



---

## Introduction

This document describes the command-line tools included with SMRT Link v4.0.0. These tools are for use by bioinformaticians working with secondary analysis results.

- The command-line tools are located in the `$SMRT_ROOT/smrtlink/smrtcmds/bin` subdirectory.

## Installation

The command-line tools are installed as an integral component of the SMRT Link software. For installation details, see **SMRT Link Software Installation (v4.0.0)**.

- To install **only** the command-line tools, use the `--smrttools-only` option with the installation command, whether for a new installation or an upgrade. Examples:

```
smrtlink-*.run --rootdir smrtlink --smrttools-only  
smrtlink-*.run --rootdir smrtlink --smrttools-only --upgrade
```

## Pacific Biosciences Command-Line Tools

Following is information on the Pacific Biosciences-supplied command-line tools included in the installation. Third-party tools installed are described at the end of the document.

**arrow** This is the `variantCaller` tool with the consensus algorithm set to `arrow`. See page 67 for details.

**bam2bam** The `bam2bam` tool reprocesses, and optionally converts, BAM files from one convention to another. For example, a BAM file containing HQ regions could be processed, adapter hits and barcodes identified, and a new subreads BAM file produced.

This tool is useful where `PostPrimary` on the instrument was used incorrectly; such as forgetting to request barcode analysis, and then reanalyzing the data with barcoding enabled.

- Both production and pulse BAM files can be processed.
- "Scraps" BAM files are always required to reconstitute the ZMW reads internally. Conversely, "scraps" BAM files will be output.
- ZMW reads are **not** allowed as input, due to the missing HQ-region annotations.

- Input read convention is determined from the `READTYPE` annotation in the `@RG: :DS` tags of the input BAM files.

`bam2bam` is installed on every Sequel™ System and is shipped with SMRT® Analysis.

## Usage

```
-o outputPrefix [options] input.(subreads|hqregion).bam
input.(scraps).bam
```

## Example

```
bam2bam in.subreads.bam in.scraps.bam -o out --barcodes bc.fasta
```

## Required Parameter

Required	Description
<code>-o STRING</code>	Prefix of the output file names.

## Optional Parameters

Options	Description
<code>-j INT</code>	Number of threads for parallel ZMW processing.
<code>-b INT</code>	Number of threads for parallel BAM compression.
<code>--silent</code>	Do <b>not</b> display progress output.

## BAM Conventions

Options	Description
<code>--zmw</code>	Create a ZMW read.
<code>--hqregion</code>	Output <code>*.hqregions.bam</code> and <code>*.scraps.bam</code> .

## Parameter for Finding Adapters

Option	Description
<code>--adapters=adapterSequences.fasta</code>	The file name of the adapter sequence(s). This signals that adapter-calling should be run.

## Parameters for Finding Barcodes

Options	Description
<code>--barcodes=barcodeSequences.fasta</code>	Specify a FASTA file of input barcode sequences to enable finding and labeling barcodes.
<code>--hotStartMode</code>	Enable searching for barcodes at the beginning of a read if no adapters are found. (Default = <code>False</code> )
<code>--scoreFirst</code>	Alternate name for <code>--hotStartMode</code> . (Default = <code>False</code> )
<code>--hotStartLength=INT</code>	Number of bases used to find the hot start barcode if the adapter is missing. (Default = <code>100</code> )

Options	Description
<code>--maxAdapters=INT</code>	Number of adapters used to determine the barcodes. (Default = 4)
<code>--scoreMode=STRING</code>	Barcode-calling mode. <i>symmetric</i> : Each barcode sequence identifies a single bin for demultiplexing reads. <i>asymmetric</i> : Barcode sequences are different on either end of an insert present in a SMRTbell™ template. (Default = <i>symmetric</i> )

### Parameters for Filtering Control Sequences

Options	Description
<code>--controls=controlSequences.fasta</code>	Enables control sequence-filtering.
<code>--maxControls=INT</code>	Number of subreads used to determine if a ZMW is a control. (Default = 3)

### Additional Output Read Types

Options	Description
<code>--fasta</code>	Output <code>fasta.gz</code> .
<code>--fastq</code>	Output <code>fastq.gz</code> .
<code>--noBam</code>	Do <b>not</b> produce BAM outputs.

### Parameters for Fine Tuning

Options	Description
<code>--minAdapterScore=int</code>	Minimum score for an adapter.
<code>--minPolyLength=INT</code>	Minimum ZMW real length. (Default = 1)
<code>--minSubLength=INT</code>	Minimum subread length. (Default = 1)
<code>--fullHQ</code>	Disable HQRf - the entire ZMW read will be deemed "HQ".

### Examples

To rerun with the barcoding option turned on because the initial run did not specify that the data was barcoded:

```
$ bam2bam --barcodes barcodes.fasta \
-o movieName.newBarcodes \
movieName.subreads.bam movieName.scrap.bam
```

To use a new adapter-finding algorithm on an older data set:

```
$ bam2bam --adapter adapters.fasta \
-o movieName.newAdapters \
movieName.subreads.bam movieName.scrap.bam
```

To perform adapter-finding and barcode-labeling from scratch:

```
$ bam2bam --barcodes barcodes.fasta \
--adapter adapters.fasta \
-o movieName.newAdapters \
movieName.subreads.bam movieName.scrap.bam
```

---

To convert subreads+scraps to ZMW reads:

```
$ bam2bam --zmw \
-o movieName.stitched \
movieName.subreads.bam movieName.scraps.bam
```

To output stitched ZMW reads additionally in FASTA.GZ format (a  
\*.zmw.BAM file is automatically created):

```
$ bam2bam --zmw \
--fasta \
-o movieName.stitched \
movieName.subreads.bam movieName.scraps.bam
```

To convert subreads+scraps to hqregions+scraps:

```
$ bam2bam --hqregion \
-o movieName.stitchedHQ \
movieName.subreads.bam movieName.scraps.bam
```

To output hqregions additionally in FASTA.GZ format (a  
\*.hqregions.BAM file is automatically created):

```
$ bam2bam --hqregion \
--fasta \
-o movieName.stitched \
movieName.subreads.bam movieName.scraps.bam
```

To output subreads in FASTA.GZ format **only**:

```
$ bam2bam --nobam \
--fasta \
-o movieName.new \
movieName.subreads.bam movieName.scraps.bam
```

To convert hqregions+scraps to subreads+scraps with adapter and  
barcodes:

```
$ bam2bam --barcodes barcodes.fasta \
--adapter adapters.fasta \
-o movieName.newVersion \
movieName.hqregions.bam movieName.scraps.bam
```

To perform a sanity check to ensure that the output is the same as the  
input, and add a new BAM header entry with the bam2bam version:

```
$ bam2bam -o movieName.sanity \
movieName.subreads.bam movieName.scraps.bam
$ samtools view movieName.subreads.bam > a1
$ samtools view movieName.sanity.subreads.bam > b1
$ diff a1 b1
$ samtools view movieName.scraps.bam > a1
$ samtools view movieName.sanity.scraps.bam > b1
$ diff a1 b1
$ rm a1 b1
```

---

To perform spike-in control filtering on a local computer because the filter controls were not specified on the instrument:

```
$ bam2bam --controls control_orig.fasta \
-o movieName.control_orig \
movieName.subreads.bam movieName.scrapshots.bam
```

To use a better reference for the spike-in controls:

```
$ bam2bam --controls control_better.fasta \
-o movieName.control_better \
movieName.subreads.bam movieName.scrapshots.bam
```

To perform a complete analysis from scratch, as the primary analysis software was released with a new set of improved algorithms: (**Note:** Only HQ regions **cannot** be computed from scratch)

```
$ bam2bam --barcodes barcodes.fasta \
--adapter adapters.fasta \
--controls control.fasta \
-o movieName.newPPAVersion \
movieName.subreads.bam movieName.scrapshots.bam
```

To treat the complete ZMW read as an HQ region and perform adapter-finding:

```
$ bam2bam --fullHQ \
--adapter adapters.fasta \
-o movieName.fullhq \
movieName.subreads.bam movieName.scrapshots.bam
```

## **bam2fasta/ bam2fastq**

The `BAM2fastx` tools convert PacBio® BAM files into gzipped FASTA and FASTQ files, including demultiplexing of barcoded data.

### **Usage**

Both tools have an identical interface and take BAM and/or Data Set files as input.

### **Examples**

```
bam2fasta -o projectName m54008_160330_053509.subreads.bam
```

```
bam2fastq -o myEcoliRuns m54008_160330_053509.subreads.bam  
m54008_160331_235636.subreads.bam
```

```
bam2fasta -o myHumanGenomem54012_160401_000001.subreadset.xml
```

### **Input Files**

- One or more \*.bam files
- \*.subreadset.xml file (Data Set file)

### **Output Files**

- \*.fasta.gz
- \*.fastq.gz

---

**bax2bam** The `bax2bam` tool converts the legacy PacBio basecall format (`bax.h5`) into the BAM basecall format.

### Usage

```
bax2bam [options] <input files...>
```

### Options

Options	Description
<code>-h, --help</code>	Display help information and exits.
<code>--version</code>	Displays program version number and exits

### Pulse feature options

These options configure pulse features in the output BAM. Supported features include:

Pulse Feature	BAM Tag	Default
DeletionQV	dq	Y
DeletionTag	dt	Y
InsertionQV	iq	Y
IPD	ip	Y
PulseWidth	pw	N
MergeQV	mq	Y
SubstitutionQV	sq	Y
SubstitutionTag	st	N

If the Pulse Feature option is used, then **only** those features listed will be included, regardless of the default state.

- `--pulsefeatures=STRING` (Comma-separated list of desired pulse features, using the names listed above.)
- `--losslessframes` (Store full, 16-bit IPD/PulseWidth data, instead of (default) downsampled, 8-bit encoding.)

### Input Files

- `movie.1.bax.h5, movie.2.bax.h5 ...` (**Note:** Input files should be from the same movie.)
- `--xml=STRING` (Data Set XML file containing a list of movie names.)
- `-f STRING, --fofn=STRING` (File-of-file-names containing a list of input files.)

### Output Files

- `-o STRING` (Prefix of output file names. Movie name will be used if no prefix is provided.)

- `--output-xml=STRING` (Explicit output XML name. If **not** provided using this option, `bax2bam` will use the `-o` prefix (`<prefix>.dataset.xml`). If that is not specified either, the output XML file name will be `<moviename>.dataset.xml`)
- Output read types: (**Note:** These are mutually exclusive.)
  - `--subread` Output subreads (default)
  - `--hqregion` Output HQ regions
  - `--polymeraseread` Output full polymerase read
  - `--ccs` Output CCS sequences
- Output BAM file type:
  - `--internal` Output BAMs in internal mode. Currently this indicates that non-sequencing ZMWs should be included in the output scraps BAM file, if applicable.

### Example

Assuming your original file is named `mydata.bas.h5`, you can produce a file `mynewbam.subreads.bam` using the following command:

```
bax2bam -o mynewbam mydata.1.bax.h5 mydata.2.bax.h5 mydata.3.bax.h5
```

**blasr** The `blasr` tool aligns long reads against a reference sequence, possibly a multi-contig reference.

`blasr` maps reads to genomes by finding the highest scoring local alignment or set of local alignments between the read and the genome. The initial set of candidate alignments is found by querying a rapidly-searched precomputed index of the reference genome, and then refining until only high scoring alignments are kept. The base assignment in alignments is optimized and scored using all available quality information, such as insertion and deletion quality values.

Because alignment approximates an exhaustive search, alignment significance is computed by comparing optimal alignment score to the distribution of all other significant alignment scores.

`blasr` also produces output in SAM and BAM format, but it **must** be built with `pbbam` libraries to generate both formats. See the Blasr wiki page for detailed installation instructions.

### Usage

```
blasr {subreads|ccs}.bam genome.fasta --bam --out aligned.bam [--options]
```

```
blasr {subreadset|consensusreadset}.xml genome.fasta --bam --out aligned.bam [--options]
```

```
blasr reads.fasta genome.fasta [--options]
```

## Input Files

- `{subreads|ccs}.bam` is in PacBio BAM format, which is the native Sequel output format of SMRT reads. PacBio BAM files carry rich quality information (such as insertion, deletion, and substitution quality values) needed for mapping, consensus calling and variant detection. For the PacBio BAM format specifications, see <http://pacbiofileformats.readthedocs.io/en/3.0/BAM.html>.
- `{subreadset|consensusreadset}.xml` is in PacBio DataSet format. For the PacBio DataSet format specifications, see: <http://pacbiofileformats.readthedocs.io/en/3.0/DataSet.html>.
- `reads.fasta` is a multi-FASTA file of reads. While any FASTA file is valid input, `bam` or `dataset` files are preferable as they contain more rich quality value information.
- `genome.fasta`: A FASTA file to which reads should map, usually containing reference sequences.

## Output Files

- `aligned.bam`: The pairwise alignments for each read, in PacBio BAM format.

## Input Options

Options	Description
<code>--sa suffixArrayFile</code>	Use the suffix array <code>sa</code> for detecting matches between the reads and the reference. (The suffix array is prepared by the <code>sawriter</code> program.)
<code>--ctab tab</code>	A table of tuple counts used to estimate match significance, created by <code>printTupleCountTable</code> . While it is quick to generate on the fly, if there are many invocations of <code>blasr</code> , it is useful to precompute the <code>ctab</code> .
<code>--regionTable table</code>	Read in a read-region table in HDF format for masking portions of reads. This may be a single table if there is just one input file, or a <code>fofn</code> (file-of-file names). When a region table is specified, any region table inside the <code>reads.plx.h5</code> or <code>reads.bax.h5</code> files is ignored. <b>Note:</b> This option works <b>only</b> with RS II HDF5 files.
<code>--noSplitSubreads</code>	Do not split subreads at adapters. This is typically only useful when the genome in an unrolled version of a known template, and contains template-adaptor-reverse-template sequences. (Default = <code>False</code> )

## Options for Aligning Output

Options	Description
<code>--bestn n</code>	Provide the top <code>n</code> alignments for the hit policy to select from. (Default = 10)
<code>--sam</code>	Write output in SAM format.
<code>--bam</code>	Write output in PacBio BAM format.
<code>--clipping</code>	Use <code>no/hard/soft</code> clipping for SAM output. (Default = <code>none</code> )
<code>--out file</code>	Write output to <code>file</code> . (Default = <code>terminal</code> )
<code>--unaligned file</code>	Output reads that are not aligned to <code>file</code> .



Options	Description
<code>--m t</code>	If <b>not</b> printing SAM, modify the output of the alignment. <ul style="list-style-type: none"> <li><code>t=0</code>: Print blast-like output with  's connecting matched nucleotides.</li> <li><code>1</code>: Print only a summary: Score and position.</li> <li><code>2</code>: Print in <code>Compare.xml</code> format.</li> <li><code>3</code>: Print in vulgar format (Deprecated).</li> <li><code>4</code>: Print a longer tabular version of the alignment.</li> <li><code>5</code>: Print in a machine-parsable format that is read by <code>compareSequences.py</code>.</li> </ul>
<code>--noSortRefinedAlignments</code>	Once candidate alignments are generated and scored via sparse dynamic programming, they are rescored using local alignment that accounts for different error profiles. Resorting based on the local alignment may change the order in which the hits are returned. (Default = <code>False</code> )
<code>--allowAdjacentIndels</code>	When specified, adjacent insertion or deletions are allowed. Otherwise, adjacent insertion and deletions are merged into one operation. Using quality values to guide pairwise alignments may dictate that the higher probability alignment contains adjacent insertions or deletions. Tools such as GATK do <b>not</b> permit this and so they are not reported by default.
<code>--header</code>	Print a header as the first line of the output file describing the contents of each column.
<code>--titleTable tab</code>	Build a table of reference sequence titles. The reference sequences are enumerated by row, <code>0, 1, ...</code> . The reference index is printed in alignment results rather than the full reference name. This makes output concise, particularly when very verbose titles exist in reference names. (Default = <code>NULL</code> )
<code>--minPctIdentity p</code>	Only report alignments if they are greater than <code>p</code> percent identity. (Default = <code>0</code> )
<code>--holeNumbers LIST</code>	When specified, only align reads whose ZMW hole numbers are in <code>LIST</code> . <code>LIST</code> is a comma-delimited string of ranges, such as <code>1, 2, 3, 10-13</code> . This option <b>only</b> works when reads are in base or pulse h5 format.
<code>--hitPolicy policy</code>	Specifies how <code>blasr</code> treats multiple hits: <ul style="list-style-type: none"> <li><code>all</code>: Reports <b>all</b> alignments.</li> <li><code>allbest</code>: Reports all equally top-scoring alignments.</li> <li><code>random</code>: Reports a single random alignment.</li> <li><code>randombest</code>: Reports a single random alignment from multiple equally top-scoring alignments.</li> <li><code>leftmost</code>: Reports an alignment which has the best alignment score and has the smallest mapping coordinates in any reference.</li> </ul>

### Options for Anchoring Alignment Regions

- These options will have the greatest effects on speed and sensitivity.

