

From sample to fastq

Outline

- How library quality affects data recovery
- Sequencing costs
- Estimate cost for your own experiment

Requirements for library prep protocol for lcWGS




- To prepare libraries for hundreds of samples, we need a protocol that is
 - Cheap
 - Efficient
 - Reliable
- Sometimes robustness to sample degradation is also important



Sequencing-Ready Fragment




Index 2 Primer



-  P5 – complementary to Illumina flow cell oligo
-  Indexing sequence 2
-  Read 1 Sequencing Primer

Index 1 Primer






-  Read 2 Sequencing Primer
-  Indexing sequence 1
-  P7 – complementary to Illumina flow cell oligo

Unique vs. combinatorial dual index barcodes






Index 2 Primer



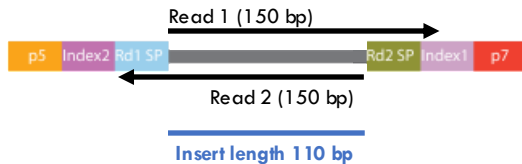
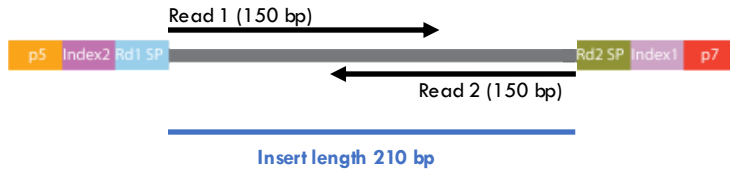
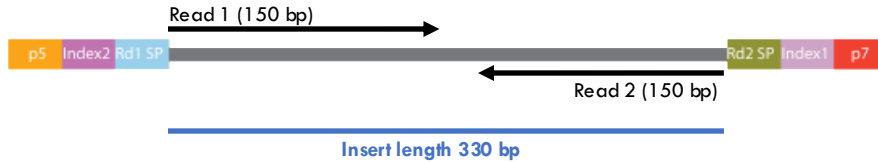
-  P5 – complementary to Illumina flow cell oligo
-  Indexing sequence 2
-  Read 1 Sequencing Primer

Index 1 Primer

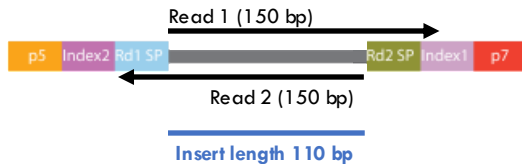
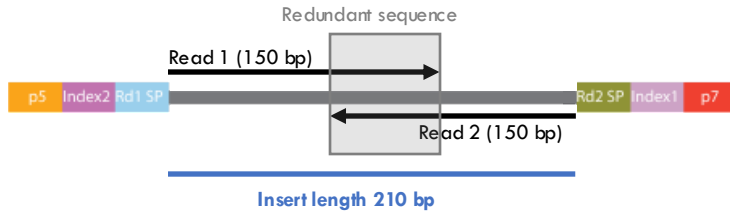
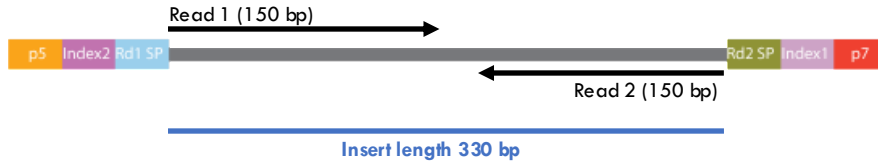


-  Read 2 Sequencing Primer
-  Indexing sequence 1
-  P7 – complementary to Illumina flow cell oligo

