

# Picard\_MarkDup-2.7.0

The [DE Quick Start tutorial](#) provides an introduction to basic DE functionality and navigation.

Please work through the tutorial and add your comments on the bottom of this page. Or send comments per email to [xwang@cshl.edu](mailto:xwang@cshl.edu). Thank you.

## ***Rationale and background:***

Picard: <http://broadinstitute.github.io/picard>

Picard is a set of command line tools for manipulating high-throughput sequencing (HTS) data and formats such as SAM/BAM/CRAM and VCF.

The MarkDuplicates tool is to locate and tag duplicate reads in a BAM or SAM file, where duplicate reads are defined as originating from a single fragment of DNA. Duplicates can arise during sample preparation e.g. library construction using PCR. See also [EstimateLibraryComplexity](#) for additional notes on PCR duplication artifacts. Duplicate reads can also result from a single amplification cluster, incorrectly detected as multiple clusters by the optical sensor of the sequencing instrument.

The program can take either coordinate-sorted or query-sorted inputs, however the behavior is slightly different. When the input is coordinate-sorted, unmapped mates of mapped records and supplementary/secondary alignments are not marked as duplicates. However, when the input is query-sorted (actually query-grouped), then unmapped mates and secondary/supplementary reads are not excluded from the duplication test and can be marked as duplicate reads.

MarkDuplicates also produces a metrics file indicating the numbers of duplicates for both single- and paired-end reads.

If desired, duplicates can be removed using the REMOVE\_DUPLICATE and REMOVE\_SEQUENCING\_DUPLICATES options.

## **Pre-Requisites**

1. A CyVerse account. (Register for an CyVerse account here - [user.cyverse.org](http://user.cyverse.org))
2. Mandatory arguments
  - a. Input directory: Directory of SAM/BAM files to analyze. The alignment files must be coordinate sorted.
3. Optional arguments:
  - a. VALIDATION\_STRINGENCY: Validation stringency for all SAM files read by this program. Setting stringency to SILENT can improve performance when processing a BAM file in which variable-length data (read, qualities, tags) do not otherwise need to be decoded. Default value: LENIENT. This option can be set to 'null' to clear the default value. Possible values: {STRICT, LENIENT, SILENT}
  - b. REMOVE\_DUPLICATES: If true do not write duplicates to the output file instead of writing them with appropriate flags set. Default value: true.

## **Sample data**

The following test data are provided for testing BWA-index-mem here `/iplant/home/xiaofei_ipiant/Sorghum_chr8/chr8_test:`

1. BWA-index-mem\_0.7.10\_Apr10\_Test/G3\_P\_H3\_chr8\_BWA\_bam and BWA-index-mem\_0.7.10\_Apr10\_Test//G3\_P\_K4me3\_chr8\_BWA\_bam

Note: These are the outputs of BWA-index-mem\_0.7.10.

## **Results**

Successful execution of the Picard\_MarkDup\_2.7.0 will create 2 directories named out for BAM files and metrics.

## **Outputs**

1. BAM files
  - G3\_P\_H3\_chr8\_BWA\_bam\_rmDup\_BAM:
    - a. G3\_P\_H3\_rep1\_chr8\_rmDup.sorted.bam
    - b. G3\_P\_H3\_rep1\_chr8\_rmDup.sorted.bam.bai
    - c. G3\_P\_H3\_rep2\_chr8\_rmDup.sorted.bam
    - d. G3\_P\_H3\_rep2\_chr8\_rmDup.sorted.bam.bai
  - G3\_P\_K4me3\_chr8\_BWA\_bam\_rmDup\_BAM:
    - a. G3\_P\_K4me3\_rep1\_chr8\_rmDup.sorted.bam
    - b. G3\_P\_K4me3\_rep1\_chr8\_rmDup.sorted.bam.bai
    - c. G3\_P\_K4me3\_rep2\_chr8\_rmDup.sorted.bam

d. G3\_P\_K4me3\_rep2\_chr8\_rmDup.sorted.bam.bai

2. Metrics files

G3\_P\_H3\_chr8\_BWA\_bam\_rmDup\_metrics:

- a. G3\_P\_H3\_rep1\_chr8\_rmDup\_metrics.txt
- b. G3\_P\_H3\_rep2\_chr8\_rmDup\_metrics.txt

G3\_P\_K4me3\_chr8\_BWA\_bam\_rmDup\_metrics:

- a. G3\_P\_K4me3\_rep1\_chr8\_rmDup\_metrics.txt
- b. G3\_P\_K4me3\_rep2\_chr8\_rmDup\_metrics.txt