# Introduction to Read Mapping
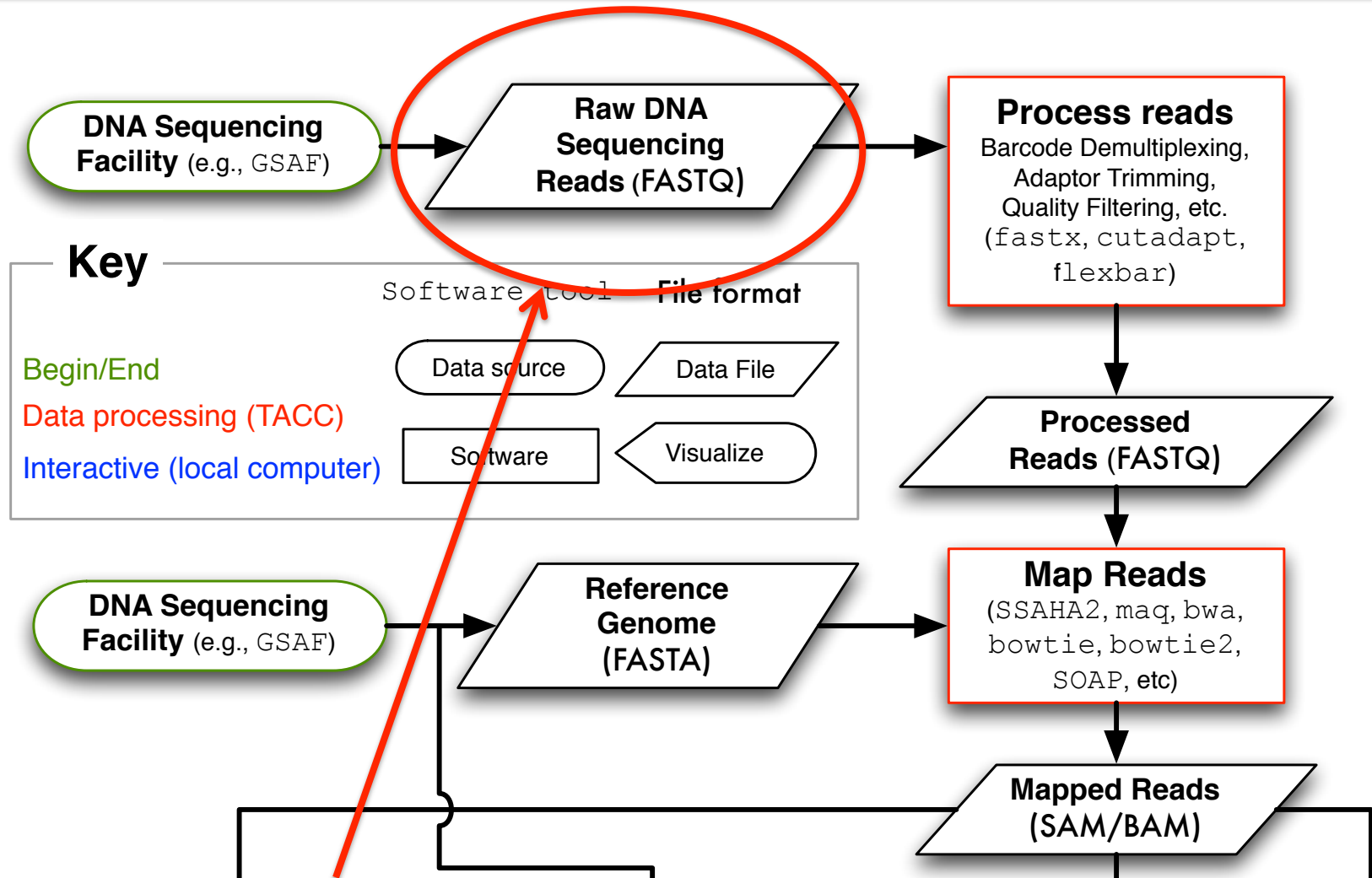