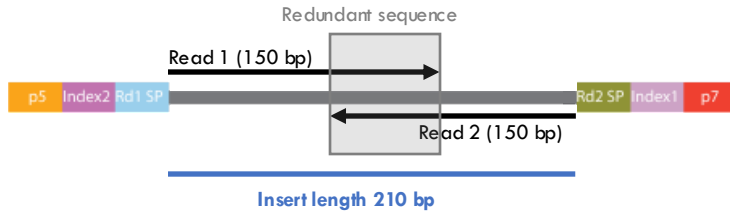
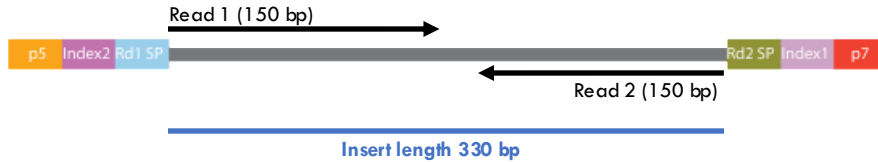
Insert length relative to read length



Insert length relative to read length

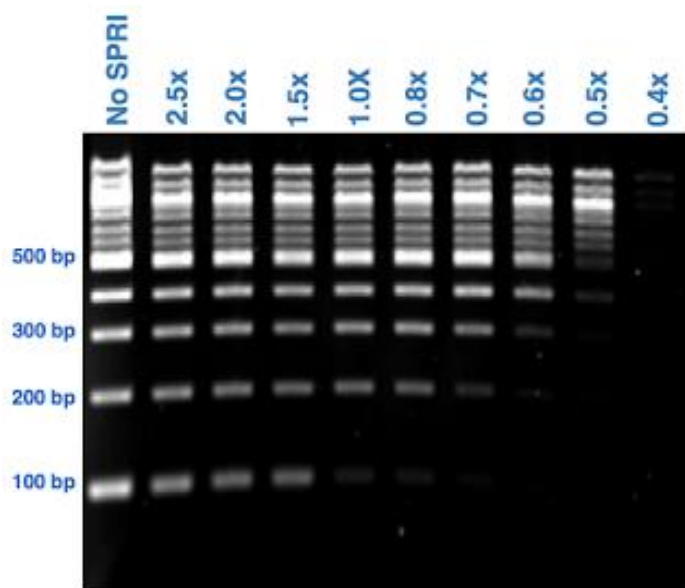


Insert length relative to read length

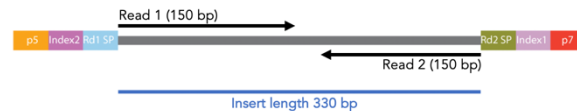


Size selection with Ampure beads

Tune the size distribution of your library fragments to minimize “waste” of sequence due to paired-end overlap and adapter read-through

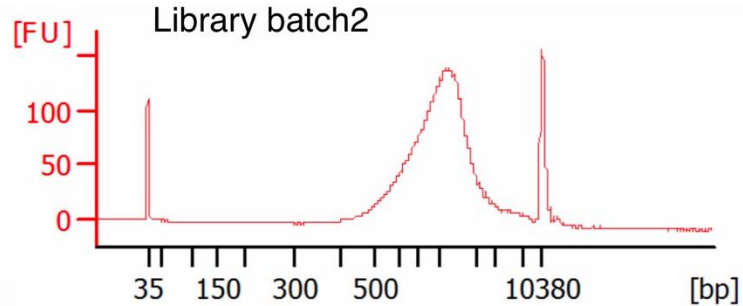
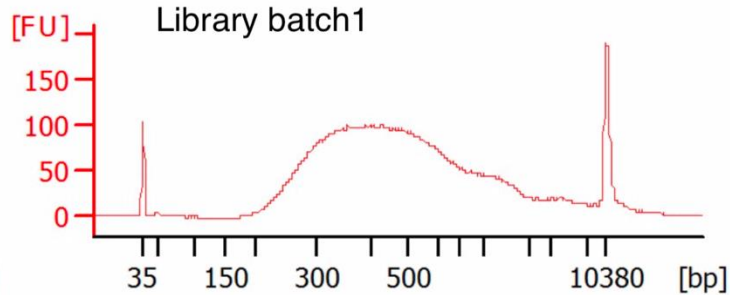


Ideally, we want all library fragments to be greater than the adapter length plus 2 x the read length (for PE)

