
Table of Contents

1. Introduction	1
1.1. What makes MOSAIK different?	1
1.2. Overview	2
1.3. Contact Information	2
2. The MOSAIK suite	3
2.1. MosaikBuild	3
2.2. MosaikAligner	5
2.3. MosaikSort	11
2.4. MosaikMerge	13
2.5. MosaikAssembler	14
3. Utilities	15
3.1. MosaikText	15
3.2. MosaikCoverage	15
3.3. MosaikDupSnoop	16
3.4. MosaikJump	16
4. Understanding Program Output	18
4.1 MosaikBuild	18
4.2 MosaikAligner	19
4.3 MosaikSort	20
5. Performance	21
5.1. Speeding Up Alignments	21
5.2. Yeast Alignments	22
5.3. Roundworm Alignments	22
5.4. Human Alignments	23
5.5. Aligner Settings By Sequencing Technology	23
6. Visualization	24
6.1. consed	24
6.2. Gambit	24
7. Data Access (MosaikTools)	25
Appendix 1: Crib Sheet	26

1. Introduction

MOSAIK is a reference-guided assembler comprising of four main modular programs (Figure 1.1):

- Mosaik**Build**
- Mosaik**Aligner**
- Mosaik**Sort**
- Mosaik**Assembler**.

Mosaik**Build** converts various sequence formats into Mosaik's native read format. Mosaik**Aligner** pairwise aligns each read to a specified series of reference sequences. Mosaik**Sort** resolves paired-end reads and sorts the alignments by the reference sequence coordinates. Finally, Mosaik**Assembler** parses the sorted alignment archive and produces a multiple sequence alignment which is then saved into an assembly file format.

Historically, the workflow consists of supplying sequences in *FASTA*, *FASTQ*, *Illumina Bustard & Gerald*, or *SRF file formats* and producing assembly files (*phrap ace and GigaBayes giga formats*) which can be viewed with utilities such as *consed*.

Our lab has been developing a new SNP caller called BamBayes which uses the new **BAM alignment file format**. Our current workflow involves producing alignment files with MosaikAligner and then exporting the alignments to the BAM format with MosaikText.

1.1. What makes MOSAIK different?

Unlike many current read aligners, MOSAIK produces gapped alignments using the Smith-Waterman algorithm. Additionally, our program goes beyond producing pairwise alignments and produces reference-guided assemblies with gapped alignments. These features make it ideal for downstream single nucleotide polymorphism (SNP) and short insertion/deletion (INDEL) discovery.

MOSAIK is written in highly portable C++ and currently targetted for the following platforms: Microsoft Windows, Apple Mac OS X, Linux/x86, Linux/Itanium2, and Sun Solaris/ UltraSPARC operating systems. Other platforms can easily be supported upon request.

MOSAIK is multithreaded. If you have a machine with 8 processors, you can use all 8 processors to align reads faster while using the same memory footprint as when using one processor.

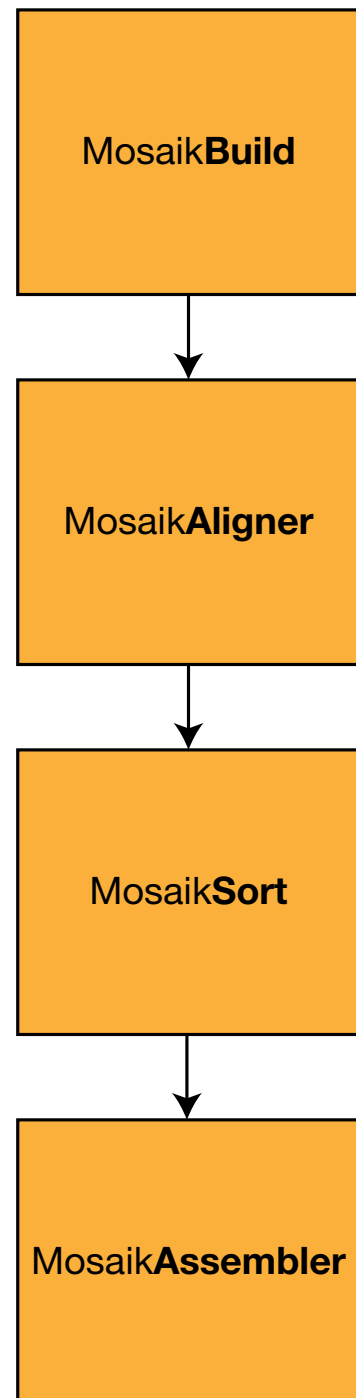


Figure 1.1 MOSAIK Pipeline.

This figure illustrates the basic MOSAIK pipeline. Mosaik**Build** converts external read formats, Mosaik**Aligner** pairwise aligns the reads, Mosaik**Sort** sorts the alignments and resolves paired-end reads, and finally Mosaik**Assembler** creates a gapped assembly file.

MOSAIK supports multiple sequencing technologies. In addition to legacy technologies such as Sanger capillary sequencing, our program supports next generation technologies such as Roche 454, Illumina, AB SOLiD, and experimental support for the Helicos Heliscope.

MOSAIK aligns AB SOLiD reads in colorspace and then converts the reads seamlessly back to basespace. This enables the user to take advantage of the SOLiD platform without all of the downstream bioinformatics headaches associated with working in colorspace.

1.2. Overview

The primary MOSAIK programs (MosaikBuild, MosaikAligner, MosaikSort, MosaikMerge, and MosaikAssembler) are discussed individually in Chapter 2.

The MOSAIK utility programs (MosaikText, MosaikCoverage, MosaikDupSnoop, and MosaikJump) are discussed individually in Chapter 3.

“MOSAIK output demystified” is probably an apt title for Chapter 4. Chapter 5 discusses how to fine tune performance for a variety of reference genomes and sequencing technologies.

The final two chapters discuss methods of visualizing and accessing the data after alignment. Chapter 6 mentions two visualization tools: consed & Gambit. Chapter 7 introduces the MOSAIK API for interacting with the alignment archives natively.

For the impatient, a one page crib sheet with the most important parameters to the primary MOSAIK programs is provided in the appendix.

1.3. Contact Information

Feel free to contact me to suggest improvements, submit bug reports, or just to offer some moral support while I write my Ph.D. dissertation:

Michael Strömberg
Biology Department
Boston College - Higgins Hall
140 Commonwealth Ave
Chestnut Hill, MA 02467
USA

email: mikaels@bc.edu

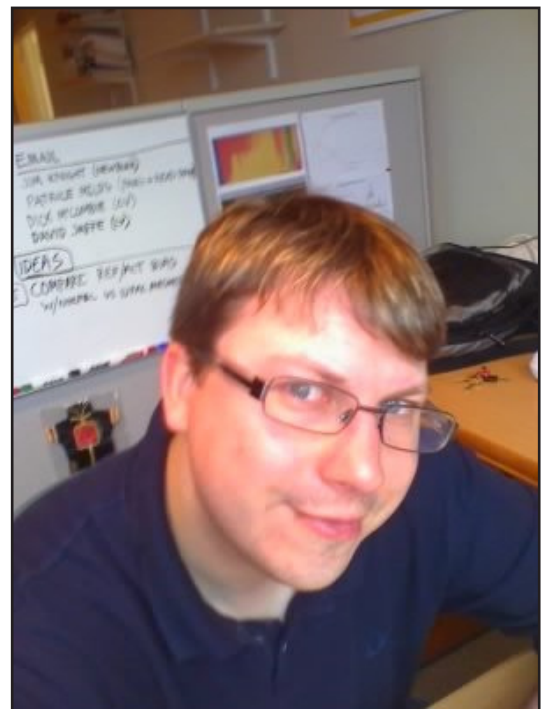


Figure 1.2 **Michael Strömberg.**

Feel free to say “Hello” if you see this character at a conference. Offer him a Duvel and he might be able to get any feature you want put into MOSAIK.

2. The MOSAIK suite

One of the author's pet peeves are programs that require way too many command-line parameters when decent default values should suffice. As a result, MOSAIK has been revised to require as few command-line parameters as possible. In essence, most programs require only an input filename and an output filename. Parameters with default values have been marked blue in the information panels.

Apparently there are more than a few system administrators out who are threatening bodily harm to faithful MOSAIK users (Figure 2.1). MosaikSort, MosaikMerge, and MosaikAssembler store temporary files in /tmp or similar platform-specific directories. These programs delete the temporary files when they finish.

If the /tmp directory is too small or if your local system administrator is lurking over your cubicle with a knife in hand, **you can specify another directory** by setting the MOSAIK_TMP environment variable. For example, Gábor might use the following command in *bash*:

```
export MOSAIK_TMP = /home/marth/tmp
```

2.1. MosaikBuild

To speed up the assembly pipeline, compressed binary file formats are used extensively throughout MOSAIK. MosaikBuild translates external read formats to a format that the aligner can readily use. In addition to processing reads, the program also converts reference sequences from a FASTA file to an efficient binary format.