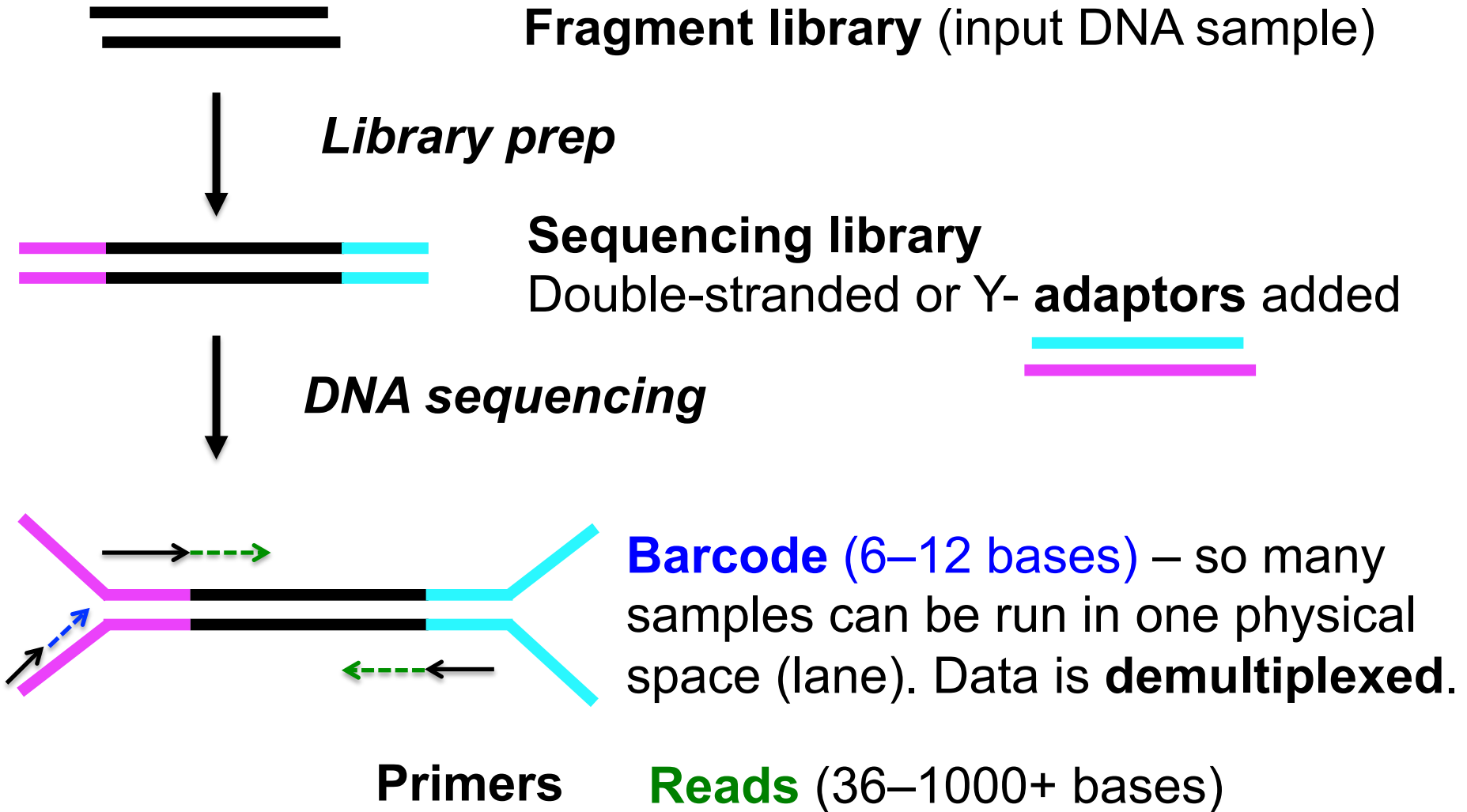
# Basic steps of mapping reads

