

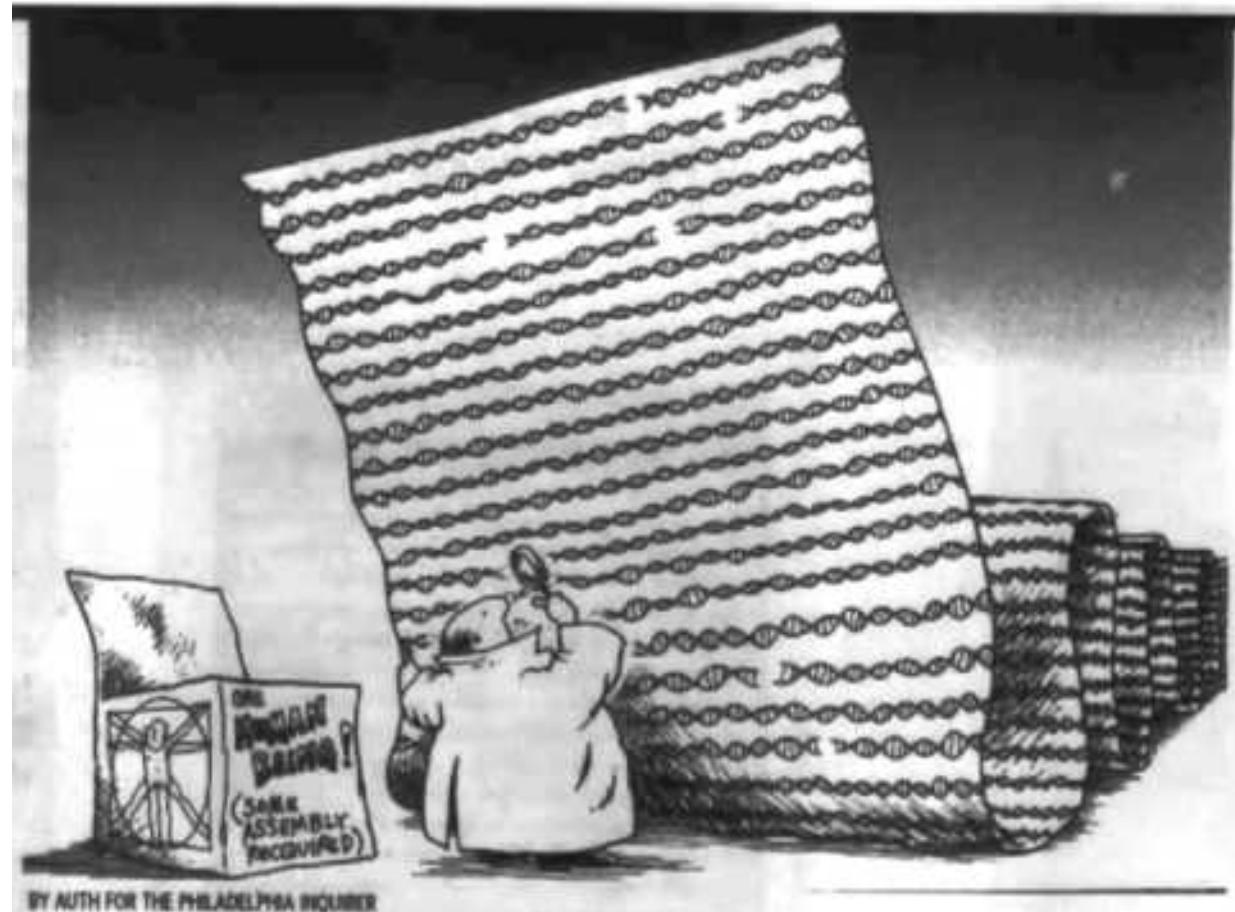
Overview of High-throughput Sequencing

Outline

- Omics
- History of NGS
- NGS platform chemistry
- Illumina platform
- File formats

NGS =
Next generation sequencing

HTS =
High throughput sequencing

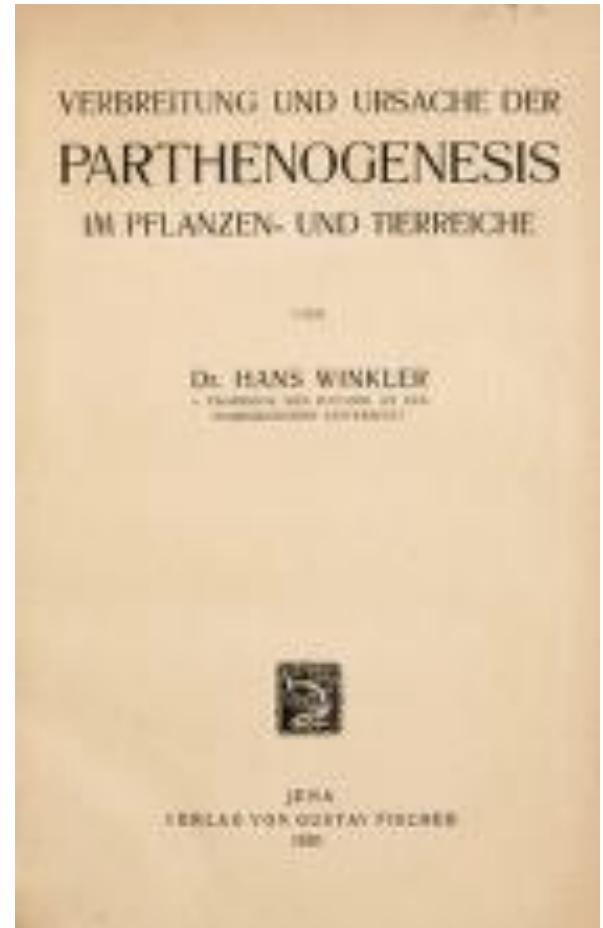


Birth of the word Genome

Translated from German:

"I propose the expression Genom
for the haploid chromosome set,
which, together with the pertinent
protoplasm, specifies the material
foundations of the species."

