



**fq2bam (BWA-MEM + GATK)**

# Table of contents

What is BWA-MEM?

---

Why BWA-MEM?

---

What is fq2bam?

---

How should I use BWA-MEM in fq2bam?

---

Quick Start

---

Compatible CPU-based BWA-MEM, GATK4 Commands

---

fq2bam Reference

---

# List of Figures

Figure 0. Fq2bam

---

Generate BAM/CRAM output given one or more pairs of FASTQ files. Can also optionally generate a BQSR report.

### **Note**

fq2bam will become an alias for [fq2bamfast](#) in the next major release. All fq2bam arguments will continue to be supported.

## What is BWA-MEM?

BWA-MEM is a fast, accurate algorithm for mapping DNA sequence reads to a reference genome, performing local alignment and producing alignment for different parts of the query sequence. It is the default algorithm in Burrows-Wheeler Aligner (BWA) for reads that are longer than 70bp and is designed for high-throughput sequencing technologies such as Illumina and Pacific Biosciences.

## Why BWA-MEM?

BWA-MEM is capable of handling longer reads and is less sensitive to errors than other alignment algorithms. It is therefore used for a variety of applications, from routine analysis of sequencing data to more advanced applications such as de novo assembly and variant calling.

Some of the advantages of using BWA-MEM over similar tools include:

1. It is faster than many other alignment algorithms, making it the ideal choice for high-throughput sequencing.
2. It has a lower false positive rate than many other alignment algorithms, which means fewer false-positive variants are reported.
3. It is memory-efficient, allowing it to be used on limited resources.

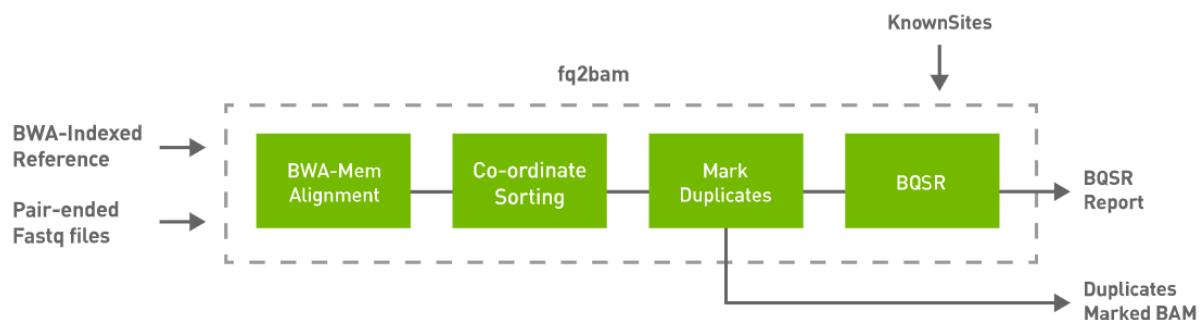
4. It is highly accurate, with a reported accuracy of over 99% on Illumina data.

## What is fq2bam?

BWA-MEM can be deployed within Parabricks, a software suite designed for accelerated secondary analysis in genomics, bringing industry standard tools and workflows from CPU to GPU, and delivering the same results at up to 60x faster runtimes. FQ2BAM is the Parabricks wrapper for BWA-MEM, which will sort the output and can mark duplicates and recalibrate base quality scores in line with GATK best practices. A 30x whole genome can be run through FQ2BAM in as little as 17 minutes on an NVIDIA DGX system, compared to 4-9 hours on a CPU instance (m5.24xlarge, 96 x vCPU).

## How should I use BWA-MEM in fq2bam?

fq2bam uses an accelerated version of BWA-MEM to generate BAM/CRAM output given one or more pairs of FASTQ files. The user can turn-off marking of duplicates by adding the `--no-markdups` option. The BQSR step is only performed if the `--knownSites input` and `--out-recal-file output` options are provided; doing so will also generate a BQSR report.



## Quick Start

```
# This command assumes all the inputs are in INPUT_DIR and all the outputs go to OUTPUT_DIR. docker run --rm --gpus all --volume INPUT_DIR:/workdir --volume OUTPUT_DIR:/outputdir \ --workdir /workdir \ nvc.io/nvidia/clara/clara-parabricks:4.3.1-1 \ pbrun fq2bam \ --ref /workdir/${REFERENCE_FILE} \ --in-fq
```

```
/workdir/${INPUT_FASTQ_1} /workdir/${INPUT_FASTQ_2} \ --knownSites  
/workdir/${KNOWN_SITES_FILE} \ --out-bam /outputdir/${OUTPUT_BAM} \ --out-  
recal-file /outputdir/${OUTPUT_RECAL_FILE}
```