MosaikBuild readily converts FASTA, FASTQ, Illumina Bustard, Illumina Gerald, and SRF files. Pyrosequences in SFF files can be converted to the FASTA format with the PyroBayes utility (<http://bioinformatics.bc.edu/marthlab/PyroBayes>). With files containing both bases and base qualities, such as FASTQ and SRF, MosaikBuild can convert an entire directory of read files into a single MOSAIK read file. This is handy, for example, when converting a run of Illumina paired-end reads that are separated by lanes and mate-pairs.

MosaikBuild automatically handles FASTA and FASTQ in both the uncompressed or compressed (gzipped) state. There's no need to uncompress the files prior to importing them into MOSAIK.

MosaikBuild stores the specified sequencing technology (**-st** parameter) and read type (single ended vs paired-end) so that it can be tracked throughout the entire pipeline.



Figure 2.1 Angry System Administrator.

This photo was taken moments before a user was savagely attacked by an angry system administrator for filling up the /tmp directory.

Metadata

MOSAİK now supports the concept of read groups that was introduced by the SAM & BAM file format specification. The following data can now be associated with each read archive:

- center name (**-cn**)
- description (**-ds**)
- read group identifier (**-id**)
- library name (**-ln**)
- median fragment length (**-mfl**)
- platform unit (**-pu**)
- sample name (**-sam**)
- sequencing technology (**-st**)

When using the **local alignment search** feature in MosaikAligner, the *median fragment length* metadata is used to define the search area.

Using MosaikBuild (reference sequences)

```
MosaikBuild -fr c_elegans_chr2_ref.fa.gz -oa
c_elegans_chr2_ref.dat
```

Here we convert the reference sequence for chromosome 2 in *C. elegans* from a FASTA file to a native MOSAİK reference archive.

In all of the MOSAİK programs, parameters can be placed in any order.

Using MosaikBuild (reads)

```
MosaikBuild -q c_elegans_chr2_mate1.fq.gz -q2
c_elegans_chr2_mate2.fq.gz -out c_elegans_chr2.dat -st
illumina -p B_
```

Here we convert the paired-end Illumina reads from a pair of FASTQ files to a native MOSAİK read archive. In this example we prepend each read name with “B_” so that we know that these reads are from the Bristol strain during visualization and SNP discovery.

Box 2.1 MosaikBuild Parameters.

FASTA

-fr, -fr2 specifies the files containing the bases for the first mate (-fr) and the second mate (-fr2).

-fq, -fq2 specifies the files containing the base qualities for the first mate (-fq) and the second mate (-fq2).

-assignQual assigns a base quality for every base.

FASTQ

-q, -q2 specifies the files for the first mate (-q) and the second mate (-q2).

SRF

-srf specifies the short read format file (currently only single-ended reads are supported).

Illumina Bustard

-bd specifies the Illumina Bustard directory.

-il specifies which lanes to use. e.g. -il 137 will select lanes 1, 3, and 7.

-split split the reads into two equal portions for paired-end reads.

Illumina Gerald

-gd specifies the Illumina Gerald directory.

-il specifies which lanes to use. e.g. -il 137 will select lanes 1, 3, and 7.

Shared Options

-cs translates the reference sequence from basespace to colorspace.

-p adds the specified prefix to each read name.

-st specifies the sequencing technology: 454, helicos, illumina, sanger, solid

-out specifies the output file when converting reads.

-oa specifies the output file when converting the reference sequence.

2.2. MosaikAligner

MosaikAligner performs pairwise alignment between every read in the read archive and a set of reference sequences (Figure 2.2). The program uses a hashing scheme similar to BLAT and BLAST and places all of the hashes (k-words or seeds) into a hash map or into a “jump database”.

When presented with a new read, MosaikAligner hashes up the read in a similar fashion and retrieves the reference positions for each hash in the hash table. These hash positions are clustered together and then evaluated with a full Smith-Waterman algorithm. Alignments are then screened according to the filter criteria specified by the user.

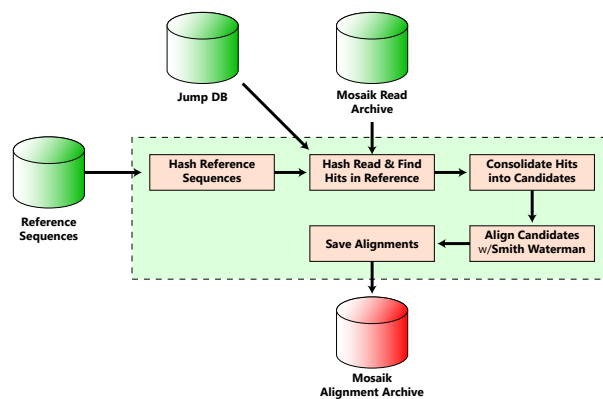


Figure 2.2 MosaikAligner Illustrated.

Hashing Strategy

The first incarnation of MOSAİK used a fast hashing strategy that stored one hash position per seed. If a seed was seen more than once in the reference sequence, it was marked non-unique and removed from subsequent use. When clustering hash positions, the longest contiguous set of hash was used to seed the full Smith-Waterman alignment. While speedy, this strategy (**FAST**) was vulnerable to microrepeat structures in reference sequences and often produced read pile-ups in repeat regions.

Subsequently the hashing strategy was revised to store *n* number of hash positions per seed and all hash position clusters were used to seed the full Smith-Waterman alignment. As the number of hash positions per seed increased, so did alignment accuracy and the memory requirements.

Three newer strategies are available:

- **SINGLE** (stores 1 hash position per seed)
- **MULTI** (stores 9 hash positions per seed)
- **ALL** (stores all hash positions per seed).

We have standardized on the use of the most accurate version (**ALL**) for all of our current lab projects.

Using Jump Databases

The aforementioned hashing strategies involve using hash maps as their primary data structure. However, putting an entire mammalian genome into a hash map is not very efficient. *MosaikAligner consumes around 57 GB RAM when using the ALL strategy and aligning against the entire mouse or human genome!*

Using a “jump database” makes full mammalian genome alignment possible (Figure 2.3). A jump database is comprised of two main files: the *keys file* and the *positions file*. The

reference sequence: **AATATGAATTAAATCAAAG**

key structure (4 ⁿ)		position structure (n _{hashes} + n _{positions})				
AAAA	AAAC	1	0	2	6	11
AAAG » 20	AAAT	1	16			
AACA	AACC					
AACG	AACT					
AAGA	AAGC					
AAGG	AAGT					
AATA » 0	AATC					
AATG	AATT » 8					

Figure 2.3 The Jump Database Illustrated.

In this example we are hashing the reference sequence with a hash size of 4. The hashes are store in 2-bit notation and thus have an equivalent integer value. e.g. AAAA=0, AAAC=1, AAAG=2, etc.

Using the integer value we can jump directly in the *keys file* to retrieve the *positions file* offset. e.g. if we look into the bucket corresponding to hash AAAG, we retrieve a file offset of 20.

Every record in the *positions file* starts with a number indicating the hash position count and is followed by the hash positions. e.g. at file offset 20, we discover one hash position at reference sequence offset 16.

positions file is a tightly packed file with all of the possible hash positions for each seed. The *keys file* is a sparse look up table that allows MOSAIK to find the relevant hash positions with only two hops in memory. This data structure guarantees no collisions, which is often prevalent in highly loaded hash maps.

What does this mean to the user? *Instead of 57 GB, we can now deliver the same data at 19 GB RAM.* Since there are no collisions, alignments are faster.

As an added bonus, the user can select whether the keys and positions files should remain on disk or be loaded into memory for additional performance. Technically, this means that a full genome alignment can be performed with as little as 4 GB RAM (used by the mammalian genome reference sequence). However the random access usage patterns in the “jump database” still make this alternative prohibitively slow. Flash-based memory such as solid state drives (SSDs) might make this feature more practical.

The Art of Being Ambiguous and Masking

Many aligners convert reads and reference sequences into an efficient 2-bit format. While this may reduce the memory footprint, it discards valuable information contained in the IUPAC ambiguity codes.

MOSAIK uses the full set of IUPAC ambiguity codes during alignment (Figure 2.4). For example, if you use a reference sequence where all of the major confirmed SNPs from dbSNP are marked with the appropriate ambiguity codes, alignment bias for one allele over the other can be avoided. i.e. if the reference genome has an ‘M’ allele, a read position with an ‘A’ or ‘C’ at that location will be scored as a match.

The only ambiguity code that is not strictly implemented is ‘N’. Many finished genomes uses large regions (often 50,000 bases) of N’s to denote where two contigs are linked but the actual genomic data is missing. In order to prevent every single read from aligning to those regions, ‘N’ is interpreted as ‘X’ (meaning this does not align to any of the four nucleotides).

MOSAIK has no problems with reference sequences that are hard masked with X’s.

Hash Size Selection (Speed vs. Sensitivity)

In general, using a larger hash size equates to a faster alignment speed. However this is often at the expense of sensitivity. e.g. using a hash size of 30 when aligning 35 bp reads will certainly be fast, but your reads will not be seeded if you have any internal sequencing errors or SNPs. In contrast,

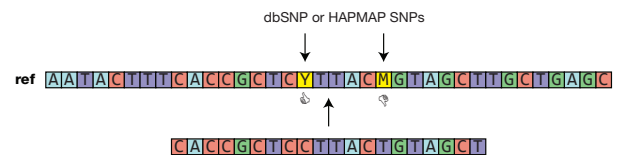


Figure 2.4 IUPAC Ambiguity Codes.