1. Read file quality control and processing
2. Build reference sequence index
3. Map DNA sequencing reads
   – Exact tool/approach depends on sequencing technology and DNA fragment library type
4. Convert result to SAM/BAM database
5. Application specific analysis...
   – These steps are common to any reference-based (opposed to *de novo*) data analysis.
   – We will use the mapped reads for variant calling.

# Input: Raw DNA Sequencing Reads

# Read terminology

**Fragment library** (input DNA sample)

*Library prep*

**Sequencing library**
Double-stranded or Y- **adaptors** added

*DNA sequencing*

**Barcode** (6–12 bases) – so many samples can be run in one physical space (lane). Data is **demultiplexed**.

**Primers**   **Reads** (36–1000+ bases)

# Types of Illumina fragment libraries

# Read file format

FASTQ Format

```
@HWI-EAS216_91209:1:2:454:192#0/1
GTTGATGAATTTCTCCAGCGCGAATTTGTGGGCT
+HWI-EAS216_91209:1:2:454:192#0/1
B@BBBBBB@BBBBAAAA>@AABA?BBBAAB??>A?
```

Line 1: @read name

Line 2: called base sequence

Line 3: +read name (optional after +)

Line 4: base quality scores
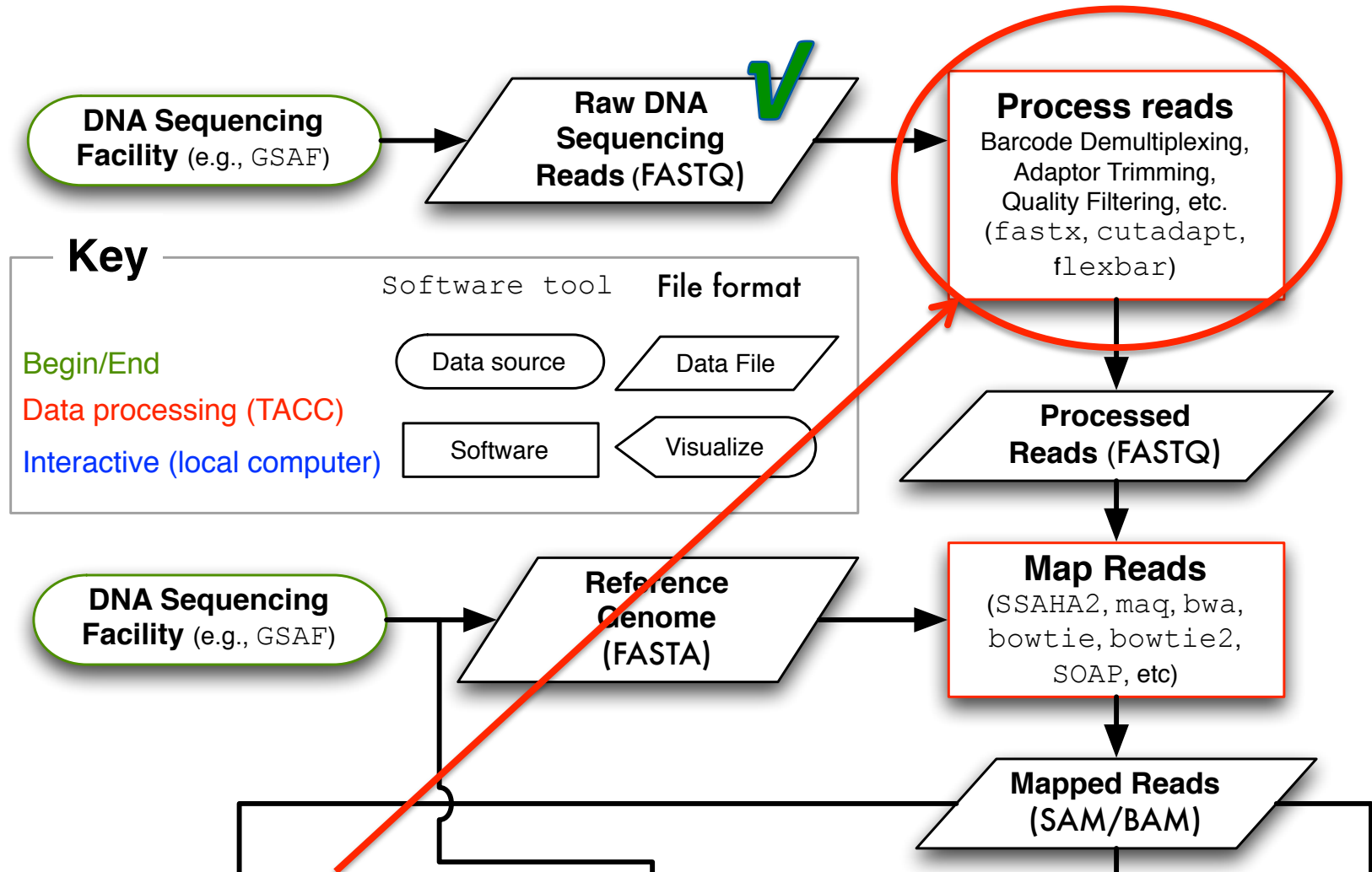
# Deciphering base quality (Q) scores