Options	Description
<code>--minMatch m</code>	Minimum seed length. A higher value will speed up alignment, but decrease sensitivity. (Default = <code>12</code> )
<code>--maxMatch m</code> <code>--maxLCPLength m</code>	Stop mapping a read to the genome when the LCP length reaches <code>m</code> . This is useful when the query is part of the reference, for example when constructing pairwise alignments for <i>de novo</i> assembly. (Both options work the same.)
<code>--maxAnchorsPerPosition m</code>	Do <b>not</b> add anchors from a position if it matches to more than <code>m</code> locations in the target.
<code>--advanceExactMatches E</code>	Another trick for speeding up alignments with <code>match -E</code> fewer anchors. Rather than finding anchors between the read and the genome at every position in the read, when an anchor is found at position <code>i</code> in a read of length <code>L</code> , the next position in a read to find an anchor is at <code>i+L-E</code> . Use this when aligning already assembled contigs. (Default = <code>0</code> )

Options	Description
<code>--nCandidates n</code>	Keep up to <i>n</i> candidates for the best alignment. A large value will slow mapping as the slower dynamic programming steps are applied to more clusters of anchors - this can be a rate-limiting step when reads are very long. (Default = 10)
<code>--concordant</code>	Map all subreads of a ZMW (hole) to where the longest full pass subread of the ZMW aligned to. This requires to use the region table and <i>hq</i> regions. This option <b>only</b> works when reads are in base or pulse <i>h5</i> format. (Default = <code>False</code> )

### Options for Refining Hits

Options	Description
<code>--sdpTupleSize K</code>	Use matches of length <i>K</i> to speed dynamic programming alignments. This controls accuracy of assigning gaps in pairwise alignments once a mapping has been found, rather than mapping sensitivity itself. (Default = 11)
<code>--scoreMatrix "score matrix string"</code>	Specify an alternative score matrix for scoring FASTA reads. The matrix is in the format ACGTN A abcde C fg hij G klmno T pqrst N uvwxy The values <i>a . . . y</i> should be input as a quoted space separated string: " <i>a b c . . . y</i> ". Lower scores are better, so matches should be less than mismatches e.g. <i>a,g,m,s</i> = -5 (match), mismatch = 6 .
<code>--affineOpen value</code>	Set the penalty for opening an affine alignment. (Default = 10)
<code>--affineExtend a</code>	Change affine (extension) gap penalty. Lower value allows more gaps. (Default = 0)

### Options for Overlap/Dynamic Programming Alignments and Pairwise Overlap for *de novo* Assembly

Options	Description
<code>--useQuality</code>	Use substitution/insertion/deletion/merge quality values to score gap and mismatch penalties in pairwise alignments. As the insertion and deletion rates are much higher than substitution, this will make many alignments favor an insertion/deletion over a substitution. Naive consensus-calling methods will then often miss substitution polymorphisms. Use this option when calling consensus using the Quiver method. <b>Note:</b> When <b>not</b> using quality values to score alignments, there will be a lower consensus accuracy in homopolymer regions. (Default = <code>False</code> )
<code>--affineAlign</code>	Refine alignment using affine guided align. (Default = <code>False</code> )

### Options for Filtering Reads

Options	Description
<code>--minReadLength l</code>	Skip reads that have a full length less than <i>l</i> . Subreads may be shorter. (Default = 50)
<code>--minSubreadLength l</code>	Do <b>not</b> align subreads of length less than <i>l</i> . (Default = 0)
<code>--maxScore m</code>	Maximum score to output; high is bad, negative is good. (Default = 0)

---

## Options for Parallel Alignment

Options	Description
<code>--nproc N</code>	Align using $N$ processes. All large data structures such as the suffix array and tuple count table are shared. (Default = 1)
<code>--start S</code>	Index of the first read to begin aligning. This is useful when multiple instances are running on the same data, for example when on a multi-rack cluster. (Default = 0)
<code>--stride S</code>	Align one read every $S$ reads. (Default = 1)

## Options for Subsampling Reads

Options	Description
<code>--subsample p</code>	Proportion $p$ of reads to randomly subsample and align; expressed as a decimal. (Default = 0)
<code>--help</code>	Displays help information and exits.
<code>--version</code>	Displays version information using the format <code>MajorVersion.Subversion.SHA1</code> (Example: <code>5.3.abcd123</code> ) and exits.

## Examples

To align reads from `reads.bam` to the `ecoli_K12` genome, and output in PacBio BAM format:

```
blasr reads.bam ecoli_K12.fasta --bam --out ecoli_aligned.bam
```

To use multiple threads:

```
blasr reads.bam ecoli_K12.fasta --bam --out ecoli_aligned.bam --proc 16
```

To include a larger minimal match, for faster but less sensitive alignments:

```
blasr reads.bam ecoli_K12.fasta --bam --out ecoli_aligned.bam --proc 16 --minMatch 15
```

To produce alignments in a pairwise human-readable format:

```
blasr reads.bam ecoli_K12.fasta -m 0
```

To use a precomputed suffix array for faster startup:

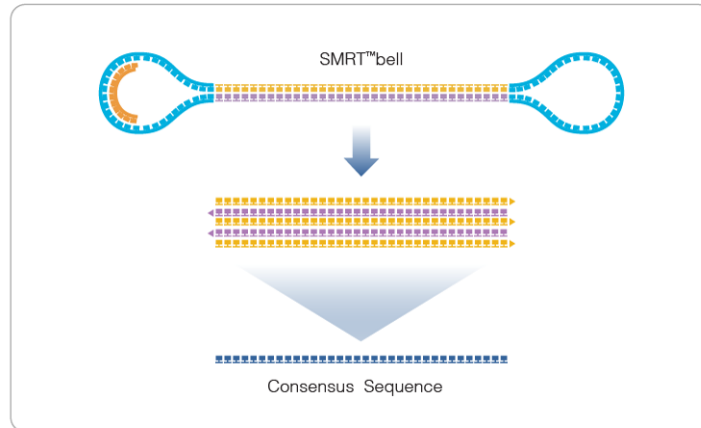
```
sawriter hg19.fasta.sa hg19.fasta #First precompute the suffix array
blasr reads.bam hg19.fasta --sa hg19.fasta.sa
```

To use a precomputed BWT-FM index for smaller runtime memory footprint, but slower alignments:

```
sa2bwt hg19.fasta hg19.fasta.sa hg19.fasta.bwt
blasr reads.bam hg19.fasta --bwt hg19.fasta.bwt
```

**ccs** Circular consensus sequencing (CCS) calculates consensus sequences from multiple “passes” around a circularized single DNA molecule (SMRTbell™ template). CCS uses the Quiver framework to achieve optimal consensus results given the number of passes available.

### Circular consensus sequencing (CCS)



### Input Files

- One `.subreads.bam` file containing the subreads for each SMRTbell™ template sequenced.

**Note:** Sequence data generated by the PacBio RS II is in `bas.h5` format, while the sequence data generated by the Sequel System is in BAM file format. If you have a `bas.h5` file, you will need to convert it into a BAM file. This can be done using the tool `bax2bam` which simply needs the name of any `bas.h5` files to convert and the prefix of the output file. See page 6 for details.

### Output Files

- A BAM file with one entry for each consensus sequence derived from a ZMW. BAM is a general file format for storing sequence data, which is described fully by the SAM/BAM working group. The CCS output format is a version of this general format, where the consensus sequence is represented by the "Query Sequence" and several tags have been added to provide additional meta information. An example BAM entry for a consensus as seen by `samtools` is shown below.

```
m141008_060349_42194_c100704972550000001823137703241586_s1_p0/63/ccs 4 * 0 255
* * 0 0 CCCGGGATCCTCTAGAATGC ~~~~~ RG:Z:83ba013f
np:i:35 rq:i:999 rs:B:i,37,0,0,1,0 sn:B:f,11.3175,6.64119,11.6261,14.5199
za:f:2.25461 zm:i:63 zs:B:f,-
1.57799,3.41424,2.96088,2.76274,3.65339,2.89414,2.446,3.04751,2.35529,3.65944,2.76774,4.119,1.679
81,1.66385,3.59421,2.32752,4.17803,-0.00353378,nan,0.531571,2.21918,3.88627,-
0.382997,0.650671,3.28113,0.798569,4.052,0.933297,3.00698,2.87132,2.66324,0.160431,1.99552,1.6935
4,1.90644,1.64448,3.13003,1.19977
```

Following are some of the common fields contained in the output BAM file:

Field	Description
Query Name	Movie Name / ZMW # /ccs
FLAG	Required by the format but meaningless in this context. Always set to 4 to indicate the read is unmapped.
Reference Name	Required by the format but meaningless in this context. Always set to *
Mapping Start	Required by the format but meaningless in this context. Always set to 0
Mapping Quality	Required by the format but meaningless in this context. Always set to 255
CIGAR	Required by the format but meaningless in this context. Always set to *
RNEXT	Required by the format but meaningless in this context. Always set to *
PNEXT	Required by the format but meaningless in this context. Always set to 0
TLEN	Required by the format but meaningless in this context. Always set to 0
Consensus Sequence	This is the consensus sequence generated.
Quality Values	The per-base parametric quality metric. For details see the "Interpreting QUAL Values" section.
RG Tag	The read group identifier.
bc Tag	A 2-entry array of upstream-provided barcode calls for this ZMW.
bq Tag	The quality of the barcode call. ( <b>Optional:</b> Depends on barcoded inputs.)
np Tag	The number of full passes that went into the subread. ( <b>Optional:</b> Depends on barcoded inputs.)
rq Tag	The predicted read quality.
rs Tag	An array of counts for the effect of adding each subread. The first element indicates the number of success and the remaining indicate the number of failures. This is a comma-separated list of the number of reads successfully added, failed to converge in likelihood, failed the Z filtering, failed to pass the pre-POA size filtering, or were excluded for another reason.
za Tag	The average Z-score for all reads successfully added.
zm Tag	The ZMW hole number.
zs Tag	This is a comma-separated list of the Z-scores for each subread when compared to the initial candidate template. A nan value indicates that the subread was <b>not</b> added.

## Usage

```
ccs [OPTIONS] OUTPUT FILES...
```

## Example

```
ccs --minLength=100 myCCS.bam myData.subreads.bam
```

Required	Description
Output File Name	The name of the output BAM file; comes after all other options listed. (Example = myResult.bam)
Input Files	The name of one or more subread.bam files to be processed. This comes at the end of the ccs command. (Example = myData.subreads.bam)

Options	Description
--version	Prints the version number.
--reportFile	Contains a result tally of the outcomes for all ZMWs that were processed. If no file name is given, the report is output to the file <code>ccs_report.txt</code> . In addition to the count of successfully-produced consensus sequences, this file lists how many ZMWs failed various data quality filters (SNR too low, not enough full passes, and so on) and is useful for diagnosing unexpected drops in yield.
--minSnr	Removes data that is likely to contain deletions. SNR is a measure of the strength of signal for all 4 channels (A, C, G, T) used to detect base pair incorporation. The SNR can vary depending on where in the ZMW a SMRTbell™ template stochastically lands when loading occurs. SMRTbell™ templates that land near the edge and away from the center of the ZMW have a less intense signal, and as a result can contain sequences with more "missed" base pairs. This value sets the threshold for minimum required SNR for any of the four channels. Data with SNR < 3.75 is typically considered lower quality. (Default = 3.75)
--minReadScore	The minimum value for the predicted quality of any subread used for CCS. Note that this filters the input to CCS (the subread quality must be above this value), whereas the <code>--minPredictedAccuracy</code> option filters the output (the predicted consensus sequence must be above a certain predicted accuracy). (Default = 0.75)
--minLength	The minimum length requirement for the median size of insert reads to generate a consensus sequence. If the targeted template is known to be a particular size range, this can filter out alternative DNA templates. (Default = 10)
--maxLength	The maximum length for both subreads that will be processed as well as the consensus sequence that will be generated. For robust results while avoiding unnecessary computation on unusual data, set to ~20% above the largest expected insert size. (Default = 7000)
--minPasses	The minimum number of passes for a ZMW to be emitted. This is the number of full passes. Full passes <b>must</b> have an adapter hit before and after the insert sequence and so do <b>not</b> include any partial passes at the start and end of the sequencing reaction. Additionally, the full pass count does <b>not</b> include any reads that were dropped by the Z-Filter. (Default = 3)
--minPredictedAccuracy	The minimum predicted accuracy of a read. CCS generates an accuracy prediction for each read, defined as the expected percentage of matches in an alignment of the consensus sequence to the true read. A value of 0.99 indicates that only reads expected to be 99% accurate are emitted. (Default = 0.9)
--minZScore	The minimum Z-Score for a subread to be included in the consensus generating process. For more information, see the "What are Z-Scores" section. (Default = 3.5)
--maxDropFraction	The maximum number of subreads that can be dropped before the entire ZMW is discarded. Subreads that appear very unlikely given the initial template (low Z-score), are discarded before generating the consensus sequence as part of an initial quality filter. Typically, very few reads should be discarded but if a high proportion are, then the entire ZMW is dropped. (Default = .34)
--zmws	If the consensus sequence for only a subset of ZMWs is required, they can be specified here. ZMWs can be specified either by range ( <code>--zmws=1-2000</code> ) by values ( <code>--zmws=5,10,20</code> ), or by both ( <code>--zmws=5-10,35,1000-2000</code> ). Use a comma-separated list with no spaces.
--numThreads	How many threads to use while processing. By default, CCS will use as many threads as there are available cores to minimize processing time, but fewer threads can be specified here.
--logFile	The name of a log file to use. If none is given, the logging information is printed to <code>STDERR</code> . (Example: <code>mylog.txt</code> )
--logLevel	Specifies verbosity of log data to produce. By setting <code>--logLevel=DEBUG</code> , you can obtain detailed information on what ZMWs were dropped during processing, as well as any errors which may have appeared. (Default = <code>INFO</code> )

Options	Description
<code>--noPolish</code>	After constructing the initial template, do <b>not</b> proceed with the polishing steps. This is significantly faster, but generates less accurate data with no RQ or QUAL values associated with each base.
<code>--byStrand</code>	Separately generate a consensus sequence from the forward and reverse strands. Useful for identifying heteroduplexes formed during sample preparation.
<code>--force</code>	Overwrite the output file when you don't care that it already exists.

## Interpreting QUAL Values

The QUAL value of a read is a measure of the posterior likelihood of an error at a particular position. Increasing QUAL values are associated with a decreasing probability of error. For indels and homopolymers, there is ambiguity as to which QUAL value is associated with the error probability. Shown below are different types of alignment errors, with a \* indicating which sequence BP should be associated with the alignment error.

### Mismatch

```

      *
ccs: ACGTATA
ref: ACATATA

```

### Deletion

```

      *
ccs: AC-TATA
ref: ACATATA

```

### Insertion

```

      *
ccs: ACGTATA
ref: AC-TATA

```

### Homopolymer Insertion or Deletion

Indels should always be left-aligned and the error probability is only given for the first base in a homopolymer.

```

      *                               *
ccs: ACGGGGTATA                     ccs: AC-GGGTATA
ref: AC-GGGTATA                       ref: ACGGGGTATA

```

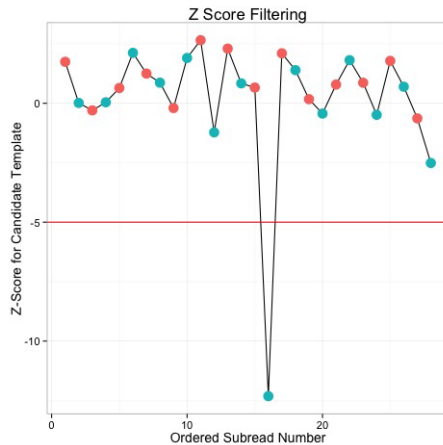
## What are Z-Scores?

Z-score filtering is a way to remove outliers and contaminating data from the CCS data prior to consensus generation, a crucial step for any analysis.

The Z-score for a subread is a metric which quantifies how it fits the model or assumptions of CCS scoring. In CCS, an initial template sequence is proposed, and then further refined using data in the templates. The initial template is usually quite close to the final consensus sequence, and at this

stage CCS will evaluate how likely each read is based on the candidate template. The likelihood of a read for a template is summarized by its Z-score, which asymptotically is normally distributed with a mean near 0.

Subreads with very low Z-scores are very **unlikely** to have been produced according to the CCS model, and so represent outliers. For example, the plot below shows the Z-scores for several subreads. With a -5 cutoff, we can see that one subread is excluded from the data.



### CCS Yield Report

The CCS Report specifies the number of ZMWs that successfully produced consensus sequences, as well as a count of how many ZMWs did not produce a consensus sequence for various reasons. The entries in this report, as well as parameters that can be set to change to increase or decrease the number of ZMWs that pass various filters, are explained in the table below.

ZMW Results	Parameters Affecting Results	Description
Below SNR threshold	<code>--minSnr</code>	ZMW had at least one channel's SNR below the minimum threshold.
No usable subreads	<code>--minReadScore,</code> <code>--minLength,</code> <code>--maxLength</code>	The ZMW had no usable subreads. Either there were no subreads, or all the subreads were below the minimum quality threshold or were above/below the specified length thresholds.
Insert size too long	<code>--maxLength</code>	The consensus sequence was above the maximum length threshold. Note that if all the input subreads were already below this threshold, they would <b>all</b> have been excluded, leading to a "No usable subreads" result.
Insert size too small	<code>--minLength</code>	The consensus sequence was below the minimum length threshold. Note that if all the input subreads were already below this threshold, they would <b>all</b> have been excluded, leading to a "No usable subreads" result.
Not enough full passes	<code>--minPasses</code>	There were not enough subreads that had an adapter at the start and end of the subread (a "full pass").



ZMW Results	Parameters Affecting Results	Description
Too many unusable subreads	<code>--minZScore,</code> <code>--maxDropFraction</code>	The ZMW had too many subreads that could not be used. A read can be unusable if it appears too unlikely given the initial template (low Z-score), or rarely, if a numerical rounding error occurs during processing.
CCS did not converge	None	The consensus sequence did not converge after the maximum number of allowed rounds of polishing.
CCS below minimum predicted accuracy	<code>--minPredictedAccuracy</code>	Each CCS read has a predicted level of accuracy associated with it, reads that are below the minimum specified threshold are removed.
Unknown error during processing	None	These should not occur.

**dataset** The `dataset` tool is a utility to create, open, manipulate and write Data Set XML files. The commands allow you to perform operations on the various types of data held by a Data Set XML: merge, split, write, and so on.

### Usage

```
dataset [-h] [--version] [--log-file LOG_FILE]
        [--log-level {DEBUG,INFO,WARNING,ERROR,CRITICAL}] | --debug | --quiet | -v]
        [--strict] [--skipCounts]
```

{create,filter,merge,split,validate,summarize,consolidate,loadstats,newuuid,loadmetadata,copyto,absolutize,relativize}

Options	Description
<code>-h, --help</code>	Displays help information and exits.
<code>&lt;Command&gt; -h</code>	Displays help for a specific command.
<code>-v, --version</code>	Displays program version number and exits.
<code>--log-file LOG_FILE</code>	Write the log to file. (Default = None, writes to stdout.)
<code>--log-level</code>	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL]. (Default = INFO)
<code>--debug</code>	Alias for setting the log level to DEBUG. (Default = False)
<code>--quiet</code>	Alias for setting the log level to CRITICAL to suppress output. (Default = False)
<code>-v</code>	Set the verbosity level. (Default = NONE)
<code>--strict</code>	Turn on strict tests and display all errors. (Default = False)
<code>--skipCounts</code>	Skip updating NumRecords and TotalLength counts. (Default = False)

**create Command:** Create an XML file from a `fofn` (file-of-file names) or BAM file. Possible types: `SubreadSet`, `AlignmentSet`, `ReferenceSet`, `HdfSubreadSet`, `BarcodeSet`, `ConsensusAlignmentSet`, `ConsensusReadSet`, `ContigSet`.

```
dataset create [-h] [--type DSTYPE] [--name DSNAME] [--generateIndices]
              [--metadata METADATA] [--novalidate] [--relative]
              outfile infile [infile ...]
```

Required	Description
outfile	The name of the XML file to create.
infile	The <code>f_ofn</code> (file-of-file-names) or BAM file(s) to convert into an XML file.

