

Short Read Sequencing Overview

GL4U: Introduction 2024

Amanda M. Saravia-Butler, Ph.D.

NASA GeneLab Science Lead

Contractor: KBR

Biological & Physical Sciences

National Aeronautics and
Space Administration





1977

1st Generation
Sanger Sequencing

Assembly
required



1990

2nd Generation
“NextGen”
Pyrosequencing

Assembly
required

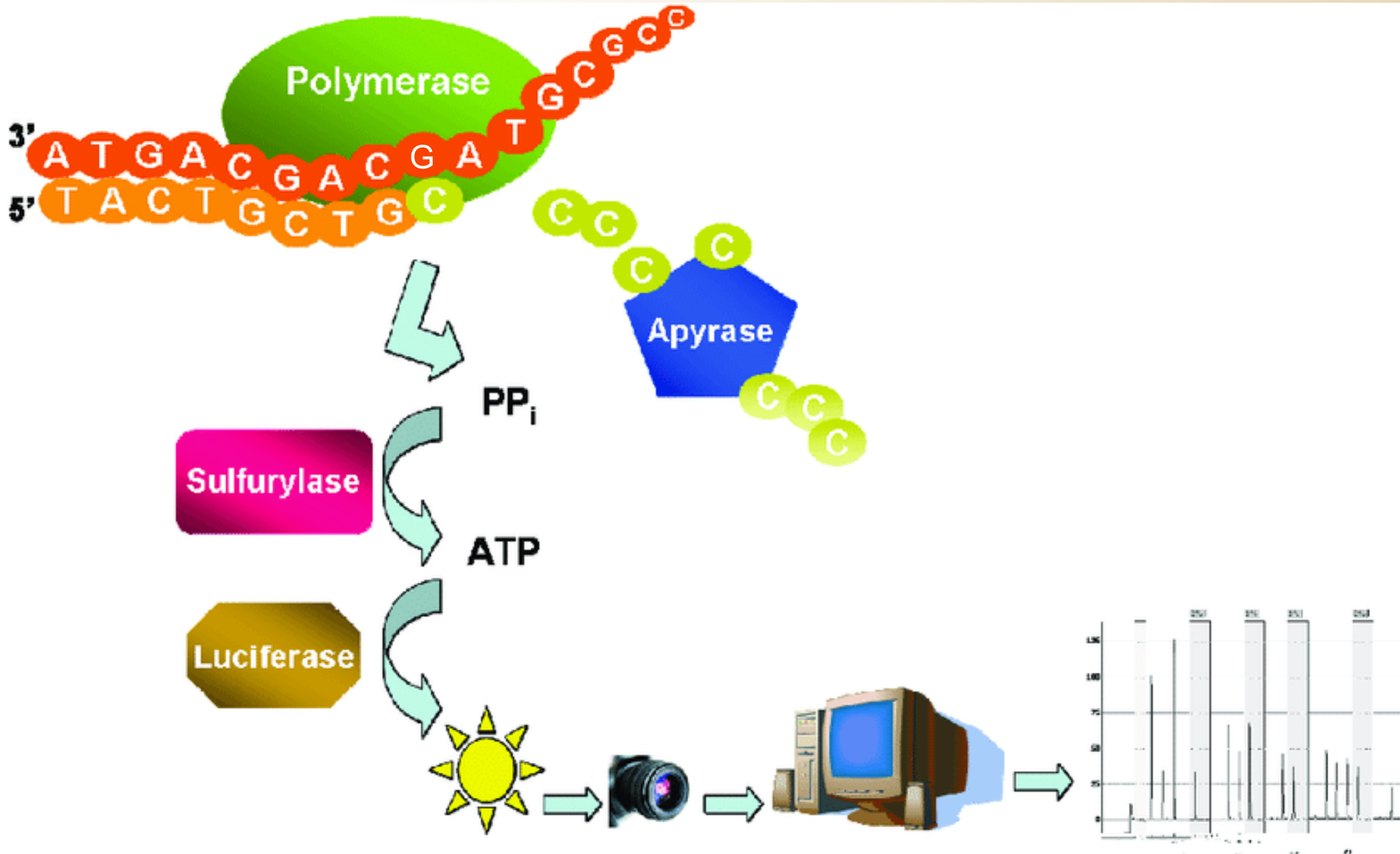


Now

3rd Generation
Single Molecule
Nanopores

Not much
assembly
required

∞



- Illumina is a very successful biotech company specializing in next generation technology that uses the pyrosequencing method
- ~90% of all sequencing worldwide is performed on an Illumina instrument (including GeneLab)
- The Illumina sequencing workflow has the following 3 steps:

➤ Library Construction

➤ Cluster Formation

➤ Sequencing



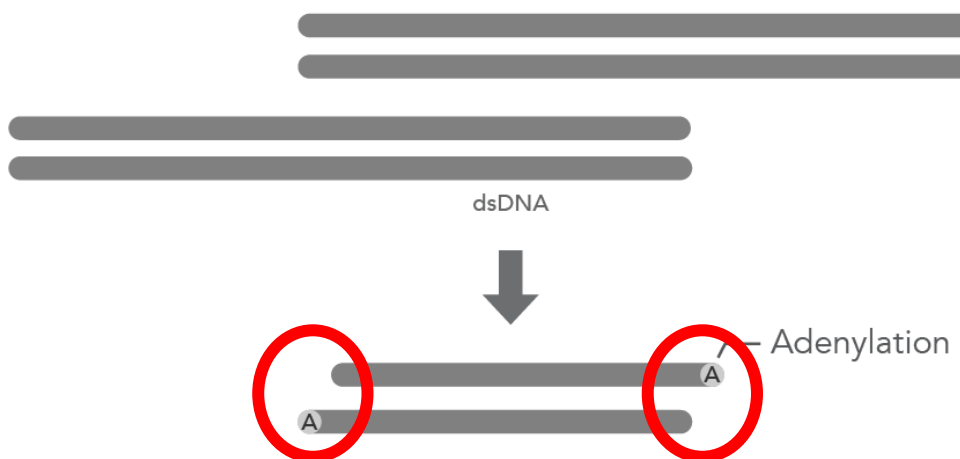
Library Preparation DNA Sequencing

Step 1: Create DNA fragments (with a means to attach adapters) from the extracted sample DNA

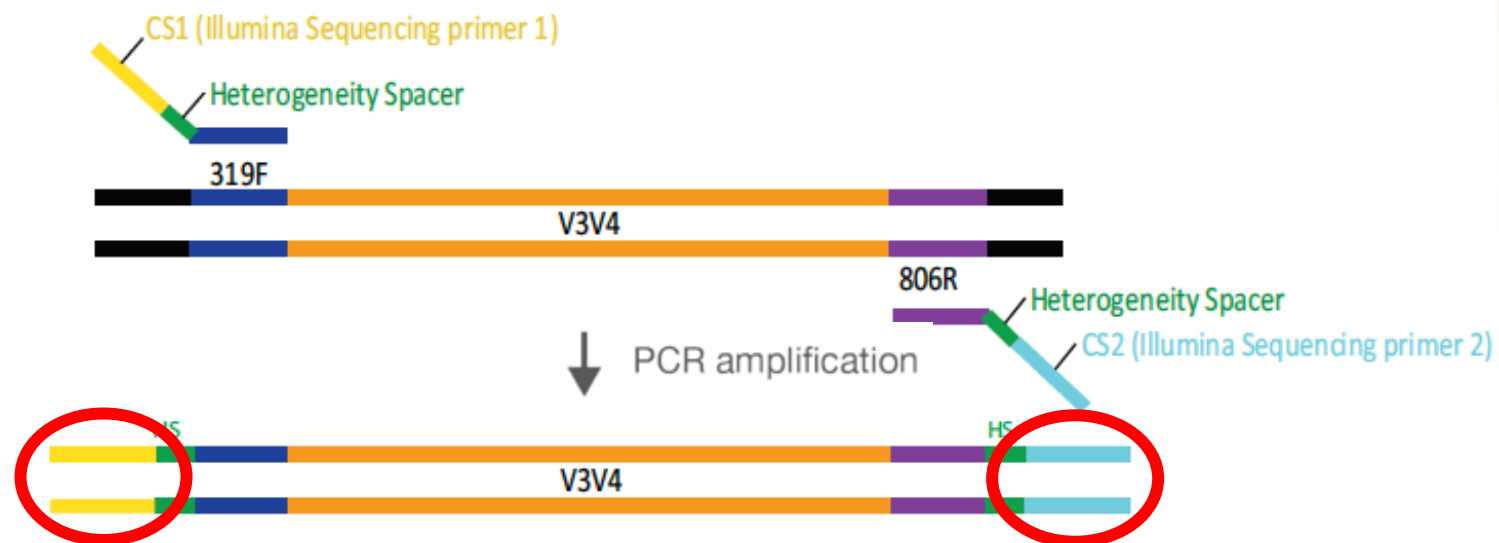
Tagmentation: Illumina DNA Prep (formerly Nextera Flex)