-Hans Winkler, German botanist,
from "Spread and cause of
pathogenesis in plant and animal
kingdoms," 1920



Genome becomes Omics

- Biome established as a word in 1916, genome in 1920
 - In molecular biology, “ome” has come to mean all constituents considered collectively
 - Surprisingly, unrelated to the word chromosome
 - All from 1990s:
 - Transcriptomics
 - Metabolomics
 - Proteomics
- “the collective characterization and quantification of pools of biological molecules that translate into the structure, function, and dynamics of an organism or organisms.”

Generations of High-Throughput Sequencing

- 1977 - 1st generation - Sanger sequencing
‘chain-termination’ or dideoxy technique
~800 base reads
- 2003 - Next generation or second generation
more data, less expensive
new sequencing by synthesis chemistries
35-300 base reads
- 2011 - Third generation
longer reads but more expensive per base
single molecule sequencing chemistries
300-90,000 bases

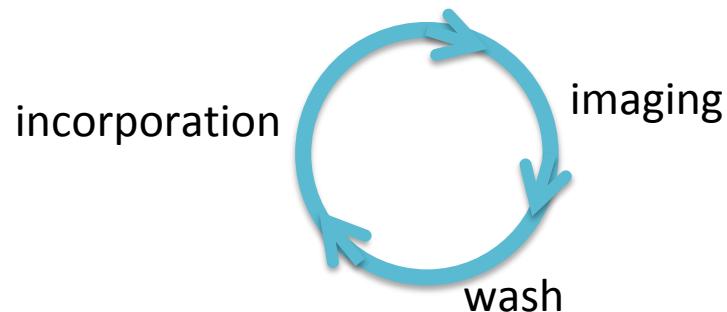
illumina

PB
PACIFIC
BIOSCIENCES®

Oxford
NANOPORE
Technologies®

Shared by most NGS

- DNA sample preparation
 - addition of defined sequences, known as “adapters,” to the ends of randomly fragmented DNA
 - Creates a “library”
- Immobilization
 - Adapters are used to anchor the individual DNA fragments to a solid surface like a glass slide
 - amplification is required to form spatially distinct and detectable sequence features (except PacBio)
- Sequencing
 - DNA polymerase synthesis with fluorescent nucleotides



illumina

Illumina Sequencing Technology

- Company was worth \$39.2 billion in June 2018
- 75%+ market share
- Why are they so popular?
 - Low price
 - High throughput
 - Paired end sequencing

Market Summary > Illumina, Inc.