Options	Description
<code>--type DSTYPE</code>	Specify the type of XML file to create. (Default = NONE)
<code>--name DSNAME</code>	The name of the new Data Set XML file.
<code>--generateIndices</code>	Generate index files ( <code>.pbi</code> and <code>.bai</code> for BAM, <code>.fai</code> for FASTA). Requires <code>samtools/pysam</code> and <code>pindex</code> . (Default = FALSE)
<code>--metadata METADATA</code>	A <code>metadata.xml</code> file (or Data Set XML) to supply metadata. (Default = NONE)
<code>--novalidate</code>	Don't validate the resulting XML. Leaves the paths as they are.
<code>--relative</code>	Make the included paths relative instead of absolute. This is <b>not</b> compatible with <code>--novalidate</code> .

`filter` Command: Filter an XML file using filters and threshold values.

- Suggested filters: [`accuracy`, `bc`, `bcf`, `bcq`, `bcr`, `bq`, `cx`, `length`, `movie`, `n_subreads`, `pos`, `qend`, `qname`, `qstart`, `readstart`, `rname`, `rq`, `tend`, `tstart`, `zm`].
- More resource-intensive filter: [`qs`]

**Note:** Multiple filters with different names will be ANDed together. Multiple filters with the **same** name will be ORed together, duplicating existing requirements.

```
dataset filter [-h] infile outfile filters [filters ...]
```

Required	Description
infile	The name of the XML file to filter.
outfile	The name of the output filtered XML file.
filters	The values to filter on. (Example: <code>rq&gt;0.85</code> )

`merge` Command: Combine XML files.

```
dataset merge [-h] outfile infiles [infiles ...]
```

Required	Description
infiles	The names of the XML files to merge.
outfile	The name of the output XML file.

`split` Command: Split a Data Set XML file.

```
dataset split [-h] [--contigs] [--barcodes] [--zmws] [--byRefLength]
              [--noCounts] [--chunks CHUNKS] [--maxChunks MAXCHUNKS]
```

```

[--targetSize TARGETSIZE] [--breakContigs]
[--subdatasets] [--outdir
infile [outfiles...]
```

Required	Description
infile	The name of the XML file to split.

Options	Description
outfiles	The names of the resulting XML files.
--contigs	Split the XML file based on contigs. (Default = FALSE)
--barcodes	Split on barcodes. (Default = FALSE)
--zmws	Split on ZMWs. (Default = FALSE)
--byRefLength	Split contigs by contig length. (Default = TRUE)
--noCounts	Update Data Set counts after split. (Default = FALSE)
--chunks x	Split contigs into x total windows. (Default = 0)
--maxChunks x	Split contig list into at most x groups. (Default = 0)
--targetSize x	Target the minimum number of records per chunk. (Default = 5000)
--breakContigs	Break contigs to get closer to maxCounts. (Default = False)
--subdatasets	Split the XML file based on subdatasets. (Default = False)
--outdir OUTDIR	Specify an output directory for the resulting XML files. (Default = <in-place>, <b>not</b> the current working directory.)

`validate` Command: Validate XML and ResourceId files. (This is an internal testing functionality that may be useful.)

**Note:** This command requires that `pyxb` (**not** distributed with SMRT Link) be installed. If **not** installed, `validate` simply checks that the files pointed to in ResourceIds exist.

```
dataset validate [-h] [--skipFiles] infile
```

Required	Description
infile	The name of the XML file to validate.

Options	Description
--skipFiles	Skip validating external resources. (Default = False)

`summarize` Command: Summarize a Data Set XML file.

```
dataset summarize [-h] infile
```

Required	Description
infile	The name of the XML file to summarize.

---

`consolidate` Command: Consolidate XML files.

```
dataset consolidate [-h] [--numFiles NUMFILES] [--noTmp]
infile datafile xmlfile
```

Required	Description
<code>infile</code>	The name of the XML file to consolidate.
<code>datafile</code>	The name of the resulting data file.
<code>xmlfile</code>	The name of the resulting XML file.

Options	Description
<code>--numFiles x</code>	The number of data files to produce. (Default = 1)
<code>--noTmp</code>	Do <b>not</b> copy to a temporary location to ensure local disk use. (Default = <code>False</code> )

`loadstats` Command: Load a `sts.xml` file containing pipeline statistics into a Data Set XML file.

```
dataset loadstats [-h] [--outfile OUTFILE] infile statsfile
```

Required	Description
<code>infile</code>	The name of the Data Set XML file to modify.
<code>statsfile</code>	The name of the <code>.sts.xml</code> file to load.

Options	Description
<code>--outfile OUTFILE</code>	The name of the XML file to output. (Default = <code>None</code> )

`newuuid` Command: Refresh a Data Set's Unique ID.

```
dataset newuuid [-h] [--random] infile
```

Required	Description
<code>infile</code>	The name of the XML file to refresh.

Options	Description
<code>--random</code>	Generate a random UUID, instead of a hash. (Default = <code>False</code> )

`loadmetadata` Command: Load a `.metadata.xml` file into a Data Set XML file.

```
dataset loadmetadata [-h] [--outfile OUTFILE] infile metadata
```

Required	Description
infile	The name of the Data Set XML file to modify.
metadata	The .metadata.xml file to load, or Data Set to borrow from.

Options	Description
--outfile OUTFILE	The XML file to output. (Default = None)

`copyto` Command: Copy a Data Set and resources to a new location.

```
dataset copyto [-h] [--relative] infile outdir
```

Required	Description
infile	The name of the XML file to copy.
outdir	The directory to copy to.

Options	Description
--relative	Make the included paths relative instead of absolute. (Default = False)

`absolutize` Command: Make the paths in an XML file absolute.

```
dataset absolutize [-h] [--outdir OUTDIR] infile
```

Required	Description
infile	The name of the XML file whose paths should be absolute.

Options	Description
--outdir OUTDIR	Specify an optional output directory. (Default = None)

`relativize` Command: Make the paths in an XML file relative.

```
dataset relativize [-h] infile
```

Required	Description
infile	The name of the XML file whose paths should be relative.

### Example - Filter Reads

To filter one or more BAM file's worth of subreads, aligned or otherwise, and then place them into a single BAM file:

```
# usage: dataset filter <in_fn.xml> <out_fn.xml> <filters>
dataset filter in_fn.subreadset.xml filtered_fn.subreadset.xml 'rq>0.85'
```

---

```
# usage: dataset consolidate <in_fn.xml> <out_data_fn.bam> <out_fn.xml>
dataset consolidate filtered_fn.subreadset.xml consolidate.subreads.bam
out_fn.subreadset.xml
```

The filtered Data Set and the consolidated Data Set should be read-for-read equivalent when used with SMRT Analysis software.

### Example - Resequencing Pipeline

- Align two movie's worth of subreads in two SubreadSets to a reference.
  - Merge the subreads together.
  - Split the subreads into Data Set chunks by contig.
  - Process using `quiver` on a chunkwise basis (in parallel).
1. Align each movie to the reference, producing a Data Set with one BAM file for each execution:

```
pballign movie1.subreadset.xml referenceset.xml movie1.alignmentset.xml
pballign movie2.subreadset.xml referenceset.xml movie2.alignmentset.xml
```

2. Merge the files into a FOFN-like Data Set; BAMs are **not** touched:

```
# dataset merge <out_fn> <in_fn> [<in_fn> <in_fn> ...]
dataset merge merged.alignmentset.xml movie1.alignmentset.xml movie2.alignmentset.xml
```

3. Split the Data Set into chunks by contig name; BAMs are **not** touched:
  - Note that supplying output files splits the Data Set into that many output files (up to the number of contigs), with multiple contigs per file.
  - **Not** supplying output files splits the Data Set into **one** output file per contig, named automatically.
  - Specifying a number of chunks instead will produce that many files, with contig or even subcontig (reference window) splitting.

```
dataset split --contigs --chunks 8 merged.alignmentset.xml
```

4. Process the chunks using `Quiver`:

```
variantCaller.py --alignmentSetRefWindows --referenceFileName referenceset.xml --
outputFilename chunk1consensus.fasta --algorithm quiver chunk1contigs.alignmentset.xml
```

```
variantCaller.py --alignmentSetRefWindows --referenceFileName referenceset.xml --
outputFilename chunk2consensus.fasta --algorithm quiver chunk2contigs.alignmentset.xml
```

The chunking works by duplicating the original merged Data Set (no BAM duplication) and adding filters to each duplicate such that only reads belonging to the appropriate contigs are emitted. The contigs are distributed among the output files in such a way that the total number of records per chunk is about even.

## fasta-to-gmap-reference

The `fasta-to-gmap-reference` tool converts a reference FASTA file to a GMAP database (including the index files required by GMAP), and creates a GmapReferenceSet XML file. The GmapReferenceSet XML file can be imported into SMRT Link and used as a reference with the **Iso-Seq™ with Mapping** workflow. (For SMRT Link v4.0.0, this is a **mandatory** step for using the application, as SMRT Link cannot generate this Data Set type itself.)

### Usage

```
fasta-to-gmap-reference [options] fasta-file output-dir name
```

Required	Description
<code>fasta-file</code>	The path to the reference FASTA file.
<code>output-dir</code>	The location for the output GMAP database and Data Set XML file.
<code>name</code>	The name of the output GmapReferenceSet XML file.

Options	Description
<code>--organism &lt;value&gt;</code>	The name of the organism.
<code>--ploidy &lt;value&gt;</code>	Ploidy.
<code>--in-place</code>	Do <b>not</b> copy the input FASTA file to the output location.
<code>--log2stdout</code>	If true, log output will be displayed to the console. (Default = <code>False</code> )
<code>--log-level &lt;value&gt;</code>	Specify the log level; values are [ <code>ERROR</code> , <code>WARN</code> , <code>DEBUG</code> , <code>INFO</code> . ] (Default = <code>ERROR</code> )
<code>--debug</code>	Same as <code>--log-level DEBUG</code> .
<code>--quiet</code>	Same as <code>--log-level ERROR</code> .
<code>--verbose</code>	Same as <code>--log-level INFO</code> .
<code>--log-file &lt;value&gt;</code>	Log output file name. (Default = <code>"."</code> )
<code>--logback &lt;value&gt;</code>	Override all logger configuration with the specified <code>logback.xml</code> file.

## fasta-to-reference

The `fasta-to-reference` tool converts a FASTA file to a ReferenceSet Data Set XML that contains the required index files:

- `samtools` index (`.fai`)
- `sawriter` index (`fasta.sa`)
- SMRT View indexes (`fasta.config.index` and `fasta.index`)

`fasta-to-reference` is provided with SMRT Link, and requires the `samtools` and `sawriter` executables.

Note that `fasta-to-reference` will run on a single CPU on the host which it is executed, and not distributed on the cluster. For human-scale references, this may take up to half a day or more to run, and consumes a significant amount of memory. The indexing step with `sawriter` can use

---

over 34 GB of memory. When running this program, make sure the process has sufficient compute resources and will not be interrupted. We suggest redirecting STDERR/STDOUT to a log file. For example:

```
fasta-to-reference hg38.fasta /opt/smrtlink/references hg38 --organism Homo_sapiens >
fasta2ref.log 2>&1
```

## Usage

```
fasta-to-reference [options] fasta-file output-dir name
```

Required	Description
fasta-file	The path to the input FASTA file.
output-dir	The path to the output PacBio Reference Dataset XML.
name	The name of the ReferenceSet.

Options	Description
--organism <value>	The name of the organism.
--ploidy <value>	Ploidy.
-d, --debug	Emit logging to stdout.
-h, --help	Displays help information and exits.

## Input File

- \*.fasta file to convert.

## Output File

- \*.dataset.xml.

## ipdSummary

The `ipdSummary` tool detects DNA base-modifications from kinetic signatures. It is part of the `kineticsTool` package.

`kineticsTool` loads IPDs observed at each position in the genome, compares those IPDs to value expected for unmodified DNA, and outputs the result of this statistical test. The expected IPD value for unmodified DNA can come from either an in-silico control or an amplified control. The in-silico control is trained by Pacific Biosciences and shipped with the package. It predicts the IPD using the local sequence context around the current position. An amplified control Data Set is generated by sequencing unmodified DNA with the same sequence as the test sample. An amplified control sample is usually generated by whole-genome amplification of the original sample.

## Modification Detection

The basic mode of `kineticsTool` does an independent comparison of IPDs at each position on the genome, for each strand, and emits various statistics to CSV and GFF files (after applying a significance filter.)



---

## Modifications Identification

`kineticsTool` also has a Modification Identification mode that can decode multi-site IPD 'fingerprints' into a reduced set of calls of specific modifications. This feature has the following benefits:

- Different modifications occurring on the same base can be distinguished; for example, m5C and m4C.
- The signal from one modification is combined into one statistic, improving sensitivity, removing extra peaks, and correctly centering the call.

## Algorithm: Synthetic Control

Studies of the relationship between IPD and sequence context reveal that most of the variation in mean IPD across a genome can be predicted from a 12-base sequence context surrounding the active site of the DNA polymerase. The bounds of the relevant context window correspond to the window of DNA in contact with the polymerase, as seen in DNA/polymerase crystal structures. To simplify the process of finding DNA modifications with PacBio data, the tool includes a pre-trained lookup table mapping 12-mer DNA sequences to mean IPDs observed in C2 chemistry.

## Algorithm: Filtering and Trimming

`kineticsTool` uses the Mapping QV generated by BLASR and stored in the `cmp.h5` or BAM file (or AlignmentSet) to ignore reads that aren't confidently mapped. The default minimum Mapping QV required is 10, implying that BLASR has 90% confidence that the read is correctly mapped. Because of the range of read lengths inherent in PacBio data, this can be changed using the `--mapQvThreshold` command-line argument.

There are a few features of PacBio data that require special attention to achieve good modification detection performance. `kineticsTool` inspects the alignment between the observed bases and the reference sequence for an IPD measurement to be included in the analysis. The PacBio read sequence **must** match the reference sequence for  $k$  around the cognate base. In the current module,  $k=1$ . The IPD distribution at some locus can be thought of as a mixture between the 'normal' incorporation process IPD, which is sensitive to the local sequence context and DNA modifications, and a contaminating 'pause' process IPD, which has a much longer duration (mean > 10 times longer than normal), but happen rarely (~1% of IPDs).

**Note:** Our current understanding is that pauses do **not** carry useful information about the methylation state of the DNA, however a more careful analysis may be warranted. Also note that modifications that drastically increase the roughly 1% of observed IPDs are generated by pause events. Capping observed IPDs at the global 99<sup>th</sup> percentile is motivated by theory from robust hypothesis testing. Some sequence contexts may have naturally longer IPDs; to avoid capping too much data at those contexts, the cap threshold is adjusted per context as follows:

```
capThreshold = max(global99, 5*modelPrediction,
percentile(ipdObservations, 75))
```

### Algorithm: Statistical Testing

We test the hypothesis that IPDs observed at a particular locus in the sample have a longer means than IPDs observed at the same locus in unmodified DNA. If we have generated a Whole Genome Amplified Data Set, which removes DNA modifications, we use a case-control, two-sample t-test. This tool also provides a pre-calibrated 'synthetic control' model which predicts the unmodified IPD, given a 12-base sequence context. In the synthetic control case we use a one-sample t-test, with an adjustment to account for error in the synthetic control model.

### Usage

To run using a BAM input, and output GFF and HDF5 files:

```
ipdSummary aligned.bam --reference ref.fasta m6A,m4C --gff basemods.gff \
--csv_h5 kinetics.h5
```

To run using `cmp.h5` input, perform methyl fraction calculation, and output GFF and CSV files:

```
ipdSummary aligned.cmp.h5 --reference ref.fasta m6A,m4C --methylFraction \
--gff basemods.gff --csv kinetics.csv
```

Output Options	Description
<code>--gff FILENAME</code>	GFF format.
<code>--csv FILENAME</code>	Comma-separated value format.
<code>--csv_h5 FILENAME</code>	Compact binary-equivalent of .csv, in HDF5 format.
<code>--bigwig FILENAME</code>	BigWig file format; mostly only useful for SMRT View.

### Input Files

- A standard PacBio alignment file - either AlignmentSet XML, BAM, or `cmp.h5` - containing alignments and IPD information. The standard `cmp.h5` file of a SMRT Portal job is `data/aligned_read.cmp.h5`.
- Reference Sequence used to perform alignments. This can be either a FASTA file or a ReferenceSet XML.

### Output Files

The tool provides results in a variety of formats suitable for in-depth statistical analysis, quick reference, and consumption by visualization tools such as SMRT View. Results are generally indexed by reference position and reference strand. In all cases the strand value refers to the strand carrying the modification in the DNA sample. Remember that the kinetic effect of the modification is observed in read sequences aligning to the opposite strand. So reads aligning to the positive strand carry information about modification on the negative strand and vice versa, but the strand containing the putative modification is always reported.

- `modifications.gff`: Compliant with the GFF Version 3 specification (<http://www.sequenceontology.org/gff3.shtml>). Each template position / strand pair whose probability value exceeds the probability value threshold appears as a row. The template position is 1-based, per the GFF specifications. The strand column refers to the strand carrying the detected modification, which is the opposite strand from those used to detect the modification. The GFF confidence column is a Phred-transformed probability value of detection.

The auxiliary data column of the GFF file contains other statistics which may be useful downstream analysis or filtering. These include the coverage level of the reads used to make the call, and +/- 20 bp sequence context surrounding the site.

- `modifications.csv`: Contains one row for each (reference position, strand) pair that appeared in the Data Set with coverage at least `x`. `x` defaults to 3, but is configurable with the `--minCoverage` flag. The reference position index is 1-based for compatibility with the GFF file in the R environment. Note that this output type scales poorly and is **not** recommended for large genomes; the HDF5 output should perform much better in these cases.

