

BIOL5488: Genomics Lab II

Table of Contents

- ❖ Review
- ❖ Assignment 2
- ❖ Q & A

Adapted from
Bio5488 Genomics 2023 TAs

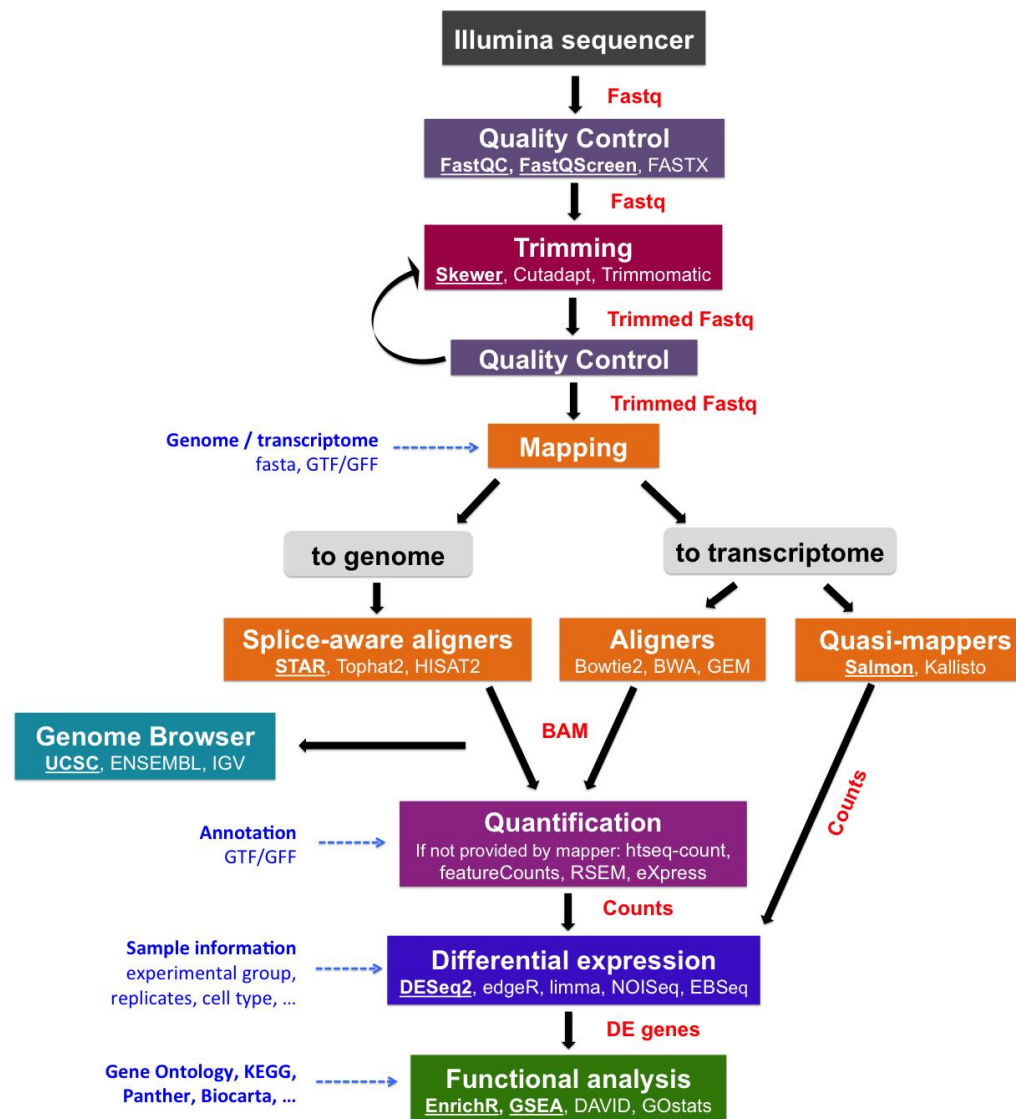
```
<?php body=
d="fb-root"></div>
t>(function(d, s, id) {
js, fjs = d.getElementsByTagName(s)[0];
(d.getElementById(id)) return;
= d.createElement(s); js.id = id;
src = "//connect.facebook.net/en_US/sdk.js#xfbml=1&vers
.parentNode.insertBefore(js, fjs);
ument, 'script', 'facebook-jssdk'));</script>
id="page" class="site">
<a class="skip-link screen-reader-text" href="#content">

<header id="masthead" class="site-header" role="banner">
  <div class="site-branding">
    <div class="navBtn pull-left">
      <?php if(is_home() && $xpanel['homepage-st
      <a href="#" id="openMenu"><i class="fa fa-
      <?php } else { ?>
      <a href="#" id="openMenu2"><i class="fa fa
      <?php } ?>
    </div>
    <div class="logo pull-left">
      <a href="<?php echo esc_url( home_url() )
        
    </div>
    <div class="search-box hidden-xs hidden-sm p
      <?php get_search_form(); ?>
    </div>
    <div class="submit-btn hidden-xs hidden-sm
      <a href="<?php echo get_page_link($xpan
    </div>
    <div class="user-info pull-right mr-10">
```