NASDAQ: ILMN

+ Follow

352.54 USD +3.17 (0.91%) ↑

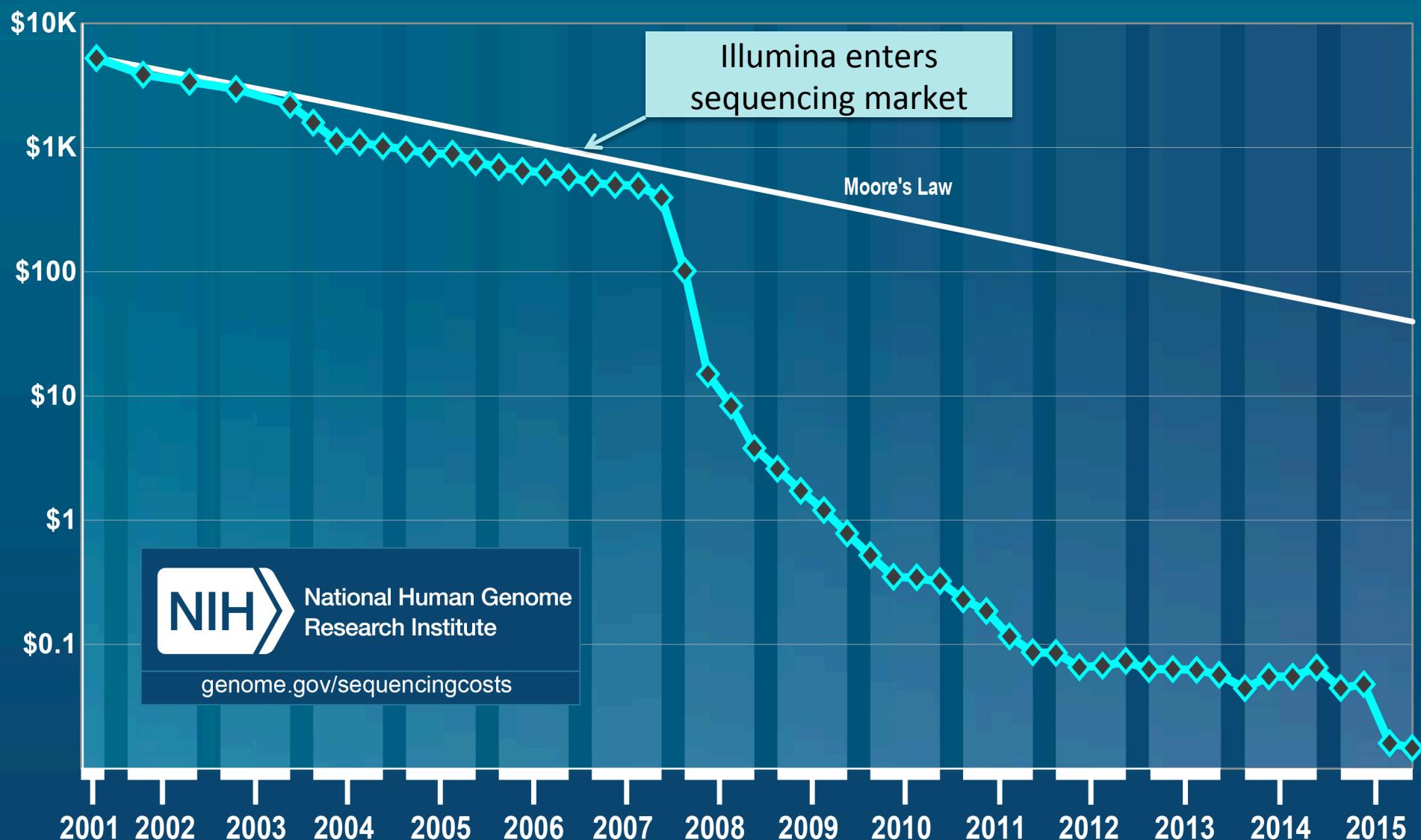
Sep 11, 12:22 PM EDT · Disclaimer

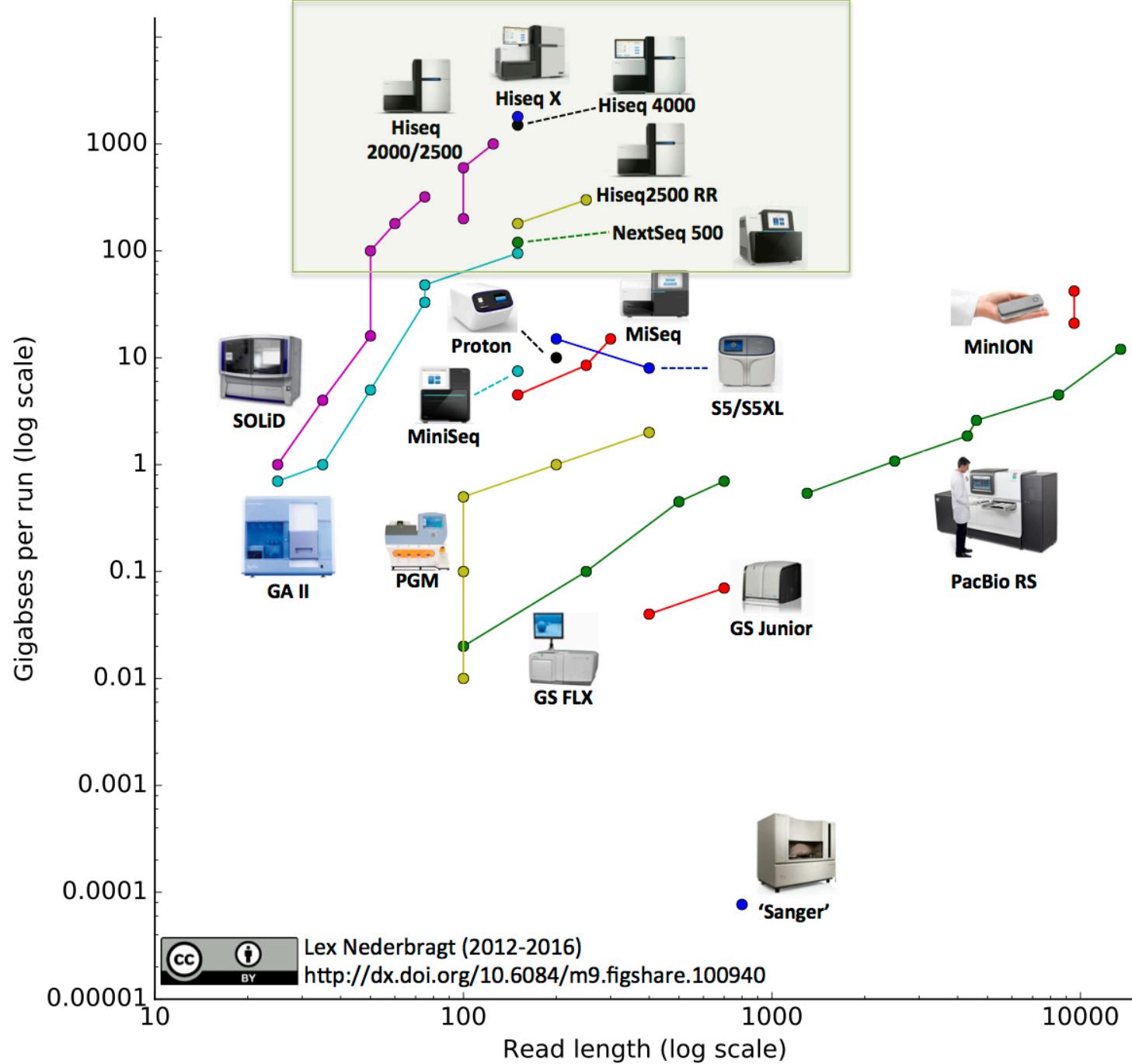
1 day 5 days 1 month 6 months YTD 1 year 5 years Max



Open	348.77	Div yield	-
High	353.16	Prev close	349.37
Low	345.09	52-wk high	357.93
Mkt cap	51.82B	52-wk low	196.00
P/E ratio	65.38		

Cost per Raw Megabase of DNA Sequence





Price and Throughput

	read type	price	Average read pair yield
MiSeq v3	Paired End (2x300)	\$1986	13.2Gb/run
HiSeq 4000	Paired End (2x150)	\$2588	100 Gb/lane
NextSeq 500	Paired End (2x150)	\$4869	100-120 Gb/run
NovaSeq 6000	Paired End (2x150)	\$9216	400-3000 Gb/cell



© 2014 Illumina, Inc. All rights reserved.