### Output Columns: In-Silico Control Mode

Column	Description
<code>refId</code>	Reference sequence ID of this observation.
<code>tpl</code>	1-based template position.
<code>strand</code>	Native sample strand where kinetics were generated. 0 is the strand of the original FASTA, 1 is opposite strand from FASTA.
<code>base</code>	The cognate base at this position in the reference.
<code>score</code>	Phred-transformed probability value that a kinetic deviation exists at this position.
<code>tMean</code>	Capped mean of normalized IPDs observed at this position.
<code>tErr</code>	Capped standard error of normalized IPDs observed at this position ( $\text{standard deviation}/\sqrt{\text{coverage}}$ ).
<code>modelPrediction</code>	Normalized mean IPD predicted by the synthetic control model for this sequence context.
<code>ipdRatio</code>	$\text{tMean}/\text{modelPrediction}$ .
<code>coverage</code>	Count of valid IPDs at this position.
<code>frac</code>	Estimate of the fraction of molecules that carry the modification.
<code>fracLow</code>	2.5% confidence bound of the <code>frac</code> estimate.
<code>fracUpp</code>	97.5% confidence bound of the <code>frac</code> estimate.

### Output Columns: Case Control Mode

Column	Description
<code>refId</code>	Reference sequence ID of this observation.
<code>tpl</code>	1-based template position.

Column	Description
strand	Native sample strand where kinetics were generated. 0 is the strand of the original FASTA, 1 is opposite strand from FASTA.
base	The cognate base at this position in the reference.
score	Phred-transformed probability value that a kinetic deviation exists at this position.
caseMean	Mean of normalized case IPDs observed at this position.
controlMean	Mean of normalized control IPDs observed at this position.
caseStd	Standard deviation of case IPDs observed at this position.
controlStd	Standard deviation of control IPDs observed at this position.
ipdRatio	tMean/modelPrediction.
testStatistic	T-test statistic.
coverage	Mean of case and control coverage.
controlCoverage	Count of valid control IPDs at this position.
caseCoverage	Count of valid case IPDs at this position.

**laa** Long Amplicon Analysis (LAA) finds phased consensus sequences from a pooled set of (possibly polyploid) amplicons sequenced with Pacific Biosciences' SMRT technology. Sometimes referred to as **LAA2**, the executable `laa` is a complete rewrite of the `AmpliconAnalysis` module from the `ConsensusTools` package included with earlier versions of SMRT Analysis, which performed a similar function in the Quiver framework. This is a computational and memory-intensive software tool that builds upon the Arrow framework for generating high quality consensus sequences. It is generally preferable to run LAA using the SMRT Link interface for efficient distribution across a compute cluster. However, it is occasionally useful to run LAA from the command-line to identify optimal parameter settings or to diagnose a problem.

### Run Modes

`AmpliconAnalysis` is a general solution for the analysis of PCR products generated with SMRT sequencing, and it can be run in multiple configurations depending on the design of the experiment.

1. **Pooled Polyploid Amplicons:** The default mode assumes that the data contains a single complex mixture of amplicons, which may come from different genes and may have multiple alleles.
2. **Barcoded Polyploid Amplicons:** If passed a file of barcoding results, `AmpliconAnalysis` will instead separate the data by barcode and run the above process on each subset.
3. **Barcoded Simple Amplicons:** Another common use case is to generate consensus sequences for a large number of simple amplicons, such as for synthetic construct validation or high-throughput screening.

---

## Input Files

Long Amplicon Analysis **only** accepts PacBio-compatible BAM files or Data Set XML files as input.

- If your data was generated on a PacBio *RS* or PacBio RS II instrument, see page 6 for details on how to convert older data to the new file formats.

In addition, the underlying files themselves now contain barcode information, so if you have multiplexed data, see page 1 for details on how to barcode your data. This document assumes that you already have a barcoded PacBio BAM file containing the data to be analyzed.

## Output Files

LAA produces two sets of FASTQ files containing a sequence for each phased template sequence in each coarse cluster, and for each barcode.

- `amplicon_analysis.fastq`: Contains all of the high-quality non-artifactual sequences found.
- `amplicon_analysis_chimeras_noise.fastq`: Contains sequences thought to be some form of PCR or sequencing artifact.

**Note:** A sequence is defined as an artifact if, in the summary CSV file, the value of either the `IsDuplicate`, `NoiseSequence` or `IsChimera` columns are `True`.

- `amplicon_analysis_summary.csv`: Contains summary information about each read. Empty fields and values of `-1` represent inapplicable columns, while fields with `1` represent `True` and `0` represents `False`. Contains the following fields:
  - `BarcodeName`: Name of the barcode the reads came from. This is set to `0` for non-barcode runs.
  - `FastaName`: Sequence ID or header string.
  - `CoarseCluster`: Number of the coarse cluster the sequence came from.
  - `Phase`: Number of the phase of the sequence in the coarse cluster.
  - `TotalCoverage`: Total number of subreads mapped to this sequence. This may be capped using the `numPhasingReads` option.
  - `SequenceLength`: Length of this consensus sequence.
  - `ConsensusConverged`: `1` if a final consensus was reached within the allotted iterations, `0` if otherwise. `0` may indicate problems with the underlying sample or data.
  - `PredictedAccuracy`: Predicted accuracy of the consensus sequence, calculated by multiplying together the QVs generated by Arrow.
  - `NoiseSequence`: `1` if the sequence has a low-quality consensus, corresponding to a predicted accuracy less than 95% indicating a probable PCR artifact; `0` if otherwise.
  - `IsDuplicate`: `1` if the sequence is a duplicate of another with more coverage, otherwise `0`.
  - `DuplicateOf`: If `IsDuplicate` is `1`, contains the name of the other sequence, otherwise empty.

- IsChimera: 1 if the sequence is tagged as a chimeric by the UCHIME-like chimera labeler, 0 if otherwise.
- ChimeraScore: UCHIME-like score for sequences tested as possible chimeras.
- ParentSequenceA: If chimeric, the name of the consensus thought to be the source of the left half.
- ParentSequenceB: If chimeric, the name of the consensus thought to be the source of the right half.
- CrossoverPosition: Position in the chimeric sequence where the junction between the parent sequences is thought to have occurred.
- amplicon\_analysis\_subreads.X.csv: Contains mapping probabilities for each subread used to call of the consensus sequences generated. A **separate** file is written for **each** barcode pair, where x is replaced with the name of that pair. Contains the following fields:
  - SubreadId: The name of a particular subread used in the current run.
  - <A Consensus Sequence Name>: The mapping probability for the subread listed in SubreadId to the particular consensus sequence named.

## Usage

laa [options] INPUT

Options	Description
-h, --help	Displays help information and exits.
--verbose, -v	Set the verbosity level.
--version	Displays program version number and exits.
--log level	Sets the logging level. (Default = INFO)
--rngSeed	RNG seed, modulates reservoir filtering of reads. (Default = 42)
--generateBamIndex	Generate PacBio indicies (*.pbi) for BAM files that don't have them.
--ignoreBamIndex	Ignore PacBio indicies (*.pbi) for BAM files if they exist.
-M, --modelPath	Path to a model file or directory containing model files.
-m, --modelSpec	Name of chemistry or model to use, overriding default selection.
--numThreads, -n	Number of threads to use; 0 means autodetection. (Default = 0)
--takeN	Report only the top N consensus sequences for each barcode. To <b>disable</b> , use a number less than 1. (Default = 0)
-t, --trimEnds	Trim N bases from each end of each consensus. (Default = 0)
--minPredictedAccuracy	Minimum predicted consensus accuracy below which a consensus is treated as noise. (Default = 0.949999988079071)
--chimeraScoreThreshold	Minimum score to consider a sequence chimeric. (Default = 1)
--ChimeraFilter	Activate the chimera filter and separate chimeric consensus outputs.
--noChimeraFilter	Deactivate the chimera filter and output all consensus.
--logFile	Output file to write logging information to.
--resultFile	Output file name for high-quality results. (Default = amplicon_analysis.fastq)

Options	Description
--junkFile	Output file name for low-quality or chimeric results. (Default = amplicon_analysis_chimeras_noise.fastq)
--reportFile	Output file name for the summary report. (Default = amplicon_analysis_summary.csv)
--inputReportFile	Output file name for the output estimates of input PCR quality, based on subread mappings. (Default = amplicon_analysis_input.csv)
--subreadsReportPrefix	Prefix for the output subreads report. (Default = amplicon_analysis_subreads)
-b, --barcodes	FASTA file name of the barcode sequences used, which <b>overwrites</b> any barcode names in the Data Set. <b>Note:</b> This is used <b>only</b> to find barcode names.
--minBarcodeScore	Minimum average barcode score required for subreads. (Default = 0)
--fullLength	Filter input reads by presence of both flanking barcodes.
--doBc	A comma-separated list of barcode pairs to analyse. This can be by name ("lbc1--lbc1") or by Index ("0--0").
--ignoreBc	Disable barcode filtering so that all data be treated as one sample.
-l, --minLength	Minimum length of input reads to use. (Default = 3000)
-L, --maxLength	Maximum length of input reads to use. To <b>disable</b> , set to 0. (Default = 0)
-s, --minReadScore	Minimum read score of input reads to use. (Default = 0.75)
--minSnr	Minimum SNR of input reads to use. (Default = 3.75)
--whitelist	A file of ReadIds, in either Text or FASTA format, to allow from the input file. (Default = NONE)
-r, --maxReads	Maximum number of input reads, per barcode, to use in analysis. (Default = 2000)
-c, --maxClusteringReads	Maximum number of input reads to use in the initial clustering step. (Default = 500)
--fullLengthPreference	Prefer full-length subreads when clustering.
--fullLengthClustering	Only use full-length subreads when clustering.
--clusterInflation	Markov Clustering inflation parameter. (Default = 2)
--clusterLoopWeight	Markov Clustering loop weight parameter. (Default = 0.00100000004749745)
--skipRate	Skip some high-scoring alignments to disperse the cluster more. (Default = 0.0)
--minClusterSize	Minimum number of reads supporting a cluster before it is reported. (Default = 20)
--doCluster	Only analyze one specified cluster. (Default = -1)
--Clustering	Enable coarse clustering.
--noClustering	Disable coarse clustering.
-i, --ignoreEnds	When splitting, ignore N bases at the end. This prevents excessive splitting caused by degenerate primers. (Default = 0)
--maxPhasingReads	Maximum number of reads to use for phasing/consensus. (Default = 500)
--minQScore	Minimum score to require of mutations used for phasing. (Default = 20)

Options	Description
<code>--minPrevalence</code>	Minimum prevalence to require of mutations used for phasing. (Default = 0.0900000035762787)
<code>--minSplitReads</code>	Minimum number reads favoring the minor phase required to split a haplotype. (Default = 20)
<code>--minSplitFraction</code>	Minimum fraction of reads favoring the minor phase required to split a haplotype. (Default = 0.100000001490116)
<code>--minSplitScore</code>	The global likelihood improvement required to split a haplotype. (Default = 500)
<code>--minZScore</code>	Minimum ZScore to allow before adding a read to a haplotype. (Default = -10)
<code>--Phasing</code>	Enable the fine phasing step.
<code>--noPhasing</code>	Disable the fine phasing step.
<code>--emit-tool-contract</code>	Emit the tool contract.
<code>--resolved-tool-contract</code>	Use arguments from the resolved tool contract.

## Algorithm Description

LAA proceeds in six main phases: Data filtering, coarse clustering, waterfall clustering, fine phasing, consensus polishing, and post-processing.

- **Data filtering** is used to separate out sequences by their barcode calls, if present, so that only reads long enough to meaningfully contribute to phasing are used.
- The **Coarse and Waterfall Clustering** steps are used to find and separate reads coming from different amplicons.
- The reads from each cluster are then put through the **phasing** step, which recursively separates full-length haplotypes using a variant of the Arrow model. Those haplotypes are then **polished** within the Arrow framework to achieve a high-quality consensus sequence.
- Finally, a **post-processing** step attempts to identify and remove spurious consensus sequences and sequences representing PCR artifacts.

## Data Filtering

In this first step, we separate sequences by barcode and then apply a series of user-selected quality filters to speed up down-stream processing and improve result quality. Filters are used primarily to remove short subreads (which may not be long enough to phase variants of interest) and subreads with low barcode scores (representing reads for whom the barcode call is uncertain and may be incorrect). A “Whitelist” option is also available so that users can specify the exact list of subreads or ZMWs to be used.

## Coarse Clustering and Waterfall Clustering

The coarse clustering step groups the number of subreads (set by the `maxClusteringReads` option) that originate from different amplicons into



---

different clusters. It works by detecting subread-to-subread similarities, building a graph of the results, and then clustering nodes (subreads) using the Markov Clustering algorithm (<http://micans.org/mcl/>). The Markov clustering step is needed to remove spurious similarities caused by chimeric reads that can arise from PCR errors or doubly-loaded ZMWs, or just by chance due to sequencing error.

Next, if the number of subreads specified with the `maxReads` option is greater than the number used in coarse clustering, any remaining subreads are aligned to a rough consensus of each cluster and added to the cluster with the greatest similarity. This “Waterfall” step allows for a larger number of reads to be used much more quickly than if all subreads had to be clustered using the normal coarse clustering process.

At the end of clustering, subreads in each cluster are then sorted for downstream analysis using the PageRank algorithm (Page, Lawrence, et al. “The PageRank citation ranking: Bringing order to the web.” (1999)). This ensures that the most representative reads of the cluster are used first in the generation of consensus sequences.

### **Phasing/Consensus**

The reads assigned to each cluster are loaded into the Arrow framework, and an initial consensus of all reads is found. SNP differences between subreads and the initial consensus are scored with the Arrow model, and combinations of high-scoring SNPs are tested for their ability to segregate the reads into multiple haplotypes. If sufficient evidence of a second haplotype is found, the template sequence is ‘split’ into two copies, one with the SNPs applied to the template and one without. This process is repeated recursively so long as new haplotypes with sufficient scores can be found with at least some minimum level of coverage.

### **Post-Processing Filters**

LAA implements a post-processing step to flag likely PCR artifacts in the set of phased output sequences. First, consensus sequences that are identical duplicates of other consensus sequences in the results are removed. Next, those with unusually low predicted accuracy are flagged as being probable sequencing artifacts and removed. We implemented a filter for consensus sequences from PCR cross-over events, which on average make up ~5 to 20% of products generated by PCR amplifications >3 kb in length.

For artifacts of PCR cross-over events, or “chimeras”, we implemented a variant of the UCHIME algorithm (Edgar, Robert C., et al. “UCHIME improves sensitivity and speed of chimera detection.” *Bioinformatics* 27.16(2011): 2194-2200). The consensus sequences are sorted in order of decreasing read coverage, and the first two sequences are accepted as non-chimeric since they have no possible parent sequences with greater coverage. The remaining sequences are evaluated in descending order, with **each** test sequence aligned to all non-chimeric sequences so far

---

processed. Cross-overs between pairs of non-chimeric sequences are checked to see if they would yield a sequence very similar to the test sequence. If one is found with a sufficient score, the test sequence is marked as chimeric. If not, the test sequence is added to the list of non-chimeric sequences.

## motifMaker

The `MotifMaker` tool identifies motifs associated with DNA modifications in prokaryotic genomes. Modified DNA in prokaryotes commonly arises from restriction-modification systems that methylate a specific base in a specific sequence motif. The canonical example is the m6A methylation of adenine in GATC contexts in *E. coli*. Prokaryotes may have a very large number of active restriction-modification systems present, leading to a complicated mixture of sequence motifs.

PacBio® SMRT sequencing is sensitive to the presence of methylated DNA at single base resolution, via shifts in the polymerase kinetics observed in the real-time sequencing traces. For more background on modification detection, see <http://nar.oxfordjournals.org/content/early/2011/12/07/nar.gkr1146.full>.

## Algorithm

Existing motif-finding algorithms such as MEME-chip and YMF are sub-optimal for this case for the following reasons:

- They search for a **single** motif, rather than attempting to identify a complicated mixture of motifs.
- They generally don't accept the notion of aligned motifs - the input to the tools is a window into the reference sequence which can contain the motif at any offset, rather than a single center position that is available with kinetic modification detection.
- Implementations generally either use a Markov model of the reference (MEME-chip), or do exact counting on the reference, but place restrictions on the size and complexity of the motifs that can be discovered.

Following is a rough overview of the algorithm used by `MotifMaker`: Define a motif as a set of tuples: Position relative to methylation, required base. Positions not listed in the motif are implicitly degenerate. Given a list of modification detections and a genome sequence, we define the following objective function on motifs:

```
Motif score(motif) = (# of detections matching motif) / (# of genome sites matching motif) * (Sum of log-pvalue of detections matching motif) = (fraction methylated) * (sum of log-pvalues of matches)
```

We search (close to exhaustively) through the space of all possible motifs, progressively testing longer motifs using a branch-and-bound search. The 'fraction methylated' term must be less than 1, so the maximum achievable score of a child node is the sum of scores of modification hits in the current

---

node, allowing us to prune all search paths whose maximum achievable score is less than the best score discovered so far.

## Usage

For command-line motif-finding, run the `find` command, and pass the reference FASTA and the `modifications.gff` (.gz) file output by the PacBio modification detection workflow.

The `reprocess` subcommand annotates the GFF file with motif information for better genome browsing.

```
MotifMaker [options] [command] [command options]
```

`find` Command: Run motif-finding.

```
find [options]
```

Options	Description
<code>-h, --help</code>	Displays help information and exits.
* <code>-f, --fasta</code>	Reference FASTA file.
* <code>-g, --gff</code>	<code>Modifications.gff</code> or <code>.gff.gz</code> file.
<code>-m, --minScore</code>	Minimum Qmod score to use in motif finding. (Default = 40.0)
* <code>-o, --output</code>	Output <code>motifs.csv</code> file.
<code>-x, --xml</code>	Output motifs XML file.

`reprocess` Command: Update a `modifications.gff` file with motif information based on new Modification QV thresholds.

```
reprocess [options]
```

Options	Description
<code>-c, --csv</code>	Raw <code>modifications.csv</code> file.
* <code>-f, --fasta</code>	Reference FASTA file.
* <code>-g, --gff</code>	<code>Modifications.gff</code> or <code>.gff.gz</code> file.
<code>-m, --minFraction</code>	Only use motifs above this methylated fraction. (Default = 0.75)
<code>-m, --motifs</code>	<code>Motifs.csv</code> file.
* <code>-o, --output</code>	Reprocessed <code>modifications.gff</code> file.

## Output Files

Using the `find` command:

- **Output CSV file:** This file has the same format as the standard "Fields included in motif\_summary.csv" described in the Methylome Analysis White Paper (<https://github.com/PacificBiosciences/Bioinformatics-Training/wiki/Methylome-Analysis-Technical-Note>).