http://www.asciitable.com/

```
Quality character    !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHI
                     |           |           |           |           |
ASCII Value          33         43          53          63          73
Base Quality (Q)     0          10          20          30          40
```

Probability of Error $= 10^{-Q/10}$

(This is a **Phred** score, also used for other types of qualities.)

\* Very low quality scores can mean something special –
   Illumina Q ≤ 3 means something like: "I'm lost, you might
   want to stop believing sequencing cycles from here on out."

\* In older FASTQ files, the formula and ASCII offset might differ.
   Consult: http://en.wikipedia.org/wiki/FASTQ_format
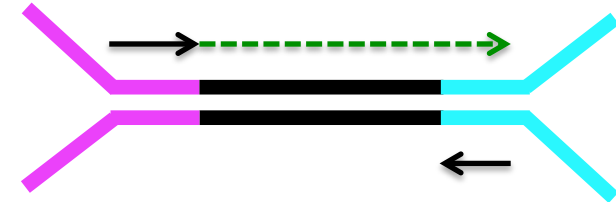
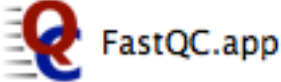# Input: Raw DNA Sequencing Reads

# Read sequence quality control
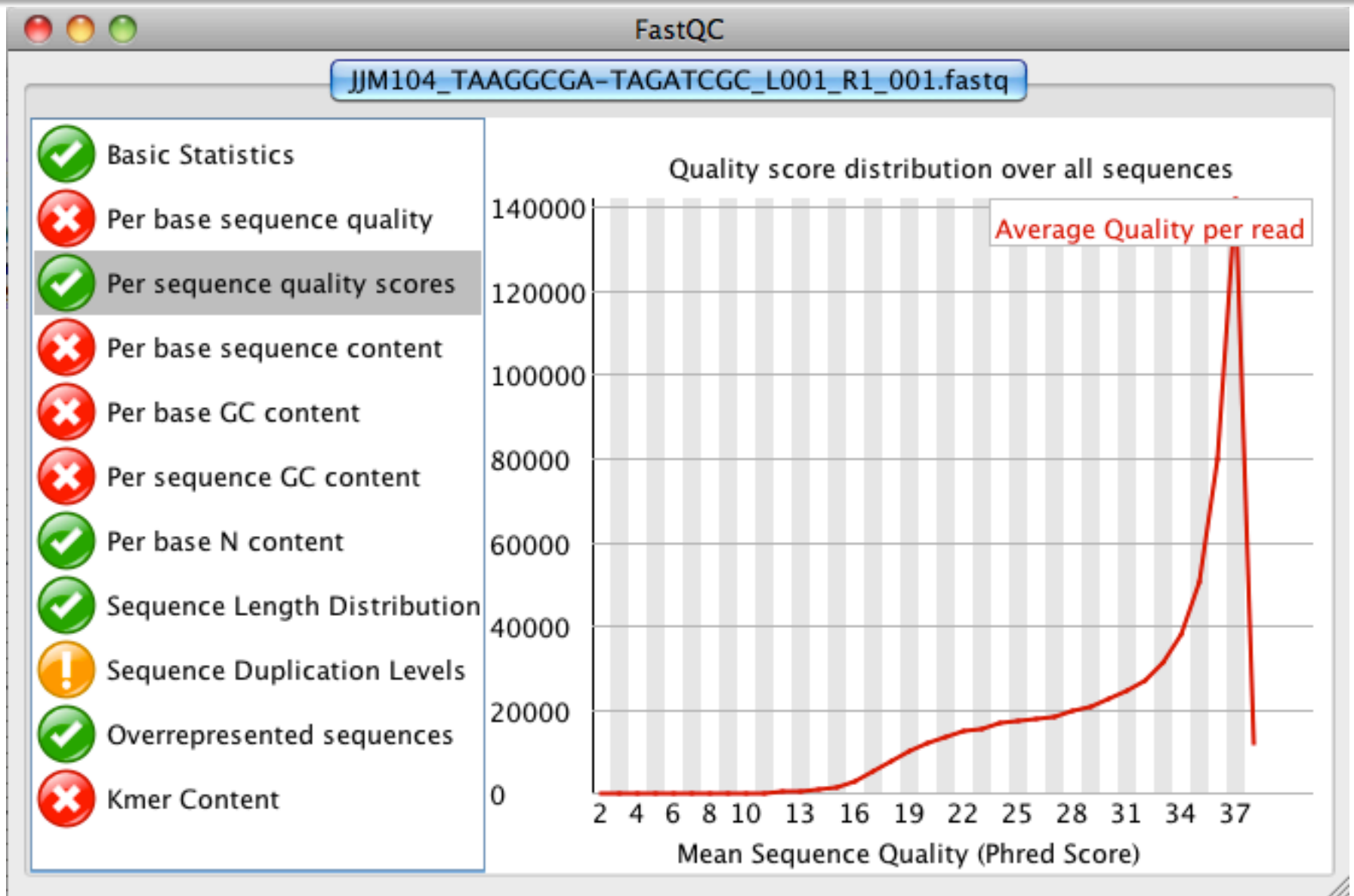
**Garbage in = garbage out**

- Contaminated with other samples?

- Sample barcodes removed?

- Adaptor/bar codes trimmed?

  – Esp. important for MiSeq data

- Trim ends of reads with poor quality?

  – Less data but higher quality data

- Know your data

  – Paired reads? Relative orientations?

  – Technology specific concerns?