Bias against known SNPs can be reduced by using a reference sequence masked with IUPAC ambiguity codes. In the example, the alignment would result in one mismatch (the Y matches the C base, but M doesn’t match the T base).

a hash size of 11 would guarantee that all reads with up to two mismatches are seeded but performance would suffer.

When aligning mammalian reads, a hash size of 15 provides a happy medium between speed and sensitivity.

Limiting Hash Positions

While scaling up to mammalian alignments, we added a feature (**-mhp**) that places a maximum number of hash positions per seed. Alignment representation bias is minimized by selecting a random subset of the hash positions.

When using a hash size of 15, each seed has an average of 5.25 hash positions in the human genome. Limiting the number of hash positions to 100 increases alignment speed significantly while having little impact on alignment accuracy.

Setting a Minimum Cluster Size

A new feature that dramatically improves alignment speed with little impact on accuracy is the alignment candidate threshold (**-act**) (Figure 2.5). Normally all clusters are submitted for Smith-Waterman alignment. Initially we imposed a criterium (**-dh**) that two consecutive hashes had to be clustered before being aligned. This double-hit mechanism ensured that fewer spurious hash hits in the reference sequence would cause a full alignment to be performed.

The *alignment candidate threshold* extends this double-hit idea. An alignment candidate is simply the set of all seeds that form a cluster. The alignment candidate size is the length from the first base in the cluster to the last base in the cluster.

e.g. If a hash size of 11 is used and two hashes form cluster separated by a SNP, the alignment candidate size is 23. If the act paramater had been set to **-act 20**, this read would be submitted for Smith-Waterman alignment.

Uniqueness and Filters

MosaikAligner has a strict, but simple definition of what makes a unique read different from a non-unique read. If the read aligns to more than one location according to the user criteria, the alignments are non-unique. If the read aligns to a single location according to the user criteria, the alignment is unique.

MOSAİK offers two different alignment modes: unique and all (Figure 2.6). When aligning in the unique mode (**-m unique**), alignment stops as soon as there is evidence that the read can be aligned to two places. Both alignments are saved to the alignment file which can be handled by downstream applications. When aligning in the all mode (**-m all**), the

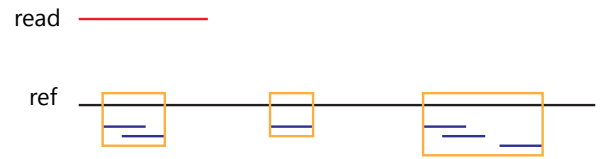


Figure 2.5 Alignment Candidate Threshold.

Before performing a pairwise alignment, MOSAİK hashes up each read and retrieves the hash positions for each seed in the reference sequence.

In this figure the hash positions are depicted as the horizontal blue lines. When we cluster these hash positions, three clusters are formed (orange boxes). Each cluster represents an alignment candidate that will be pairwise aligned.

Perhaps the middle cluster represents a spurious hit due to the repetitive nature of the reference sequence. To prevent spurious hits from eating up processing cycles, we can enforce that each cluster must have a specified length (the alignment candidate threshold) before being pairwise aligned.

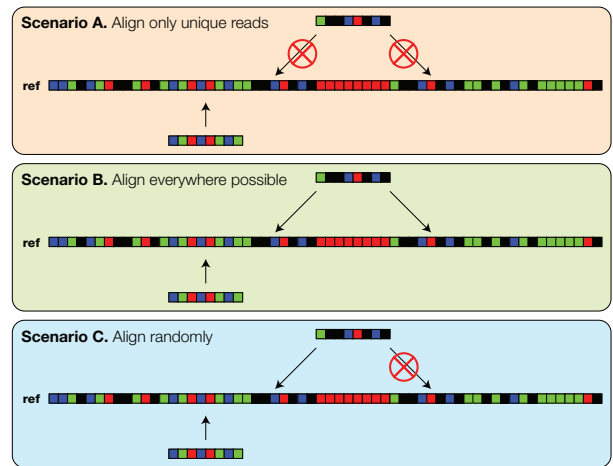


Figure 2.6 Alignment Mode.

MOSAİK has the ability to either place only the uniquely aligned reads (**-m unique**) or to place all reads (**-m all**). Other aligners choose a random location when placing non-unique reads, this benefit of this scenario confuses the author and therefore is not supported.

aligner finds all possible alignments. The all mode is suggested when trying to resolve paired-end reads.

When filtering reads, the mismatch threshold can be either a maximum number of mismatches allowed (**-mm**) or a maximum percentage of mismatches (**-mmp**) (w.r.t. read length).

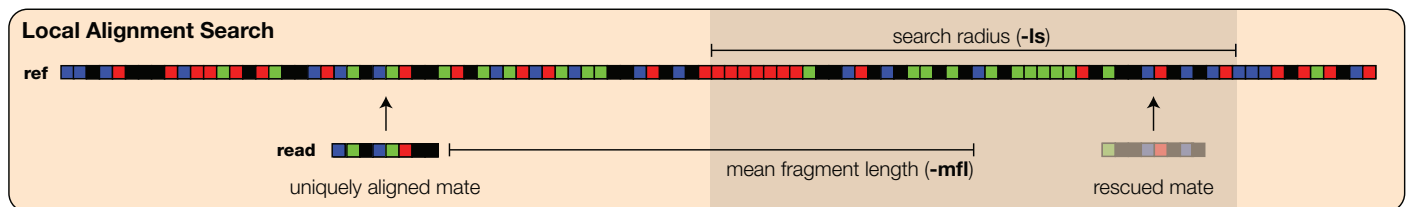
Since Smith-Waterman is a local alignment algorithm, undesirable subsequences may be aligned. Normally, unaligned portions of the read are counted as mismatches, but this behavior can be disabled with the **-mmal** parameter. When using this parameter, it may be useful to force the alignments to have a minimum length with respect to the original read length. This threshold can be set with the **-minp** parameter.

We regularly allow four mismatches (**-mm 4**) for 36 bp Illumina reads and allow up to 5 % mismatches (**-mmp .05**) with variable length read technologies such as Helicos and 454.

Local Alignment Search

When aligning paired-end/mate-pair libraries, sometimes one mate will align uniquely, but the other mate will end up in a highly repetitive region (like an ALU) resulting in thousands of potential locations. Depending on the chosen alignment parameters, the proper location that resolves with the unique mate might not be found.

To handle this situation, MOSAİK includes a local alignment search (**-ls**) parameter (Figure 2.7). This parameter uses the mean fragment length (**-mfl**) setting that was recorded when the read archive was created. The **-ls** parameter specifies a search radius relative to the mean fragment length.



For example, if the mean fragment length was set to 200 bp and the search radius is set to 100 bp - the alignment algorithm will search 100 - 300 bp from the unique mate for an alignment that conforms to the proper paired-end/mate-pair orientation and ordering.

Figure 2.7 Local Alignment Search.

When aligning paired-end/mate-pair reads, MOSAİK now has the ability to locally search for a missing mate within a user-specified search radius.

Alignment Qualities

Quality scores are calculated for each alignment in MOSAIK (Figure 2.8). Similar to base qualities, alignment qualities give the probability that a read has been misaligned.

Alignment qualities occur on a logarithmic scale of 0 to 99. An alignment quality of 20 indicates that there is a 1 % chance that the alignment was misaligned, whereas an alignment quality of 30 indicates that there is a 0.1 % chance of misalignment.

In an effort to create **highly accurate alignment qualities**, we used logistic regression to characterize reads that have a substitution error model (Illumina/AB SOLiD) and reads that have an insertion/deletion error model (Roche 454). These models use reference sequence length, read length, a weighted mismatch metric, and information content as predictors.

A high complexity 36 bp Illumina read aligned to the human genome with zero mismatches, will receive an alignment quality of around 57. The predictors affect the alignment quality in the following manner:

- As reference sequence length ↑, alignment quality ↓
- As # of mismatches ↑, alignment quality ↓
- As read length ↑, alignment quality ↑
- As information content ↑, alignment quality ↑.

At the moment, alignment qualities are not adjusted in resolved paired-end/mate-pair reads (some aligners sum up the two alignment qualities). However, the qualities are penalized if the chosen alignment parameters prevent the aligner from discovering all of the possible alignments for any given read.

Handling Colourspace

Despite the virtues of colourspace, the fact remains that many tools work only in basespace. Previously, MOSAIK has always stored the alignments in colourspace. In this version, alignments are converted seamlessly back to basespace directly after the Smith-Waterman algorithm.

Currently the algorithm used to translate dibase qualities involves taking the minimum of the two qualities that overlap a nucleotide in basespace.

There are a couple of advantages to handling the conversion immediately after alignment:

1. **Conventional alignment parameters** can be used. Allowing two mismatches means two bases that differ from the reference in basespace. In colourspace, two mismatches could mean one SNP or two sequencing errors.