Using the `reprocess` command:

- **Output GFF file:** The format of the output file is the same as the input file, and is described in the Methylome Analysis White Paper (<https://github.com/PacificBiosciences/Bioinformatics-Training/wiki/Methylome-Analysis-White-Paper>) under "Fields included in the modifications.gff file".

**pballign** The `pballign` tool aligns PacBio reads to reference sequences, filters aligned reads according to user-specific filtering criteria, and converts the output to PacBio BAM, SAM, or PacBio DataSet format.

### Input Files

The `pballign` tool distinguishes input and output file formats by file extensions. The tool supports the following input formats:

- **BAM:** `.bam`
- **DataSet:** `.subreadset.xml` OR `.consensusreadset.xml`
- **FASTA:** `.fa` or `.fasta`
- **File-Of-File-Names:** `.fofn`

The input reference sequences can be in a FASTA file or a reference dataset created by `fasta-to-reference`, a PacBio tool for converting references in a FASTA file to PacBio reference dataset. See page 23 for details.

### Output Files

The tool supports the following output formats:

- **BAM:** `.bam`
- **DataSet:** `.xml`
- **SAM:** `.sam`

### Usage

```
pballign [-h] [--verbose] [--version] [--profile] [--debug]
          [--regionTable REGIONTABLE] [--configFile CONFIGFILE]
          [--algorithm {blasr,bowtie}] [--maxHits MAXHITS]
          [--minAnchorSize MINANCHORSIZE]
          [--maxMatch MAXMATCH]
          [--useccs {useccs,useccsall,useccsdenovo}]
          [--noSplitSubreads] [--nproc NPROC]
          [--algorithmOptions ALGORITHMOPTIONS]
          [--maxDivergence MAXDIVERGENCE] [--minAccuracy MINACCURACY]
          [--minLength MINLENGTH]
```

```

[--scoreFunction {alignerscore,editdist,blasrscore}]
[--scoreCutoff SCORECUTOFF]
[--hitPolicy {randidbest,allbest,random,all}] [--forQuiver]
[--seed SEED] [--tmpDir TMPDIR]
inputFileName referencePath outputFileName

```

Required	Description
inputFileName	The input file of PacBio reads. Can be a BAM, Data Set, FASTA file, or a fofn (File-Of-File-Names).
referencePath	Either a reference FASTA file or a PacBio reference dataset file.
outputFileName	The output .bam, .xml or .sam file.

Options	Description
-h, --help	Displays help information and exits.
--verbose, -v	Set the verbosity level.
--version	Displays program version number and exits.
--profile	Print runtime profile at exit.
--debug	Run within a debugger session.
--configFile	Specify a set of user-defined argument values.
--algorithm	Select an algorithm from blasr or bowtie. (Default = blasr)
--maxHits	The maximum number of matches of each read to the reference sequence that will be evaluated. (Default = 10)
--minAnchorSize	The length of the read that must match against the reference sequence. (Default = 12)
--maxMatch	Stop extending an anchor between the read and the reference sequence when its length reaches this value. Bypasses the blasr maxMatch option. (Default = 30)
--noSplitSubreads	Do <b>not</b> split reads into subreads even if subread regions are available. (Default = False)
--nproc NPROC	Number of threads. (Default = 8)
--algorithmOptions	Pass alignment options through.
--maxDivergence	The maximum allowed percentage divergence of a read from the reference sequence. (Default = 30)
--minAccuracy	The minimum percentage accuracy of alignments that will be evaluated. (Default = 70)
--minLength	The minimum aligned read length of alignments that will be evaluated. (Default = 50)
--scoreFunction	Specify a score function for evaluating alignments. <ul style="list-style-type: none"> <li>alignerscore: Aligner's score in the SAM tag as.</li> <li>editdist: Edit distance between read and reference.</li> <li>blasrscore: The blasr default score function.</li> </ul> (Default = alignerscore)
--scoreCutoff	The worst score to output an alignment.

Options	Description
<code>--hitPolicy</code>	Specify a policy for how to treat multiple hits. <ul style="list-style-type: none"> <li><code>random</code>: Selects a random hit.</li> <li><code>all</code>: Selects <b>all</b> hits.</li> <li><code>allbest</code>: Selects all the best score hits.</li> <li><code>randombest</code>: Selects a random hit from all best alignment score hits.</li> </ul> (Default = <code>randombest</code> )
<code>--seed</code>	Initialize the random number generator with a non-zero integer. Zero means that current system time is used. (Default = 1)
<code>--tmpDir</code>	Specify a directory for saving temporary files. (Default = <code>/scratch</code> )

## Examples

### Basic usage:

```
$ pbalign tests/data/example/read.bam \
  tests/data/example/ref.fasta \
  tests/data/example/example.bam
```

### Basic usage with optional arguments:

```
$ pbalign --maxHits 10 --hitPolicy all \
  tests/data/example_read.fasta \
  tests/data/example_ref.fasta \
  example.sam
```

### Advanced usage - To import predefined options from a configuration file:

```
$ pbalign --configFile=tests/data/1.config \
  tests/data/example/read.fasta \
  tests/data/example/ref.fasta \
  example.sam
```

### Advanced usage - To pass options through to the Aligner:

```
$ pbalign --algorithmOptions='-nCandidates 10 -sdpTupleSize 12'\
  tests/data/example/read.fasta \
  tests/data/example/ref.fasta \
  example.sam
```

### Advanced usage - To use pbalign as a library using the Python API:

```
$ python
>>> from pbalign.pbalignrunner import PBAAlignRunner
>>> # Specify arguments in a list.
>>> args = ['--maxHits', '20', 'tests/data/example/read.fasta', \
...       'tests/data/example/ref.fasta', 'example.sam']
>>> # Create a PBAAlignRunner object.
>>> a = PBAAlignRunner(args)
>>> # Execute.
>>> exitCode = a.start()
>>> # Show all files used.
>>> print a.fileNames
```

---

**pbdagcon** The `pbdagcon` tool implements DAGCon (Directed Acyclic Graph Consensus) which is a sequence consensus algorithm based on using directed acyclic graphs to encode multiple sequence alignments.

`pbdagcon` uses the alignment information from `blasr` to align sequence reads to a "backbone" sequence. Based on the underlying alignment directed acyclic graph (DAG), it will be able to use the new information from the reads to find the discrepancies between the reads and the "backbone" sequences. A dynamic programming process is then applied to the DAG to find the optimum sequence of bases as the consensus. The new consensus can be used as a new backbone sequence to iteratively improve the consensus quality.

While the code is developed for processing Pacific Biosciences raw sequence data, the algorithm can be used for general consensus purpose. Currently, it only takes FASTA input. For shorter read sequences, one might need to adjust the `blasr` alignment parameters to get the alignment string properly.

**Note:** This code is **not** an official Pacific Biosciences software release.

### Examples

To generate consensus from `blasr` alignments:

This is the most basic use case where one can generate a consensus from a set of alignments using the `pbdagcon` executable directly.

At the most basic level, `pbdagcon` takes information from `BLASR` alignments sorted by target and generates FASTA-formatted corrected target sequences. The alignments from `blasr` can be formatted with either `-m 4` or `-m 5`. For `-m 4` format, the alignments **must** be run through a format adapter, `m4topre.py`, to generate suitable input to `pbdagcon`.

The following example shows the simplest way to generate a consensus for one target using `blasr -m 5` alignments as input:

```
blasr queries.fasta target.fasta -bestn 1 -m 5 -out mapped.m5
pbdagcon mapped.m5 > consensus.fasta
```

To generate corrected reads from `daligner` alignments:

Support for generating consensus from `daligner` output was added in the form of a new executable `dazcon`. Note that `dazcon` is sensitive to the version of `daligner` used and may crash if using inputs generated by versions other than what is referenced in the submodules.

```
dazcon -ox -j 4 -s subreads.db -a subreads.las > corrected.fasta
```

To correct PacBio reads using HGAP:

Describes how to use `pbdagcon` to correct PacBio reads. This example demonstrates how correction is performed in PacBio's "Hierarchical Genome Assembly Process" (HGAP) workflow. HGAP uses `blasr -m 4` output.

This example makes use of the `filterm4.py` and `m4topre.py` scripts:

```
# First filter the m4 file to help remove chimeras:
filterm4.py mapped.m4 > mapped.m4.filt

# Next run the m4 adapter script, generating 'pre-alignments':
m4topre.py mapped.m4.filt mapped.m4.filt reads.fasta 24 > mapped.pre

# Finally, correct using pbdagcon, typically using multiple consensus threads:
pbdagcon -j 4 -a mapped.pre > corrected.fasta
```

**pbindex** The `pbindex` tool creates an index file that enables random-access to PacBio-specific data in BAM files.

### Usage

```
pbindex <input>
```

Options	Description
<code>-h, --help</code>	Displays help information and exits.
<code>--version</code>	Displays program version number and exits.

### Input File

- `*.bam` file containing PacBio data.

### Output File

- `*.pbi` index file, with the same prefix as the input file name.

**pbservice** The `pbservice` tool performs a variety of useful tasks within SMRT Link.

- To get help for `pbservice`, use `pbservice -h`.
- To get help for a specific `pbservice` command, use `pbservice <command> -h`.

`status` Command: Use to get system status.

```
pbservice status [-h] [--host HOST] [--port PORT]
                [--log-file LOG_FILE]
                [--log-level INFO]
                [--debug] [--quiet]
```

Options	Description
<code>--host=http://localhost</code>	The server host. Override the default with the environmental variable <code>PB_SERVICE_HOST</code> .



Options	Description
--port=8070	The server port. Override the default with the environmental variable PB_SERVICE_PORT.
--log-file LOG_FILE	Write the log to file. (Default = None, writes to stdout.)
--log-level=INFO	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = INFO)
--debug=False	Alias for setting the log level to DEBUG. (Default = False)
--quiet=False	Alias for setting the log level to CRITICAL to suppress output. (Default = False)

**import-dataset Command: Import Local Data Set XML.** The file location **must** be accessible from the host where the services are running; often on a shared file system

```
pbserve import-dataset [-h] [--host HOST] [--port PORT]
                    [--log-file LOG_FILE]
                    [--log-level INFO]
                    [--debug] [--quiet]
                    xml_or_dir
```

Required	Description
xml_or_dir	Specify a directory or XML file for the Data Set.

Options	Description
--host=http://localhost	The server host. Override the default with the environmental variable PB_SERVICE_HOST.
--port=8070	The server port. Override the default with the environmental variable PB_SERVICE_PORT.
--log-file LOG_FILE	Write the log to file. (Default = None, writes to stdout.)
--log-level=INFO	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = INFO)
--debug=False	Alias for setting the log level to DEBUG. (Default = False)
--quiet=False	Alias for setting the log level to CRITICAL to suppress output. (Default = False)

**import-fasta Command: Import a FASTA file and convert to a ReferenceSet file.** The file location **must** be accessible from the host where the services are running; often on a shared file system.

```
pbserve import-fasta [-h] --name NAME --organism ORGANISM --ploidy
                    PLOIDY [--block] [--host HOST] [--port PORT]
                    [--log-file LOG_FILE]
                    [--log-level INFO]
                    [--debug] [--quiet]
```

fasta\_path

Required	Description
fasta_path	Path to the FASTA file to import.

Options	Description
--name	Name of the ReferenceSet to convert the FASTA file to.
--organism	Name of the organism.
--ploidy	Ploidy.
--block=False	Block during importing process.
--host=http://localhost	The server host. Override the default with the environmental variable PB_SERVICE_HOST.
--port=8070	The server port. Override the default with the environmental variable PB_SERVICE_PORT.
--log-file LOG_FILE	Write the log to file. (Default = None, writes to stdout.)
--log-level=INFO	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = INFO)
--debug=False	Alias for setting the log level to DEBUG. (Default = False)
--quiet=False	Alias for setting the log level to CRITICAL to suppress output. (Default = False)

run-analysis Command: Run a secondary analysis pipeline using an analysis.json file.

```
pbserve run-analysis [-h] [--host HOST] [--port PORT]
                    [--log-file LOG_FILE]
                    [--log-level INFO]
                    [--debug] [--quiet] [--block]
                    json_path
```

Required	Description
json_path	Path to the analysis.json file.

Options	Description
--host=http://localhost	The server host. Override the default with the environmental variable PB_SERVICE_HOST.
--port=8070	The server port. Override the default with the environmental variable PB_SERVICE_PORT.
--log-file LOG_FILE	Write the log to file. (Default = None, writes to stdout.)
--log-level=INFO	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = INFO)
--debug=False	Alias for setting the log level to DEBUG. (Default = False)
--quiet=False	Alias for setting the log level to CRITICAL to suppress output. (Default = False)
--block=False	Block during importing process.

Options	Description
json_path	Path to the analysis.json file. <b>(Required)</b> .

`emit-analysis-template` **Command:** Output an analysis.json template to stdout that can be run using the `run-analysis` command.

```
pbservice emit-analysis-template [-h] [--log-file LOG_FILE]
                                [--log-level INFO]
                                [--debug] [--quiet]
```

Options	Description
--log-file LOG_FILE	Write the log to file. (Default = None, writes to stdout.)
--log-level=INFO	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = INFO)
--debug=False	Alias for setting the log level to DEBUG. (Default = False)
--quiet=False	Alias for setting the log level to CRITICAL to suppress output. (Default = False)

`get-job` **Command:** Get a Job Summary by Job Id.

```
pbservice get-job [-h] [--host HOST] [--port PORT]
                 [--log-file LOG_FILE]
                 [--log-level INFO]
                 [--debug] [--quiet]
                 job_id
```

Required	Description
job_id	Job id or UUID.

Options	Description
--host=http://localhost	The server host. Override the default with the environmental variable PB_SERVICE_HOST.
--port=8070	The server port. Override the default with the environmental variable PB_SERVICE_PORT.
--log-file LOG_FILE	Write the log to file. (Default = None, writes to stdout.)
--log-level=INFO	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = INFO)
--debug=False	Alias for setting the log level to DEBUG. (Default = False)
--quiet=False	Alias for setting the log level to CRITICAL to suppress output. (Default = False)

`get-jobs` **Command:** Get Job Summaries by Job Id.

```
pbservice get-jobs [-h] [-m MAX_ITEMS] [--host HOST] [--port PORT]
                  [--log-file LOG_FILE]
                  [--log-level INFO]
                  [--debug] [--quiet]
```

Options	Description
-m=25, --max-items=25	Maximum number of jobs to get.
--host=http://localhost	The server host. Override the default with the environmental variable PB_SERVICE_HOST.
--port=8070	The server port. Override the default with the environmental variable PB_SERVICE_PORT.
--log-file LOG_FILE	Write the log to file. (Default = None, writes to stdout.)
--log-level=INFO	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = INFO)
--debug=False	Alias for setting the log level to DEBUG. (Default = False)
--quiet=False	Alias for setting the log level to CRITICAL to suppress output. (Default = False)

**get-dataset Command:** Get a Data Set summary by Data Set Id or UUID.

```
pbservice get-dataset [-h] [--host HOST] [--port PORT]
                    [--log-file LOG_FILE]
                    [--log-level INFO]
                    [--debug] [--quiet]
                    id_or_uuid
```

Required	Description
id_or_uuid	Data Set Id or UUID.

Options	Description
--host=http://localhost	The server host. Override the default with the environmental variable PB_SERVICE_HOST.
--port=8070	The server port. Override the default with the environmental variable PB_SERVICE_PORT.
--log-file LOG_FILE	Write the log to file. (Default = None, writes to stdout.)
--log-level=INFO	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = INFO)
--debug=False	Alias for setting the log level to DEBUG. (Default = False)
--quiet=False	Alias for setting the log level to CRITICAL to suppress output. (Default = False)

**get-datasets Command:** Get a Data Set list summary by Data Set type.

```
pbservice get-datasets [-h] [--host HOST] [--port PORT]
                    [--log-file LOG_FILE]
                    [--log-level INFO]
                    [--debug] [--quiet] [-m MAX_ITEMS]
                    [-t DATASET_TYPE]
```

Required	Description
-t=subreads, --dataset-type=subreads	The type of Data Set to retrieve: subreads, alignments, references, barcodes.

Options	Description
<code>--host=http://localhost</code>	The server host. Override the default with the environmental variable <code>PB_SERVICE_HOST</code> .
<code>--port=8070</code>	The server port. Override the default with the environmental variable <code>PB_SERVICE_PORT</code> .
<code>--log-file LOG_FILE</code>	Write the log to file. (Default = None, writes to stdout.)
<code>--log-level=INFO</code>	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = INFO)
<code>--debug=False</code>	Alias for setting the log level to DEBUG. (Default = False)
<code>--quiet=False</code>	Alias for setting the log level to CRITICAL to suppress output. (Default = False)
<code>-m=25, --max-items=25</code>	Maximum number of jobs to get.

## Examples

To obtain system status, the Data Set summary, and the job summary:

```
pbservice status --host smrtlink-release --port 9091
```

To import a Data Set XML:

```
pbservice import-dataset --host smrtlink-release --port 9091 \
path/to/subreadset.xml
```

To obtain a job summary using the Job Id:

```
pbservice get-job --host smrtlink-release --port 9091 \
--log-level CRITICAL 1
```

To obtain Data Sets by using the Data Set Type subreads:

```
pbservice get-datasets --host smrtlink-alpha --port 8081 \
--quiet --max-items 1 -t subreads
```

To obtain Data Sets by using the Data Set Type alignments:

```
pbservice get-datasets --host smrtlink-alpha --port 8081 \
--quiet --max-items 1 -t alignments
```

To obtain Data Sets by using the Data Set Type references:

```
pbservice get-datasets --host smrtlink-alpha --port 8081 \
--quiet --max-items 1 -t references
```

To obtain Data Sets by using the Data Set Type barcodes:

```
pbservice get-datasets --host smrtlink-alpha --port 8081 \
--quiet --max-items 1 -t barcodes
```

To obtain Data Sets by using the Data Set UUID:

```
pbservice get-dataset --host smrtlink-alpha --port 8081 \
--quiet 43156b3a-3974-4ddb-2548-bb0ec95270ee
```

---

**pbsmrtpipe** The `pbsmrtpipe` tool is the secondary analysis workflow engine of Pacific Biosciences' SMRT Analysis software. `pbsmrtpipe` is easily extensible, and supports logging, distributed computing, error handling, analysis parameters, and temporary files.