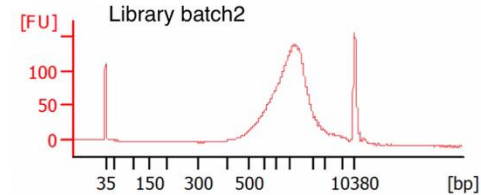
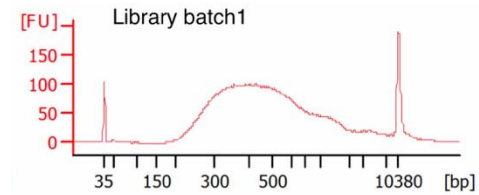
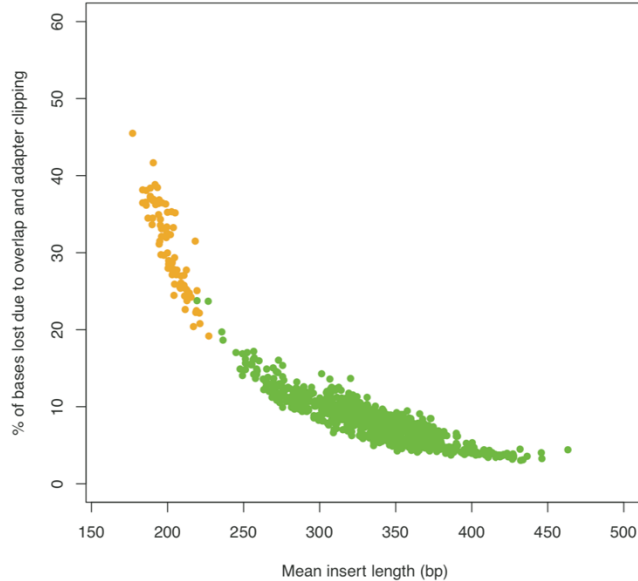


Ideal minimum fragment length

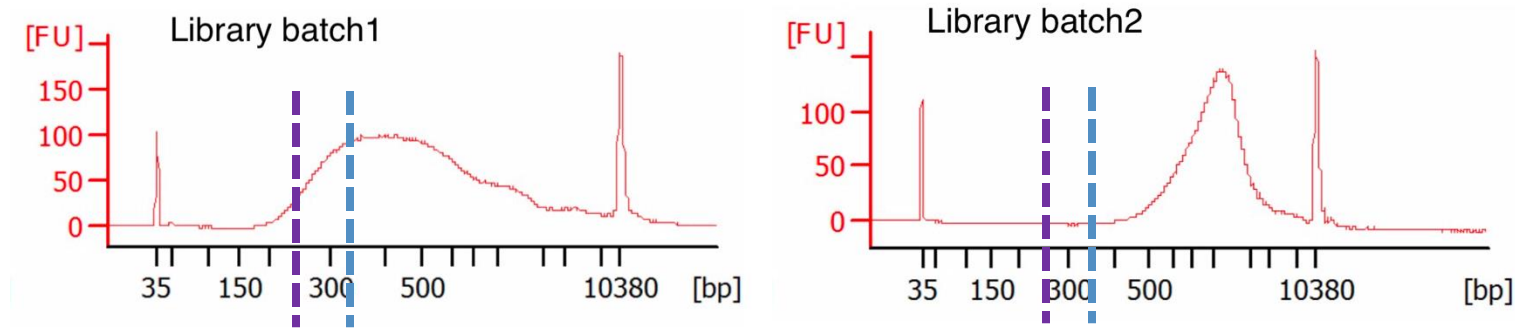
Two examples of our library pools



The library fragment size distribution can substantially influence the amount of data lost in data QC steps



