

Devin D  
Bi 624  
De-Duper Part1

PCR duplicates are two or more reads that come from the same original sequence.

AACCTTGG-original seq  
AACCTTGG –duplicate  
AACCTTGG-duplicate

Duplicates arise during the library amplification steps of preparation. If certain sequences happen to amplify really well you can get uneven amplification. Uneven amplification can lead to skewed sequence data. Removing PCR duplicates can be useful for reducing over-amplification data bias.

Algorithm design / pseudocode

- Write a function to check sequence file UMI's against all known UMI's
  - Remove any sequences with unknown UMI
  - Def UMICheck()
    - Compare sequence umi against list of known UMI's
    - return lines that pass UMI check to new .sam file
- Run UMI check function
- Write a function to search the cigar string for softclipping
  - Adjust the alignment start position if Softclipping occurs
  - Def SoftCheck()
    - Grab cigar string determine if softclipping ie(2S99M)
    - Adjust start position to reflect presence of soft clipping
      - If SPOS= 12 and 2S99M new SPOS = 10
    - Return adjusted start position
- Write a function that will compare duplicates for best possible alignment ie least softclipping
  - Def RetainRead()
    - If no soft clipping retain read
    - If soft clipping store clipping amount in variable for comparison to other reads with same start and soft clipping
    - Lowest amount is retained.
- Use Sam tools sort to align file by chromosome and alignment position
- Open the sorted file
  - for line in sortedFile
    - Split the line into parts
    - Run soft clipping function
    - Retain chromosome #, Alignment start Position and UMI as a tuple(1)
  - Read the next line
    - Split into parts
    - Run soft clipping function
    - Retain chromosome #, Alignment start Position and UMI as a tuple(2)
  - If tuple(1) == tuple(2)
    - Match = Duplicate found

[illegible]

### Example Sam Out

[illegible]