In a typical installation of the SMRT Analysis Software, SMRT Link's SMRT Analysis module calls `pbsmrtpipe` when an analysis is started. SMRT Link's SMRT Analysis module provides a convenient and user-friendly way to analyze Pacific Biosciences' sequencing data through `pbsmrtpipe`.

For power users, there is more flexibility and customization available by instead running `pbsmrtpipe` analyses from the command line.

The `pbsmrtpipe` command is normally run in one of several modes, which are specified as a positional argument.

For details about a specific pipeline, specify the ID (the last field in each item in the output of `show-templates`) using the `show-template-details` command:

```
$ pbsmrtpipe show-template-details pbsmrtpipe.pipelines.sa3_ds_resequencing
```

Note that if you are starting from PacBio's `bax.h5` basecalling files, you will need to do an initial conversion step.

## Pipelines

Following are the available pipelines, their purpose, and their outputs.

**Note:** All pipeline names are prefixed with `pbsmrtpipe.pipelines`; this is omitted from the table.

Pipeline Name	Description/Common Outputs
<code>sa3_sat</code>	<ul style="list-style-type: none"><li>Site Acceptance Test run on all new PacBio installations.</li><li>variants GFF, SAT report.</li></ul>
<code>sa3_ds_resequencing</code>	<ul style="list-style-type: none"><li>Map subreads to reference genome and determine consensus sequence with Arrow.</li><li>AlignmentSet, consensus ContigSet, variants GFF.</li></ul>
<code>sa3_ds_ccs</code>	<ul style="list-style-type: none"><li>Generate high-accuracy Circular Consensus Sequences from subreads.</li><li>ConsensusReadSet, FASTA and FastQ files.</li></ul>
<code>sa3_ds_ccs_mapping</code>	<ul style="list-style-type: none"><li>ConsensusRead (CCS) + mapping to reference genome, starting from subreads.</li><li>ConsensusReadSet, FASTA and FastQ files, ConsensusAlignmentSet.</li></ul>
<code>sa3_ds_isoseq_classify</code>	<ul style="list-style-type: none"><li>IsoSeq transcript classification, starting from subreads.</li><li>ContigSets of classified transcripts.</li></ul>
<code>sa3_ds_isoseq</code>	<ul style="list-style-type: none"><li>Full Iso-Seq analysis with clustering and Quiver polishing. (This is much slower.)</li><li>ContigSets of classified transcripts plus polished isoform ContigSet.</li></ul>

Pipeline Name	Description/Common Outputs
ds_modification_motif_analysis	<ul style="list-style-type: none"> <li>Base modification detection and motif-finding, starting from subreads.</li> <li>Resequencing output plus basemods GFF, motifs CSV.</li> </ul>
sa3_hdfsubread_to_subread	<ul style="list-style-type: none"> <li>Convert HdfSubreadSet to SubreadSet (import bax.h5 basecalling files).</li> <li>SubreadSet</li> </ul>
sa3_ds_laa	<ul style="list-style-type: none"> <li>Basic Long Amplicon Analysis (LAA) pipeline, from barcoded subreads.</li> <li>Consensus FASTA</li> </ul>
polished_falcon_fat	<ul style="list-style-type: none"> <li>HGAP 4 assembly pipeline starting from subreads and a configuration file.</li> <li>ContigSet of assembled contigs</li> </ul>

## Parallelization

The algorithms used to analyze PacBio data are computationally intensive but also intrinsically highly parallel. `pbsmrtpipe` can scale to at least hundreds of processors on multi-core systems and/or managed clusters. This is handled by two distinct but complementary methods:

- **Multiprocessing** is implemented in the underlying tasks, all of which are generally shared-memory programs. This is effectively always turned on unless the `max_nchunk` parameter is set to 1. (See the Examples section for a description of how to modify parameter values.) For most compute node configurations, a value between 8 and 16 is appropriate.
- **Parallelization (chunking)** is implemented by `pbsmrtpipe` and works by applying filters to the input Data Sets, which direct tasks to operate on a subset (“chunk”) of the data. These chunks are most commonly either a contiguous subset of reads or windows in the reference genome sequence.

Note that at present, the task-level output directories (and the locations of the final result files) may be slightly different depending on whether chunking is used, since an intermediate “gather” step is required to join chunked results.

## Usage

```
pbsmrtpipe [-h] [--version]
{pipeline, pipeline-id, task, show-templates, show-template-details, show-tasks, show-
task-details, show-workflow-options, run-diagnostic, show-chunk-operators}
...
```

Options	Description
<code>--help</code>	Displays information about command-line usage and options, and then exits. <code>pipeline-id --help</code> : Displays information about a specific pipeline.
<code>--version</code>	Displays program version number and exits.

`pipeline` Command: Run a pipeline using a pipeline template or with explicit Bindings and EntryPoints.

```

pbsmrtpipe pipeline [-h] [--debug] -e ENTRY_POINTS [ENTRY_POINTS ...]
                    [-o OUTPUT_DIR] [--preset-xml PRESET_XML]
                    [--preset-json PRESET_JSON]
                    [--preset-rc-xml PRESET_RC_XML]
                    [--service-uri SERVICE_URI]
                    [--force-distributed | --local-only]
                    [--force-chunk-mode | --disable-chunk-mode]
                    pipeline_template_xml

```

Required	Description
pipeline_template_xml	Path to a pipeline template XML file.

Options	Description
--debug=False	Alias for setting log level to DEBUG.
-e, --entry	Entry Points using entry_idx:/path/to/file.txt format.
-o=, --output-dir=	Path to the job output directory. The directory will be created if it does not exist.
--preset-xml=[]	Preset/Option XML file. This option may be repeated if you have multiple preset files.
--preset-json=[]	Preset/Option JSON file. This option may be repeated if you have multiple preset files.
--preset-rc-xml	Skips loading preset from the environmental variable PB_SMRTPPIPE_XML_PRESET and explicitly loads the supplied preset.xml.
--service-uri	Remote Web services to send update and log status to. (This is a JSON file containing the host name and port number.)
--force-distributed	Override XML settings to enable distributed mode, if a cluster manager is provided.
--local-only	Override XML settings to disable distributed mode.
--force-chunk-mode	Override to enable Chunk mode.
--disable-chunk-mode	Override to disable Chunk mode.

**pipeline-id Command:** Run a registered pipeline by specifying the pipeline id.

```

pbsmrtpipe pipeline-id [-h] [--debug] -e ENTRY_POINTS
                       [ENTRY_POINTS ...] [-o OUTPUT_DIR]
                       [--preset-xml PRESET_XML]
                       [--preset-json PRESET_JSON]
                       [--preset-rc-xml PRESET_RC_XML]
                       [--service-uri SERVICE_URI]
                       [--force-distributed | --local-only]
                       [--force-chunk-mode | --disable-chunk-mode]
                       pipeline_id

```

Required	Description
pipeline_id	Registered pipeline id. Run show-templates to display a list of the registered pipelines.

Options	Description
--debug=False	Alias for setting log level to DEBUG.



Options	Description
<code>-e, --entry</code>	Entry Points using <code>entry_idX:/path/to/file.txt</code> format.
<code>-o=, --output-dir=</code>	Path to the job output directory. The directory will be created if it does not exist.
<code>--preset-xml=[]</code>	Preset/Option XML file. This option may be repeated if you have multiple preset files.
<code>--preset-json=[]</code>	Preset/Option JSON file. This option may be repeated if you have multiple preset files.
<code>--preset-rc-xml</code>	Skips loading preset from the environmental variable <code>PB_SMRTPIPE_XML_PRESET</code> and explicitly loads the supplied <code>preset.xml</code> .
<code>--service-uri</code>	Remote Web services to send update and log status to. (This is a JSON file containing the host name and port number.)
<code>--force-distributed</code>	Override XML settings to enable distributed mode, if a cluster manager is provided.
<code>--local-only</code>	Override XML settings to disable distributed mode.
<code>--force-chunk-mode</code>	Override to enable Chunk mode.
<code>--disable-chunk-mode</code>	Override to disable Chunk mode.

`task` Command: Run a task, such as a ToolContract, by id.

```
pbsmrtpipe task [-h] [--debug] -e ENTRY_POINTS [ENTRY_POINTS ...]
                [-o OUTPUT_DIR]
                [--preset-xml PRESET_XML]
                [--preset-json PRESET_JSON]
                [--preset-rc-xml PRESET_RC_XML]
                [--service-uri SERVICE_URI]
                [--force-distributed | --local-only]
                [--force-chunk-mode | --disable-chunk-mode]
                task_id
```

Required	Description
<code>task_id</code>	Show details of a registered task by id.

Options	Description
<code>--debug=False</code>	Alias for setting log level to <code>DEBUG</code> .
<code>-e, --entry</code>	Entry Points using <code>entry_idX:/path/to/file.txt</code> format.
<code>-o=, --output-dir=</code>	Path to the job output directory. The directory will be created if it does not exist.
<code>--preset-xml=[]</code>	Preset/Option XML file. This option may be repeated if you have multiple preset files.
<code>--preset-rc-xml</code>	Skips loading preset from the environmental variable <code>PB_SMRTPIPE_XML_PRESET</code> and explicitly loads the supplied <code>preset.xml</code> .
<code>--service-uri</code>	Remote Web services to send update and log status to. (This is a JSON file containing the host name and port number.)
<code>--force-distributed</code>	Override XML settings to enable distributed mode, if a cluster manager is provided.
<code>--local-only</code>	Override XML settings to disable distributed mode.
<code>--force-chunk-mode</code>	Override to enable Chunk mode.
<code>--disable-chunk-mode</code>	Override to disable Chunk mode.

`show-templates` Command: List all pipeline templates.

A pipeline 'id' can be referenced in your `my_pipeline.xml` file using `<import-template id="pbsmrtpipe.pipelines.my_pipeline_id" />`. This can replace the explicit listing of EntryPoints and Bindings.

```
pbsmrtpipe show-templates [-h]
                        [--log-level {DEBUG,INFO,WARNING,ERROR,CRITICAL}]
                        [--output-templates-avro OUTPUT_TEMPLATES_AVRO]
                        [--output-templates-json OUTPUT_TEMPLATES_JSON]
```

Options	Description
<code>--log-level</code>	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.]
<code>--output-templates-avro</code>	Resolve, validate and output Registered pipeline templates to AVRO files in <code>output-dir</code> .
<code>--output-templates-json</code>	Resolve, validate and output Registered pipeline templates to JSON files in <code>output-dir</code> .

`show-template-details` Command: Displays information about a specific pipeline.

This command lists the entry points required for the pipeline. These are usually PacBio Data Set XML files, although single raw data files (BAM or FASTA format) may be acceptable for some use cases. The most common input will be `eid_subread`, a SubreadSet XML Data Set, which contains one or more BAM files containing the raw unaligned subreads. Also common is `eid_ref_dataset`, for a ReferenceSet or genomic FASTA file.

```
pbsmrtpipe show-template-details [-h] [-o OUTPUT_PRESET_XML]
                                template_id
```

Required	Description
<code>template_id</code>	Show details of a registered Template by id.

Options	Description
<code>-o, --output-preset-xml</code>	Write pipeline/task <code>preset.xml</code> of options.

`show-tasks` Command: Show completed list of tasks by id.

Use the environmental variable `PB_TOOL_CONTRACT_DIR` to define a custom directory of tool contracts. These tool contracts will override the installed tool contracts, such as `PB_TOOL_CONTRACT_DIR=/path/to/my-tc-dir/`.

```
pbsmrtpipe show-tasks [-h]
```

`show-task-details` Command: Show details of a particular task by id, such as `pbsmrtpipe.tasks.filter_report`.

- Use `show-tasks` to get a complete list of registered tasks.

```
pbsmrtpipe show-task-details [-h] [-o OUTPUT_PRESET_XML] task_id
```

Required	Description
task_id	Show details of a registered task by id.

Options	Description
-o, --output-preset-xml	Write pipeline/task preset.xml of options.

`show-workflow-options` Command: Display all workflow-level options that can be set in `<options />` for `preset.xml`.

```
pbsmrtpipe show-workflow-options [-h] [-o OUTPUT_PRESET_XML]
```

Options	Description
-o, --output-preset-xml	Write pipeline/task preset.xml of options.

`run-diagnostic` Command: Performs diagnostic tests of `preset.xml` and the cluster configuration.

```
pbsmrtpipe run-diagnostic [-h] [--debug] [-o OUTPUT_DIR] [--simple]
                        preset_xml
```

Required	Description
preset_xml	Path to Preset XML file.

Options	Description
--debug=False	Alias for setting log level to <code>DEBUG</code> .
-o=, --output-dir=	Path to the job output directory. The directory will be created if it does not exist.
--simple=False	Perform full diagnostics tests; submit a test job to the cluster.

`show-chunk-operators` Command: Show a list of loaded chunk operators for Scatter/Gather Tasks. Extend resource loading by exporting the environmental variable `PB_CHUNK_OPERATOR_DIR`.

**Example:** `export PB_CHUNK_OPERATOR_DIR=/path/to/chunk-operators-xml-dir.`

```
pbsmrtpipe show-chunk-operators [-h]
```

### Example - Basic Resequencing

This pipeline uses `pbalign` to map reads to a reference genome, and `quiver` to determine the consensus sequence. The example uses the `sa3_ds_resequencing` pipeline:

```
$ pbsmrtpipe show-template-details pbsmrtpipe.pipelines.sa3_ds_resequencing
```

---

This requires two entry points: a SubreadSet and a ReferenceSet. A typical invocation might look like this (for a hypothetical lambda virus genome):

```
$ pbsmrtpipe pipeline-id pbsmrtpipe.pipelines.sa3_ds_resequencing \  
-e eid_subread:/data/smrt/2372215/0007/Analysis_Results/\  
m150404_101626_42267_c100807920800000001823174110291514_s1_p0.all.subreadset.xml \  
-e eid_ref_dataset:/data/references/lambdaNEB/lambdaNEB.referenceset.xml
```

This will run for a while and emit several directories, including tasks, logs, and workflow. The tasks directory is the most useful, as it stores the intermediate results and resolved tool contracts (how the task was executed) for each task. The directory names (`task_ids`) should be somewhat self-explanatory. To direct the output to a subdirectory in the current working directory, use the `-o` flag: `-o job_output_1`.

Other pipelines related to resequencing, such as the basemods detection and motif-finding, have nearly identical command-line arguments except for the pipeline ID.

For a general overview of the resequencing results, the GFF file written by `summarizeConsensus` is the most useful:

```
job_output_2/tasks/genomicconsensus.tasks.summarize_consensus-0/  
alignment_summary_variants.gff
```

The GFF file contains records for a complete set of sequence regions in the reference genome, including coverage statistics and the number of gaps, substitutions, insertions or deletions. For example:

```
lambda_NEB3011 .      region 1      50      0.00  +      .  
cov=116,190,190;cov2=183.000,14.633;gaps=0,0;cQv=20,20,20;del=0;ins=0;sub=0
```

### Example - Quiver (Genomic Consensus)

If you already have an AlignmentSet on which you just want to run quiver, the `sa3_ds_genomic_consensus` pipeline will be faster:

```
$ pbsmrtpipe pipeline-id pbsmrtpipe.pipelines.sa3_ds_genomic_consensus \  
-e eid_bam_alignment:/data/project/my_lambda_genome.alignmentset.xml \  
-e eid_ref_dataset:/data/references/lambda.referenceset.xml \  
--preset-xml=preset.xml
```

### Example - Circular Consensus Sequences

To obtain high-quality consensus sequences (also known as CCS reads) for individual SMRT Cell ZMWs from high-coverage subreads:

```
$ pbsmrtpipe pipeline-id pbsmrtpipe.pipelines.sa3_ds_ccs \  
-e eid_subread:/data/smrt/2372215/0007/Analysis_Results/\  
m150404_101626_42267_c100807920800000001823174110291514_s1_p0.all.subreadset.xml \  
--preset-xml preset.xml -o job_output
```

---

This pipeline is relatively simple and parallelizes especially well. The essential outputs are a ConsensusRead Data Set (composed of one or more unmapped BAM files) and corresponding FASTA and FASTQ files:

```
job_output/tasks/pbccs.tasks.ccs-0/ccs.consensusreadset.xml
job_output/tasks/pbsmrtpipe.tasks.bam2fasta_ccs-0/file.fasta
job_output/tasks/pbsmrtpipe.tasks.bam2fastq_ccs-0/file.fastq
```

The `pbccs.tasks.ccs-0` task directory will also contain a JSON report with basic metrics for the run such as number of reads passed and rejected for various reasons. (Note, as explained below, that the location of the final ConsensusRead XML - and JSON report - will be different in chunk mode.)

As the full resequencing workflow operates directly on subreads to produce a genomic consensus, it is not applicable to CCS reads. However, a CCS pipeline is available that incorporates the `blasr` mapping step:

```
$ pbsmrtpipe pipeline-id pbsmrtpipe.pipelines.sa3_ds_ccs_align \
  -e eid_subread:/data/smrt/2372215/0007/Analysis_Results/ \
m150404_101626_42267_c100807920800000001823174110291514_s1_p0.all.subreadset.xml \
  -e eid_ref_dataset:/data/references/lambda.referenceset.xml \
  --preset-xml preset.xml -o job_output
```

### Example - Iso-Seq™ Transcriptome Analysis

The Iso-Seq Transcriptome Analysis workflows automate use of the `pbtranscript` package for investigating mRNA transcript isoforms. The transcript analysis uses CCS reads where possible, and the pipeline incorporates the CCS pipeline with looser settings. The starting point is therefore still a SubreadSet. The simpler of the two pipelines is `sa3_ds_iseq_classify`, which runs CCS and classifies the reads as full-length or not:

```
$ pbsmrtpipe pipeline-id pbsmrtpipe.pipelines.sa3_ds_iseq_classify \
  -e eid_subread:/data/smrt/2372215/0007/Analysis_Results/ \
m150404_101626_42267_c100807920800000001823174110291514_s1_p0.all.subreadset.xml \
  --preset-xml preset.xml -o job_output
```

The output files from the CCS pipeline will again be present. Note however that the sequences will be lower-quality as the pipeline tries to use as many reads as possible. The output task folder `pbtranscript.tasks.classify-0` (or gathered equivalent; see below) contains the classified transcripts in various ContigSet Data Sets (or underlying FASTA files).

A more thorough analysis yielding Quiver-polished, high-quality isoforms is the `pbsmrtpipe.pipelines.sa3_ds_iseq` pipeline, which is invoked identically to the classify-only pipeline. Note that this is significantly slower, as the clustering step may take days to run for large Data Sets.