Illumina Sequencing at UTK

MiSeq in Center for
Environmental
Biotechnology

[https://ceb.utk.edu/
dna-sequencing/](https://ceb.utk.edu/dna-sequencing/)

Bioinformatics
Resource Center



Illumina Limitations

- Short read length
 - MiSeq (smaller throughput instrument)
 - 2x300
 - HiSeq/NextSeq/NovaSeq
 - 2x150
- Bias against sequencing through GC-rich regions or AT-rich regions
- Errors are likely to be SNPs and are likely to cluster at the ends of sequences

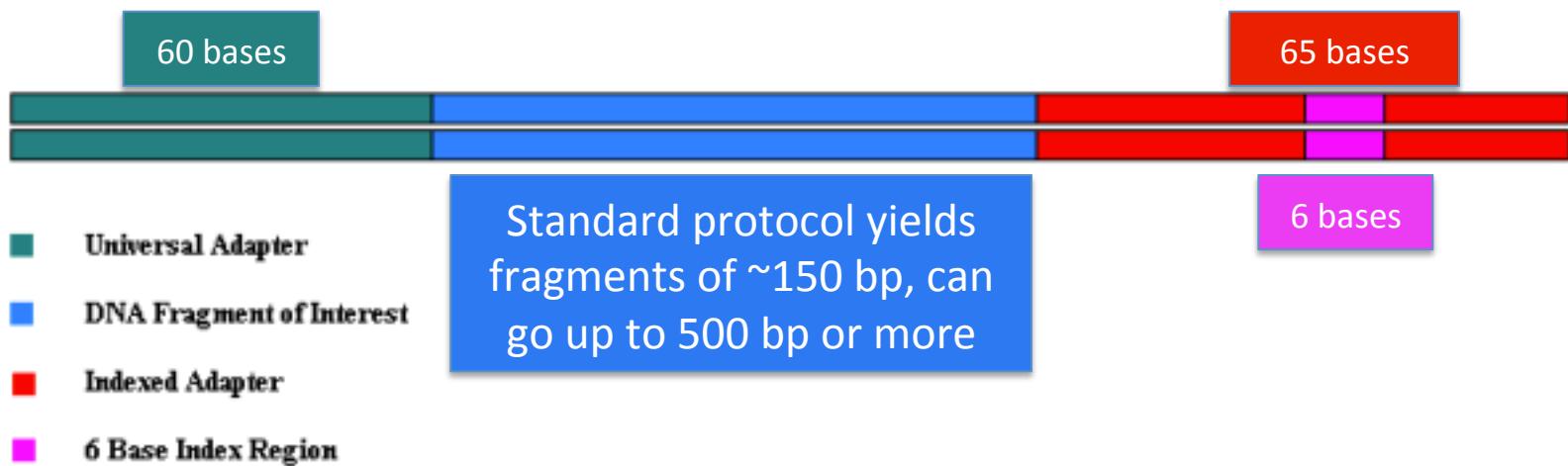
Illumina Sequencing Technology

- Video of Illumina Sequencing Technique w Nextera:
 - <https://www.youtube.com/watch?v=womKfikWlxM>
- Newest sequencing technologies (MiniSeq and NextSeq) use only 2 dyes instead of 4
 - <http://www.illumina.com/technology/next-generation-sequencing/sequencing-technology/2-channel-sbs.html>

How does it work?

Library construction can vary by kit

TruSeq Example:

