilent 4200 TapeStation HS DNA1000 and Kapa Illumina Library Quantification qPCR assays. The QC revealed that it was likely the "amp_Block" samples were not amplified during the 2° PCR reactions. A second pool of the normalized PCR products was made, including only products from the AOA and AOB samples. This pool was QC'd as above.

The 16S and AOA/AOB library pools were combined in proportion to the number of libraries in each (16/28). This combined pool was loaded onto an Illumina MiSeq v3 flow cell and sequencing was performed in a 2x300bp paired end format using a MiSeq v3 600 cycle reagent cartridge. Custom sequencing and index primers complementary to the 515f/806 as well as primers complementary to the Fluidigm CS1/CS2 oligos were added to appropriate wells of the reagent cartridge. Base calling was done by Illumina Real Time Analysis (RTA) v1.18.54 and output of RTA was demultiplexed and converted to FastQ format with Illumina Bcl2fastq v2.19.1. A summary of the run output is attached below. Basic QC information about your sequence data is provided by the accompanying FastQC reports. Please see the FastQC Tutorial and FAQ (<https://rtsf.natsci.msu.edu/genomics/tech-notes/fastqc-tutorial-and-faq>) for information regarding interpretation of these reports.

The number of reads per sample is on average 2X more for the 16S-V4 libraries than for the AOA/AOB amplicons. It is not surprising that different amplicons perform better or worse when sequencing on an Illumina instrument.

Overall the read quality is poor, especially so for Read 2. Poor quality for amplicon libraries sequenced on a MiSeq v3 PE300 run was a concern raised when first discussing this project.

Data may be downloaded using the Shade lab account on the Genomics FTP server as before. Data for this project is in subdirectory 20190920_Amplicon_PE300. You must use secure FTP (FTPS) when connecting to the Genomics server. See the Genomics FAQ for general instructions (<https://rtsf.natsci.msu.edu/genomics/data-retrieval>). Sequence data typically remain available on the FTP server for 60 days. It is the responsibility of the researcher to download and store their data long term, including a safe backup copy. The RTSF Genomics Core only guarantees retention of sequence data for one year from the date of availability.

Regards,

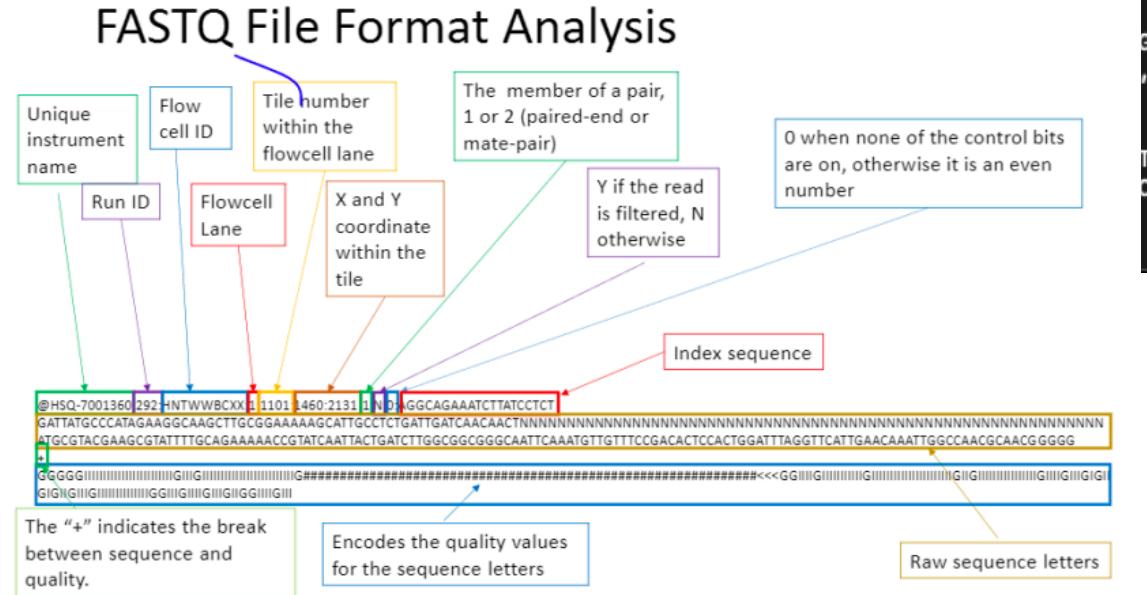
Kevin M. Carr

Download the FASTQ files to
/mnt/research/ShadeLab/Sequences
/raw_sequences/ along with the
email content and the QC file that is
attached in email

```
[(base) -bash-4.2$ cd ShadeLab/
[(base) -bash-4.2$ cd Sequence/
[(base) -bash-4.2$ cd raw_sequence/
[(base) -bash-4.2$ ls -lth
total 120K
drwxr-s--- 2 bintarti ShadeLab 8.0K Oct 15 14:08 its.tutorial
drwxr-s--- 4 stopnise ShadeLab 8.0K Oct  8 15:33 Soil_aggregates
drwxr-s--- 4 bintarti ShadeLab 8.0K Sep 13 12:20 Bean_variation
drwxr-sr-x 4 bowsher1 ShadeLab 8.0K Jun 19 22:37 Bowsher
drwxrws---. 5 stopnise ShadeLab 8.0K Jun 19 11:11 Bean_development
drwxrwsrwx. 4 kearnspa ShadeLab 8.0K Jun 13 14:58 Bean_biology
drwxr-sr-x 3 stopnise ShadeLab 8.0K Jun 13 14:41 Phyllo
drwxrwsrwx 6 stopnise ShadeLab 8.0K Jun 13 13:21 GLBRC
drwxrwsrwx 8 gradykea ShadeLab 8.0K May  3 10:35 AppleReps
drwxr-sr-x 2 dooleys1 ShadeLab 8.0K Sep  5 2018 GLBRC_analysis
drwxr-sr-x 6 kearnspa ShadeLab 8.0K Apr 23 2018 PRI_Kearns
drwxr-sr-x 4 stopnise ShadeLab 8.0K Aug 17 2017 Soil
drwxr-sr-x 2 shadeash ShadeLab 8.0K Jul 11 2017 JGI_iTags_1
drwxrwsrwx 3 sorens75 ShadeLab 8.0K Jun  8 2017 Centralia
drwx----- 2 shadeash ShadeLab 8.0K Aug  4 2016 160701_Jack
(base) -bash-4.2$
```

FASTQ file

If PE option, you will receive 2 files per sample, containing *R1*.fastq and *R2* labels



Amplicon processing tools



<https://www.drive5.com/usearch/>
(alternative: <https://github.com/torognes/vsearch>)



<https://www.mothur.org/>



USEARCH 64-bit version