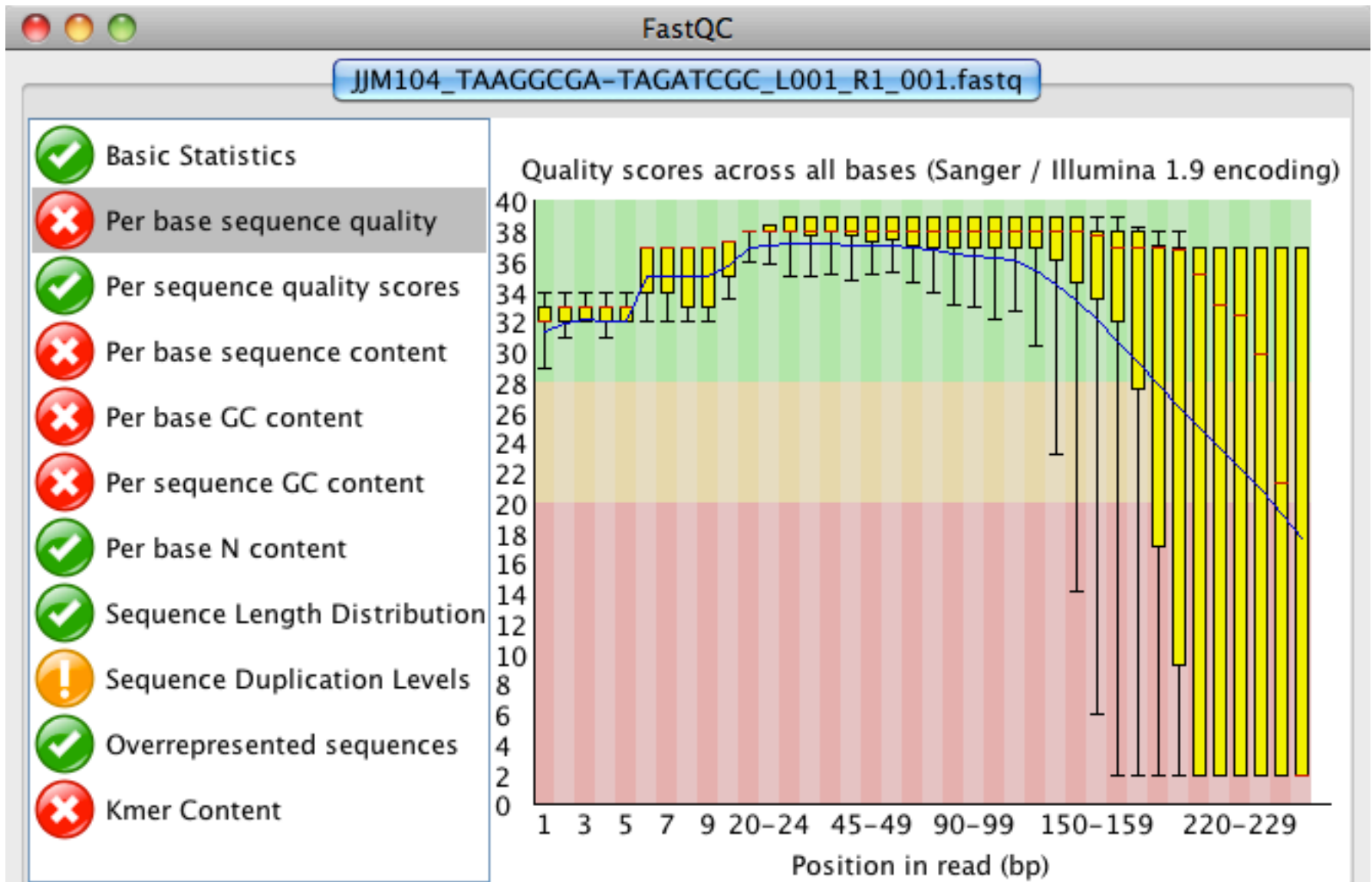    • Error hotspots (e.g., 454 Indels, Illumina GGT)

# Read quality control software

- FastQC is pretty much the only game in town
  - TACC module or run on your own computer
  - Generates nice graphical output  FastQC.app

- Do not be surprised if some criteria "fail" even for really good FASTQ data !!!