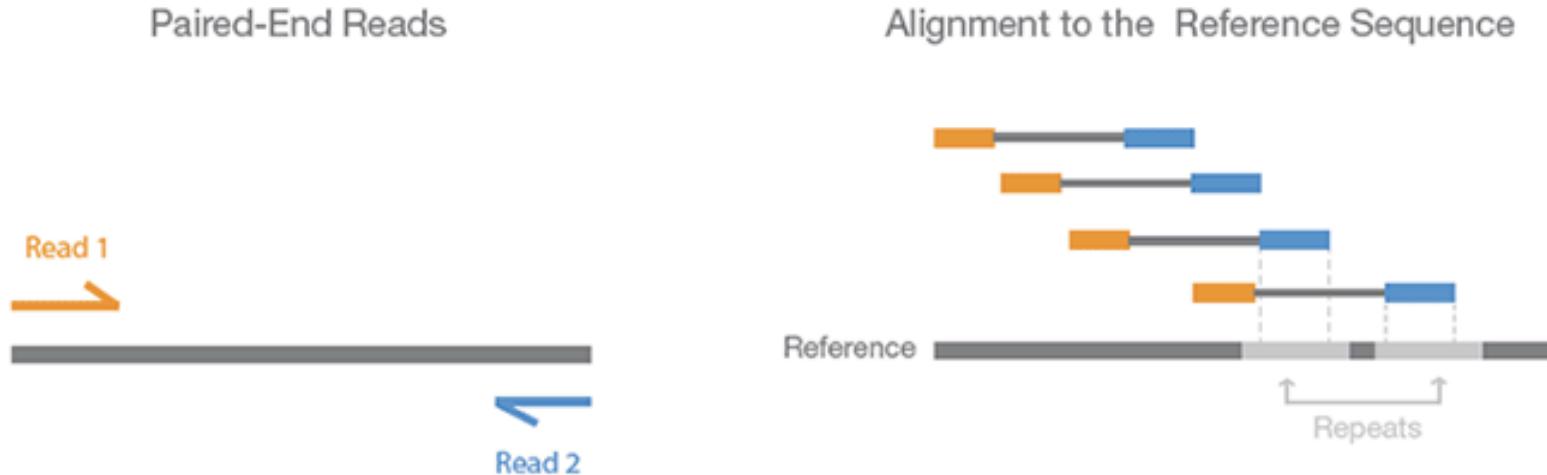


You will need the adapter sequences and a good understanding of adapter locations to later trim them out of your data

Paired End Sequencing

Overcome lack of length.

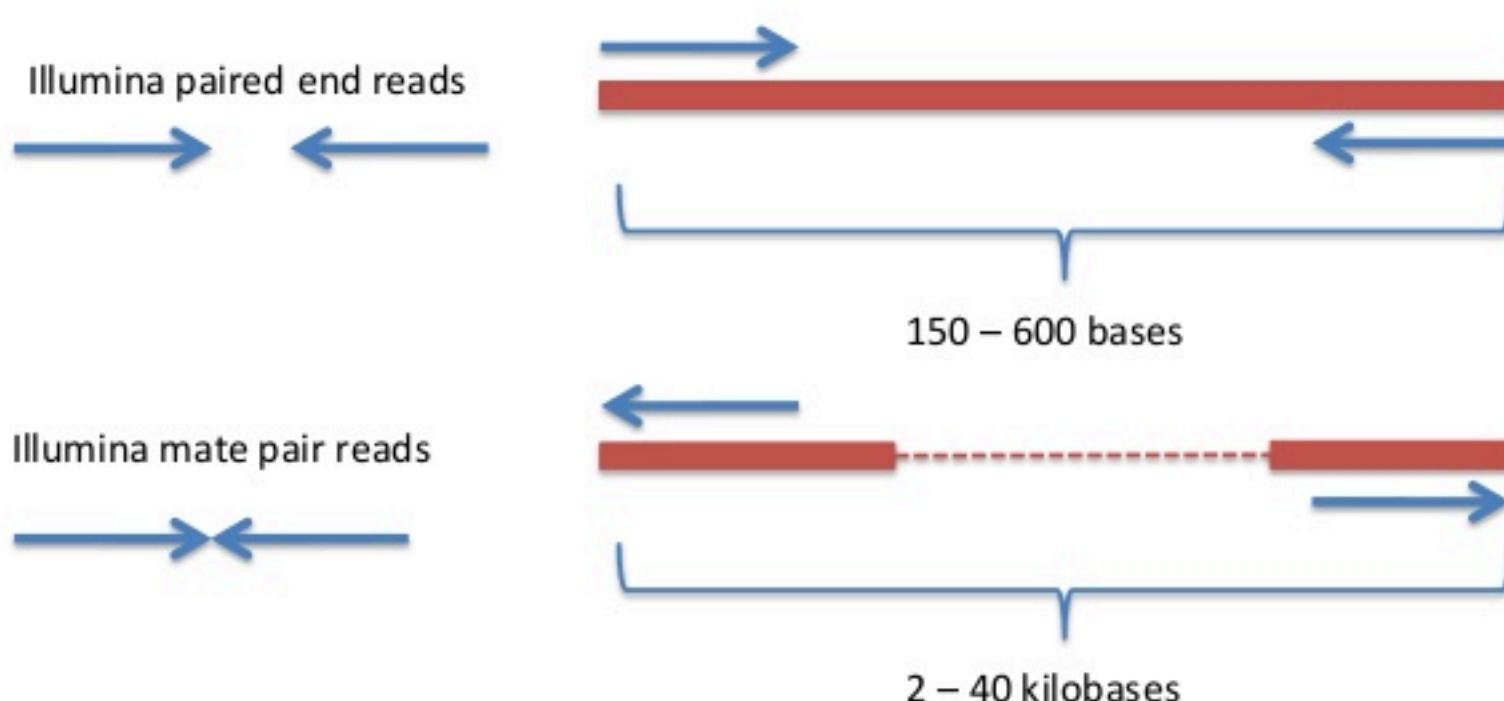
Figure 4. Paired-End Sequencing and Alignment



Paired-end sequencing enables both ends of the DNA fragment to be sequenced. Because the distance between each paired read is known, alignment algorithms can use this information to map the reads over repetitive regions more precisely. This results in much better alignment of the reads, especially across difficult-to-sequence, repetitive regions of the genome.

