
MinHash Alignment Process (MHAP) Documentation

Release 2.0

Sergey Koren and Konstantin Berlin

March 02, 2016

1	Overview	1
1.1	Installation	1
1.2	Quick Start	2
1.3	Utilities	5
1.4	Contact	6

Overview

MHAP (pronounced MAP) is a reference implementation of a probabilistic sequence overlapping algorithm. Designed to efficiently detect all overlaps between noisy long-read sequence data. It efficiently estimates Jaccard similarity by compressing sequences to their representative fingerprints composed on min-mers (minimum k-mer).

MHAP is included within the [Canu](#) assembler. Canu can be downloaded [here](#).

Contents:

1.1 Installation

1.1.1 Before your start

MHAP requires a recent version of the [JVM](#) (1.8u6+). JDK 1.7 or earlier will not work. If you would like to build the code from source, you need to have the [JDK](#) and the [Maven](#) build system available.

1.1.2 Prerequisites

- java (1.8u6+)
- maven (3.0+)

If you have not already installed the dependencies using maven, you will need an internet connection to do so during maven installation.

Here is a list of currently supported Operating Systems:

1. Mac OSX (10.7 or newer)
2. Linux 64-bit (tested on CentOS, Fedora, RedHat, OpenSUSE and Ubuntu)
3. Windows (XP or newer)

1.1.3 Installation

Pre-compiled

The pre-compiled version is recommended to users who want to run MHAP, without doing development. To download a pre-compiled tar run:

```
$ wget https://github.com/marbl/MHAP/releases/download/v2.0/mhap-2.0.tar.gz
```

And if wget not available, you can use curl instead:

```
$ curl -L https://github.com/marbl/MHAP/releases/download/v2.0/mhap-2.0.tar.gz > mhap-2.0.tar.gz
```

Then run

```
$ tar xvzf mhap-2.0.tar.gz
```

Source

To build the code from the release:

```
$ wget https://github.com/marbl/MHAP/archive/v2.0.zip
```

If you see a certificate not trusted error, you can add the following option to wget:

```
$ --no-check-certificate
```

And if wget not available, you can use curl instead:

```
$ curl -L https://github.com/marbl/MHAP/archive/v2.0.zip > v2.0.zip
```

You can also browse the <https://github.com/marbl/MHAP/tree/v2.0> and click on Downloads.

Once downloaded, extract to unpack:

```
$ unzip v2.0.zip
```

Change to MetAMOS directory:

```
$ cd MHAP-2.0
```

Once inside the MetAMOS directory, run:

```
$ maven install
```

This will compile the program and create a target/mhap-2.0.jar file which you can use to run MHAP. The quick-start instructions assume you are in the target directory when running the program. You can also use the target/mhap-2.0.jar file to copy MHAP to a different system or directory. If you would like to run the validation utilities you must also download and build the [SSW Library](#). Follow the instructions on the utilities page.

1.2 Quick Start

1.2.1 Running MHAP

Running MHAP provides command-line documentation if you run it without parameters. Assuming you have followed the installation instructions, you can run:

```
$ java -jar mhap-2.0.jar
```

MHAP has two main usage modes, the main finds all overlaps between the input sequences. The second only constructs an index which can be subsequently reused.

1.2.2 Finding overlaps

```
$ java -Xmx32g -server -jar mhap-2.0.jar -s<fasta/dat from/self file> [-q<fasta/dat to file or directory>]
```

Both the `-s` and `-q` options can accept either FastA sequences or binary dat files (generated as described below). The `-q` option can accept either a file or a directory, in which case all FastA/dat files in the specified directory will be used. By default, only the sequences specified by `-s` are indexed and the sequences in `-q` are streamed against the constructed index. Since MHAP is written in Java, the memory usage can be high. Generally, 32GB of RAM is sufficient to index 40K sequences. If you have more sequences, you can partition your data and run MHAP on the partitions. You can also increase the memory MHAP is allowed to use by changing the `Xmx` parameter to a larger limit.

The optional `-f` flag provides a file of repetitive k-mers which should not be selected as min-mers. The file is a two-column tab-delimited input specifying the kmer and the fraction of total kmers the k-mer comprises. For example:

```
$ head kmers.ignore
GGGGGGGGGGGGG      0.0005
```

means the k-mer GGGGGGGGGGGG represents 0.05% of the k-mers in the dataset (so if there are 100,000 total k-mers, it occurs 50 times).

1.2.3 Constructing binary index

```
$ java -Xmx32g -server -jar mhap-2.0.jar -p<directory of fasta files> -q <output directory> [-f<kmer file>]
```

In this use case, files in the `-p` directory will be converted to binary sketch files in the `-q` directory. Subsequent runs using these files (instead of FastA files) will be faster as the sequences no longer need to be sketched, only loaded into memory.

1.2.4 Output

MHAP outputs overlaps in a format similar to BLASR's M4 format. Example output:

```
[A ID] [B ID] [% error] [# shared min-mers] [0=A fwd, 1=A rc] [A start] [A end] [A length] [0=B fwd, 1=B rc]
```

An example of output from a small dataset is below:

```
155 11 0.164156 206 0 69 1693 1704 0 1208 2831 5871
155 15 0.157788 163 0 16 1041 1704 1 67 1088 2935
155 27 0.185483 159 0 455 1678 1704 0 0 1225 1862
```

In this case sequence 155 overlaps 11, 15, and 27. The error percent is computed from the Jaccard estimate using [mash distance](#).

1.2.5 Options

The full list of options is available via command-line help (`-help` or `-h`). Below is a list of commonly used options.