Fragmentation/A-tailing: IDT

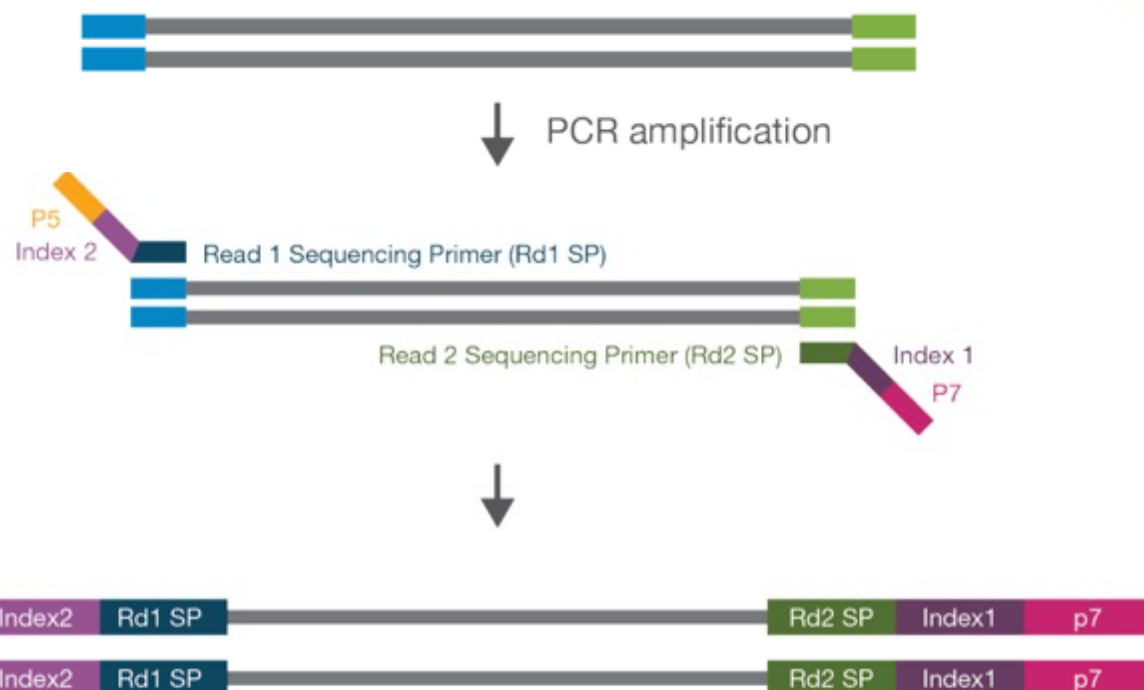


Target PCR: Illumina Amplicon

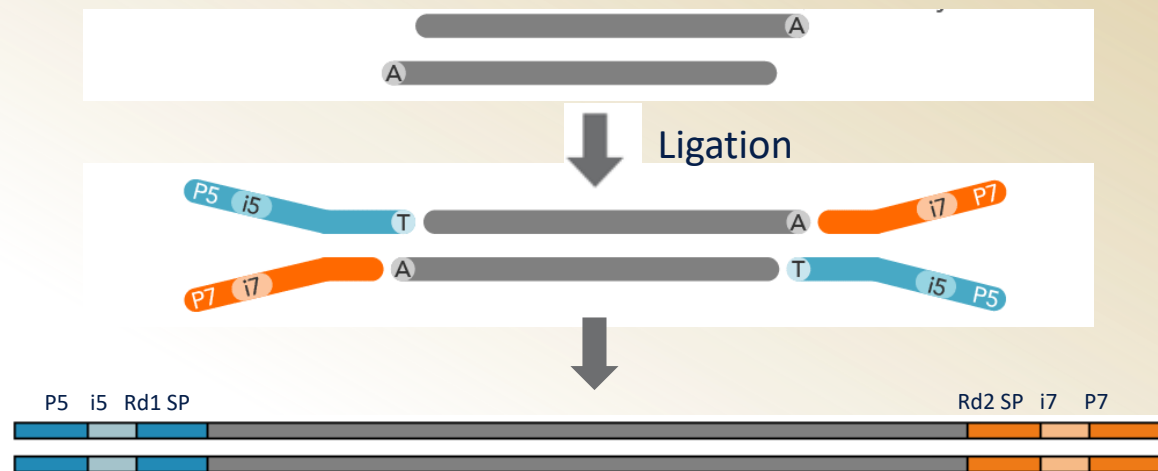


Step 2: Attach adapters

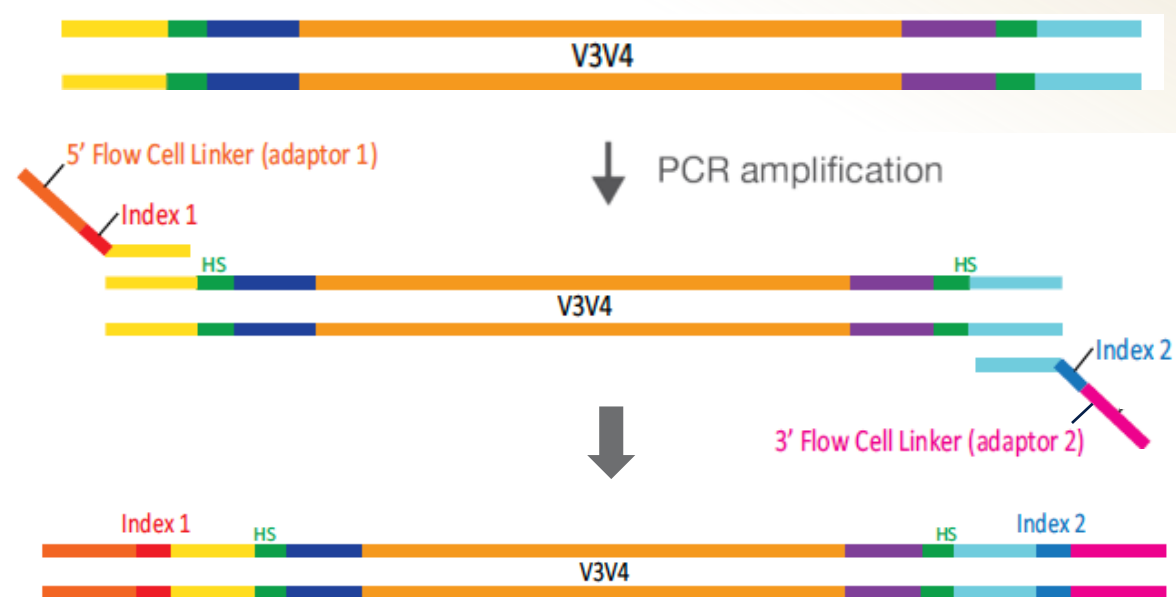
Tagmentation: Illumina DNA Prep



Fragmentation/A-tailing: IDT



Target PCR: Illumina Amplicon

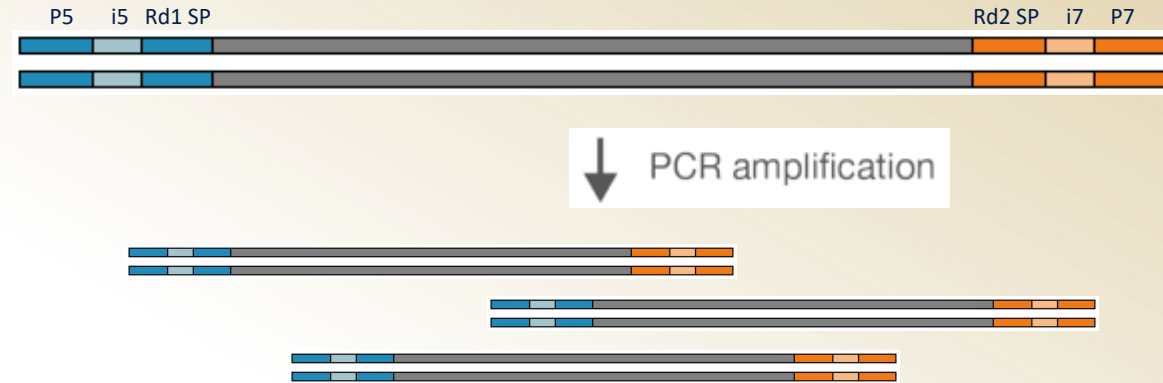


Step 3: Amplify libraries

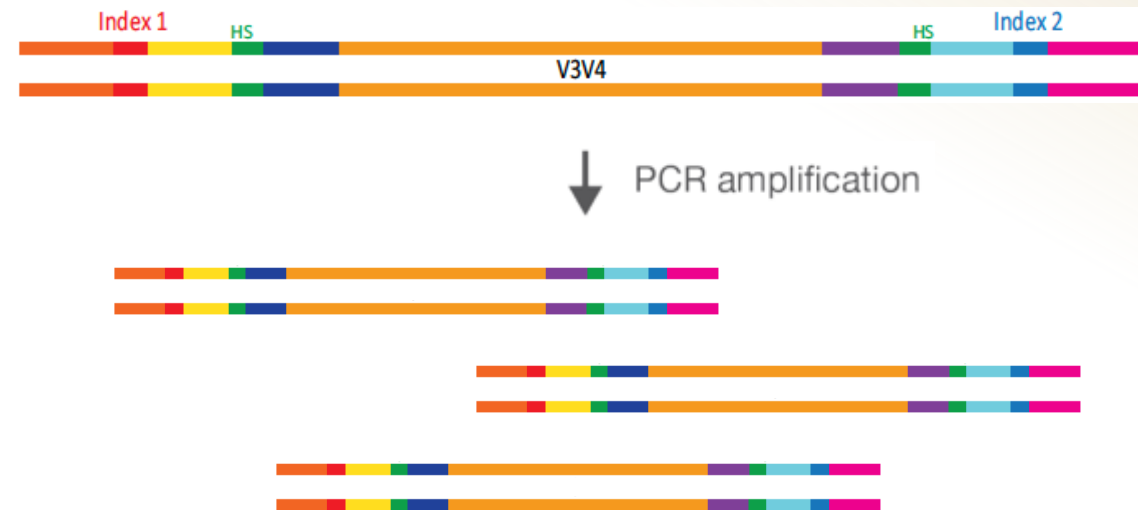
Tagmentation: Illumina DNA Prep



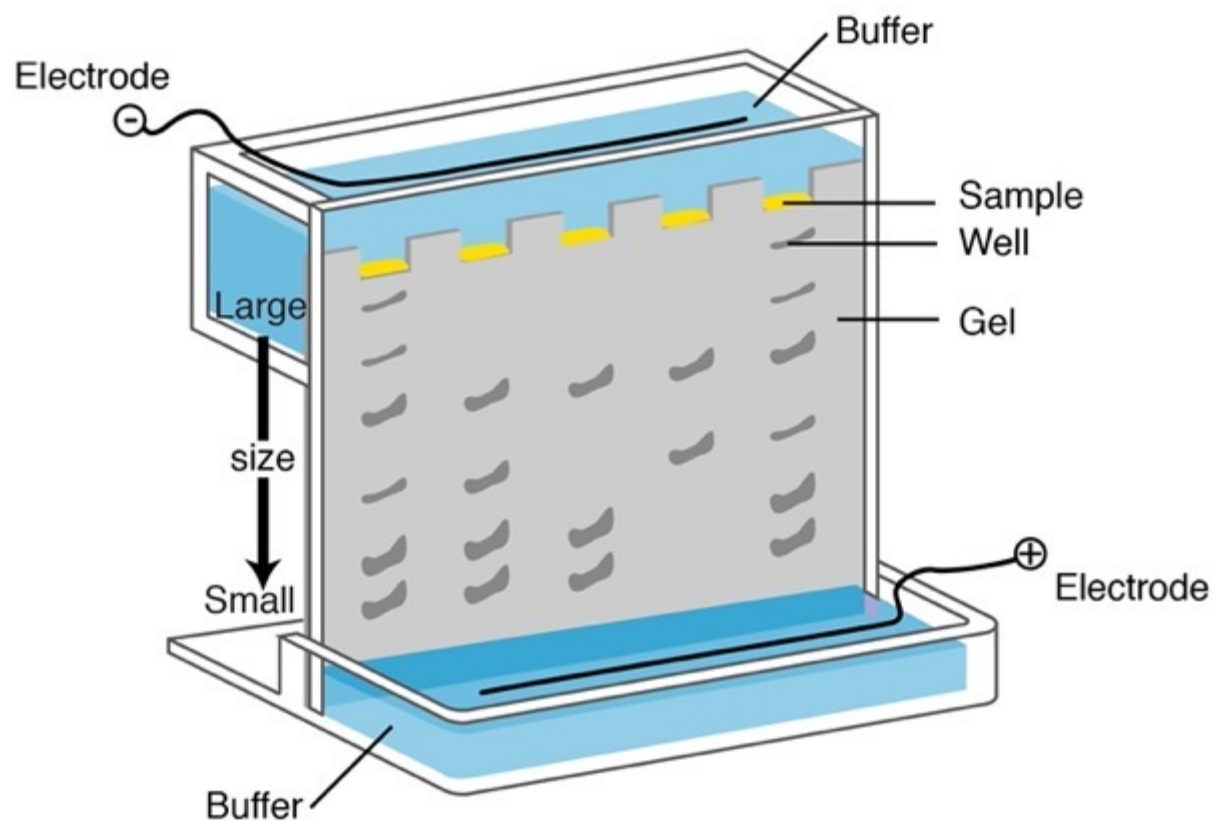
Fragmentation/A-tailing: IDT