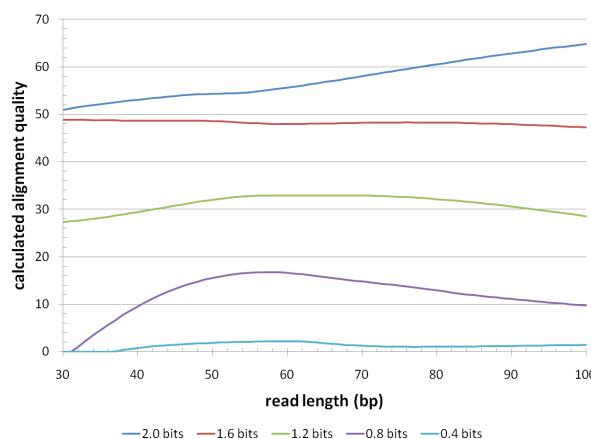
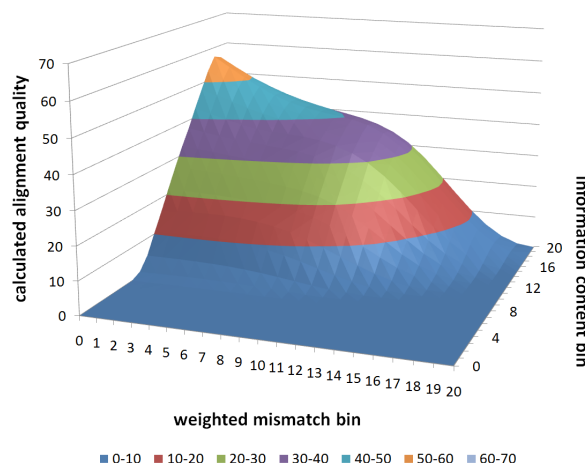


Figure 2.8 Alignment Qualities.

The graph above shows how alignment qualities vary in the Illumina logistic regression model when aligning to the full human genome. The number of mismatches was locked at zero, but the information content varied from 0.4 bits to 2.0 bits. For each information content curve, the alignment quality is plotted on a read length sweep from 30 - 100 bp.

The three dimensional surface plot below depicts the alignment qualities in the Illumina logistic regression model when aligning to the full human genome. The read length was locked at 36 bp, but the information content and mismatch metric were varied.



2. Post-alignment tasks are all **consistently performed in basespace**. Data can be **easily exported** into various file formats. For example, the BAM file format cannot handle colorspace sequences natively.

How do you use MOSAİK in colorspace?

1. Use MosaikBuild to import your reads into a read archive.
2. Use MosaikBuild to create a reference sequence archive (this will naturally be in basespace).
3. Use MosaikBuild to create a reference sequence archive in colorspace using the **-cs** parameter.
4. Use MosaikAligner to align the read archive using the colorspace reference archive with the **-ia** parameter and the basespace reference archive with the **-ibs** parameter.

Using MosaikAligner:

```
MosaikAligner -in myreads.dat -out h_sapiens_aligned.
dat -ia h.sapiens.dat -hs 15 -mm 4 -mhp 100 -act 20 -j
h.sapiens_15 -p 10
```

Here we specify an input read file (**-in**) and an output alignment file (**-out**) that stores all of the alignments. Additionally we specify a binary reference sequence file (**-ia**). All of the aforementioned parameters are required.

A hash size of 15 was specified (**-hs**) and a maximum of 4 mismatches is allowed (**-mm**).

A jump database (**-j**) will be used instead of the normal hash map. All hash positions are initially stored by the database, but only 100 random hash positions will kept for each seed (**-mhp**). In each seed cluster, a minimum length of 20 bp is required (**-act**).

A total of 10 processors (**-p**) will be used to increase alignment speed.

For some additional performance, try out **the new banded Smith-Waterman algorithm** (**-bw**) designed by Wan-Ping Lee. A bandwidth of 13 works well for 36 bp Illumina reads, 29 works well for 76 bp Illumina reads, and 51 works well for Roche 454 Titanium reads.

Box 2.2 MosaikAligner Parameters.

Input & Output

-in	specifies the input read archive.
-out	specifies the output alignment archive.
-ia	specifies the input reference sequence archive.
-ibs	specifies the basespace input reference sequence archive when aligning SOLiD reads.
-rur	stores unaligned reads in a FASTQ file.

Essential Parameters

-m	specifies the alignment mode: unique or all. Default: all.
-hs	specifies the hash size [4 - 32]. Default: 15.
-p	uses the specified number of processors. Default: 1.
-bw	uses the banded Smith-Waterman algorithm for increased performance.

Filtering

-act	specifies the alignment candidate threshold.
-mm	specifies the number of mismatches allowed. Default: 4.
-mmp	specifies the maximum percentage of the read length that are allowed to be errors. [0.0 - 1.0]
-mmal	uses the aligned read length instead of the original read length when counting errors.
-minp	specifies what minimum percentage of the read length should be aligned. [0.0 - 1.0]
-mhp	specifies the maximum number of hash positions to be used per seed.

Jump Database

-j	specifies the jump database filename stub.
-----------	--------------------------------------------

2.3. MosaikSort

MosaikSort takes the alignment output and prepares it for multiple sequence alignment. For single-ended reads, MosaikSort simply resorts the reads in the order they occur on each reference sequence.

For mate-pair/paired-end reads, MOSAİK resolves the reads according to user-specified criteria before resorting the reads in the order they occur on each reference sequence.

Paired-End Read Resolution

When both mate-pairs are available, MosaikSort first samples the fragment lengths from uniquely aligned reads. Using this information, a minimum and maximum fragment length is calculated using the empirical 99.73% confidence interval.

When one mate is unique and the other is non-unique, the paired-end read is resolved if the program finds at most one alignment that occurs within the desired confidence interval (Figure 2.9). Similarly when both mate-pairs are non-unique, the paired-end read is resolved if the program finds at most one combination of alignments that fit the confidence interval.

Box 2.2 MosaikAligner Parameters Continued.

Pairwise Alignment Scores

- ms** the match score. Default: 10.
- mms** the mismatch score. Default: -9.
- gop** the gap open penalty. Default: 15.
- gep** the gap extend penalty. Default: 6.66.
- hgop** the gap open penalty used in homopolymer stretches when aligning with 454 reads. Default: 4.

Paired-end Terminology

We use the term **mate** to mean one of the reads at each end of the fragment. In MosaikAligner, each mate is aligned separately. Sometimes one mate aligns, but the other does not. We use the term **orphaned reads** to describe such an event.

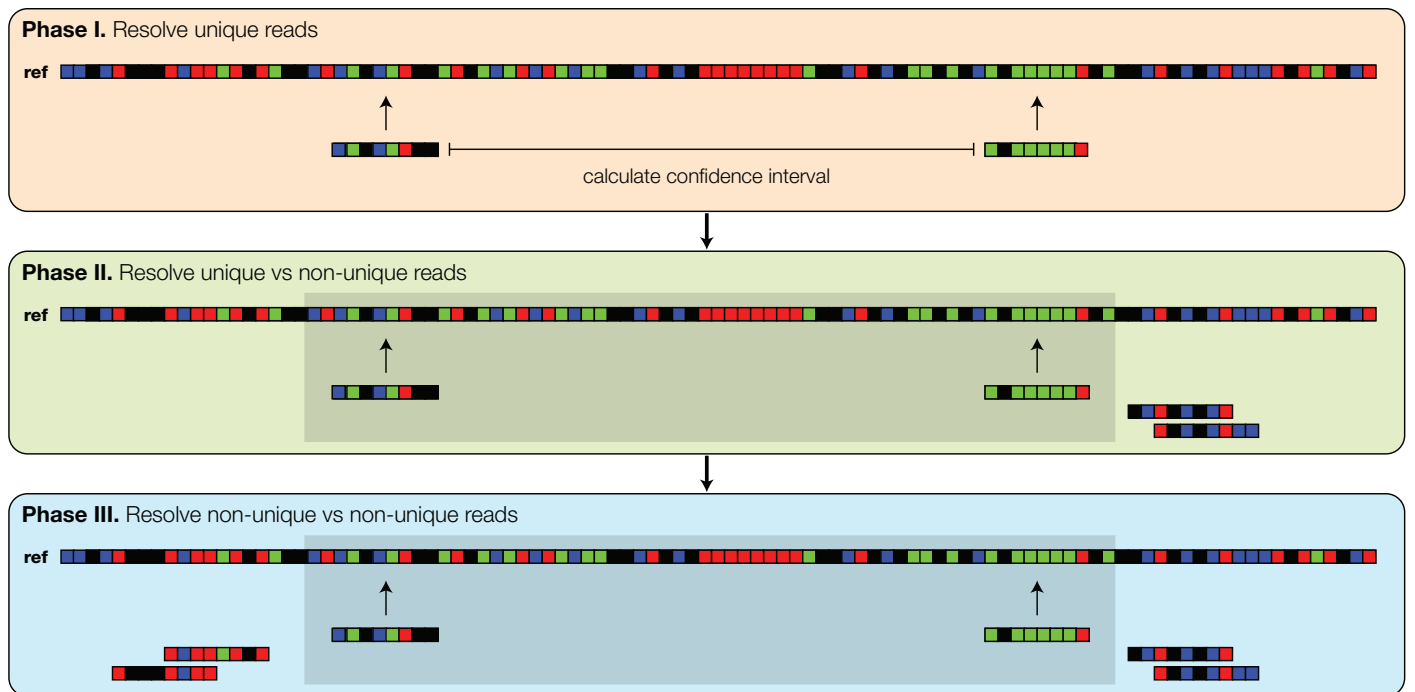


Figure 2.9 Paired-End Read Resolution

Read Name Mangling

Some downstream utilities (such as `consed`) require that each alignment has a unique name. When the `-consed` option is used, the following adjustments are made to read names:

- when a single-ended read aligns to multiple locations, a period is appended along with the current alignment number. e.g. `IL7_1_1_203_187.27` would indicate the 27th alignment for this read
- when a paired-end/mate-pair read is resolved, a forward slash is appended followed by the mate number. e.g. `IL7_1_1_203_187/1` indicates the first mate
- orphaned paired-end/mate-pair read names are not modified.