-h	Displays the help menu.
--version	Displays the version.
--pachio-fast	Set all the parameters for the PacBio fast setting. This is the current best guidance, and could change at any time without warning, default = false.

--pachio-sensitive	Set all the parameters for the PacBio sensitive settings. This is the current best guidance, and could change at any time without warning, default = false.
--nanopore	Set all the parameters for the Nanopore settings. This is the current best guidance, and could change at any time without warning, default = false.
-k	[int], k-mer size used for MinHashing. The k-mer size for second stage filter is separate, default = 16.
--num-hashes	[int], number of min-mers to be used in MinHashing, default = 512.
--num-min-matches	[int], minimum # min-mer that must be shared before computing second stage filter. Any sequences below that value are considered non-overlapping, default = 3.
--threshold	[double], the threshold cutoff for the second stage sort-merge filter. This is based on the identity score computed from the Jaccard distance of k-mers (size given by ordered-kmer-size) in the overlapping regions, default = 0.78.
--filter-threshold	[double], the cutoff at which the k-mer in the k-mer filter file is considered repetitive. This value for a specific k-mer is specified in the second column in the filter file. If no filter file is provided, this option is ignored, default = 1.0E-5.
--weighted	Perform weighted MinHashing using tf-idf scaling which biases repetitive k-mers to higher hash values. default=false.
--max-shift	[double], region size to the left and right of the estimated overlap, as derived from the median shift and sequence length, where a k-mer matches are still considered valid. Second stage filter only, default = 0.2.
--min-store-length	[int], The minimum length of the read that is stored in the box. Used to filter out short reads from FASTA file, default = 0.
--no-self	Do not compute the overlaps between sequences inside a box. Should be used when the to and from sequences are coming from different files, default = false.
--num-threads	[int], number of threads to use for computation. Typically set to #cores, , default = 8.
--ordered-kmer-size	[int], The size of k-mers used in the ordered second stage filter, , default = 12.
--ordered-sketch-size	[int], The sketch size for second stage filter, default = 1536.
--store-full-id	Store full IDs as seen in FASTA file, rather than storing just the sequence position in the file. Some FASTA files have long IDs, slowing output of results. This options is ignored when using compressed file format, default = false.
-f	[string], k-mer filter file used for filtering out highly repetitive k-mers. Must be sorted in descending order of frequency (second column), default = "".

1.3 Utilities

1.3.1 Using MHAP extras

In addition to the main overlapping algorithm, MHAP includes several utilities for validating overlaps and simulating data.

1.3.2 Validating overlaps

Assuming you have a mapping of sequences to a truth (such as a reference genome) in BLASR's M4 format, you can validate overlaps using MHAP's EstimateROC utility which will compute PPV/Sensitivity/Specificity:

```
$ java -cp mhap-2.0.jar edu.umd.marbl.mhap.main.EstimateROC <reference mapping M4> <overlaps M4/MHAP>
```

The default minimum overlap length is 2000 and default number of trials is 10000. This will estimate sensitivity/specificity to within 1%. It can be increased at the expense of runtime. Specifying 0 will examine all possible N^2 overlap pairs.

If the dynamic programming is turned on (by typing true for the parameter), overlaps not present in the reference mapping will be confirmed if a Smith-Waterman alignment can identify the overlap specified. This step requires the [SSW Library](#) to be separately compiled and installed:

```
$ wget https://github.com/mengyao/Complete-Striped-Smith-Waterman-Library/archive/master.zip
$ unzip master.gip && cd Complete-Striped-Smith-Waterman-Library-master/src
$ make all
$ cd /full/path/to/mhap/target/lib
$ ln -s /full/path/to/Complete-Striped-Smith-Waterman-Library-master/src/libsswjni.so
```

Now, you can run the EstimateROC command above.

1.3.3 Simulating Data

MHAP includes a tool to simulate sequencing data with random error as well as estimate Jaccard similarity for the simulated data.

```
$ java -cp mhap-2.0.jar edu.umd.marbl.mhap.main.KmerStatSimulator <# sequences> <sequence length (bp)>
```

The error rates must be between 0 and 1 and are additive. Specifying 10% insertion, 2% deletion, and 1% substitution will result in sequences with a 13% error rate. If no reference sequence is given, completely random sequences are generated and errors added. Otherwise, random sequences are drawn from the reference and errors added. Errors are added randomly with no bias.

```
$ java -cp mhap-2.0.jar edu.umd.marbl.mhap.main.KmerStatSimulator <# trials> <kmer size> <sequence length (bp)>
```

This usage will output a distribution of Jaccard similarity between a pair of overlapping sequences with the specified error rate (when using the specified k-mer size) and two random sequences of the same length. If no reference sequence is given, completely random sequences are generated and errors added, otherwise sequences are drawn from the reference. When one-sided error is specified (by typing true for the parameter), only one of the two sequences will have error simulated, matching a mapping of a noisy sequence to a reference. If a set of k-mers for filtering is given, they are excluded when computing Jaccard similarity, both between random and overlapping sequences.

1.4 Contact

1.4.1 Bugs, feature requests, comments:

If you encounter any problems/bugs, please check the known issues pages:

<https://github.com/marbl/MHAP/issues>

If not, please report the issue either using the contact information below or by submitting a new issue online.

Please include information on your run:

- 1) any output produced by MHAP
 - 3) sample data, if possible, to reproduce the issue

Who to contact to report bugs, forward complaints, feature requests:

Konstantin Berlin: kberlin@gmail.com

Sergey Koren: sergek@umd.edu