Review





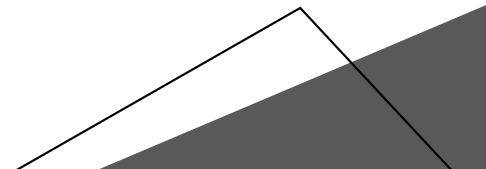
← We are here



FASTA format

1. Represent nucleotide sequences or amino acid sequences
2. The filename ends with .fa/.fna/.fasta
3. The first line start with “>” as the title & description
(sometimes “;” will use to start a comment line)
4. The following lines are data

```
>NG_008679.1:5001-38170 Homo sapiens paired box 6 (PAX6)
ACCCTCTTTTCTTATCATTGACATTTAAACTCTGGGGCAGGTCCTCGCGTAGAACGCGGCTGTCAGATCT
GCCACTTCCCCTGCCGAGCGGCGGTGAGAAGTGTGGGAACCGGCGCTGCCAGGCTCACCTGCCTCCCCGC
CCTCCGCTCCCAGGTAACCGCCCCGGGCTCCGGCCCCGGCCCGGCTCGGGGCCCCGCGGGGCCTCTCCGCTG
CCAGCGACTGCTGTCCCCAAATCAAAGCCCCGCCCAAGTGGCCCCGGGGCTTGATTTTTGCTTTTAAAAG
GAGGCATACAAAGATGGAAGCGAGTTACTGAGGGAGGGATAGGAAGGGGGGTGGAGGAGGGACTTGTCTT
TGCCGAGTGTGCTCTTCTGCAAAAGTAGCAAAATGTTCCACTCCTAAGAGTGGACTTCCAGTCCGGCCCT
GAGCTGGGAGTAGGGGGCGGGAGTCTGCTGCTGTCTGCTAAAGCCACTCGCGACCGCGAAAAATGCA
GGAGGTGGGGACGCACTTTGCATCCAGACCTCCTCTGCATCGCAGTTCACGACATCCACGCTTGGGAAAG
TCCGTACCCGCGCCTGGAGCGCTTAAAGACACCCTGCCGCGGGTTCGGGCGAGGTGCAGCAGAAGTTTCCC
GCGGTTGCAAAGTGCAGATGGCTGGACCGCAACAAAGTCTAGAGATGGGGTTCGTTTCTCAGAAAGACGC
```





FASTQ format

1. Storing both nucleotide sequence and corresponding quality scores
-> usually handle reads from the sequencing machine

2. The filename ends with .fq/.fastq

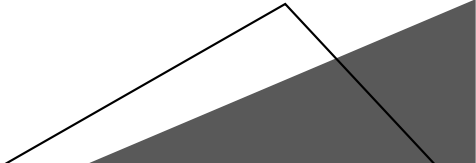
3. Each sequence has four lines:

Line 1: begins with '@', the following context is about sequence identifier
(sequence tech, flow cell IDs, info about read pairs etc)

Line 2: the sequence

Line 3: begins with '+', the following context could be the same as Line 1

Line 4: encodes the quality of the sequence in line 2 by ASCII. The letters in Line 4 should be the same as the letters in Line 2.





FASTQ format: example

Identifier | @HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Sequence | TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTTNNNNNNNNNNNTAGTTTCTTGAGA
+ sign & identifier | +HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Quality scores | efcfffffcfeefffcfffffddf`feed]`_]_Ba_^__[YBBBBBBBBBBRTT\]][] dddd`