Duplicate Removal

When used in conjunction with `MosaikDupSnoop`, `MosaikSort` can remove duplicates (`-dup`) with respect to the originating sequencing library.

Using MosaikSort:

```
MosaikSort -in yeast_aligned.dat -out yeast_sorted.dat
```

In the example above we specify an input alignment file (`-in`) containing paired-end reads and a sorted output alignment file (`-out`).

Box 2.3 MosaikSort Parameters.

Input & Output

- `-in` specifies the input alignment archive.
- `-out` specifies the output alignment archive.
- `-dup` enables duplicate filtering with databases in the specified directory.
- `-mem` specifies how many alignments to cache. Default: 6,000,000.

Single-end Options

- `-nu` include non-unique reads.

Paired-end Options

- `-afl` allows all fragments lengths when evaluating unique read pairs.
- `-ci` sets the fragment length confidence interval. Default: 0.9973.
- `-sa` samples fragment lengths from all unique read pairs.

Paired-end Resolution

- `-iuo` ignore unique orphaned reads.
- `-iuu` ignore unique vs unique read pairs.
- `-ium` ignore unique vs multiple read pairs.
- `-rmm` resolve multiple vs multiple read pairs.

2.4. MosaikMerge

MosaikMerge is not normally a part of the MOSAİK pipeline. It takes multiple sorted alignment archives from MosaikSort and merges them into a single alignment archive that can be processed by MosaikAssembler, MosaikText, or MosaikCoverage.

When aligning multiple runs belonging to different sequence libraries, the best practice is to use MosaikSort on each lane or run and then combine them with MosaikMerge. This practice avoids the unintended consequences of attempting to resolve paired-end reads when multiple fragment length distributions are present.

Another common use of MosaikMerge is to combine runs from different read technologies. This allows the researcher to combine 454, Illumina, Helicos, SOLiD, and Sanger capillary technologies into one co-assembly.

Using MosaikMerge:

```
MosaikMerge -in 454_aligned.dat -in helicos_aligned.dat  
-in illumina_alignments/ -out coassembly.dat
```

In this example we combine the alignments from the two files, *454_aligned.dat* and *helicos_aligned.dat*, with all of the alignments contained in the directory *illumina_alignments*. The merged output will be saved in *coassembly.dat*.

Box 2.4 MosaikMerge Parameters.

Options

- | | |
|-------------|---------------------------------------------------------------------------------------------------------------------------------------------|
| -in | specifies either an input alignment archive or a directory containing alignment archives. <i>This parameter can be used multiple times.</i> |
| -out | specifies the merged output alignment archive. |
| -mem | specifies how many alignments to cache. Default: 6,000,000. |

2.5. MosaikAssembler

MosaikAssembler takes the sorted alignment file and produces a multiple sequence alignment which is saved in an assembly file format. At the moment, MosaikAssembler saves the assembly in the phrap ace format and the GigaBayes gig format.

MosaikAssembler has been completely rewritten to take advantage of the sorted alignment files. As a result, assembly file creation is now orders of magnitude faster than before and limited only by sequential hard disk transfer speeds rather than slower random-access transfer speeds.

By default MosaikAssembler will assemble each reference sequence where reads aligned. Since the sorted alignment archives incorporate an index, a specific reference sequence can be assembled quickly with the region of interest (**-roi**) parameter.

Using MosaikAssembler:

```
MosaikAssembler -in yeast_sorted.dat -ia yeast.dat -out  
yeast_assembly -roi chrX
```

In this example, a sorted alignment file (**-in**) and a binary reference sequence file (**-ia**) are specified as input. MosaikAssembler will use the provided filename stub (**-out**) when generating an assembly specifically for the X chromosome (**-roi**).

Box 2.5 MosaikAssembler Parameters.

Input & Output

-in	specifies the input alignment archive.
-out	specifies the assembly output filename stub.
-ia	specifies the input reference sequence archive.
Options	
-f	specifies the assembly file format: ace or gig. Default: ace.
-roi	specifies the name of the reference sequence to assemble.

3. Utilities

MOSAİK contains four additional utility programs to facilitate data analysis: MosaikText, MosaikCoverage, MosaikDupSnoop, and MosaikJump.

3.1. MosaikText

MosaikText converts alignments to different text-based formats. Currently it supports the BLAT axt (**-axt**) format, the BAM format (**-bam**), the UCSC Genome Browser bed format (**-bed**), the SAM format (**-sam**), and the Illumina ELAND (**-eland**) format. Alternatively alignments can be dumped directly to the screen (**-screen**) in an axt-like format.

In addition to examining alignment archives, MosaikText supports dumping the contents of read archives.

Using MosaikText:

```
MosaikText -in yeast_aligned.dat -bam yeast_aligned.bam
```

In this example, the supplied alignment archive (**-in**) is exported to a BAM file (**-bam**).

3.2. MosaikCoverage

MosaikCoverage is a handy program for investigating representational bias. This utility parses the alignment file and produces a base-accurate coverage plot (no binning) for each reference sequence that has coverage. In addition to the coverage plot, a simple space delimited coverage file is produced (Figure 3.1). MosaikCoverage uses *gnuplot* to generate the coverage graphs in PostScript. If available, it will call *ps2pdf* to convert the graphs into more portable pdf files.

When running MosaikCoverage, output similar to this will be seen:

```
- calculating coverage statistics:
* coverage statistics for 22 (34851311 bp): 33062466 bp
(94.9 %), mean: 11.3x
```

In this example, the coverage for chromosome 22 is given. Chromosome 22 has 34 Mbp of sequence containing { A, C, G, T }. 33 Mbp (94.9 %) of the chromosome has at least 1x coverage and the mean coverage is 11.3x.

Box 3.1 MosaikText Parameters.

Read Archive

- ir** specifies the input read archive.
- fastq** stores the data in the specified FASTQ file.
- screen** displays the reads on the screen.

Alignment Archive

- in** specifies the input alignment archive.
- u** limits the output to show only uniquely aligned reads.
- axt** stores the data in the specified BLAT AXT file.
- bam** stores the data in the specified BAM file.
- bed** stores the data in the specified UCSC Genome Browser bed file.
- eland** stores the data in the specified Illumina ELAND file.
- sam** stores the data in the specified SAM file.
- screen** displays the alignments on the screen in an AXT-like format.

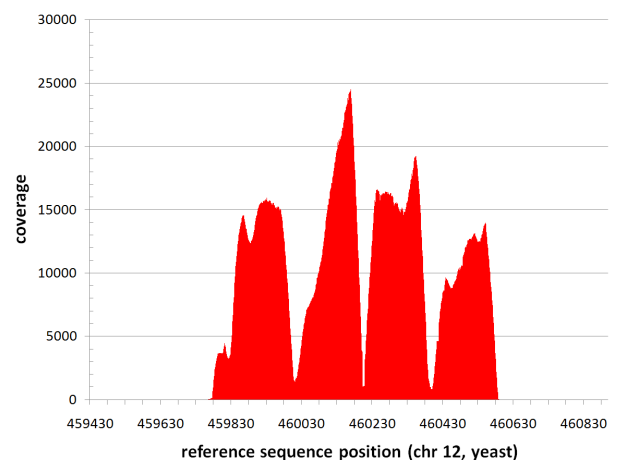


Figure 3.1 Using MosaikCoverage in Excel.

This graph depicts the observed coverage in *S. cerevisiae* chromosome 12 from an Illumina data set.

Using MosaikCoverage:

```
MosaikCoverage -in h_sapiens_aligned.dat -ia h.sapiens.dat
-u -od graphs -cg
```

In the example above, the coverage is calculated from all of the reads in the alignment file *h_sapiens_aligned.dat* (**-in**) for each sequence specified in the reference sequence file *h.sapiens.dat* (**-ia**). By enabling the **-u** parameter, only unique alignments will be used to calculate coverage. All resulting files will be placed in the output directory *graphs* (**-od**) and pdf graphs will be generated (**-cg**).