- Example FASTQ stats on the next two slides are for the 1st read of a paired-end 250-cycle MiSeq run of *E. coli* DNA.
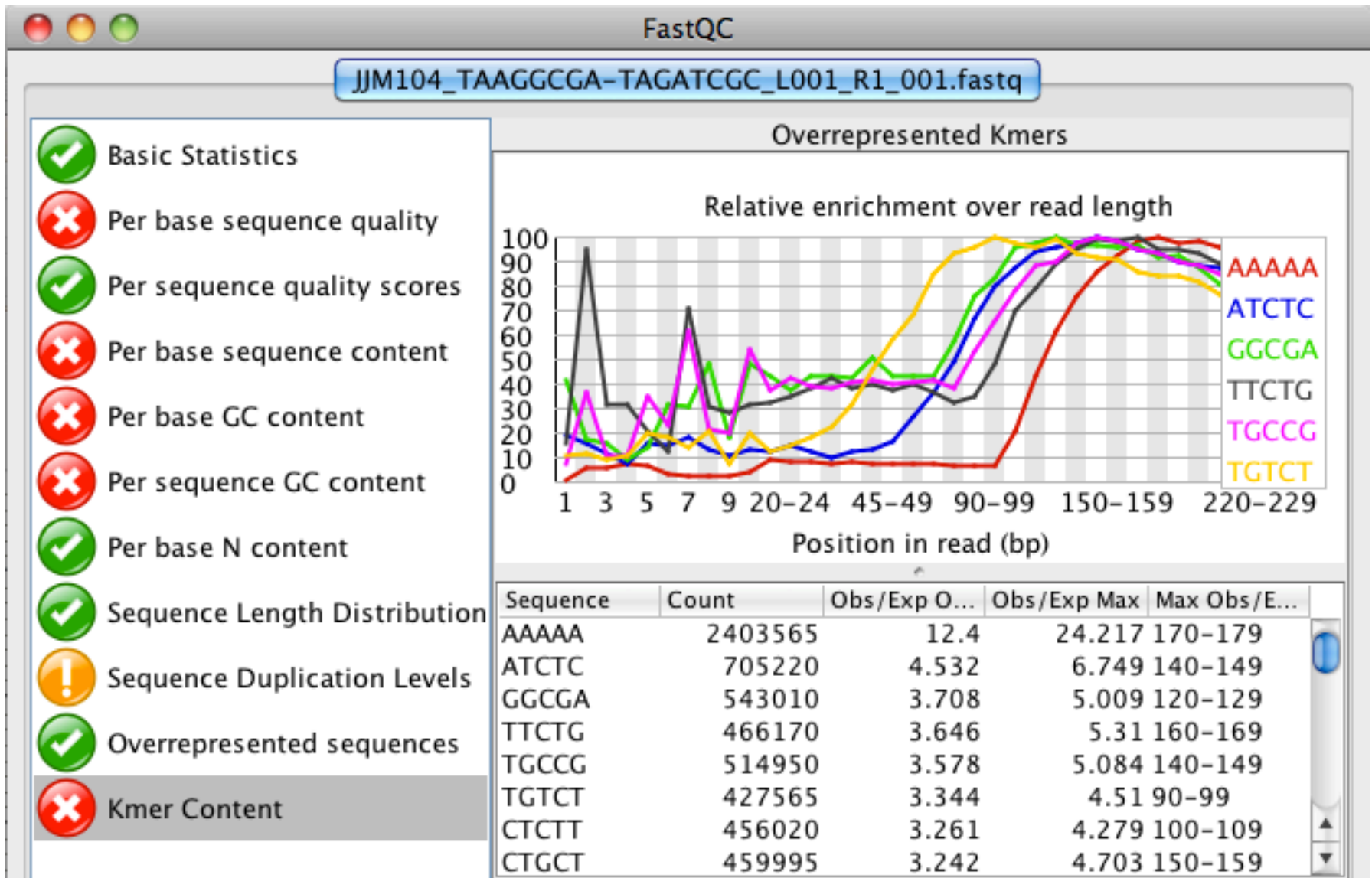
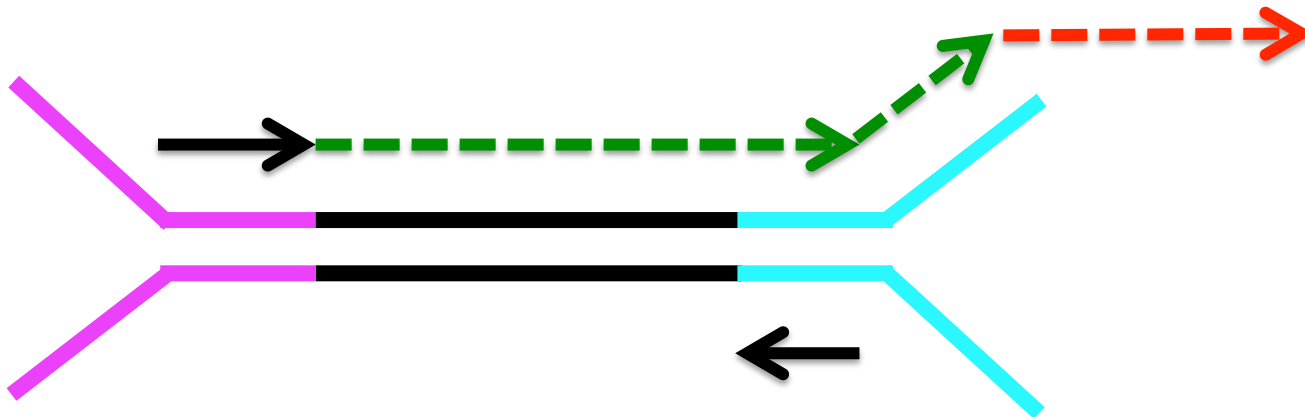# Illumina data example