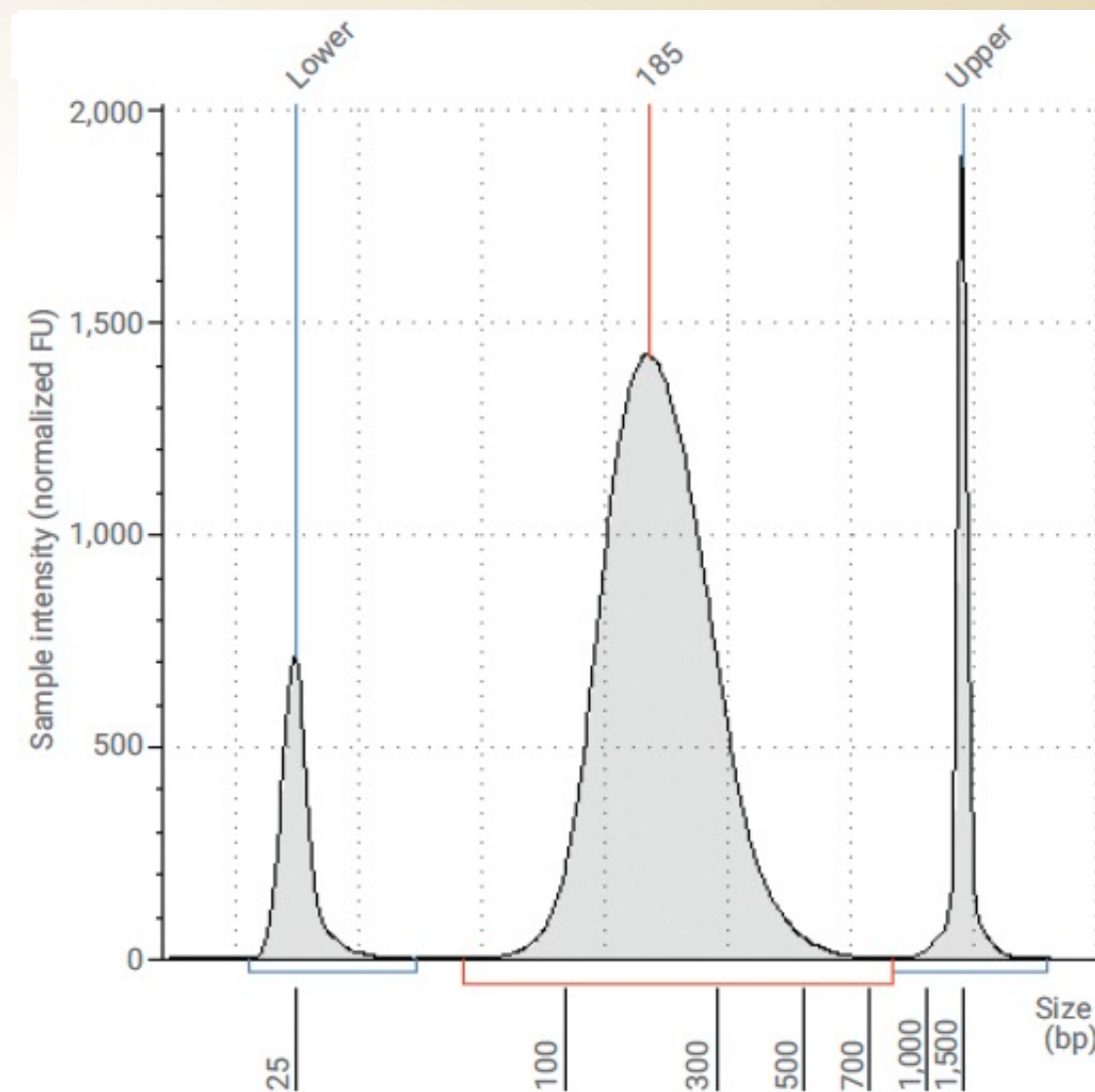
Target PCR: Illumina Amplicon



Electrophoresis

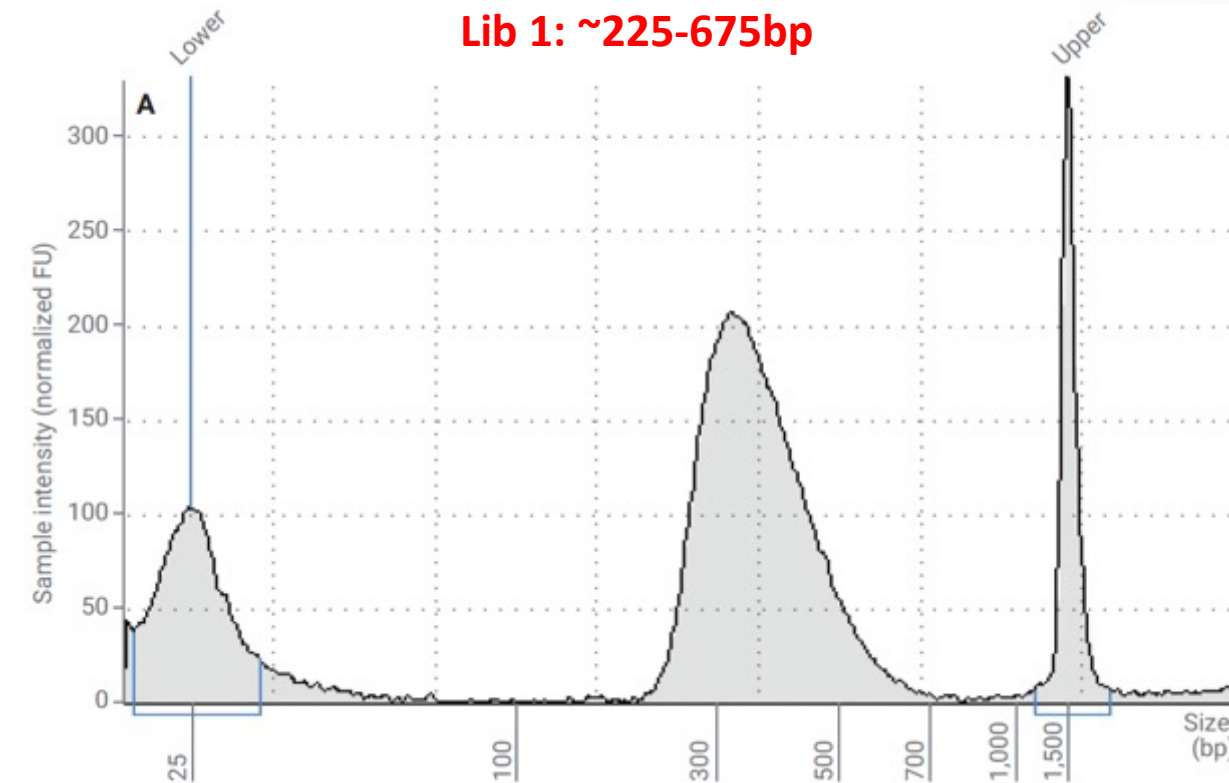


Electropherogram

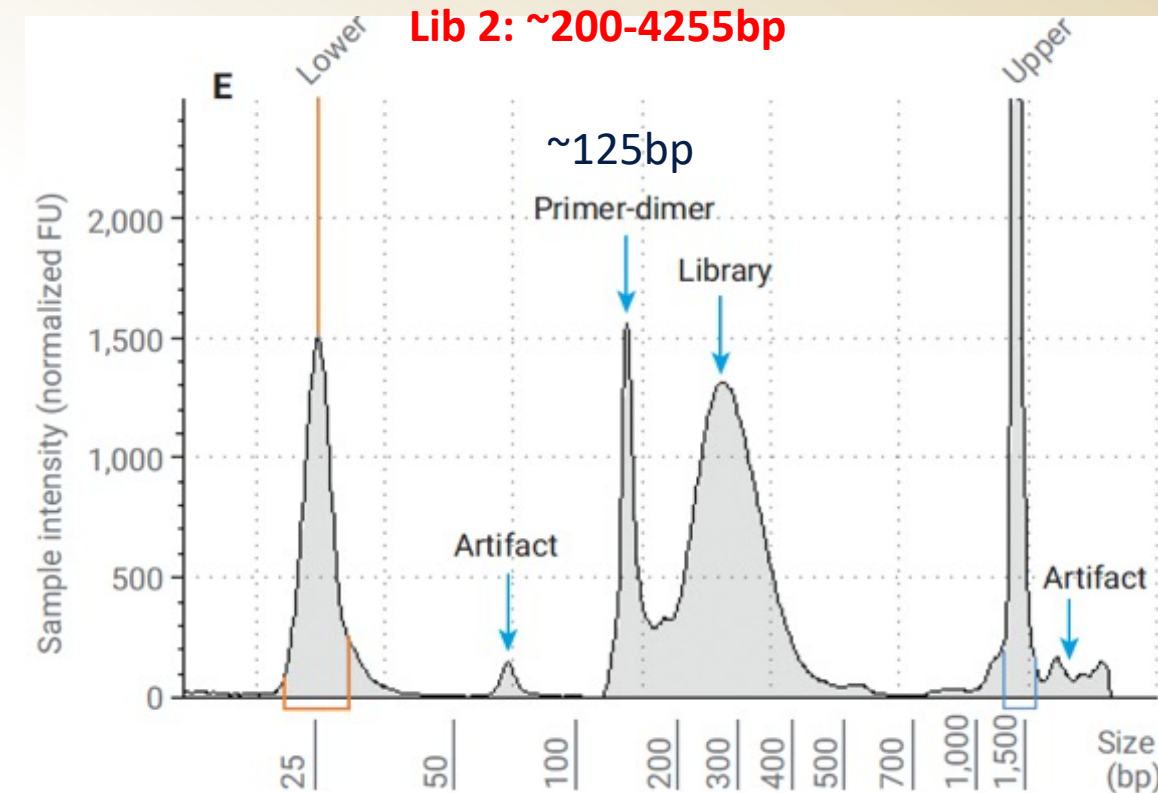


- Libraries are evaluated using a bioanalyzer or a tape station to create an electropherogram to assess quality

Good Library - 1



Bad Library - 2

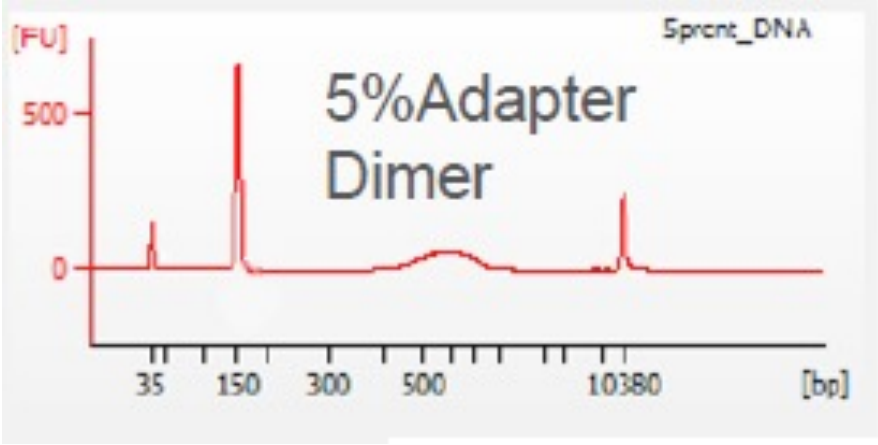
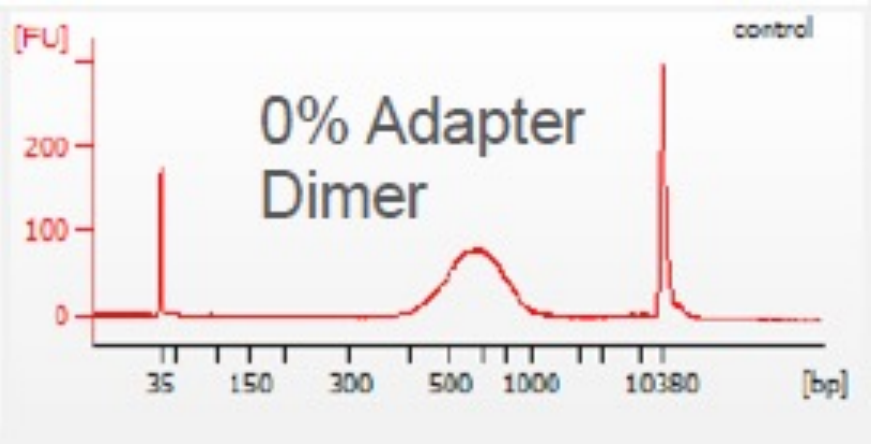
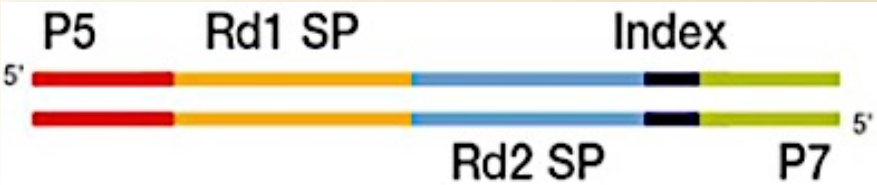


When assessing library quality look for the following:

- NO adapter dimers! – Why?
- Library size is consistent with the number of desired sequencing cycles
 - If you're sequencing at PE 250, what is a good library size? (hint: library size = insert length + adapter length)

Assuming adapters are ~65bp each (130bp)
~600bp that would give an insert length of ~470,
allowing a ~30bp overlap between R1 and R2

- Assess library quality (bioanalyzer, TapeStation)
 - Adapter dimers
 - Fragment size
- Determine library quantity (Qubit, qPCR)



% AD	% PF	% AD Reads
Control	69.54	0.24
10%	10.87	84.25
5%	21.39	60.44
1%	51.88	6.46



