# Quick Manual: Local-Haplotype Variant Calling Software (LocHap -ver1.0)

Subhajit Sengupta[1], Kamalakar Gulukota[2], Yitan Zhu[1], Yuan Ji[1,3,*]

[1] *Center for Biomedical Research Informatics, NorthShore University HealthSystem, Evanston, IL, USA.*

[2] *Center for Molecular Medicine, NorthShore University HealthSystem, Evanston, IL, USA.*

[3] *Department of Health Studies, University of Chicago, Chicago, IL, USA.*

---

[*] Email for Correspondence: koaeraser@gmail.com

# 1 Ver-1.0 Release Note

We have tested **LocHap** on mac OSX 10.8.4 and 10.9, CentOs 6.5, and Ubuntu 12.04.4 LTS. Presently we have not tested our software on any Windows operating system.

# 2 Prerequisites

Installation of **LocHap** requires g++ compiler and GNU Make. If users don't have those already, then for UNIX system they can download them from here: `ftp://ftp.gnu.org/gnu/`. For mac OSX, installing Xcode (IDE for Apple's OSX and iOS) would automatically install those tools.

# 3 Installation

**LocHap** is an open-source self-contained computational tool. Standard installation of **LocHap** places the main functions as well as all the required libraries on the user's computer. It also includes a small example for illustration purposes. In the small example, a small *bam* and a *vcf* file are also copied to the user's computer.

> For format specifications of *bam* and *vcf* files, see `http://samtools.github.io/hts-specs/` `SAMv1.pdf` and `http://samtools.github.io/hts-specs/VCFv4.1.pdf`, respectively.

Installation steps are as follows.

1. Download and extract the "LocHap-release-v1.tar.gz" file from the download section of

   `http://www.compgenome.org/`

   and by default all the files would be extracted inside the directory named "LocHap-release-v1".

2. In the extracted directory, in a terminal prompt, run the script "install.sh" which is an executable file inside the directory. It would take approximately 5 to 6 minutes to finish the installation process.

3. If successful, **LocHap** binary executable would be generated in the "bin" subdirectory under the directory "LocHap-release-v1".

A screenshot of running the installer is provided below in Fig. 1.

Figure 1: Screenshot of installation process

# 4  Basic Usage

For general applications of **LocHap**, a user must provide the following two types of input files for biological samples to be analyzed.

- A sorted (coordinate-based) *bam* file and the corresponding bam index file (*bai*).

- A *vcf* file containing SNP information for the sample that the *bam* file characterizes and the one that user wants to run LH analysis.

When multiple samples are available for analysis, the *vcf* file must contain all the corresponding *bam* files, and LocHap can be executed in parallel, one per sample.

**The main command for executing LocHap for one particular sample is'**

```
./LocHap --vcf ⟨full name of the vcf file⟩ --bam ⟨full name of the bam file⟩ --sample ⟨sample name⟩
```

> The ⟨sample name⟩ is a standard field in any vcf file. Users can find the sample name therein. The header line in any *vcf* has 8 fixed, mandatory columns. They are: #CHROM,POS, ID, REF, ALT, QUAL, FILTER, INFO. If genotype data is present in the file, these are followed by a FORMAT column header, then an arbitrary number of sample IDs. The header line is tab-delimited.

Screenshot of this main command is provided below in Fig. 2.



Figure 2: Screenshot for executing the main command of **LocHap**

There are few optional flags that can be added to this main command.

- out: [string] – output file prefix which has default value same as sample name. Otherwise if user provides "output" as the prefix "output.hcf" will be generated.

- sig: [boolean] – significant flag which sets a flag for printing only DNA segments (see main paper for definition of a DNA "segment") with at least one significant haplotype called. The default value is currently set to 0, so all segments will be printed.

- size: [integer] – Max base pairs between two adjacent SNVs. This is an important value and is denoted as $K$ in the paper. See §2.1 in Online Methods (default is 500 NTs and must be between 50 and 1000).

- qual: [integer] LocHap ignores all reads with Phred-scaled Mapping Quality less than qual. LocHap Ignores all bases with Base Quality less than qual (default is set to 30 and it must be greater than 15)

- igv: [boolean] – IGV compatibility flag which sets a flag for printing in format suitable for loading into IGV software from Broad Institute. By default the output is written in native *hcf* format and the output file name would be ⟨sample name⟩.hcf or ⟨sample name⟩.igv, based on the choice of the flag.

> **LocHap** handles one bam at a time and that sample has to be inside the *vcf*; it does not care if other samples are inside that *vcf*. For analyzing multiple samples, one can parallelize the analyses, one for each sample.

## 4.1 Demo Example

In the "Example" subdirectory under the directory "LocHap-release-v1", a test example is prepared to demonstrate the main functions of **LocHap**, which can be executed by the following steps:

- Go to the "Example" subdirectory.

- Execute the script "run-example.sh" (Linux command "./run-example.sh").

- If successful, an output *hcf* file will be produced (called "tiny.hcf") with a message on the screen (see Fig. 3). The "tiny.hcf" file is an ASCII tab-delimited text file (see Online Methods for detail) containing all the local haplotype variants produced from the proposed statistical inference. See section 5 for detail of a *hcf* file.

A screenshot is provided below in Fig. 3 for the commands above.