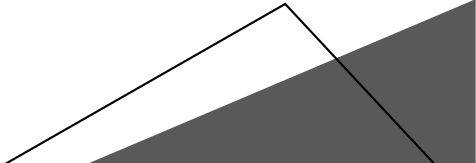
Base T
phred Quality] = 29

Each letter in line 4 represents a *Qphred value*

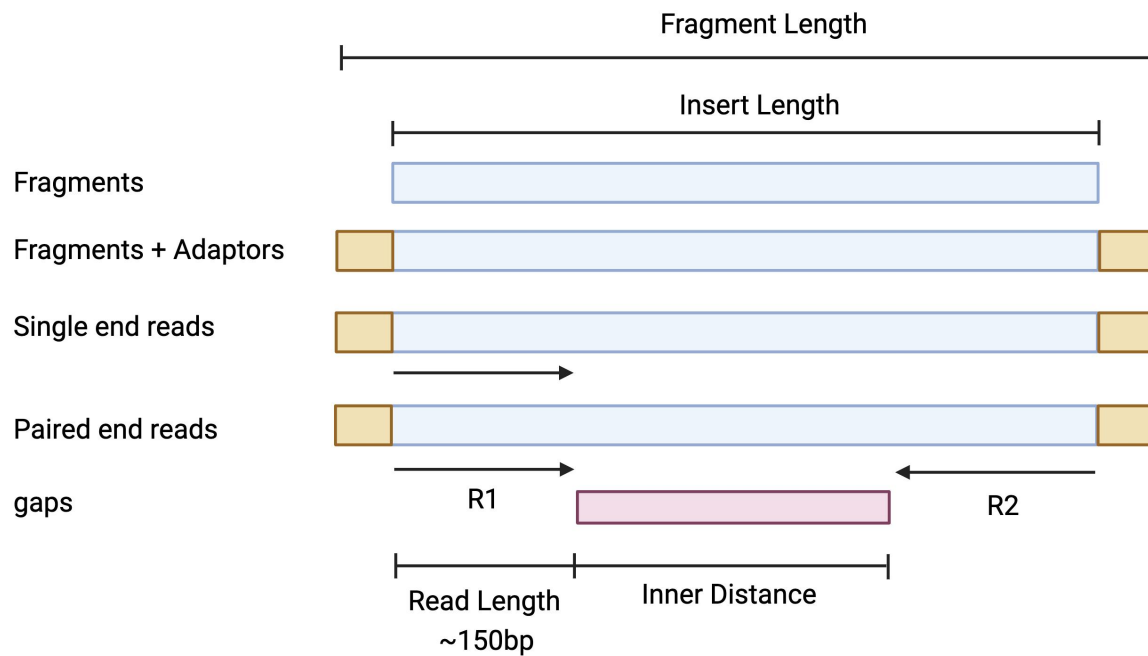
$$Q_{phred} = -10 \log_{10} p$$

Where p is the estimated probability of a base being wrong

Higher *Qphred*, better quality



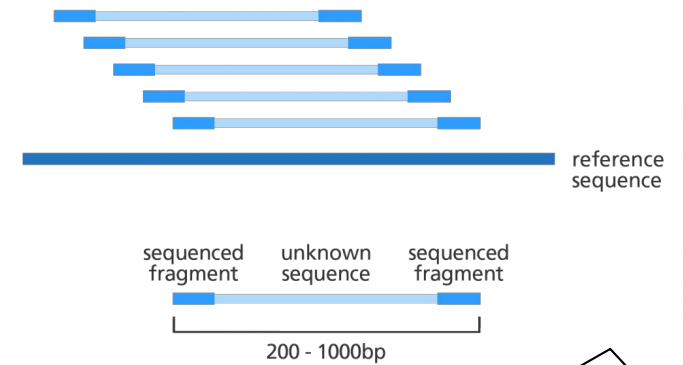
Single-end vs. Paired-end



Single-end reads



Paired-end reads

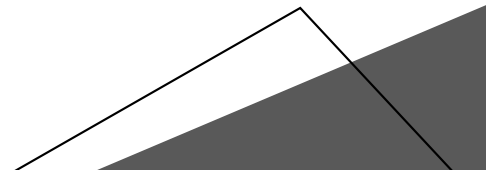




Single-end vs. Paired-end

Single-end: 1. Economic friendly
2. when the fragment is relatively short
e.g., small RNA sequencing, CHIP-seq

Paired-end: 1. Better mapping ability (especially in repetitive sequence)
2. Detect splice isoforms





Choose your aligners

1. What is the data source?
 - Do the reads contain splice junctions (i.e., RNA-seq)
 - Map the spliced reads against genome using STAR instead of Bowtie2
 - Map the non-spliced DNA-seq reads using Bowtie2
2. What is the output?
 - Typically SAM file. Make sure the out put file matches downstream analysis.
3. Trade-off between speed and completeness