Library Preparation RNA Sequencing

mRNA makes up only ~2-5% of a total RNA sample

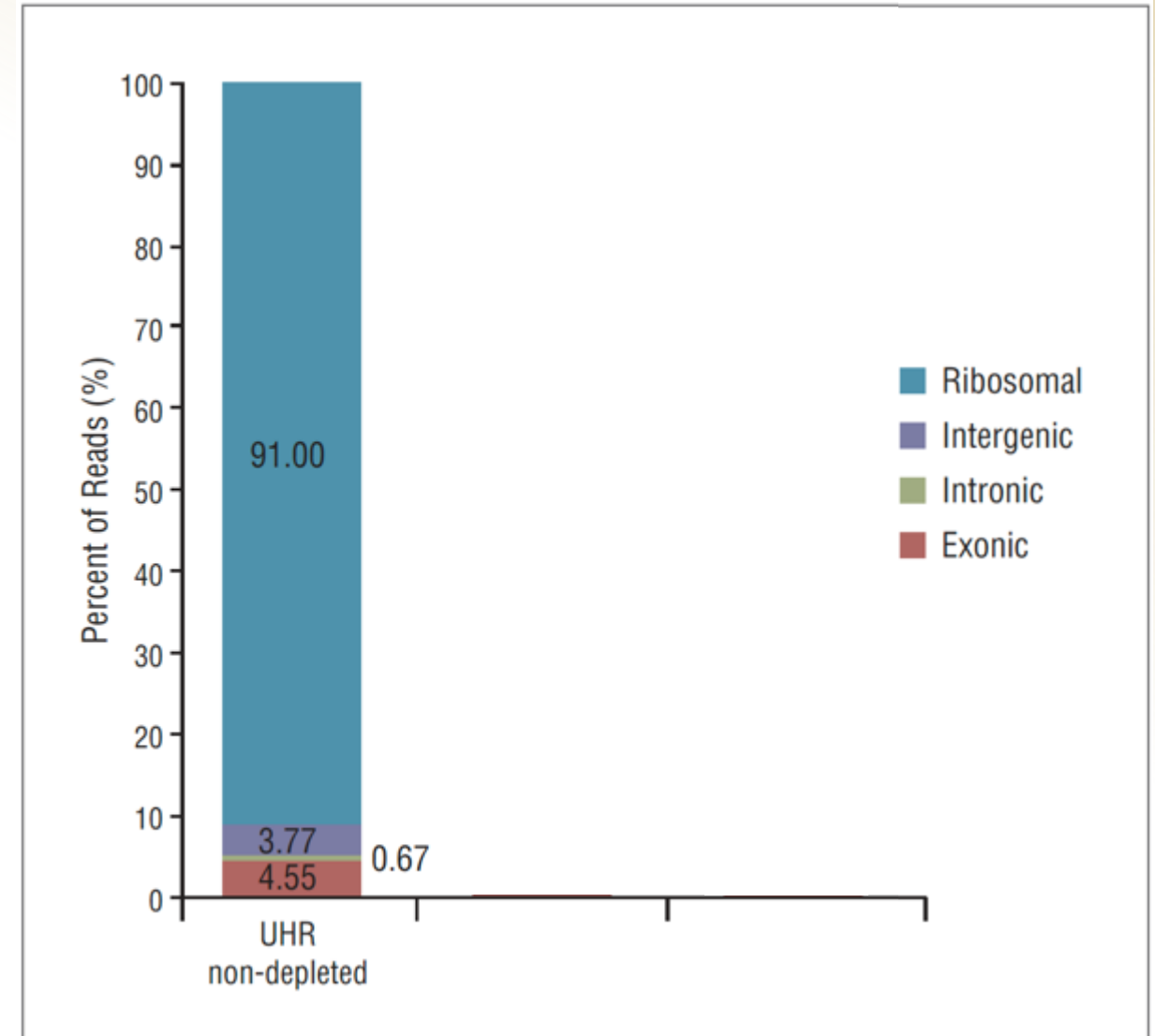
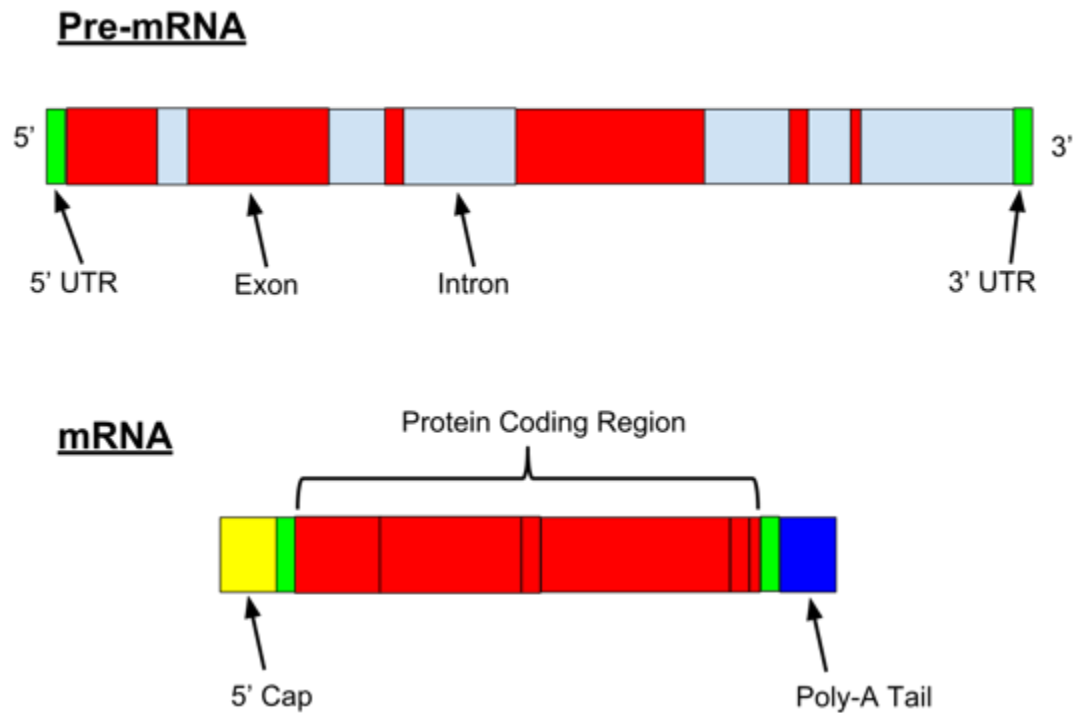


Figure 1 Ribo-Zero Depleting and Fragmenting RNA



Figure 2 Synthesizing First Strand cDNA



Figure 3 Synthesizing Second Strand cDNA



Figure 4 Adenylating 3' Ends*

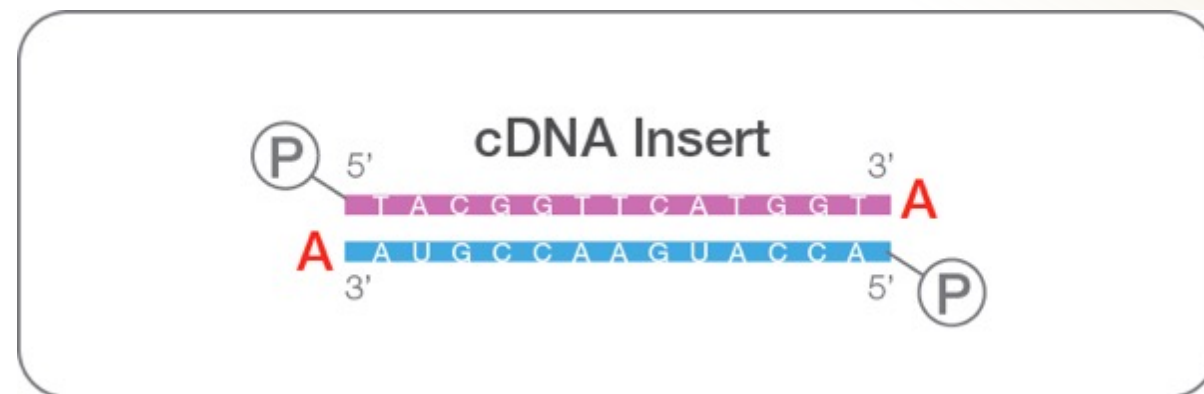
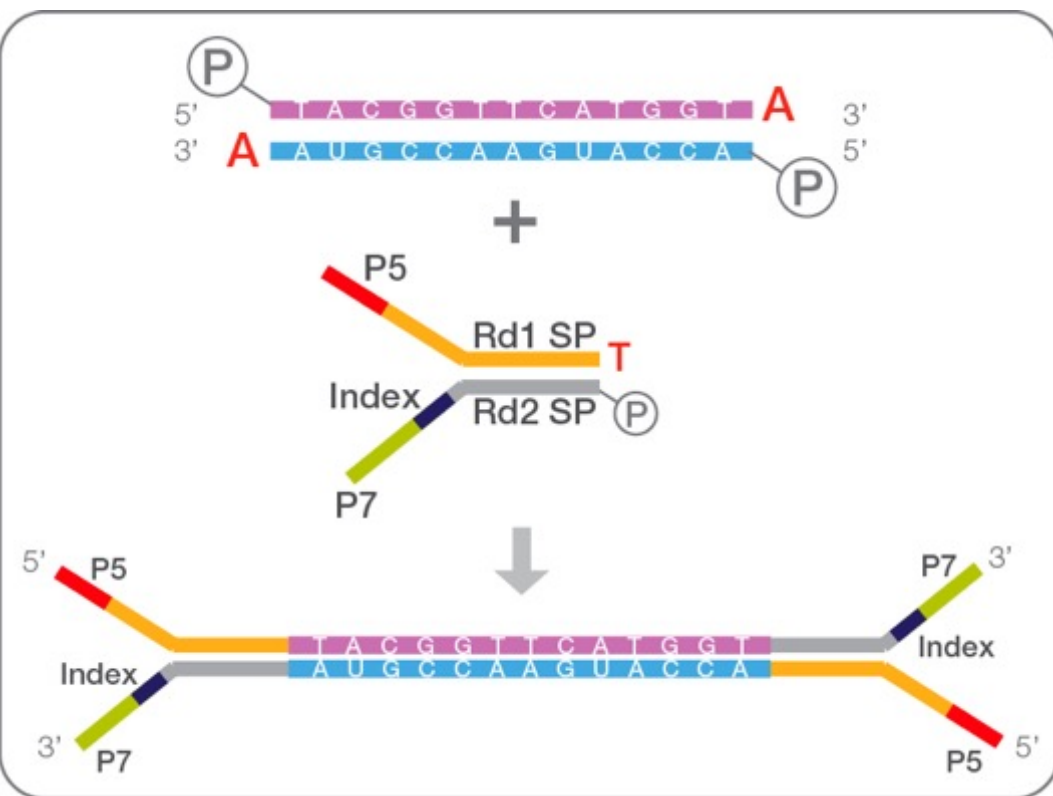


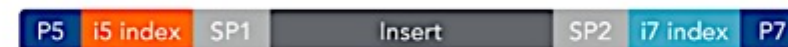
Figure 5 Ligating Adapters



Single index



Unique dual index



Dual index UMI



- Flow cell binding sequence:** Platform-specific sequences for library binding to instrument
- Sequencing primer sites:** Binding sites for general sequencing primers
- Sample indexes:** Short sequences specific to a given sample library
- Molecular index/barcode:** Short sequence used to uniquely tag each molecule in a given sample library
- Insert:** Target DNA or RNA fragment from a given sample library