3.3. MosaikDupSnoop

MosaikDupSnoop inspects a specified set of alignment archives and stores the aligned fragment locations with respect to the sequencing library. If duplicates are encountered, the fragment with the highest alignment quality is recorded. When used in conjunction with MosaikSort (**-dup**), all of the **duplicate fragments with a lower alignment quality will be discarded**.

MosaikDupSnoop treats paired-end/mate-pair reads differently from single-ended reads. For paired-end reads, all reads that share the same start coordinate but have end coordinates that **differ by up to 2 bp** will be treated as duplicates. Likewise reads that share the same end coordinate but have start coordinates that differ by up to 2 bp are also considered duplicates. For single-ended reads, only reads that share the exact start and end coordinates are considered duplicates.

Using MosaikDupSnoop:

```
MosaikDupSnoop -in yeast_aligned.dat -od fragData/
```

In the example above, the fragments in *yeast_aligned.dat* (**-in**) will be recorded in sequencing library databases located in the *fragData* directory (**-od**).

3.4. MosaikJump

MosaikJump is a tool that converts a reference sequence archive into a “jump database”. As discussed earlier in section 2.2, the jump database is a custom data structure that separates seeds and hash positions in a manner that is fast and collision-free. As such, it works as a drop-in replacement for the traditional hash maps that MOSAİK normally uses.

Box 3.2 MosaikCoverage Parameters.

Input & Output

-in	specifies the input alignment archive.
-ia	specifies the input reference sequence archive.
-u	limits the output to show only uniquely aligned reads.
-od	specifies the output directory.
-cg	creates coverage graphs if <i>gnuplot</i> is installed.

Box 3.3 MosaikDupSnoop Parameters.

Input & Output

-in	specifies the input alignment archive.
-od	specifies the output directory.

Paired-end Options

-afl	allows all fragments lengths when evaluating unique read pairs.
-ci	sets the fragment length confidence interval. Default: 0.9973.
-iuo	ignore unique orphaned reads.
-iuu	ignore unique vs unique read pairs.
-ium	ignore unique vs multiple read pairs.
-rmm	resolve multiple vs multiple read pairs.

An added benefit of using the jump database is that the user can choose which components should remain on disk and which should be loaded into memory. Maximum performance is realized when the entire jump database is loaded into memory.

The *positions file* is directly proportional to the number of hash positions in the reference sequences as well as the number of seeds present. e.g. Using a hash size of 15, there are roughly 2.9 billion hash positions in the human genome and 550,000 unique seeds. Assuming an integer is used to store each seed grouping and hash position, the positions file will be 12.7 GB.

The *keys file* stores the offsets for each seed grouping in the *positions file*. The *keys file* grows exponentially with each increasing hash size. e.g. Using a hash size of 15, the required file size is 5 GB ($5 * 4^{15}$). However if a hash size of 13 is used, the resulting file is only 313 MB. Due to the exponentially growing file size, the upper practical bounds for the jump database is a hash size of 15.

Using MosaikJump:

```
MosaikJump -ia h.sapiens.dat -out h.sapiens_15 -hs 15
```

In the example above, MosaikJump will use the reference sequence file (**-in**) to generate the jump database files for hash size 15 (**-hs**) that start with the specified filename stub (**-out**). To improve performance during jump database creation, the keys file will be kept in memory (**-km**).

Box 3.4 MosaikJump Parameters.

Input & Output

- ia** specifies the input reference sequence archive.
- out** specifies the output filename stub for the jump database.

Options

- mem** specifies the amount of RAM (GB) to use when sorting the hashes.
- hs** specifies the hash size [4 - 32].
- mhp** specifies the maximum number of hash positions to be used per seed.

4. Understanding Program Output

A lot of effort has been put into making the program output as easy to understand as possible. Sometimes the items that seem simple and obvious to us, may seem foreign and vague to others. To help bridge that gap, we will discuss what the values actually mean in this section.

4.1 MosaikBuild

```

mikaels@snp:~
-----
MosaikBuild 1.0.1307                               2009-10-14
Michael Stromberg                                 Marth Lab, Boston College Biology Department
-----
- setting center name to: WTSI
- setting read group ID to: ERR001719
- setting sample name to: NA12878
- setting sequencing technology to: illumina
- trimming leading and lagging N's. Mates with >4 interior N's will be deleted.

- parsing paired-end/mate-pair FASTQ files:
reads: 6,563,762 (52,682.9 reads/s)

Filtering statistics:
-----
# mates deleted:      152714 ( 1.2 %)
# lagging N's trimmed: 2974
-----
# reads written:      6487405
# bases written:      467090186

MosaikBuild CPU time: 124.310 s, wall time: 124.620 s
[mikaels@snp ~]$

```

Handling N's

By default MosaikBuild will trim off the N's that occur in the beginning or at the end of each read.

If more than a specified number of N's occur internally within the read, that read will be discarded. During alignment, an N base will always count as a mismatch. Therefore it's wise to filter out any reads with more N's than your maximum number of allowed mismatches.

Keeping Time

All MOSAİK tools will report both a **CPU time** and a **wall time** when finished. CPU time reflects the aggregate time spent on all processors while wall time reflects time passed in the real world. e.g. a disk intensive tool will likely have a longer wall time than CPU time. In contrast, a CPU intensive task using 8 processors for 1 minute will return a CPU time of 480 seconds and a wall time of 60 seconds.

4.2 MosaikAligner

```

mikaels@snp:~/UnifiedRelease
-----
MosaikAligner 1.0.1307                               2009-10-14
Michael Stromberg & Wan-Ping Lee  Marth Lab, Boston College Biology Department
-----
- Using the following alignment algorithm: all positions
- Using the following alignment mode: aligning reads to all possible locations
- Using a maximum mismatch threshold of 4
- Using 8 processors
- Using an alignment candidate threshold of 20bp.
- Setting hash position threshold to 100

Hashing reference sequence:
100%[=====] 3,122,914 ref bases/s      in 1 s

- Loading reference sequence... finished.

Aligning read library (1217876):
100%[=====] 5,998.3 reads/s      in 03:23

Alignment statistics (mates):
-----
# failed hash:      33020 ( 1.4 %)
# filtered out:    140027 ( 5.7 %)
# unique:          1930894 ( 79.3 %)
# non-unique:      331811 ( 13.6 %)
-----
total:              2435752
total aligned:     2262705 ( 92.9 %)

Alignment statistics (reads):
-----
# unaligned:       14513 ( 1.2 %)
# orphaned:        144021 ( 11.8 %)
# both mates unique: 854007 ( 70.1 %)
# one mate non-unique: 133466 ( 11.0 %)
# both mates non-unique: 71869 ( 5.9 %)
-----
total reads:       1217876
total reads aligned: 1203363 ( 98.8 %)

Miscellaneous statistics
-----
aligned mate bp:   79194675
alignment candidates/s: 125683.8

MosaikAligner CPU time: 1560.880 s, wall time: 207.049 s
    
```

Alignment Statistics (mates)

This section of the program output displays the statistics with respect to the individual sequences. Each paired-end read will usually have two mates sequences associated with it, whereas single-ended reads only have one sequence per read.

Failed Hash

A failed hashing attempt simply means that **a hash size of contiguous bases that match the reference couldn't be found**. e.g. if the user selects a hash size of 15, but the read has a sequencing error every 12 bases; the aligner will fail to successfully seed the alignment. *If this percentage is too high, consider reducing the hash size.*

Filtered Out

The user has a lot of flexibility in specifying how to **filter out good alignments from bad alignments**. The most common used filter is the maximum number of allowed mismatches. e.g. If the user specifies that up to 4 mismatches should be allowed, but the best alignment results in 5 mismatches; the read will be filtered out. *If this percentage is too high, consider relaxing the filter criteria.*

Unique vs Non-Unique

An read is considered to be aligned uniquely if it can be **aligned to only one location in the reference given the current alignment parameters**. Everything else is considered non-unique.

Alignment Statistics (reads)

This section of the program output is shown when aligning paired-end/mate-pair reads. Orphaned reads occur when only one of the mate sequences align to the reference. In addition to the unaligned and orphaned reads, statistics are also shown for various uniqueness conditions. *Reads where both of the mate sequences are unique have a higher efficiency at being resolved than reads where one or both mates are non-unique.*

Total Reads Aligned

You may have noticed that the “total reads aligned” statistic in the reads section is different than the one in the mates section. *In the paired-end section, a read is considered to be aligned if at least one of the mate sequences is aligned to the reference.* In the mates section, the statistic reflects all of the mate sequences in the entire data set.