Figure 3: Screenshot for testing an example

# 5 Haplotype Call Format (*hcf*)

The output segments from **LocHap** are written in the Haplotype Call Format (*hcf*). See Figure 4. Similar to *vcf*, an *hcf* file is a tab-delimited text file. There are a few header fields that contain the field names and their descriptions. After the initial header fields, each line in the *hcf* file represents a local haplotype (might not be a variant) and has seven column fields, displayed from left to right as: chromosome name (CHROM), positions on the chromosome (POS), nucleotides at those positions in the reference genome (REF), number of significant haplotypes (NumSig), called haplotypes (HAP-Call), all the possible haplotypes (All-HAP), data for the sample (DataForSample=⟨sample-name⟩).

Most of the fields are self-explanatory except the following – in the "Hap-Call" field, we include the posterior probability of each haplotype variant that is deemed statistically significant. The "All-HAP" field contains the posterior probability (before semicolon) and corresponding posterior false discovery rate (FDR) (after semicolon) for each possible haplotype. Haplotype variants in the "Hap-Call" field

are generated from those in "All-HAP" by using an FDR threshold of 0.01.

In the last field, a few basic statistics about the input-data are given, which include total number of SNPs (nSNP), total number of reads (nTot), number of reads having at least one entry in one of the SNP position (nACGT), number of blank reads (nBlank), number of discrepant reads (nDisc), number of reads with no missing entries, missing entries in one position or two positions or three positions (nM0, nM1, nM2, nM3) and number of clusters directly observed from the data (nClus). Note that, number of unique groups of data that has no missing SNP defines the number of clusters.



Figure 4: Screenshot of a sample *hcf* file.

# 6 Integrated Genome Viewer (IGV) Compatibility

We have an optional "igv" flag to generate an output file (with extension .igv). The ".igv" output files can be visualized using IGV. When loaded, SNPs belonging to an LHV are shown in *red* bars while those not in a LHV are shown in *blue* bars. Screenshot of this command is provided below in Fig. 5.

```
administrators-iMac-2:bin subhajit$
administrators-iMac-2:bin subhajit$ ls -l
total 1592
-rwxr-xr-x  1 subhajit  staff  524444 Jul 25 11:08 LocHap
-rwxr-xr-x  1 subhajit  staff  284072 Jul 25 11:09 filter
administrators-iMac-2:bin subhajit$
administrators-iMac-2:bin subhajit$
administrators-iMac-2:bin subhajit$ ./LocHap --vcf ../example/tiny.vcf --bam ../example/NA12878_tiny.bam --sample NA12878 -igv

Analyzing VCF (../example/tiny.vcf) for sample (NA12878) using bam (../example/NA12878_tiny.bam).
Output to => (NA12878.igv)
Block size = (500).
Min Quality = (30).
Sig Flag = (0)


Performing Local Haplotype Analysis and generating NA12878.igv file .....

administrators-iMac-2:bin subhajit$
administrators-iMac-2:bin subhajit$ ls -l
total 1608
-rwxr-xr-x  1 subhajit  staff  524444 Jul 25 11:08 LocHap
-rw-r--r--  1 subhajit  staff    6574 Jul 25 11:44 NA12878.igv
-rwxr-xr-x  1 subhajit  staff  284072 Jul 25 11:09 filter
administrators-iMac-2:bin subhajit$
administrators-iMac-2:bin subhajit$
```

Figure 5: Screenshot for executing the command in order to generate IGV compatible format

When multiple samples are analyzed by **LocHap**, for each sample an *hcf* file will be produced. IGV can open multiple *hcf* files as shown in Fig 6.

Figure 6: Screenshot of outputs from multiple samples opened in the IGV viewer

# 7 Optional Filtering

In order to perform type I or type II or type III filtering as mentioned in §6.1 of our Online Methods, a binary executable file named **filter** is also generated at the time of the installation inside the "bin" subdirectory under the directory "LocHap-release-v1". A user needs to provide the following files (in the strict order required by the program) in order to carry out this optional filtering:

- An *hcf* file generated by **LocHap** for one particular sample.

- A sorted (coordinate-based) *bam* file and the corresponding bam index file (*bai*) for that sample.

- A *vcf* file containing SNP information for that sample.

The main command for executing the filter for one particular sample is:

```
./filter ⟨full name of the hcf file⟩ ⟨full name of the bam file⟩ ⟨full name of the vcf
file⟩ ⟨sample name⟩ ⟨full name of output directory⟩ ⟨filter type⟩
```

User can run the filter according to the choice of the stringency level for the filtering. He or she can
enter 1 or 2 or 3 in ⟨filter type⟩ in order to run type I or type II or type III filter, respectively.

In Figure 7, screenshot of the command for executing type I filter. For type II or type III filter, the last
field of the command needs to be replaced by 2 or 3, respectively while other fields remain exactly the
same.



Figure 7: Screenshot of executing the command for optional filtering

Note that, in the provided example, type I filter would remove all the segments from the *hcf* file that
is generated by **LocHap**. The output file generated after filtering is in the same format as the original
*hcf* file.

# 8   Tips

- Any *bam* file should be indexed and the index file (*.bai*) should be in the same directory with the
  original *bam* file.

  - In case there is no index file, it can be easily generated by executing the command below:

10

```
samtools index ⟨name of the bam file⟩
```

See http://samtools.sourceforge.net/samtools.shtml for command details.

- An input *vcf* file may consist of variant calls for multiple samples. To analyze all the samples, run **LocHap** in parallel using multiple nodes, one per sample.