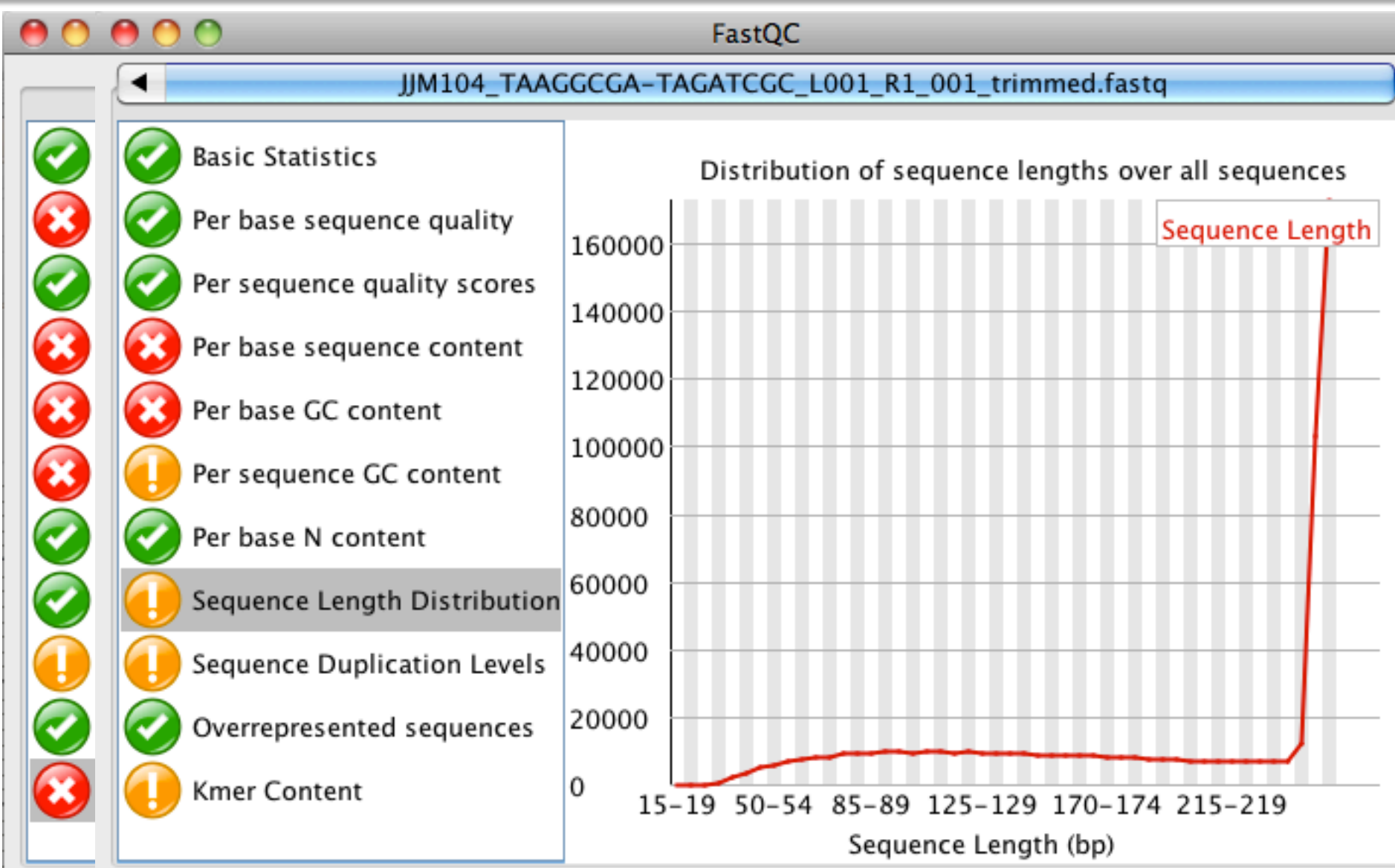# Illumina data example

# Illumina data example

# Problem in this data set?

- Adaptor/bar codes trimmed?
  – Esp. important for MiSeq data
- DNA was sheared to smaller than the read length, so many reads extend past the end. They need their 3' ends trimmed of the adaptor and junk sequence.

# Processed Illumina data example

# Processed Illumina data example