Figure 6 Enriching DNA Fragments*

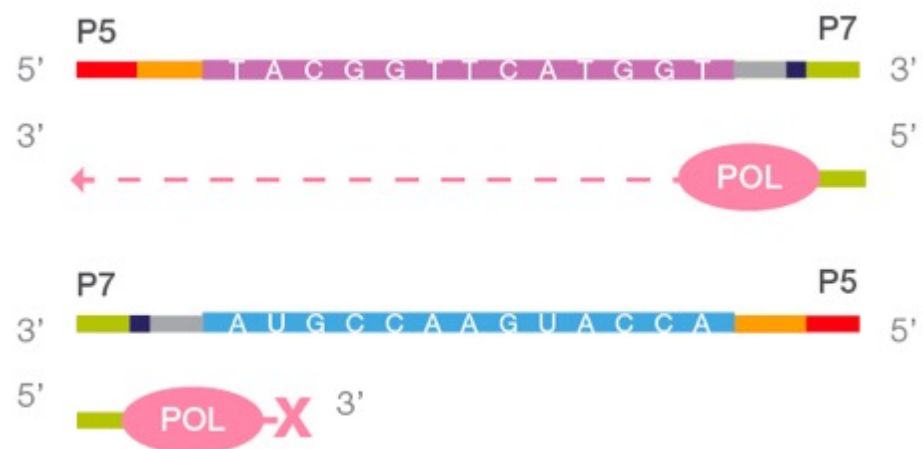
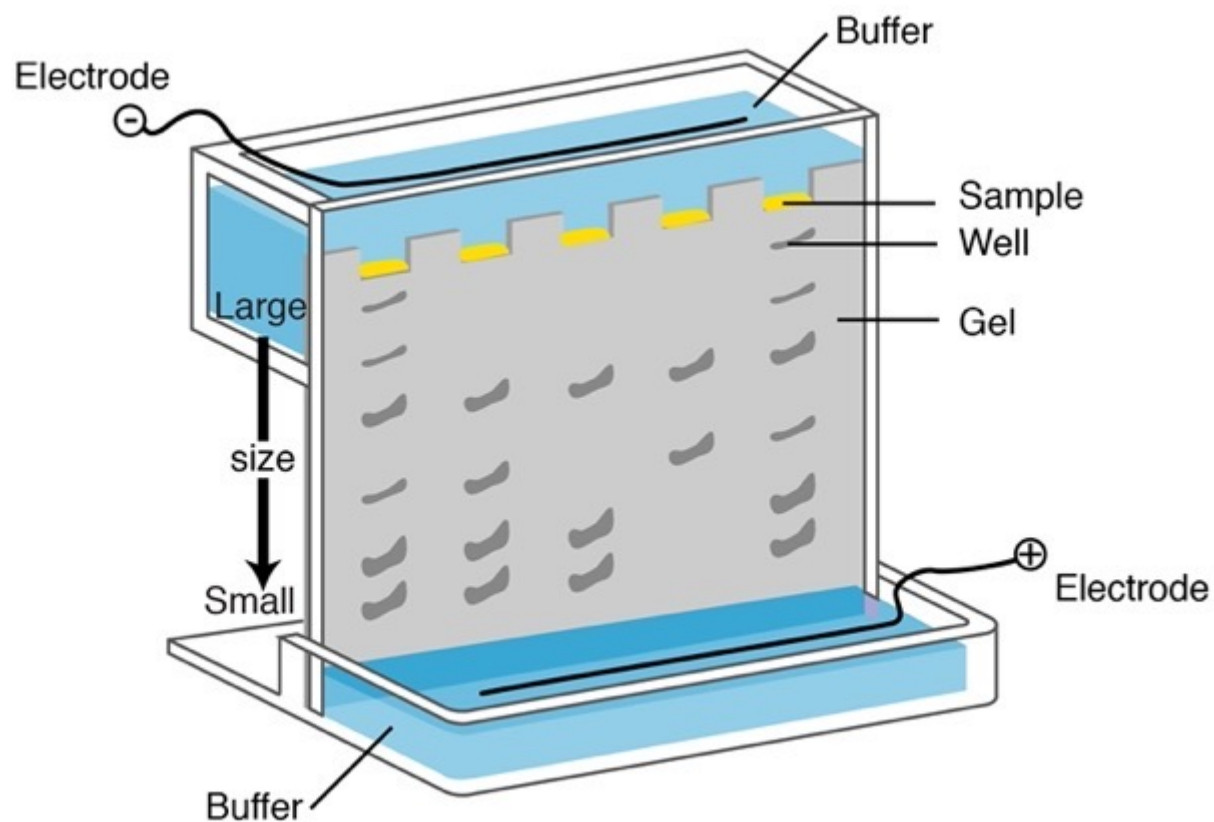