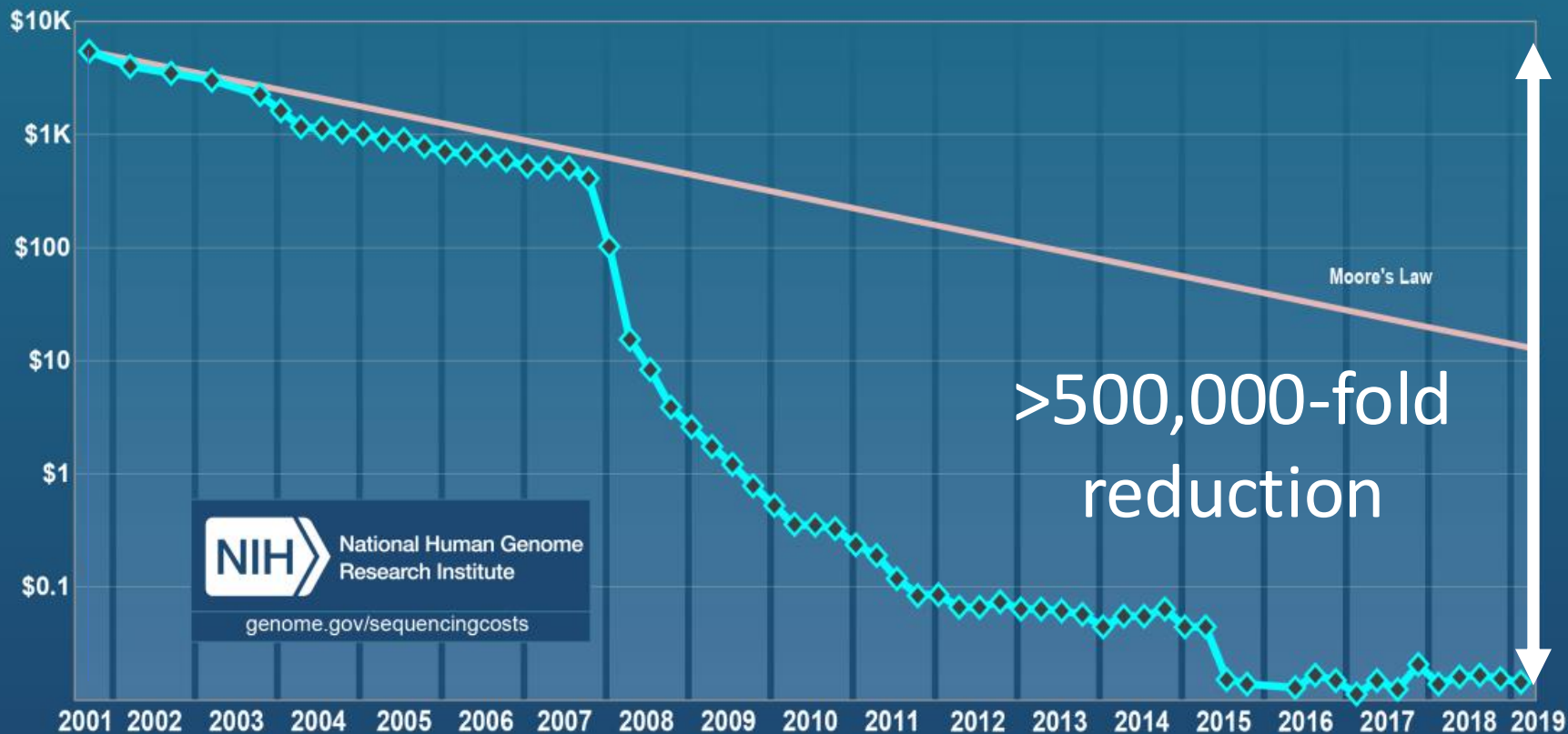
Two examples of our library pools



The length of Nextera adapters is 138 bp and libraries were sequenced with 2*125bp reads

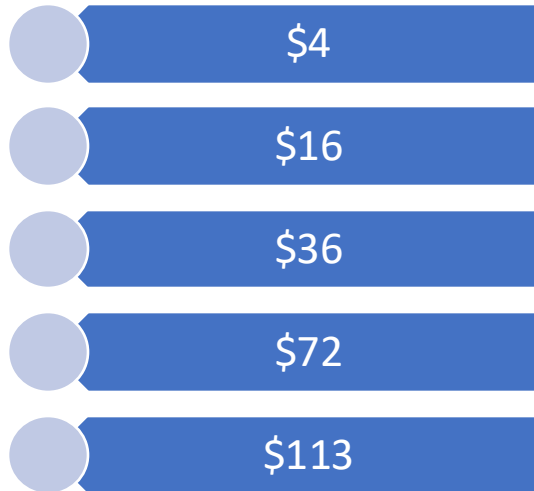
- Minimum fragment length to avoid overlap 388bp
- Minimum fragment length to avoid adapter read-through 263bp

Cost per Raw Megabase of DNA Sequence



What is the current price for 2x sequencing of an Atlantic silverside (including library preparation)?