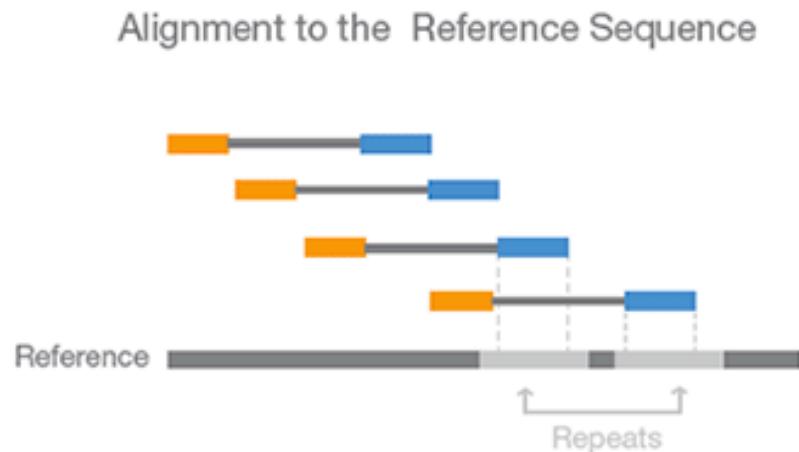
Mate Pair Sequencing

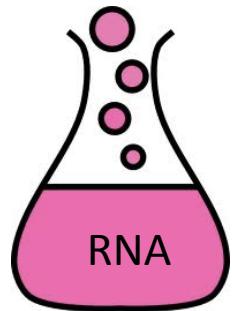
- Even longer jump distances



Utility of PE and MP for Repeats

- Can be used to sort out the number of repeats and the sequence of repeats
- Optimally, one end is in a unique region, the other end is in the repetitive region

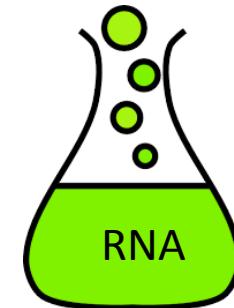
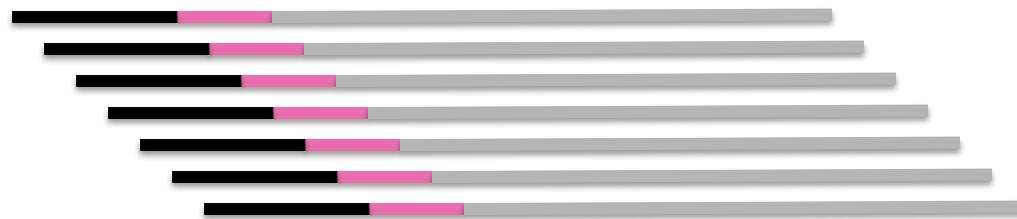




Multiplexing

Loading many samples into one lane.