---

## Example - Exporting Subreads to FASTA/FASTQ

Converting a PacBio SubreadSet to FASTA or FASTQ format for use with external software can be done as a standalone pipeline. Unlike most of the other pipelines, this one has no task-specific options and no chunking, so the invocation is very simple:

```
$ pbsmrtpipe pipeline-id pbsmrtpipe.pipelines.sa3_ds_subreads_to_fastx \  
-e eid_subread:/data/smrt/2372215/0007/Analysis_Results/  
m150404_101626_42267_c100807920800000001823174110291514_s1_p0.all.subreadset.xml \  
-o job_output
```

The result files will be here:

```
job_output/tasks/pbsmrtpipe.tasks.bam2fasta-0/file.fasta  
job_output/tasks/pbsmrtpipe.tasks.bam2fastq-0/file.fastq
```

Both are also available gzipped in the same directories.

## Chunking

To take advantage of `pbsmrtpipe`'s parallelization, we need an XML configuration file for global `pbsmrtpipe` options, which can be generated by the following command:

```
$ pbsmrtpipe show-workflow-options -o preset.xml
```

The output `preset.xml` will have this format:

```
<?xml version="1.0" encoding="utf-8" ?>  
<pipeline-preset-template>  
  <options>  
    <option id="pbsmrtpipe.options.max_nproc">  
      <value>16</value>  
    </option>  
    <option id="pbsmrtpipe.options.chunk_mode">  
      <value>False</value>  
    </option>  
    <!-- MANY MORE OPTIONS OMITTED -->  
  </options>  
</pipeline-preset-template>
```

The appropriate types should be clear; quotes are unnecessary, and boolean values should have initial capitals (`True`, `False`). To enable chunk mode, change the value of option `pbsmrtpipe.options.chunk_mode` to `True`. Several additional options may also need to be modified:

- `pbsmrtpipe.options.distributed_mode` enables execution of most tasks on a managed cluster such as Sun Grid Engine. Use this for chunk mode if available.
- `pbsmrtpipe.options.max_nchunks` sets the upper limit on the number of jobs per task in chunked mode. Note that more chunks is not always better, as there is some overhead to chunking, especially in distributed mode.

- `pbsmrtpipe.options.max_nproc` sets the upper limit on the number of processors per job (including individual chunk jobs). This should be set to a value appropriate for your compute environment.

You can adjust `max_nproc` and `max_nchunks` in the `preset.xml` to consume as many queue slots as you desire, but note that the number of slots consumed will be the product of the two numbers. For some shorter jobs (typically with low-volume input data), it may make more sense to run the job unchunked but still distribute tasks to the cluster (where they will still use multiple cores if allowed).

Once you are satisfied with the settings, add it to your command like this:

```
$ pbsmrtpipe pipeline-id pbsmrtpipe.pipelines.sa3_ds_resequencing \
  --preset-xml preset.xml \
  -e eid_subread:/data/smrt/2372215/0007/Analysis_Results/
m150404_101626_42267_c100807920800000001823174110291514_s1_p0.all.subreadset.xml \
  -e eid_ref_dataset:/data/references/lambda.referenceset.xml
```

Alternately, the flags `--force-chunk-mode`, `--force-distributed`, `--disable-chunk-mode`, and `--local-only` can be used to toggle the chunk/distributed mode settings on the command line (but this will not affect the values of `max_nproc` or `max_nchunks`).

If the pipeline runs correctly, you should see an expansion of task folders. The final results for certain steps (alignment, variantCaller, and so on), should end up in the appropriate “gather” directory. For instance, the final gathered FASTA file from quiver should be in `pbsmrtpipe.tasks.gather_contigset-1`. Note that for many Data Set types, the gathered Data Set XML file will often encapsulate multiple BAM files in multiple directories.

### HdfSubreadSet to SubreadSet Conversion

If you have existing `bax.h5` files to process with `pbsmrtpipe`, you need to convert them to a SubreadSet before continuing. Bare `bax.h5` files are **not** directly compatible with `pbsmrtpipe`, but an HdfSubreadSet XML file can be easily generated from a `fofn` (file-of-file-names) or folder of `bax.h5` files using the `dataset` tool. (See “dataset” on page 17.)

From a `fofn`, `allTheBaxFiles.fofn`:

```
$ dataset create --type HdfSubreadSet allTheBaxFiles.hdfsubreadset.xml
allTheBaxFiles.fofn
```

Or a directory with all the `bax` files:

```
$ dataset create --type HdfSubreadSet allTheBaxFiles.hdfsubreadset.xml allTheBaxFiles/
*.bax.h5
```

This can be used as an entry point to the conversion pipeline. (We recommend using chunked mode if there is more than one `bax.h5` file, so include the appropriate `preset.xml`):

---

```
$ pbsmrtpipe pipeline-id pbsmrtpipe.pipelines.sa3_hdfsubread_to_subread \  
--preset-xml preset.xml -e eid_hdfsubread:allTheBaxFiles.hdfsubreadset.xml
```

And use the gathered output XML file as an entry point to the resequencing pipeline from earlier:

```
$ pbsmrtpipe pipeline-id pbsmrtpipe.pipelines.sa3_ds_resequencing \  
--preset-xml preset.xml \  
-e eid_subread:tasks/pbsmrtpipe.tasks.gather_subreadset-0/gathered.xml \  
-e eid_ref_dataset:/data/references/lambda.referenceset.xml
```

## Working with Data Sets

Data Sets can also be created for one or more existing `subreads.bam` files or `alignedsubreads.bam` files for use with the pipeline:

```
$ dataset create --type SubreadSet allTheSubreads.subreadset.xml \  
mySubreadBams/*.bam
```

or:

```
$ dataset create --type AlignmentSet allTheMappedSubreads.alignmentset.xml \  
myMappedSubreadBams/*.bam
```

Make sure that all `.bam` files have corresponding `.bai` and `.pbi` index files before generating the Data Set, as these make some operations significantly faster and are required by many programs. You can create indices with `samtools` and `pbindex`, both included in the distribution:

```
$ samtools index subreads.bam  
$ pbindex subreads.bam
```

In addition to the BAM-based Data Sets, and `HdfSubreadSet`, `pbsmrtpipe` also works with two Data Set types based on FASTA format: `ContigSet` (used for both *de novo* assemblies and other collections of contiguous sequences such as transcripts in the IsoSeq workflows) and `ReferenceSet` (a reference genome). These are created in the same way as BAM Data Sets:

```
$ dataset create --type ReferenceSet  
human_genome.referenceset.xml \  
genome/chr*.fasta
```

FASTA files can also be indexed for increased speed using `samtools`, and this is again recommended before creating the Data Set:

```
$ samtools faidx chr1.fasta
```

Note that PacBio's specifications for BAM and FASTA files impose additional restrictions on content and formatting; files produce by non-PacBio software are **not** guaranteed to work as input. The `pbvalidate` tool can be used to check for format compliance. (See page 62 for details.)



---

## pbtranscript

The `pbtranscript` tool is part of the Iso-Seq™ analysis pipeline, and it is used for the Classify and Cluster/polish steps, as well as post-polish analysis.

Using the command-line, the Iso-Seq analysis is performed in 3 steps:

1. Run CCS on your subreads, generating a CCS BAM file. Then generate an XML file from the BAM file.
2. Run Classify on your CCSs with the XML as input, generating a FASTA file of annotated sequences.
3. Run Cluster on the FASTA file produced by Classify, generating polished isoforms.

### Step 1: CCS

Convert the subreads to circular consensus sequences, using the following command:

```
ccs --noPolish --minLength=300 --minPasses=1 --minZScore=-999 --maxDropFraction=0.8 --minPredictedAccuracy=0.8 --minSnr=4 subreads.bam ccs.bam
```

Where:

- `ccs.bam` is where the CCSs will be output.
- `subreads.bam` is the file containing your subreads.

If you think that you have transcripts of interest that are less than 300 base pairs in length, be sure to adjust the `minLength` parameter. Next, you generate an XML file from your CCSs, using the following command:

```
dataset create --type ConsensusReadSet ccs.xml ccs.bam
```

Where:

- `ccs.xml` is the name of the XML file you are generating.
- `ccs.bam` is the name of the BAM file you generated previously using the `ccs` command.

### Step 2: Classify

Iso-Seq Classify classifies reads into full-length or non-full-length reads, artificial-concatemer chimeric, or non-chimeric reads.

To classify a read as full-length or non-full-length, we search for primers and poly-A within reads. If and **only** if both primers and poly-As are seen in a read, it is classified as a **full-length read**. Otherwise, the read is classified as **non-full-length**. We also remove primers and poly-As from reads and identify read-strandedness based on this information.

Next, full-length reads are classified into artificial-concatemer chimeric reads or non-chimeric reads by locating primer hits within reads.

- HMMER: `phmmer` in the HMMER package is used to detect locations of primer hits within reads and classify reads which have primer hits in the middle of sequences as artificial-concatemer chimeric.

### Classify - Input File

- `ccs.xml`: Circular consensus sequences generated from the CCS step.

### Classify - Output Files

- `isoseq_flnc.fasta`: Contains all full-length, non-artificial-concatemer reads.
- `isoseq_nfl.fasta`: Contains all non-full-length reads.
- `isoseq_draft.fasta`: An intermediate file in order to get full-length reads, which you can ignore.

Reads in these FASTA output files look like the following:

```
>m140121_100730_42141_c100626750070000001823119808061462_s1_p0/119/30_1067_CCS
strand=+;fiveseen=1;polyAseen=1;threeseen=1;fiveend=30;polyAend=1067;threeend=1096;primer=1;chimera=0
ATAAGACGACGCTATATG
```

These lines have the format:

```
<movie_name>/<ZMW>/<start>_<end>_CCS INFO
```

The `INFO` fields are:

- `strand`: Either + or -, whether a read is forward or reverse-complement cDNA.
  - `fiveseen`: Whether or not 5' prime is seen in this read, 1 is yes, 0 is no.
  - `polyAseen`: Whether or not poly-A tail is seen, 1 is yes, 0 is no.
  - `threeseen`: Whether or not 3' prime is seen, 1 is yes, 0 is no.
  - `fiveend`: Start position of 5' in the read.
  - `threeend`: Start position of 3' in the read.
  - `polyAend`: Start position of poly-A in the read.
  - `primer`: Index of primer seen in this read.
  - `chimera`: Whether or not this read is classified as a chimeric cDNA.
- `classify_summary.txt`: This file contains the following statistics:
    - Number of reads of insert
    - Number of five prime reads
    - Number of three prime reads
    - Number of poly-A reads
    - Number of filtered short reads
    - Number of non-full-length reads
    - Number of full-length reads

- Number of full-length non-chimeric reads
- Average full-length non-chimeric read length

**Note:** By seeing that the number of full-length, non-chimeric (`flnc`) reads is only a little less than the number of full-length reads, we can confirm that the number of artificial concatemers is very low. This indicates a successful SMRTbell library preparation.

### Classify - Usage

```
pbtranscript classify [OPTIONS] ccs.xml isoseq_draft.fasta --flnc=isoseq_flnc.fasta --nfl=isoseq_nfl.fasta
```

- Where `ccs.xml` is the XML file you generated in Step 1.
- `isoseq_flnc.fasta` contains only the full-length, non-chimeric reads.
- `isoseq_nfl.fasta` contains all non-full-length reads.

Or you can run Classify creating XML files instead of FASTA files as follows:

```
pbtranscript classify [OPTIONS] ccs.xml isoseq_draft.fasta --flnc=isoseq_flnc.contigset.xml --nfl=isoseq_nfl.contigset.xml
```

- Where `ccs.xml` is the XML file you generated in Step 1.
- `isoseq_flnc.contigset.xml` contains only the full-length, non-chimeric reads.
- `isoseq_nfl.contigset.xml` contains all non-full-length reads.

**Note:** One can always use `pbtranscript subset` to further subset `isoseq_draft.fasta` if `--flnc` and `--nfl` are **not** specified when you run `pbtranscript classify`. For example:

```
pbtranscript subset isoseq_draft.fasta isoseq_flnc.fasta --FL --nonChimeric
```

### Classify Options

- To view Classify options, enter `pbtranscript classify --help`.

Required	Description
<code>readsFN</code>	First positional argument. Input CCS reads in BAM, Data Set XML, or FASTA format. Example: <code>ccs.bam,xml,fasta</code> .
<code>outReadsFN</code>	Second positional argument. Output file which contains all classified reads in FASTA or contigset XML format. Example: <code>isoseq_draft.fasta,contigset.xml</code>

Options	Description
<code>--flnc</code>	Outputs full-length non-chimeric reads in FASTA or contigset XML format. Example: <code>FLNC_FA.fasta,contigset.xml</code>
<code>-d OUTDIR, --outDir OUTDIR</code>	Directory to store HMMER output. (Default = <code>output/</code> )
<code>-summary</code>	Text file to output classify summary. (Default = <code>out.classify_summary.txt</code> ).

Options	Description
<code>-p primers.fa, --primer primers.fa</code>	Primer FASTA file. (Default = <code>primers.fa</code> )
<code>--report</code>	CSV file of primer information. Contains the same information found in the description lines of the output FASTA. (Default = <code>out.primer_info.csv</code> )
<code>-cpus CPUS</code>	Number of CPUs to run HMMER. (Default = 8)
<code>--min_seq_len MIN_SEQ_LEN</code>	Minimum CCS length to be analyzed. Fragments shorter than the minimum sequence length are excluded from analysis. (Default = 300)
<code>--min_score MIN_SCORE</code>	Minimum <code>phmmer</code> score for primer hit. (Default = 10)
<code>--detect_chimera_nfl</code>	Detect chimeric reads among non-full-length reads. Non-full-length non-chimeric/chimeric reads are saved to <code>outDir/nflnc.fasta</code> and <code>outDir/nflc.fasta</code> .
<code>--ignore_polyA</code>	Full-length criteria does <b>not</b> require poly-A tail. By default this is off, which means that poly-A tails are required for a sequence to be considered full length. When it is turned on, sequences do <b>not</b> need poly-A tails to be considered full length.

### Step 3: Cluster and Polish

Iso-Seq Cluster performs isoform-level clustering using the Iterative Clustering and Error correction (ICE) algorithm, which iteratively classifies full-length non-chimeric CCS reads into clusters and builds consensus sequences of clusters using `pbdagcon`.

ICE is customized to work well on alternative isoforms and alternative polyadenylation sites, but **not** on SNP analysis and SNP-based highly complex gene families.

Iso-Seq Polish further polishes consensus sequences of clusters (i.e., `pbdagcon` output) taking into account all the QV information. Full-length non-chimeric CCS reads and non-full-length CCS reads are assigned into clusters based on similarity. Then for each cluster, we align raw subreads of its assigned ZMWs towards its consensus sequence. Finally, we load quality values to these alignments and polish the consensus sequence using `quiver` or `Arrow`.

#### Cluster - Input Files

- A file of non-full length reads output by `Classify`.
- A file of full-length non-chimeric reads.

#### Cluster - Output Files

- A file of polished, high-quality consensus sequences.
- A file of polished, low-quality consensus sequences.
- `cluster_summary.txt`, which contains the following statistics:
  - Number of consensus isoforms.
  - Average read length of consensus isoforms.
- `cluster_report.csv`; each line contains the following fields:

- cluster\_id: ID of a consensus isoforms from ICE.
- read\_id: ID of a read which supports the consensus isoform.
- read\_type: Type of the supportive read.

### Cluster - Usage

```
pbtranscript cluster [OPTIONS] isoseq_flnc.fasta polished_clustered.fasta --quiver --nfl=isoseq_nfl.fasta --bas_fofn=my.subreadset.xml
```

Or

```
pbtranscript cluster [OPTIONS] isoseq_flnc.contigset.xml polished_clustered.contigset.xml --quiver --nfl=isoseq_nfl.contigset.xml --bas_fofn=my.subreadset.xml
```

**Note:** `--quiver --nfl=isoseq_nfl.fasta|contigset.xml` **must** be specified to get Quiver/Arrow-polished consensus isoforms.

Optionally, you may call the following command to run ICE and create unpolished consensus isoforms only:

```
pbtranscript cluster [OPTIONS] isoseq_flnc.fasta unpolished_clustered.fasta
```

### Cluster Options

- To view Cluster options, use `pbtranscript cluster --help`.

Options	Description
Input reads	Input full-length non-chimeric reads in FASTA or contigset XML format. Used for clustering consensus isoforms. Example: <code>isoseq_flnc.fasta, contigset.xml</code> <b>(Required)</b>
Output Isoforms	Output predicted (unpolished) consensus isoforms in FASTA file. Example: <code>out.fasta, congigtset.xml</code> <b>(Required)</b>
<code>--nfl_fa isoseq_nfl.fasta</code>	Input non-full-length reads in FASTA format, used for polishing consensus isoforms.
<code>--ccs_fofn ccs.fofn</code>	A <code>ccs.fofn</code> , <code>ccs.bam</code> or <code>ccs.xml</code> file. If not given, Cluster assumes there is no QV information available.
<code>--bas_fofn my.subreadset.xml</code>	A file which provides quality values of raw reads and subreads. Can be either a <code>fofn</code> (file-of-file-names) of BAM or BAX files, or a Data Set XML.
<code>-d output/, --outDir output/</code>	Directory to store temporary and output cluster files. (Default = <code>output/</code> )
<code>--tmp_dir tmp/</code>	Directory to store temporary files. (Default = <code>tmp/</code> )
<code>--summary my.cluster_summary.txt</code>	TXT file to output cluster summary. (Default = <code>my.cluster_summary.txt</code> )
<code>--report report.csv</code>	A CSV file, each line containing a cluster, an associated read of the cluster, and the read type.
<code>--pickle_fn PICKLE_FN</code>	Developers' option from which all clusters can be reconstructed.
<code>--cDNA_size</code>	Estimated cDNA size. Values = [ <code>under1k</code> , <code>between1k2k</code> , <code>between2k3k</code> , <code>above3k</code> ]
<code>--quiver</code>	Call Quiver or Arrow to polish consensus isoforms using non-full-length non-chimeric CCS reads.