Genome size ~650 Mb

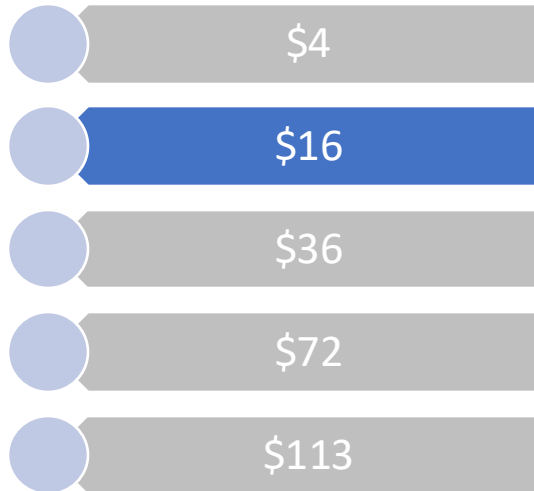


1 USD \approx 1 EURO



What is the current price for 2x sequencing of an Atlantic silverside?

Genome size ~650 Mb



1 USD \approx 1 EURO



Example costs for other genome sizes

Incl. library preparation and sequencing to 2x genome coverage*

Genome size (Gb)	Cost per sample (USD) ^a		Example organisms
	1× coverage	2× coverage	
0.2	11 (3)	13 (5)	Fruit fly, honeybee, arabidopsis
0.65	16 (8)	25 (17)	Atlantic silverside, stickleback, eastern oyster
1	21 (13)	34 (26)	Zebra finch, chicken, purple sea urchin
3	47 (39)	86 (78)	Human, Atlantic salmon, African clawed frog

*Cost estimates do not include labor and assume sequencing costs ~13 USD per Gb in shared S4 lanes on an Illumina NovaSeq and 8 USD per sample for library preparation

Example costs for other genome sizes

Incl. library preparation and sequencing to 2x genome coverage*

Genome size (Gb)	Cost per sample (USD) ^a	
	1× coverage	2× coverage
0.2	11 (3)	13 (5)
0.65	16 (8)	25 (17)
1	21 (13)	34 (26)
3	47 (39)	86 (78)

Compare to:

\$30 per sample for RADseq

\$15 per sample for RADcapture

Meek and Larson. 2019. Mol Ecol Res

*Cost estimates do not include labor and assume sequencing costs ~13 USD per Gb in shared S4 lanes on an Illumina NovaSeq and 8 USD per sample for library preparation

Exercise – how much will your experiment cost?

- Assumed costs:
 - Library preparation: \$8 per sample
 - Sequencing: \$4 per Gb (assuming a full NovaSeqX 25B lane)
 - Target coverage per sample: Expect to lose at least 30-50% of your data in filtering

Exercise – how much will your experiment cost?

- Assumed costs:
 - Library preparation: \$8 per sample
 - Sequencing: \$4 per Gb (assuming a full NovaSeqX 25B lane)
 - Target coverage per sample: Expect to lose at least 30-50% of your data in filtering
- **Example:** I would like to have 1x coverage for downstream analysis for 40 individuals from each of 5 populations (200 individuals total) of my favorite animal with a genome size of ~800 Mb
- **Calculation:** I will target 2x coverage raw sequencing. This means
$$2 * 800 \text{ Mb/individual} * 200 \text{ individuals} = 320,000 \text{ Mb (320 Gb)}$$
My total cost is thus $(320 \text{ Gb} * \$4/\text{Gb}) + (200 \text{ libraries} * \$8 \text{ per library}) = \textbf{\$2,880}$