Alignment Candidates/s

In MOSAİK, an alignment candidate refers to any region that has been clustered after hashing that meets the minimum criteria for pairwise alignment. This statistic gives the user a sense of *how many Smith-Waterman alignments are occurring per second.*

4.3 MosaikSort

```

mikaels@snp:~/UnifiedRelease
-----
MosaikSort 1.0.1307                               2009-10-14
Michael Stromberg                               Marth Lab, Boston College Biology Department
-----
- resolving the following types of read pairs: [unique orphans] [unique vs unique] [unique vs multiple]
- resolving paired-end alignments

- phase 1 of 3: building fragment length distribution:
samples: 1,000,000 (60,510.7 samples/s)

- phase 2 of 3: resolve read pairs:
100%[=====] 47,853.1 reads/s           in 25 s

- phase 3 of 3: sort resolved read pairs:
100%[=====] 592,239.7 alignments/s     in 3 s

Paired-end read statistics:
=====
                        original          resolved
-----
# orphaned:             144021 (12.0 %)    126353 (10.5 %)
# both mates unique:   854007 (71.0 %)    850518 (70.7 %)
# one mate non-unique: 133466 (11.1 %)    125390 (10.4 %)
# both mates non-unique: 71869 ( 6.0 %)      0 ( 0.0 %)
-----
total:                  1203363                1102261 (91.6 %)

Fragment statistics:
=====
min target frag len:    90
median target frag len: 149
max target frag len:    209

MosaikSort CPU time: 60.880 s, wall time: 62.154 s
    
```

Fragment Statistics

Many aligners use a user-specified 3σ standard deviation or a median absolute deviation approach when calculating the endpoints for the fragment length distribution. Often these approaches rely on approximate values and are prone to outlier bias. *MOSAİK calculates the values in realtime based on the empirical fragment length distribution of unique vs unique read pairs.*

Paired-end Resolution Efficiency

For each of the paired-end categories (orphaned, unique vs unique, unique vs multiple, and multiple vs multiple), the original number of reads will be displayed as well as the number that were successfully resolved.

During paired-end resolution, **reads are checked for the proper order, orientation, and that one and only one mate is located within the fragment length confidence interval.** Orphaned reads are resolved if the proper order and orientation is consistent.

For example, 71.0 % of the aligned paired-end reads were in the category “both mates unique”. After paired-end resolution, that percentage drops to 70.7 %. *i.e. 0.3 % of the reads failed the aforementioned filtering criteria.*

0.7 % of the “both mates non-unique” failed the filtering criteria.

By default, paired-end resolution of reads in the category “both mates non-unique” is disabled due to the potential source for misalignment error. *This category can be enabled using the **-rmm** option.*

5. Performance

All statistics were measured on a lightly loaded computer with two 3.0 GHz Intel Xeon X5450 quadcore processors and 64 GB memory.

There is a **direct correlation between MOSAİK alignment speed and the aggregate length of the reference sequences** (the target genome). MOSAİK is extremely fast when processing yeast or roundworm-sized genomes and is faster than many aligners when processing mammalian genomes.

5.1. Speeding Up Alignments

There are a number of steps you can take to increase alignment speed.

Algorithmically speaking, the parameters that will help increase alignment speed are the following: hash size (**-hs**), alignment candidate threshold (**-act**), the maximum hash position threshold (**-mhp**), and the banded Smith-Waterman bandwidth (**-bw**). The crux is that you want to maximize the improvement in speed while minimizing any negative impact on alignment accuracy and the percentage of reads aligned.

The larger the hash size, the more likely that hash will be unique in the reference sequence. Therefore increasing the hash size, reduces the number of spurious hash hits in the reference sequence. The consequence of increasing the hash size is that it also increases the number of bases that must exactly match the reference sequence. Our lab uses a hash size of 15 for most of our current analysis projects.

The alignment candidate threshold effectively dictates the minimum size of hash clusters before being submitted for pairwise alignment. Increasing this threshold reduces the number of spurious alignments being aligned. The consequence of increasing the threshold is that a read with a combination of sequencing errors and polymorphisms may also be filtered if the threshold is not met.

The maximum hash position threshold speeds up the alignment by reducing the number of hash positions that need to be clustered. In the human genome, each seed has an average of 5.25 hash positions when using hash size 15. Setting the threshold to 100 in such a case has minimal affect on the alignments while increasing the alignment speed considerably. The consequence of setting the threshold too low is that alignments in highly repetitive regions might no longer be seeded.

One of the new improvements in MOSAİK is a **banded Smith-Waterman alignment algorithm**. As the bandwidth is reduced, alignment performance will improve - especially in longer Roche 454 or Sanger capillary reads. If the bandwidth is reduced too much, alignment artifacts may result. In Box 5.1 we have recorded the bandwidths used in our own projects.

Perhaps the most obvious way to speed up alignments is to use more processors. MosaikAligner is fully multi-threaded and can handle any number of processors. The caveat is that the input/output and memory bandwidth may become saturated. e.g. On some of our 16 core systems, we have discovered that using 12 cores is usually faster than using 16 processor cores. You may have to perform a bioinformatics titration experiment where you align a subset of reads with an increasing number of processor cores in order to choose the configuration that works best for your system.

5.2. Yeast Alignments

MosaikBuild converted a SOLiD run of 61,516,412 reads (35 bp) in 7.4 minutes (139,000 reads/s).

MosaikAligner aligned the run in **9.5 minutes (109,000 reads/s)** to the entire *pichia stipitis* genome (15.4 Mbp) using 8 processor cores. We used the hash size 13 jump database for this experiment, so our total memory usage was 468 MB. 57.1 % of the reads aligned in this unfiltered data set.

MosaikSort sorted the uniquely aligned reads in **3.0 minutes**.

MosaikAssembler produced a GigaBayes assembly file with 27,592,793 Illumina alignments in **64 seconds (431,000 reads/s)**.

5.3. Roundworm Alignments

MosaikBuild converted an Illumina run of 11,386,260 reads (35 bp) in 38 seconds (308,000 reads/s).

MosaikAligner aligned the run in **4.8 minutes (40,000 reads/s)** to the entire *C. elegans* genome (100 Mbp) using 8 processor cores. We used the hash size 15 jump database for this experiment, so our total memory usage was 5.8 GB. 91.6 % of the reads aligned in this data set.

MosaikSort sorted the output in **1.6 minutes**.

MosaikAssembler produced a GigaBayes assembly file with 12,323,293 Illumina alignments in **31 seconds (398,000 reads/s)**.

Box 5.1 Banded Smith-Waterman Parameters

Illumina (36 - 43 bp)	-bw=13
Illumina (44 - 63 bp)	-bw=17
Illumina (63+ bp)	-bw=29
Roche 454 (Titanium)	-bw=51

Box 5.2 Yeast Alignment Parameters

```
-m unique -hs 13 -act 20 -mm 6 -mhp 100 -p 8
-j jumpdb/p.stipitis_13cs
```

Box 5.3 Roundworm Alignment Parameters

```
-m unique -act 20 -mm 4 -mhp 100 -p 8
-j jumpdb/c.elegans_15
```

5.4. Human Alignments

MosaikBuild converted a lane of 6,563,762 Illumina paired-end reads (36 bp) in 2.9 minutes (53,000 reads/s).

MosaikAligner aligned the lane in **45.9 minutes (2,390 paired-end reads/s)** to the entire human genome (2.9 Gbp) using 8 processor cores. We used the hash size 15 jump database for this experiment, so our total memory usage was 20 GB. 95.1 % of the mates aligned in this data set.

MosaikSort resolved the paired-end reads and sorted the output in **37 minutes** on a slower computer (2.0 GHz AMD Opteron 270).

MosaikAssembler produced a GigaBayes assembly file with 2,589,439,852 Illumina mates in **4.1 hours (174,000 mates/s)**.

Box 5.4 Human Alignment Parameters

```
-act 20 -mm 4 -mhp 100 -p 8  
-j jumpdb/h.sapiens_15
```

5.5. Aligner Settings By Sequencing Technology

By analyzing the performance characteristics and alignment accuracy, we have compiled a list of aligner settings (Box 5.5) that we typically use for each sequencing technology.

Box 5.5 Aligner Settings.