## Compatible CPU-based BWA-MEM, GATK4 Commands

The commands below are the bwa-0.7.15 and GATK4 counterpart of the Parabricks command above. The output from these commands will be identical to the output from the above command. See the [Output Comparison](#) page for comparing the results.

```
# Run bwa-mem and pipe the output to create a sorted BAM. $ bwa mem \ -t 32 \ -K  
10000000 \ -R '@RG\tID:sample_rg1\tLB:lib1\tPL:bar\tSM:sample\tPU:sample_rg1' \  
<INPUT_DIR>/${REFERENCE_FILE} <INPUT_DIR>/${INPUT_FASTQ_1}  
<INPUT_DIR>/${INPUT_FASTQ_2} | \ gatk SortSam \ --java-options -Xmx30g \ --  
MAX_RECORDS_IN_RAM 5000000 \ -I /dev/stdin \ -O cpu.bam \ --SORT_ORDER  
coordinate # Mark duplicates. $ gatk MarkDuplicates \ --java-options -Xmx30g \ -I  
cpu.bam \ -O mark_dups_cpu.bam \ -M metrics.txt # Generate a BQSR report. $ gatk  
BaseRecalibrator \ --java-options -Xmx30g \ --input mark_dups_cpu.bam \ --output  
<OUTPUT_DIR>/${OUTPUT_RECAL_FILE} \ --known-sites  
<INPUT_DIR>/${KNOWN_SITES_FILE} \ --reference <INPUT_DIR>/${REFERENCE_FILE}
```

## fq2bam Reference

Run GPU-bwa mem, co-ordinate sorting, marking duplicates, and Base Quality Score Recalibration to convert FASTQ to BAM/CRAM.

### Input/Output file options

--ref REF

Path to the reference file. (default: None)

Option is required.

--in-fq [IN\_FQ ...]

Path to the pair-ended FASTQ files followed by optional read groups with quotes (Example: "@RG\tID:foo\tLB:lib1\tPL:bar\tSM:sample\tPU:foo"). The files must be in fastq or fastq.gz format. All sets of inputs should have a read group; otherwise, none should have a read group, and it will be automatically added by the pipeline. This option can be repeated multiple times. Example 1: --in-fq sampleX\_1\_1.fastq.gz sampleX\_1\_2.fastq.gz --in-fq sampleX\_2\_1.fastq.gz sampleX\_2\_2.fastq.gz. Example 2: --in-fq sampleX\_1\_1.fastq.gz sampleX\_1\_2.fastq.gz "@RG\tID:foo\tLB:lib1\tPL:bar\tSM:sample\tPU:unit1" --in-fq sampleX\_2\_1.fastq.gz sampleX\_2\_2.fastq.gz "@RG\tID:foo2\tLB:lib1\tPL:bar\tSM:sample\tPU:unit2". For the same sample, Read Groups should have the same sample name (SM) and a different ID and PU. (default: None)

--in-se-fq [IN\_SE\_FQ ...]

Path to the single-ended FASTQ file followed by optional read group with quotes (Example: "@RG\tID:foo\tLB:lib1\tPL:bar\tSM:sample\tPU:foo"). The file must be in fastq or fastq.gz format. Either all sets of inputs have a read group, or none should have one, and it will be automatically added by the pipeline. This option can be repeated multiple times. Example 1: --in-se-fq sampleX\_1.fastq.gz --in-se-fq sampleX\_2.fastq.gz . Example 2: --in-se-fq sampleX\_1.fastq.gz "@RG\tID:foo\tLB:lib1\tPL:bar\tSM:sample\tPU:unit1" --in-se-fq sampleX\_2.fastq.gz "@RG\tID:foo2\tLB:lib1\tPL:bar\tSM:sample\tPU:unit2" . For the same sample, Read Groups should have the same sample name (SM) and a different ID and PU. (default: None)

--in-fq-list IN\_FQ\_LIST

Path to a file that contains the locations of pair-ended FASTQ files. Each line must contain the location of two FASTQ files followed by a read group, each separated by a space. Each set of files (and associated read group) must be on a separate line. Files must be in fastq/fastq.gz format. Line syntax: <fastq\_1> <fastq\_2> <read group> (default: None)

--knownSites KNOWN\_SITES

Path to a known indels file. The file must be in vcf.gz format. This option can be used multiple times. (default: None)

--interval-file INTERVAL\_FILE

Path to an interval file in one of these formats: Picard-style (.interval\_list or .picard), GATK-style (.list or .intervals), or BED file (.bed). This option can be used multiple times. (default: None)

`--out-recal-file OUT_RECAL_FILE`

Path of a report file after Base Quality Score Recalibration. (default: None)

`--out-bam OUT_BAM`

Path of a BAM/CRAM file. (default: None)

Option is required.