Options	Description
<code>--use_finer_qv</code>	Use finer classes of QV information from CCS input instead of a single QV from FASTQ. This option is slower and consumes more memory.
<code>--use_sge</code>	Specify that the cluster uses SGE.
<code>--max_sge_jobs MAX_SGE_JOBS</code>	The maximum number of jobs that will be submitted to SGE concurrently.
<code>--unique_id UNIQUE_ID</code>	Unique ID for submitting SGE jobs.
<code>--blasr_nproc BLASR_NPROC</code>	Number of cores for each BLASR job.
<code>--quiver_nproc QUIVER_NPROC</code>	Number of CPUs that each quiver/Arrow job uses.
<code>--hq_quiver_min_accuracy HQ_QUIVER_MIN_ACCURACY</code>	Minimum allowed quiver accuracy to classify an isoform as high-quality.
<code>-qv_trim_5 QV_TRIM_5</code>	Ignore QV of n bases in the 5' end.
<code>--qv_trim_3 QV_TRIM_3</code>	Ignore QV of n bases in the 3' end.
<code>--hq_isoforms_fa output/all_quivered_hq.fa</code>	Quiver/Arrow-polished, high-quality isoforms in FASTA. (Default = output/all_quivered_hq.fa)
<code>--hq_isoforms_fq output/all_quivered_hq.fq</code>	Quiver/Arrow-polished, high-quality isoforms in FASTQ. (Default = output/all_quivered_hq.fq)
<code>--lq_isoforms_fa output/all_quivered_lq.fa</code>	Quiver/Arrow-polished, low-quality isoforms in FASTA. (Default = output/all_quivered_lq.fa)
<code>--lq_isoforms_fq output/all_quivered_lq.fq</code>	Quiver/Arrow-polished, low-quality isoforms in FASTQ. (Default = output/all_quivered_lq.fq)

### Subset Options

Subset is an optional program which can be used to subset the output files for particular classes of sequences, such as non-chimeric reads, or non-full-length reads.

- To view Subset options, enter `pbtranscript subset --help`.

Options	Description
Input sequences	Input FASTA file. Example: <code>isoseq_draft.fasta</code> ( <b>Required</b> )
Output sequences	Output FASTA file. Example: <code>isoseq_subset.fasta</code> ( <b>Required</b> )
<code>--FL</code>	Output only full-length reads, with 3' primer and 5' primer and poly-A tail seen.
<code>--nonFL</code>	Output only non-full-length reads.
<code>--nonChimeric</code>	Output only non-chimeric reads.
<code>--printReadLengthOnly</code>	Only print read lengths, with no read names and sequences.
<code>--ignore_polyA</code>	Full-Length criteria does <b>not</b> require poly-A tail. By default this is off, which means that poly-A tails are required for a sequence to be considered full length. When it is turned on, sequences do <b>not</b> need poly-A tails to be considered full length.

**pbvalidate** The `pbvalidate` tool validates that files produced by PacBio software are compliant with Pacific Biosciences' own internal specifications.

## Input Files

pbvalidate supports the following input formats:

- BAM
- FASTA
- Data Set XML

See <http://pacbiofileformats.readthedocs.org/en/3.0/> for further information about each format's requirements.

## Usage

```
pbvalidate [-h] [--version] [--log-file LOG_FILE]
           [--log-level {DEBUG,INFO,WARNING,ERROR,CRITICAL} | --debug | --quiet | -v]
           [-c] [--quick] [--max MAX_ERRORS]
           [--max-records MAX_RECORDS]
           [--type
           {BAM,Fasta,AlignmentSet,ConsensusSet,ConsensusAlignmentSet,SubreadSet,BarcodeSet,ContigSet,ReferenceSet,GmapReferenceSet,HdfSubreadSet}]
           [--index] [--strict] [-x XUNIT_OUT] [--unaligned]
           [--unmapped] [--aligned] [--mapped]
           [--contents {SUBREAD,CCS}] [--reference REFERENCE]
           [--permissive-headers]
           file
```

Required	Description
file	BAM, FASTA, or Data Set XML file to validate.

Options	Description
-h, --help	Displays help information and exits.
--version	Displays program version number and exits.
--log-file LOG_FILE	Write the log to file. Default (None) will write to stdout.
--log-level	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL.] (Default = CRITICAL)
--debug=False	Alias for setting the log level to DEBUG. (Default = False)
--quiet	Alias for setting the log level to CRITICAL to suppress output. (Default = False)
--verbose, -v	Set the verbosity level. (Default = None)
--quick	Limits validation to the first 100 records (plus file header); equivalent to --max-records=100. (Default = False)
--max MAX_ERRORS	Exit after MAX_ERRORS have been recorded. (Default = None; check the entire file.)
--max-records MAX_RECORDS	Exit after MAX_RECORDS have been inspected. (Default = None; check the entire file.)
--type	Use the specified file type instead of guessing. [BAM, Fasta, AlignmentSet, ConsensusSet, ConsensusAlignmentSet, SubreadSet, BarcodeSet, ContigSet, ReferenceSet, GmapReferenceSet, HdfSubreadSet] (Default = None)

Options	Description
--index	Require index files: .fai or .pbi. (Default = False)
--strict	Turn on additional validation, primarily for Data Set XML. (Default = False)

BAM Options	Description
--unaligned	Specify that the file should contain <b>only</b> unmapped alignments. (Default = None, no requirement.)
--unmapped	Alias for --unaligned. (Default = None)
--aligned	Specify that the file should contain <b>only</b> mapped alignments. (Default = None, no requirement.)
--mapped	Alias for --aligned. (Default = None)
--contents	Enforce the read type: [SUBREAD, CCS] (Default = None)
--reference REFERENCE	Path to optional reference FASTA file, used for additional validation of mapped BAM records. (Default = None)
--permissive-headers	Don't check chemistry/basecaller versions. (Default = False)

## Examples

To validate a BAM file:

```
$ pbvalidate in.subreads.bam
```

To validate a FASTA file:

```
$ pbvalidate in.fasta
```

To validate a Data Set XML file:

```
$ pbvalidate in.subreadset.xml
```

To validate a BAM file and its index file (.pbi):

```
$ pbvalidate --index in.subreads.bam
```

To validate a BAM file and exit after 10 errors are detected:

```
$ pbvalidate --max 10 in.subreads.bam
```

To validate up to 100 records in a BAM file:

```
$ pbvalidate --max-records 100 in.subreads.bam
```

To validate up to 100 records in a BAM file (equivalent to --max-records=100):

```
$ pbvalidate --quick in.subreads.bam
```



---

To validate a BAM file, using a specified log level:

```
$ pbvalidate --log-level=INFO in.subreads.bam
```

To validate a BAM file and write log messages to a file rather than to stdout:

```
$ pbvalidate --log-file validation_results.log in.subreads.bam
```

**quiver** Quiver is the legacy consensus model based on a conditional random field approach. Quiver enables consensus accuracies on genome assemblies at accuracies approaching or even exceeding Q60 (one error per million bases). If you use the HGAP assembly protocol in SMRT Portal v2.0 or later, Quiver runs automatically as the final "assembly polishing" step.

Quiver identifies haploid SNPs and single-base indels by comparing a multiple sequence alignment of mapped reads against a reference sequence.

Over the years Quiver has proven difficult to train and develop, so we are phasing it out in favor of the new model, Arrow. Arrow is an improved consensus model based on a more straightforward hidden Markov model approach.

- Quiver is supported for PacBio *RS* and PacBio *RS II* data.
- Arrow is supported for PacBio *Sequel* data and PacBio *RS* data with the P6-C4 chemistry.

### Usage

```
% quiver aligned_reads{.cmp.h5, .bam, .fofn, or .xml} \
> -r reference{.fasta or .xml} -o variants.gff \
> -o consensus.fasta -o consensus.fastq
```

In this example we perform haploid consensus and variant calling on the mapped reads in the `aligned_reads.bam` which was aligned to `reference.fasta`. The `reference.fasta` is only used for designating variant calls, not for computing the consensus. The consensus quality score for every position can be found in the output FASTQ file.

**sawriter** The `sawriter` tool generates a suffix array file from an input FASTA file. It is used to prebuild suffix array files for reference sequences which can later be used in resequencing workflows. `sawriter` comes with `blasr`, and is independent of `python`.

### Usage

```
sawriter saOut fastaIn [fastaIn2 fastaIn3 ...] [-blt p] [-larsson] [-4bit] [-manmy]
[-kar]
or
sawriter fastaIn (writes to fastIn.sa)
```

Options	Description
-blt p	Build a lookup table on prefixes of length p. This speeds up lookups considerably (more than the LCP table), but misses matches less than p when searching.
-4bit	Read in one FASTA file as a compressed sequence file.
-larsson	Uses the Larsson and Sadakane method to build the array. (Default)
-mamy	Uses the MAnber and MYers method to build the array. This is slower than the Larsson method, and produces the same result. This is mainly for double-checking the correctness of the Larsson method.
-kark	Uses the Karkkainen DS3 method for building the suffix array. This is probably slower than the Larsson method, but takes only $N/(\sqrt{3})$ extra space.
-welter	Use lightweight (sort of light) suffix array construction. This is a bit slower than the normal Larsson method.
-welterweight N	Use a difference cover of size N for building the suffix array. Valid values are 7, 32, 64, 111, and 2281.

## summarize Modifications

The `summarizeModifications` tool generates a GFF summary file (`alignment_summary.gff`) from the output of base modification analysis (i.e. `ipdSummary`) combined with the coverage summary GFF generated by resequencing pipelines. This is also part of the standard Base Modification pipelines in `pbsmrtpipe`, and is useful for power users running custom workflows.

### Usage

```
summarizeModifications [-h] [--version] [--emit-tool-contract]
                        [--resolved-tool-contract RESOLVED_TOOL_CONTRACT]
                        [--log-file LOG_FILE]
                        [--log-level {DEBUG,INFO,WARNING,ERROR,CRITICAL}] | --debug
                        | --quiet | -v
                        modifications alignmentSummary gff_out
```

### Input Files

- `modifications`: Base Modification GFF file.
- `alignmentSummary`: Alignment Summary GFF file.

### Output Files

- `gff_out`: Coverage summary for regions (bins) spanning the reference with Base Modification results for each region.

Options	Description
-h, --help	Displays help information and exits.
--version	Displays program version number and exits.
--emit-tool-contract	Emit tool contract to <code>stdout</code> . (Default = <code>False</code> )
--resolved-tool-contract RESOLVED_TOOL_CONTRACT	Run the tool directly from a PacBio Resolved tool contract. (Default = <code>None</code> )

Options	Description
<code>--log-file LOG_FILE</code>	Write the log to file. Default (None) will write to stdout.
<code>--log-level</code>	Specify the log level; values are [DEBUG, INFO, WARNING, ERROR, CRITICAL] (Default = INFO)
<code>--debug</code>	Alias for setting the log level to DEBUG. (Default = False)
<code>--quiet</code>	Alias for setting the log level to CRITICAL to suppress output. (Default = False)
<code>--verbose, -v</code>	Set the verbosity level. (Default = None)

**variantCaller** `variantCaller.py` is a variant-calling tool provided by the `GenomicConsensus` package which provides several variant-calling algorithms for PacBio sequencing data.

### Usage

```
variantCaller.py -j8 --algorithm=arrow \
-r lambdaNEB.fa \
-o variants.gff \
aligned_subreads.bam
```

This example requests variant calling, using 8 worker processes, using the Arrow algorithm, taking input from the file `aligned_subreads.bam`, using the FASTA file `lambdaNEB.fa` as the reference, and writing output to `variants.gff`.

A particularly useful option is `--referenceWindow/-w`; which allows the user variant calling to be performed exclusively on a \*window\* of the reference genome.

### Input Files

- A sorted file of reference-aligned reads in Pacific Biosciences' standard BAM format.
- A FASTA file that follows Pacific Biosciences' FASTA file convention.

**Note:** The `quiver` and `arrow` algorithms require that certain metrics are in place in the input BAM file.

- `quiver`, which operates on PacBio RS II data **only**, requires the basecaller-computed "pulse features" `InsertionQV`, `SubstitutionQV`, `DeletionQV`, and `DeletionTag`. These features are populated in BAM tags by the `bax2bam` conversion program. See page 6 for details.
- `arrow`, which operates on PacBio RS II P6-C4 data and all Sequel data, requires per-read SNR metrics, and the per-base `PulseWidth` metric for Sequel data (but **not** for PacBio RS II P6-C4). These metrics are populated by Sequel instrument software or the `bax2bam` converter (for PacBio RS II data).

The selected algorithm will halt with an error message if any features that it requires are unavailable.

---

## Output Files

Output files are specified as arguments to the `-o` flag. The file name extension provided to the `-o` flag is meaningful, as it determines the output file format. For example:

```
variantCaller aligned_subreads.bam -r lambda.fa -o myVariants.gff -o myConsensus.fasta
```

will read input from `aligned_subreads.bam`, using the reference `lambda.fa`, and send variant call output to the file `myVariants.gff`, and consensus output to `myConsensus.fasta`.

The file formats presently supported, by extension, are:

- `.gff`: PacBio GFFv3 variants format; convertible to VCF or BED.
- `.fasta`: FASTA file recording the consensus sequence calculated for each reference contig.
- `.fastq`: FASTQ file recording the consensus sequence calculated for each reference contig, as well as per-base confidence scores

Options	Description
<code>-j</code>	The number of worker processes to use.
<code>--algorithm=</code>	The variant-calling algorithm to use; values are <code>plurality</code> , <code>quiver</code> , and <code>arrow</code> .
<code>-r</code>	The FASTA reference file to use.
<code>-o</code>	The output file format; values are <code>.gff</code> , <code>.fasta</code> , and <code>.fastq</code> .

## Available Algorithms

At this time there are three algorithms available for variant calling: `plurality`, `quiver`, and `arrow`.

- `Plurality` is a simple and very fast procedure that merely tallies the most frequent read base or bases found in alignment with each reference base, and reports deviations from the reference as potential variants. This is a very insensitive and flawed approach for PacBio sequence data, which is prone to insertion and deletion errors.
- `Quiver` is a more complex procedure based on algorithms originally developed for CCS. `Quiver` leverages the quality values (QVs) provided by upstream processing tools, which provide insight into whether insertions/deletions/substitutions were deemed likely at a given read position. Use of `quiver` requires the ConsensusCore library.

- 
- `Arrow` is the successor to `quiver`; it uses a more principled HMM model approach. It does not require basecaller quality value metrics; rather, it uses the per-read SNR metric and the per-pulse `pulsewidth` metric as part of its likelihood model. Beyond the model specifics, other aspects of the `Arrow` algorithm are similar to `quiver`. Use of `arrow` requires the `ConsensusCore2` library, which is provided by the `unanimity` codebase.

### Confidence Values

The `Arrow`, `Quiver`, and `Plurality` algorithms make a confidence metric available for every position of the consensus sequence. The confidence should be interpreted as a phred-transformed posterior probability that the consensus call is incorrect; i.e.

$$QV = -10\log_{10}(p_{err})$$

`variantCaller` clips reported QV values at 93---larger values cannot be encoded in a standard FASTQ file.

### Chemistry Specificity

The `Quiver` and `Arrow` algorithm parameters are trained per-chemistry. `Quiver` and `Arrow` identify the sequencing chemistry used for each run by looking at metadata contained in the data file (the input BAM or `cmp.h5` file). This behavior can be overridden by a command-line flag.

When multiple chemistries are represented in the reads in the input file, `Quiver/Arrow` will model reads appropriately using the parameter set for its chemistry, thus yielding optimal results.

---

## Third Party Command-Line Tools

Following is information on the third-party command-line tools included in the `smrtcnds/bin` subdirectory.

- bamtools**
- A C++ API and toolkit for reading, writing, and manipulating BAM files.
  - See <https://sourceforge.net/projects/bamtools/> for details.

- daligner,  
LASort,  
LAmerge,  
HPC.daligner**
- Finds all significant local alignments between reads.
  - See <https://dazzlerblog.wordpress.com/command-guides/daligner-command-reference-guide/> for details.

- datander**
- Finds all local self-alignment between long, noisy DNA reads.
  - See <https://github.com/thegenemyers/DAMASKER> for details.

**DB2fasta,  
DBdump,  
DBdust, DBrm,  
DBshow,  
DBsplit,  
DBstats,  
Fasta2DB**

Utilities that work with Dazzler databases:

- **DB2fasta**: Converts database files to FASTS format.
- **DBdust**: Runs the DUST algorithm over the reads in the untrimmed database, producing a track that marks all intervals of low complexity sequence.
- **DBdump/DBshow**: Displays a subset of the reads in the database; selects the information to show about the reads, including any mask tracks.
- **DBrm**: Deletes all the files in a given database.
- **DBsplit**: Divides a database conceptually into a series of blocks.
- **DBstats**: Shows overview statistics for all the reads in the trimmed database.
- **Fasta2DB**: Builds an initial database, or adds to an existing database, using a list of `.fasta` files.
- See [https://dazzlerblog.wordpress.com/command-guides/dazz\\_db-command-guide/](https://dazzlerblog.wordpress.com/command-guides/dazz_db-command-guide/) for details.

- gmap,  
gmap\_build,  
gmapl**
- A genomic mapping and alignment program for mRNA and EST Sequences.
  - See <http://research-pub.gene.com/gmap/> for details.

- ipython**
- An interactive shell for using the Pacific Biosciences API.
  - See <https://ipython.org/> for details.

- python**
- An object-oriented programming language.
  - See <https://www.python.org/> for details.

---

**REPmask,  
TANmask,  
HPC.REPmask,  
HPC.TANmask**

- A set of programs to soft-mask all tandem and interspersed repeats in Dazzler databases when computing overlaps.
- See <https://github.com/thegenemyers/DAMASKER> for details.

**samtools**

- A set of programs for interacting with high-throughput sequencing data in SAM/BAM/VCF formats.
- See <http://www.htslib.org/> for details.

For Research Use Only. Not for use in diagnostic procedures. © Copyright 2017, Pacific Biosciences of California, Inc. All rights reserved. Information in this document is subject to change without notice. Pacific Biosciences assumes no responsibility for any errors or omissions in this document. Certain notices, terms, conditions and/or use restrictions may pertain to your use of Pacific Biosciences products and/or third party products. Please refer to the applicable Pacific Biosciences Terms and Conditions of Sale and to the applicable license terms at <http://www.pacb.com/legal-and-trademarks/product-license-and-use-restrictions/>.

Pacific Biosciences, the Pacific Biosciences logo, PacBio, SMRT, SMRTbell, Iso-Seq and Sequel are trademarks of Pacific Biosciences. BluePippin and SageELF are trademarks of Sage Science, Inc. NGS-go and NGSengine are trademarks of GenDx. All other trademarks are the sole property of their respective owners.

P/N 100-939-900-01