454 GS20 & FLX

```
-hs 15 -mm 0.05 -act 26
```

454 Titanium

```
-hs 15 -mm 0.05 -act 55
```

Illumina GA1 36 bp

```
-hs 15 -mm 4 -act 20
```

Illumina GA2 51 bp

```
-hs 15 -mm 6 -act 25
```

Illumina GA2 76 bp

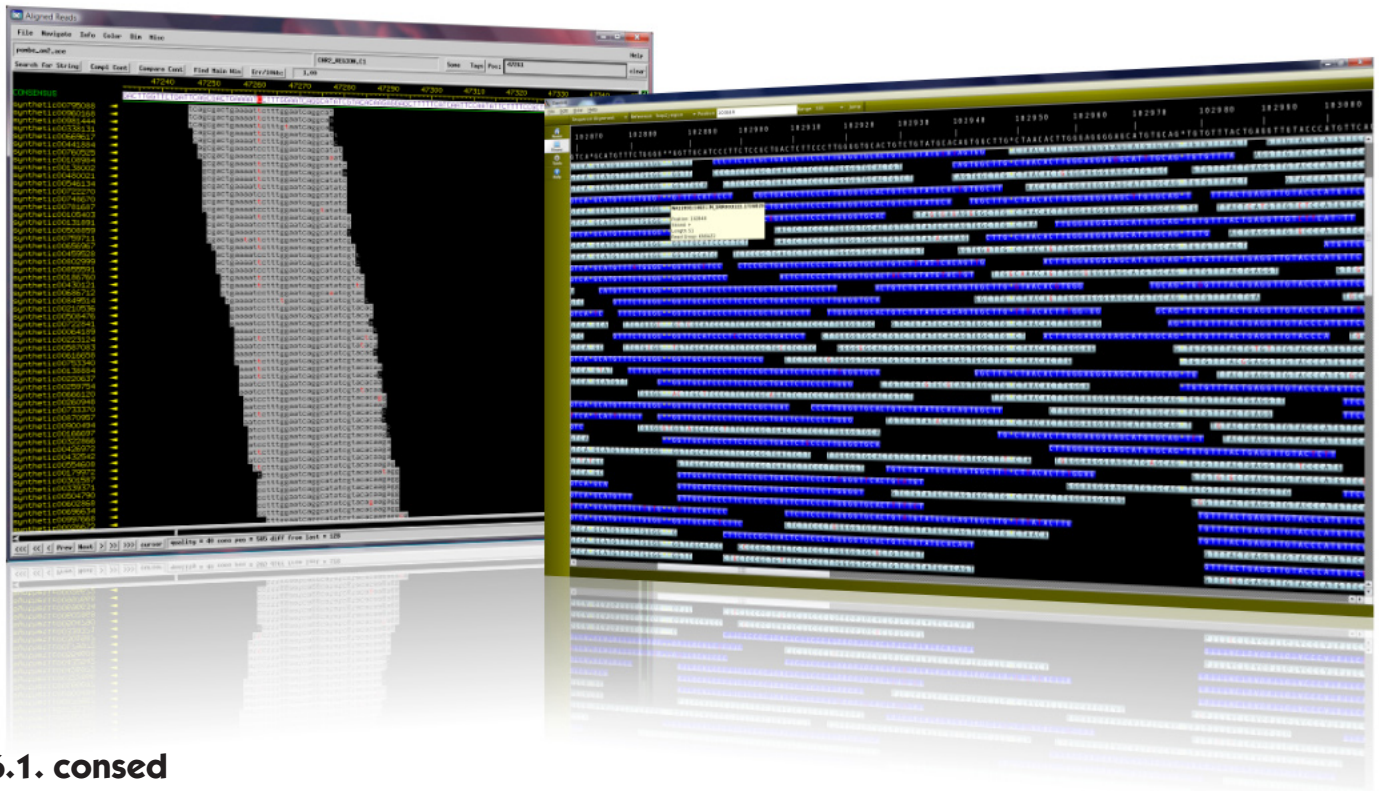
```
-hs 15 -mm 12 -act 35
```

SOLiD 35 bp

```
-hs 15 -mm 4 -act 20
```

6. Visualization

MOSAIK works well with two visualization programs: *consed* and *Gambit*. *Consed* reads ACE assembly files while *Gambit* reads BAM files. An ace file is easily created by running the MosaikAssembler with the *-face* parameter on your sorted alignment archive. BAM files can be exported from MosaikText using the *-bam* parameter.



6.1. consed

Consed (*left*) is available from Phil Green's webpage:

<http://bozeman.mbt.washington.edu/consed/consed.html>

When using *consed* with ace files created in MOSAIK, you should run it with the *-nophd* switch to prevent it from loading the non-existent phred sequence files.

6.2. Gambit

Gambit (*right*) is being developed in our lab by Derek Barnett. Compared with *consed*, *Gambit* shows a more compact view of the read assembly. *Gambit*'s strength is that it can show any number of annotation tracks and features a robust analysis plugin API to help provide context to the alignment data.

7. Data Access (MosaikTools)

A common problem when evaluating bioinformatics tools is data access. Users are often forced to write parsers for various file formats or output log files. MosaikText provides alignment data in several known text formats, but the MosaikTools API provides even faster and easier access to the aligner results.

The API is written in C++ and can be found in the **MosaikTools directory** contained in the MOSAIK package. By using the SWIG (<http://www.swig.org>) package, a version that is accessible in Perl is also available.

Here is a sample program that loads each read in a MOSAIK alignment archive:

```
#include <iostream>
#include "AlignmentReader.h"

using namespace std;

int main(int argc, char* argv[]) {

    // open the MOSAIK alignments file
    Mosaik::CAlignmentReader reader;
    reader.Open("myreads_aligned.dat");

    // get some basic statistics
    uint64_t numBases = reader.GetNumBases();
    uint64_t numReads = reader.GetNumReads();

    cout << "# of bases: " << numBases << endl;
    cout << "# of reads: " << numReads << endl;

    // keep reading all of the sequences
    Mosaik::AlignedRead ar;
    while(reader.LoadNextRead(ar)) {

        // do something with each read
        cout << "read name: " << ar.Name << endl;
    }

    // close the alignments file
    reader.Close();

    return 0;
}
```

The MosaikTools directory contains two sample programs: *MosaikReaderMain.cpp* & *MosaikConversionMain.cpp*. *MosaikReaderMain.cpp* loads and displays information about each alignment in an alignment archive. *MosaikConversionMain.cpp* loads each read from an alignment archive and saves those reads into a new alignment archive.

Box 7.1 Read Data Structures

The *AlignedRead* data structure contains all of the mate 1 and mate 2 alignments for one read. If the read is single-ended, only the *Mate1Alignments* vector will contain alignments.

```
struct AlignedRead {
    string Name;
    vector<Alignment> Mate1Alignments;
    vector<Alignment> Mate2Alignments;
};
```

In the *Alignment* data structure, *ReferenceBegin* and *ReferenceEnd* contain the start and end locations on the reference sequence. *ReferenceName* contains the name of the reference sequence.

QueryBegin and *QueryEnd* contain the start and end locations on the read.

The gapped pairwise alignment is contained in the strings *Reference* and *Query*.

Each character in the *BaseQualities* string contains a base quality for each pairwise aligned base in the same orientation as the alignment. Extra base qualities have NOT been introduced for gaps in the alignment. There is no offset between each character and the intended base quality.

Quality contains the alignment quality which measures the probability that an alignment has been misaligned.

IsReverseStrand is set to true if the alignment is on the 3' or Crick strand.

```
struct Alignment {
    unsigned int ReferenceBegin;
    unsigned int ReferenceEnd;
    unsigned short QueryBegin;
    unsigned short QueryEnd;
    unsigned char Quality;
    bool IsReverseStrand;
    char* ReferenceName;
    string Reference;
    string Query;
    string BaseQualities;
};
```

Appendix 1: Crib Sheet

MosaikBuild

MosaikBuild **-fr** myreads.fasta **-fq** myreads.fasta.qual **-out** myreads.dat

MosaikBuild **-fr** myreference.fasta **-oa** myreference.dat

- fr, -fr2** specifies the FASTA files containing the bases for the first mate (-fr) and the second mate (-fr2).
- fq, -fq2** specifies the FASTA files containing the base qualities for the first mate (-fq) and the second mate (-fq2).
- q, -q2** specifies the FASTQ files for the first mate (-q) and the second mate (-q2).
- cs** translates the reference sequence from base space to colorspace.
- st** specifies the sequencing technology: 454, helicob, illumina, sanger, solid.
- out** specifies the output file when converting reads.
- oa** specifies the output file when converting the reference sequence.

MosaikAligner

MosaikAligner **-in** myreads.dat **-out** myreads_aligned.dat **-ia** myreference.dat **-hs** 15 **-mm** 4 **-m** all
-mhp 100 **-act** 20 **-j** myjumpdb **-p** 8

- in** specifies the input read archive.
- out** specifies the output alignment archive.
- ia** specifies the input reference sequence archive.
- rur** stores unaligned reads in a FASTQ file.
- m** specifies the alignment mode: unique or all.
- hs** specifies the hash size [4 - 32]. Default: 15.
- p** uses the specified number of processors.
- act** specifies the alignment candidate threshold.
- mm** specifies the number of mismatches allowed. Default: 4.
- mmp** specifies the maximum percentage of the read length are allowed to be errors. [0.0 - 1.0]
- minp** specifies what minimum percentage of the read length should be aligned. [0.0 - 1.0]
- mmal** when enabled, unaligned portions of the read will not count as a mismatch.
- mhp** specifies the maximum number of hash positions to be used per seed.
- j** specifies the jump database filename stub.
- ls** enables local alignment search for paired-end/mate-pair reads.
- bw** uses the banded Smith-Waterman algorithm for increased performance.

MosaikSort

MosaikSort **-in** myreads_aligned.dat **-out** myreads_sorted.dat

- in** specifies the input alignment archive.
- out** specifies the output alignment archive.
- ci** sets the fragment length confidence interval. Default: 0.9973.
- dup** enables duplicate filtering with databases in the specified directory.

MosaikAssembler

MosaikAssembler **-in** myreads_sorted.dat **-ia** myreference.dat **-out** myassembly

- in** specifies the input alignment archive.
- out** specifies the assembly output filename stub.
- ia** specifies the input reference sequence archive.
- f** specifies the assembly file format: ace or gig. Default: ace.
- roi** specifies the name of the reference sequence to assemble.