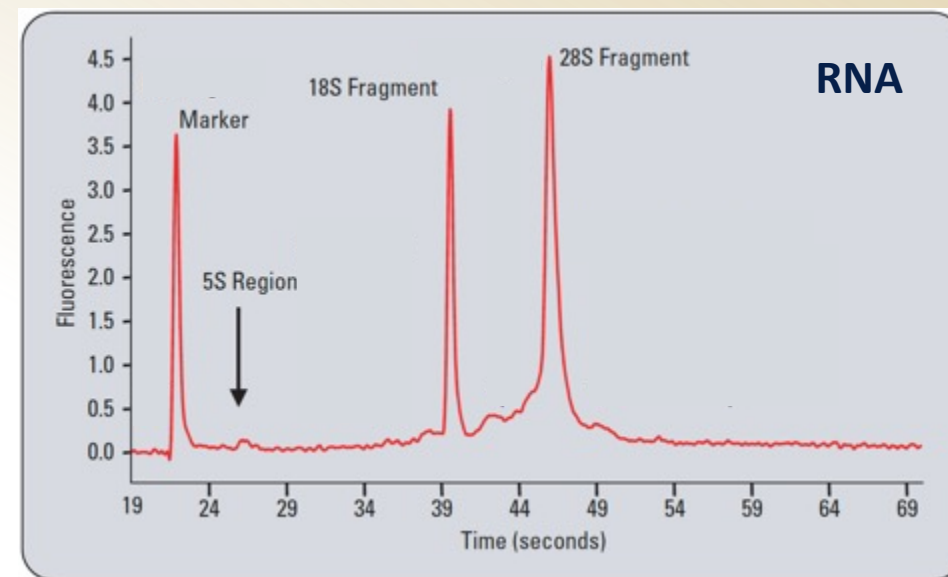
Figure 7 LS Final Library



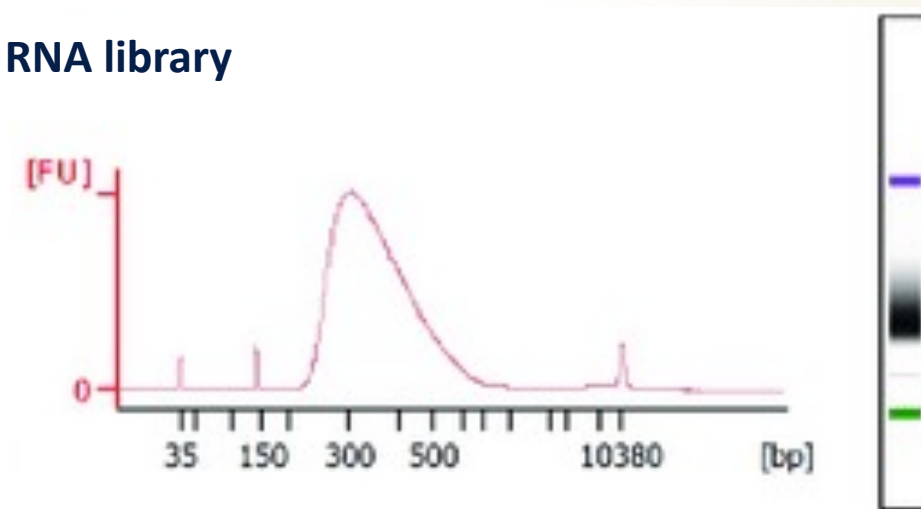
Electrophoresis



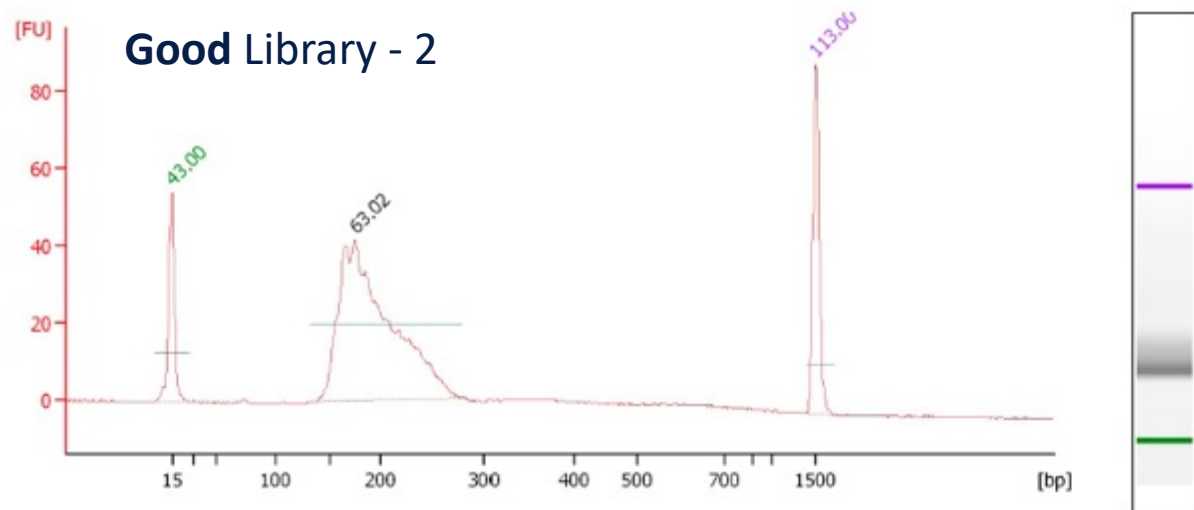
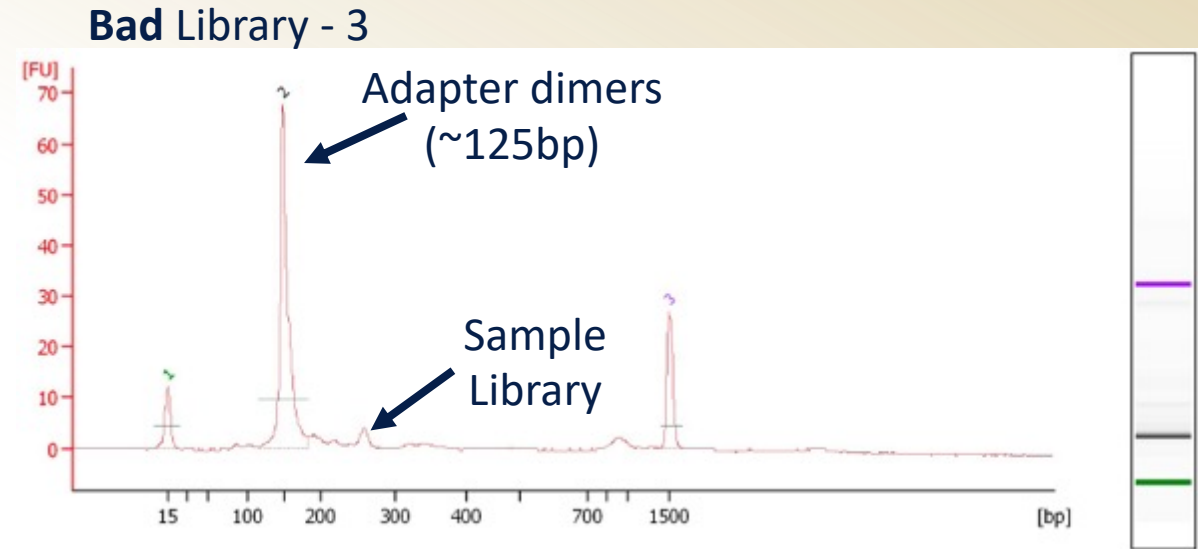
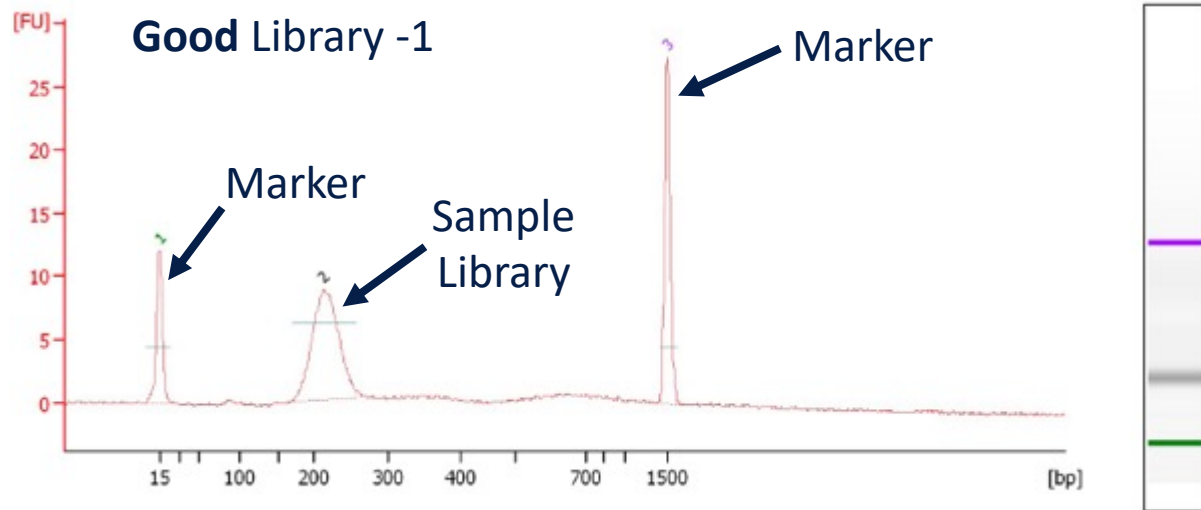
Electropherogram



RNA library



- Libraries are evaluated using a bioanalyzer or a tape station to create an electropherogram to assess quality



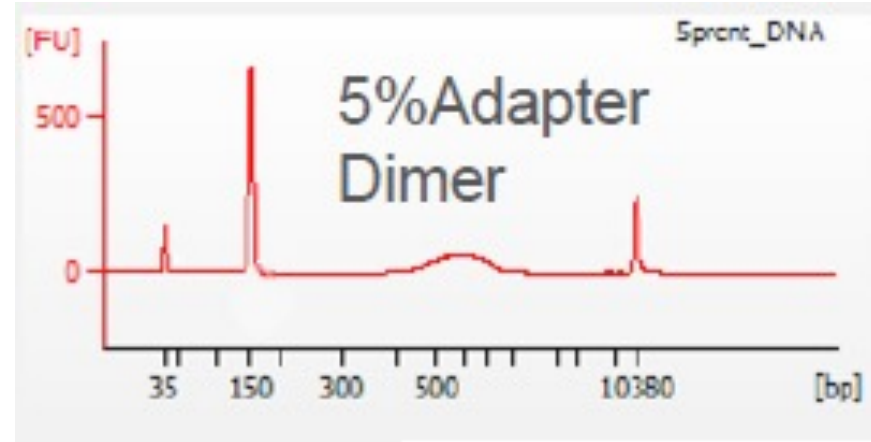
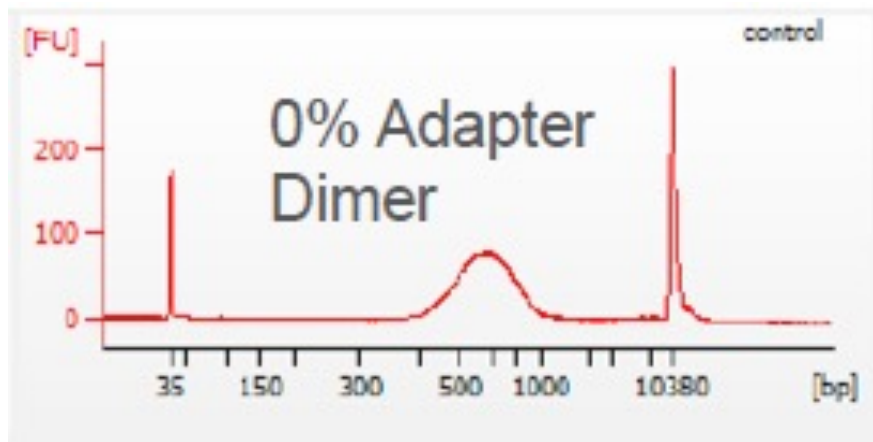
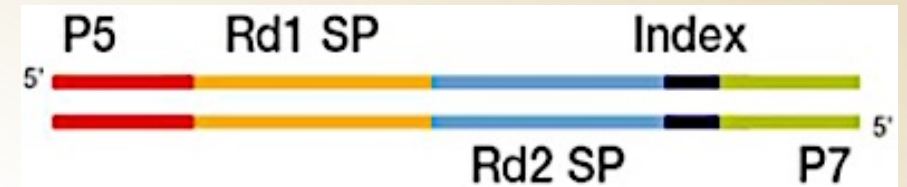
When assessing library quality look for the following:

- NO adapter dimers! – Why?
- Library size is consistent with the number of desired sequencing cycles
 - If you're sequencing at PE 100, what is a good library size? (hint: library size = insert length + adapter length)

Assuming adapters are ~65bp each (130bp): ~300bp that would give an insert length of ~170, allowing a ~30bp overlap between R1 and R2

- What is the size range of library 1? Library 2?
~175-250bp; ~150-275bp

- Assess library quality (bioanalyzer, TapeStation)
 - Adapter dimers
 - Fragment size
- Determine library quantity (Qubit, qPCR)



% AD	% PF	% AD Reads
Control	69.54	0.24
10%	10.87	84.25
5%	21.39	60.44
1%	51.88	6.46