`--out-duplicate-metrics OUT_DUPLICATE_METRICS`

Path of duplicate metrics file after Marking Duplicates. (default: None)

`--out-qc-metrics-dir OUT_QC_METRICS_DIR`

Path of the directory where QC metrics will be generated. (default: None)

## **Tool Options:**

`-L INTERVAL, --interval INTERVAL`

Interval within which to call bqsr from the input reads. All intervals will have a padding of 100 to get read records, and overlapping intervals will be combined. Interval files should be passed using the `--interval-file` option. This option can be used multiple times e.g. "`-L chr1 -L chr2:10000 -L chr3:20000+ -L chr4:10000-20000`". (default: None)

`--bwa-options BWA_OPTIONS`

Pass supported bwa mem options as one string. The current original bwa mem supported options are `-M, -Y` and `-T` e.g. `--bwa-options="-M -Y"` (default: None)

`--no-warnings`

Suppress warning messages about system thread and memory usage. (default: None)

`--filter-flag FILTER_FLAG`

Don't generate SAM entries in the output if the entry's flag's meet this criteria. Criteria: `(flag & filter != 0)` (default: 0)



--skip-multiple-hits

Filter SAM entries whose length of SA is not 0. (default: None)

--min-read-length MIN\_READ\_LENGTH

Skip reads below minimum read length. They will not be part of the output. (default: None)

--align-only

Generate output BAM after bwa-mem. The output will not be co-ordinate sorted or duplicates will not be marked. (default: None)

--no-markdups

Do not perform the Mark Duplicates step. Return BAM after sorting. (default: None)

--fix-mate

Add mate cigar (MC) and mate quality (MQ) tags to the output file. (default: None)

--markdups-assume-sortorder-queryname

Assume the reads are sorted by queryname for Marking Duplicates. This will mark secondary, supplementary, and unmapped reads as duplicates as well. This flag will not impact variant calling while increasing processing times. (default: None)

--markdups-picard-version-2182

Assume marking duplicates to be similar to Picard version 2.18.2. (default: None)

--monitor-usage

Monitor approximate CPU utilization and host memory usage during execution. (default: None)

--optical-duplicate-pixel-distance OPTICAL\_DUPLICATE\_PIXEL\_DISTANCE

The maximum offset between two duplicate clusters in order to consider them optical duplicates. Ignored if --out-duplicate-metrics is not passed. (default: None)

`--read-group-sm READ_GROUP_SM`

SM tag for read groups in this run. (default: None)

`--read-group-lb READ_GROUP_LB`

LB tag for read groups in this run. (default: None)

`--read-group-pl READ_GROUP_PL`

PL tag for read groups in this run. (default: None)

`--read-group-id-prefix READ_GROUP_ID_PREFIX`

Prefix for the ID and PU tags for read groups in this run. This prefix will be used for all pairs of fastq files in this run. The ID and PU tags will consist of this prefix and an identifier, that will be unique for a pair of fastq files. (default: None)

`-ip INTERVAL_PADDING, --interval-padding INTERVAL_PADDING`

Amount of padding (in base pairs) to add to each interval you are including. (default: None)

`--standalone-bqsr`

Run standalone BQSR. (default: None)

## **Performance Options:**

`--gpuwrite`

Use one GPU to accelerate writing final BAM. (default: None)

`--gpuwrite-deflate-algo GPUWRITE_DEFLATE_ALGO`

Choose the nvCOMP DEFLATE algorithm to use with `--gpuwrite`. Note these options do not correspond to CPU DEFLATE options. Valid options are 0 and 3. Option 0 is faster while option 3 provides a better compression ratio. (default=0) (default: None)

`--gpusort`

Use GPUs to accelerate sorting and marking. (default: None)

`--use-gds`

Use GPUDirect Storage (GDS) to enable a direct data path for direct memory access (DMA) transfers between GPU memory and storage. Must be used concurrently with `--gpuwrite`. Please refer to Parabricks Documentation > Best Performance for information on how to set up and use GPUDirect Storage. (default: None)

`--memory-limit MEMORY_LIMIT`

System memory limit in GBs during sorting and postsorting. By default, the limit is half of the total system memory. (default: 62)

`--low-memory`

Use low memory mode (default: None)

`--num-cpu-threads-per-stage NUM_CPU_THREADS_PER_STAGE`

Number of CPU threads to use per stage. (default: 8)

## **Common options:**

`--logfile LOGFILE`

Path to the log file. If not specified, messages will only be written to the standard error output. (default: None)

`--tmp-dir TMP_DIR`

Full path to the directory where temporary files will be stored.

`--with-petabase-dir WITH_PETAGENE_DIR`

Full path to the PetaGene installation directory. By default, this should have been installed at `/opt/petabase`. Use of this option also requires that the PetaLink library has been preloaded by setting the `LD_PRELOAD` environment variable. Optionally set the `PETASUITE_REFPATH` and `PGCLOUD_CREDPATH` environment variables that are used for data and credentials (default: None)

--keep-tmp

Do not delete the directory storing temporary files after completion.

--no-seccomp-override

Do not override seccomp options for docker (default: None).

--version

View compatible software versions.

### **GPU options:**

--num-gpus NUM\_GPUS

Number of GPUs to use for a run. GPUs 0..(NUM\_GPUS-1) will be used.

#### **Note**

The *--in-fq* option takes the names of two FASTQ files, optionally followed by a quoted read group. The FASTQ filenames must not start with a hyphen.

#### **Note**

When using the *--in-fq-list* option a read group is required on each line of the input file.

© Copyright 2024, Nvidia.. PDF Generated on 06/05/2024