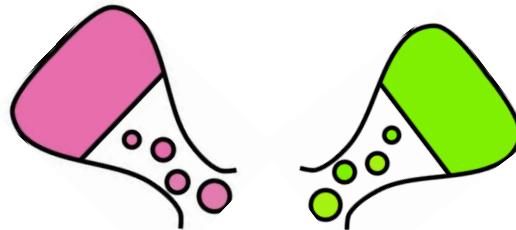
Pink Sample With **CGATGT**



Green Sample with **TGACCA**



CGATGT



Sequencing



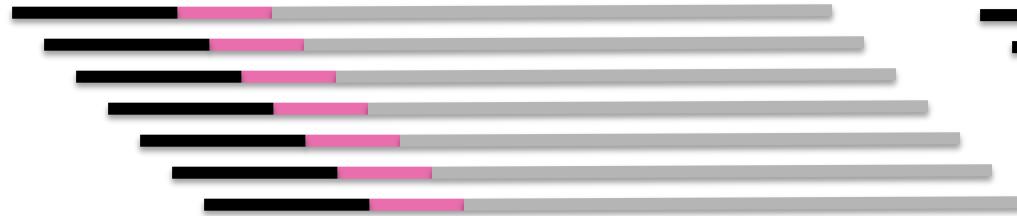
TGACCA



Software for De-multiplexing

Pink Sample File

Green Sample File



Multiplexing

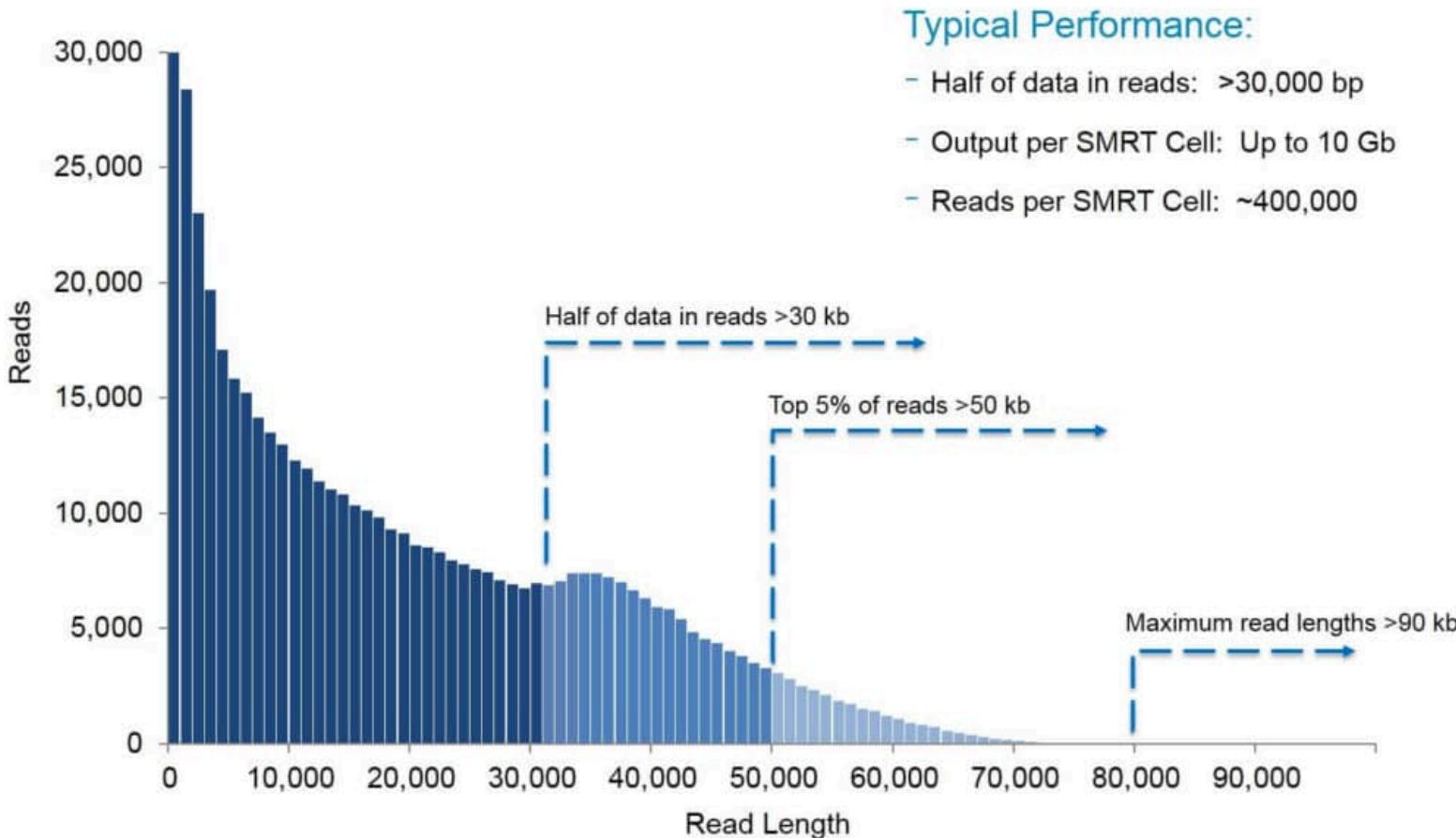
- Up to 384 barcodes available in some commercial kits
- Can be beneficial to select your barcodes wisely
 - Want an even distribution of nucleotides per cycle
 - If barcodes differ from each other at more than one position, you may be able to assign reads to a sample even with a sequencing error



- “PacBio”
- Commercially released 2011
- SMRT = single molecule real time
- No amplification needed
- More expensive
- Very long reads
- 15% error
 - Truly random error
 - Indels instead of SNP errors
- No GC bias
- Can detect methylation of nucleotides without alterations to the DNA
- Requires a large amount of high quality DNA



SEQUEL SYSTEM PERFORMANCE: GENOMIC LIBRARY





- MinION
- Extremely different mode for sequencing
 - Tiny (handheld) instruments and disposable chips
 - No fluorescence and no polymerase
 - uses ion current disruptions while “unzipping” the DNA
 - No fixed run time
 - Long reads, can read each strand twice