Cluster Formation and Sequencing by Synthesis

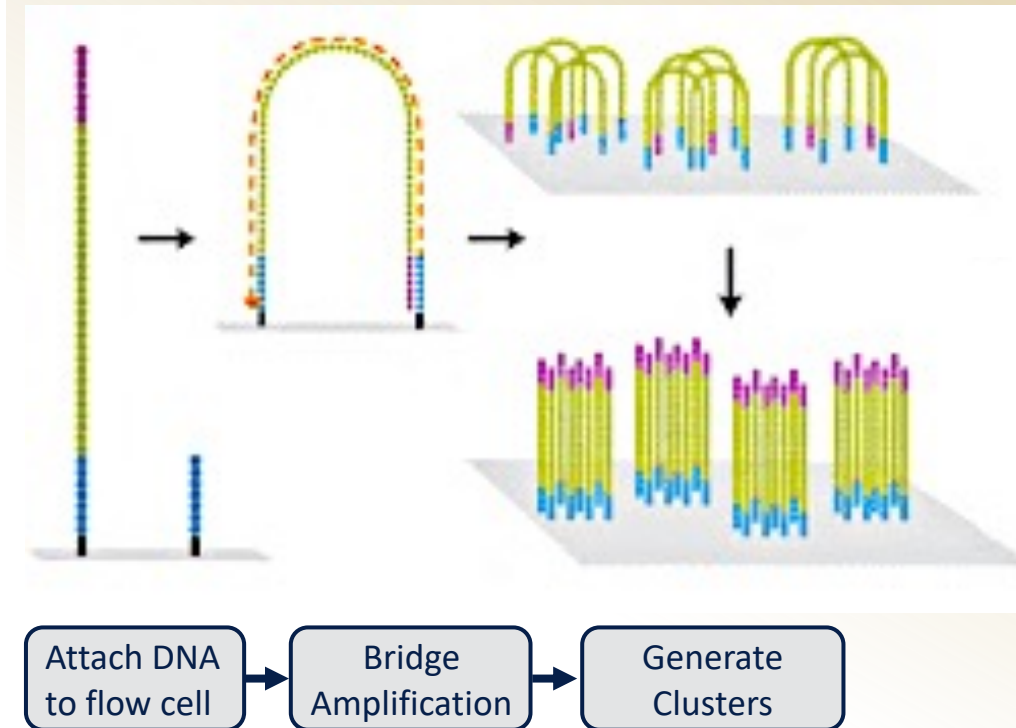
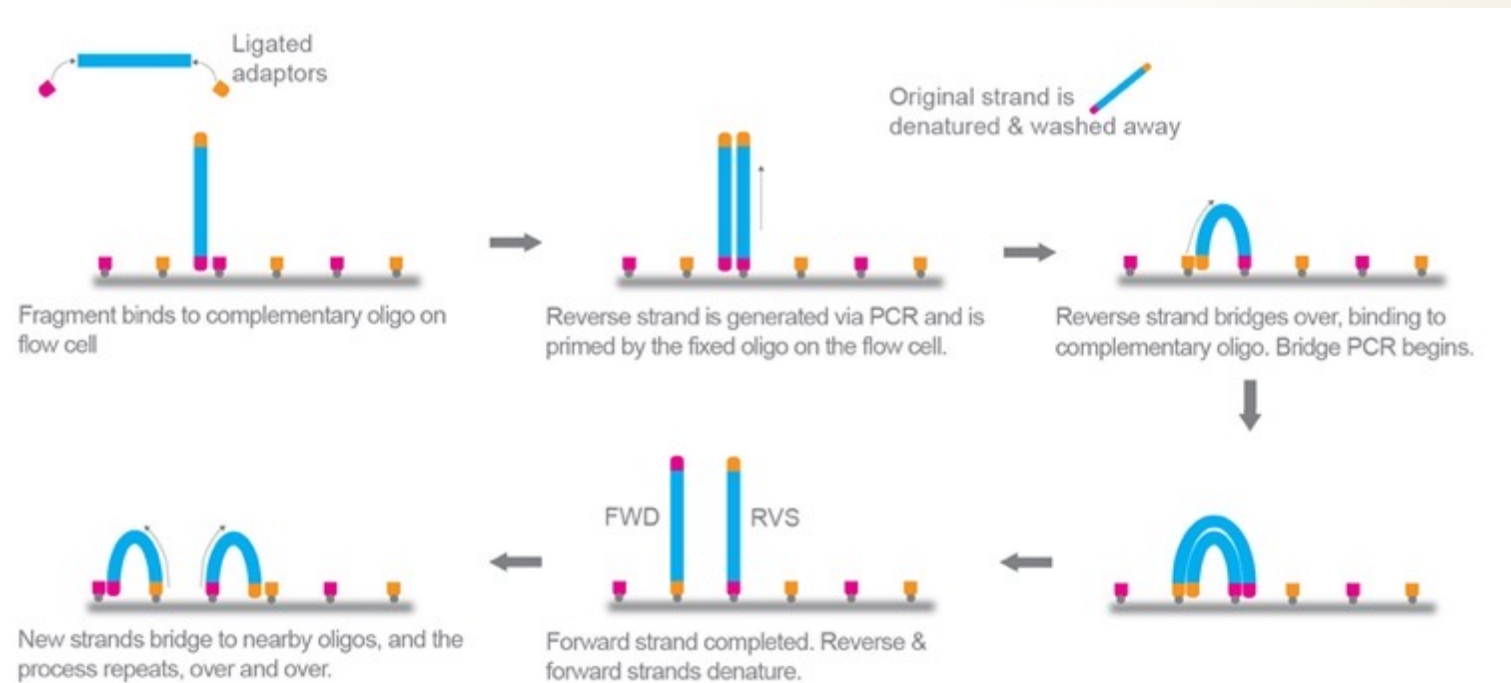
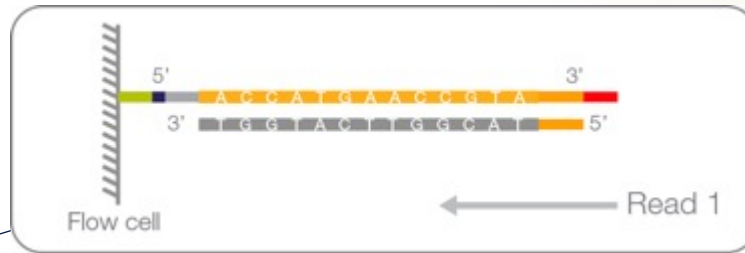
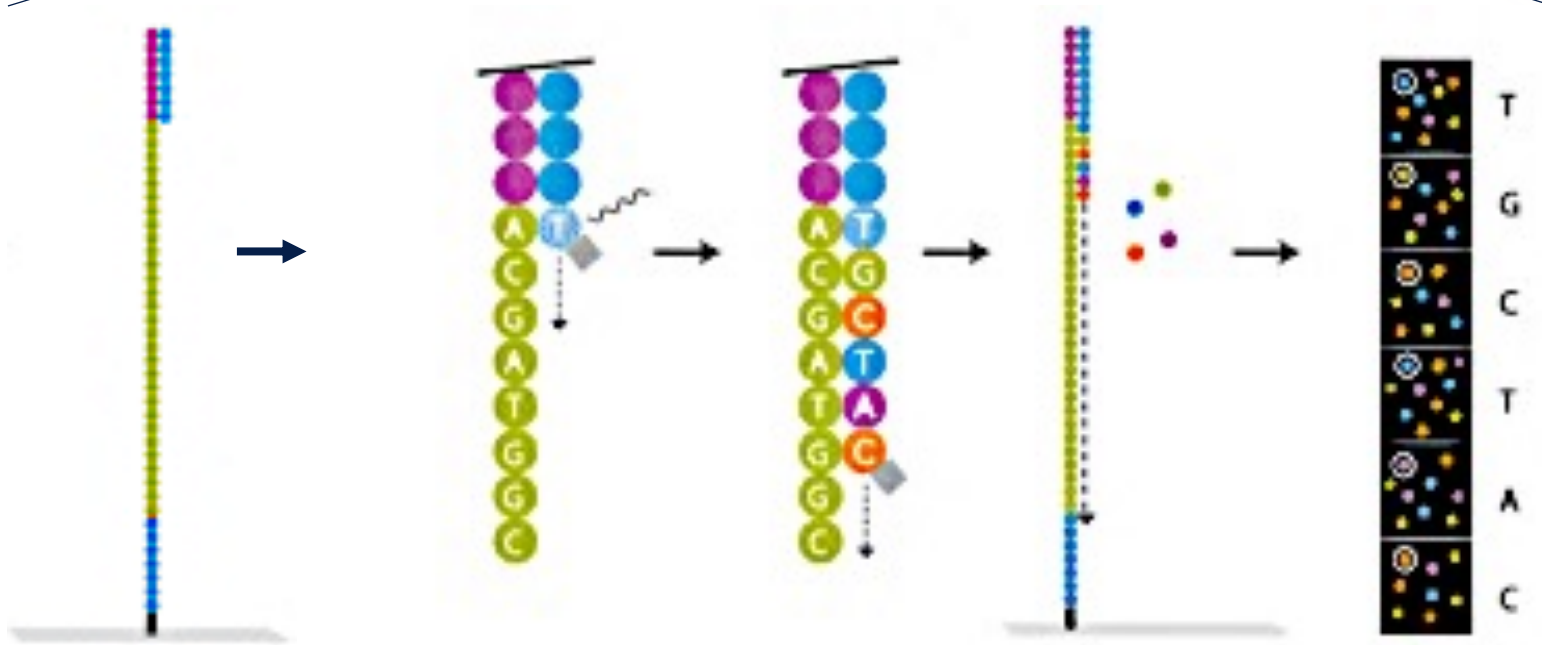


Figure 4. Bridge PCR - a PCR method used to amplify samples for sequencing.

Figure 8 Read 1 Sequencing



- For paired-end (PE) sequencing, after read 1 is sequenced, forward strand reagents are washed
- The index read(s) are sequenced next
- Sequencing of the reverse read (read 2) is initiated after the index read(s)



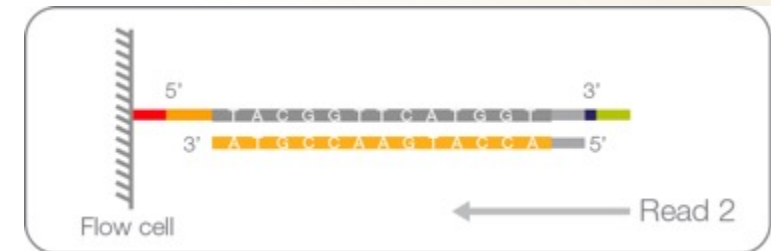
Anneal Sequencing
Primer

Extend first base,
read, and deblock

Repeat previous step
to extend strand

Generate
base calls

Figure 9 Read 2 Sequencing



	Read length	Sequencing depth	Paired-end (PE) or Single-end (SE)
DNAseq	<ul style="list-style-type: none"> ➤ Longer reads enable greater confidence in taxonomic classifications and functional annotations <ul style="list-style-type: none"> ➤ “Assembly” is often performed with short reads to facilitate this ➤ GL standard is 2x250bp 	<ul style="list-style-type: none"> ➤ Greater depth increases the likelihood of sequencing low-abundance organisms (if metagenomics) and detecting things like single-nucleotide variants and genetic rearrangements with greater confidence ➤ GL standards: <ul style="list-style-type: none"> ➤ Single organism or tissue: <ul style="list-style-type: none"> ➤ Re-sequencing (reference available): 10X minimum ➤ De novo sequencing (no reference available): 50X minimum ➤ Metagenomics (mixed community): 10M per sample, minimum 	<ul style="list-style-type: none"> ➤ PE is generally preferred
RNAseq	<ul style="list-style-type: none"> ➤ Longer reads increase gene ID confidence ➤ GL standard is 2x150bp for bulk RNAseq 	<ul style="list-style-type: none"> ➤ Greater depth increases the likelihood of sequencing low-abundant transcripts, detecting novel transcripts, and quantifying isoforms ➤ Greater depth is necessary for ribo-depleted samples (vs. poly-A enriched samples) – for RNAseq <ul style="list-style-type: none"> ➤ GL RNAseq standard for mammals prepared with ripo-depletion is 40-60M reads/sample ➤ More replicates is usually preferred over greater depth 	<ul style="list-style-type: none"> ➤ PE is preferred



Questions?