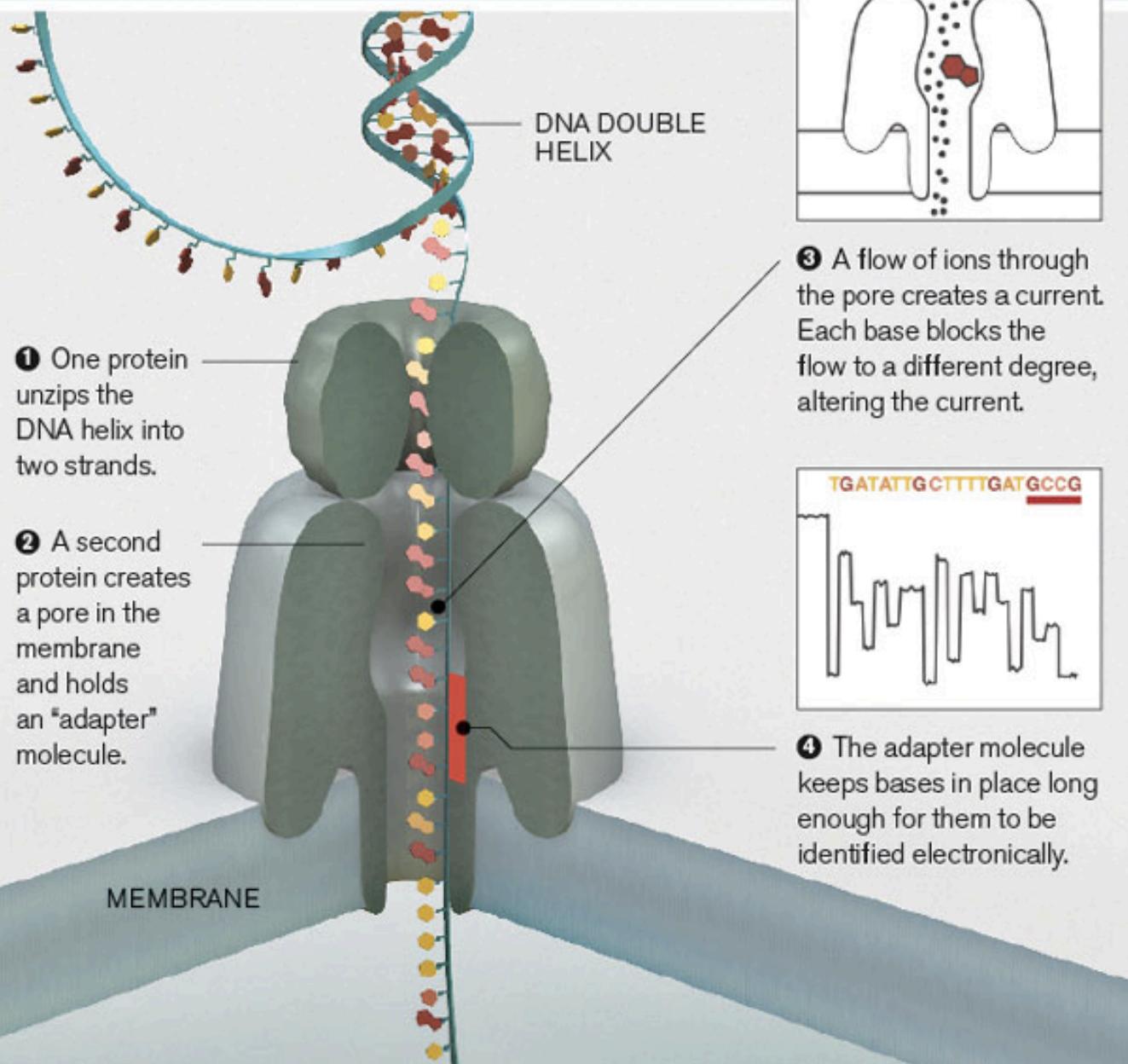


Does still require library prep, but field based applications are very exciting!



Smidglon
Coming soon!

DNA can be sequenced by threading it through a microscopic pore in a membrane. Bases are identified by the way they affect ions flowing through the pore from one side of the membrane to the other.



Video:

[https://
www.youtube.com/
watch?
v=3UHw22hBpAk](https://www.youtube.com/watch?v=3UHw22hBpAk)

First Step to Data Analysis

Preprocessing

- Turning the raw signal of the instrument into base calls and quality scores
- Images are converted to text and numbers
- Supplied by vendor-provided software and usually done by sequencing lab
- Image files are usually not kept because of their size
- You usually will not have to deal with this, but best to be aware that it is happening

Additional Steps in Data Analysis

- Why would we want to learn about the sequencing chemistry to do bioinformatics?
 - Understand biases
 - Spot nonsensical results
 - Plan for robust statistical analysis
 - Good experimental design
- Always get the sequences for your adapters and barcodes for later analysis!
 - Illumina has an open letter with most of their adapters:
 - https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry_documentation/experiment-design/illumina-customer-sequence-letter.pdf
 - Kits should come with additional documentation

File Formats

- We already learned about Fasta
- Raw reads come in Fastq format
 - Q for quality
 - This format includes quality values for each individual base

Fastq Format

```
@SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
```

```
CAGCACTAATGCACCGGATCCCATCAGAACTCCGCAGTTAA
```

```
+SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
```

```
BACBC9BCC@.>C>96;CB@?:?BB7@5>BA=:4.:B9>BB@
```

```
@SRR070570.2 HWUSI-EAS455:3:1:1:1785 length=41
```

```
CCAGAACACAAAGCTCATGACACGTTCACCTCCTGGAAAGTT
```

```
+SRR070570.2 HWUSI-EAS455:3:1:1:1785 length=41
```

```
>AB@ACBB<BCA:>B;AA;@<B=;-=;<?@?<?=1-?B<8A
```

```
@SRR070570.3 HWUSI-EAS455:3:1:1:1679 length=41
```

```
ATCGATGAAGAACGTAGCGAAATGCGATACTGGTGTGAAT
```

```
+SRR070570.3 HWUSI-EAS455:3:1:1:1679 length=41
```

```
BA==:=4?:8>A:8:>6:4:;2<07,<:@582+22'-';@>
```

Fastq Format

Sequence Identifier

Optional Description

```
@SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
```

```
CAGCACTAATGCACCGGATCCCATCAGAACTCCGCAGTTAA
```

```
+SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
```

```
BACBC9BCC@.>C>96;CB@:?BB7@5>BA=:4.:B9>BB@
```

Fastq Format

The Sequence

```
@SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
CAGCACTAATGCACCGGATCCCATCAGAACTCCGCAGTTAA
+SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
BACBC9BCC@.>C>96;CB@?:?BB7@5>BA=:4.:B9>BB@
```

Fastq Format

Totally useless line that begins with a + but does not need anything else; id and description are sometimes repeated.

```
@SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
CAGCACTAATGCACCGGATCCCATCAGAACTCCGCAGTTAA
+SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
BACBC9BCC@.>C>96;CB@?:?BB7@5>BA=:4.:B9>BB@
```

Fastq Format

Quality values for each base.

```
@SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
CAGCACTAATGCACCGGATCCCATCAGAACTCCGCAGTTAA
+SRR070570.1 HWUSI-EAS455:3:1:1:1388 length=41
BACBC9BCC@.>C>96;CB@:?BB7@5>BA=:4.:B9>BB@
```

FASTQ Quality Values

- Based on phred quality scoring system
(developed in the 1990s)

Phred quality scores are logarithmically linked to error probabilities

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%

https://en.wikipedia.org/wiki/Phred_quality_score

Ewing et al, 1998

FASTQ Quality Values

- Storing it based on the numbers 0-60 takes up too much space
- Sequence: ACTGATC
- Quality: 10 15 25 15 17 32 35

- Instead, assign individual letters, numbers and symbols to represent numeric quality values
- No need for space
- New Quality: JOXOLb

The Quality Value Debacle

This has been solved for about 6 years now, but is still a good cautionary tale and an excellent example of how NOT to create a standard format.

See wiki page for converters. http://en.wikipedia.org/wiki/FASTQ_format