

Pipeliner v 0.2.0

Evaluating the performance of bioinformatics' pipelines for Next Generation re-Sequencing

Bruno Nevado and Miguel Perez-Enciso

January 13, 2014

Contents

1	Introduction	3
2	Obtaining and installing Pipeliner	4
2.1	Pipeliner	4
2.2	Additional software	4
2.3	OSes	4
3	Pipeliner GUI	5
3.1	Pipeline	5
3.1.1	Input	5
3.1.2	NGS simulation	8
3.1.3	Alignment	10
3.1.4	SNP calling	12
3.1.5	Statistics	14
3.2	Run settings	17
3.3	Bash file	19
3.4	Plots	19
4	Pipeliner command line tools	21
5	A guided run	22
6	User-defined steps	23
6.1	User-defined Input step	24
6.2	User-defined NGS simulation step (from fasta to fastq files)	25
6.3	User-defined alignment step (from fastq to bam files)	26
6.4	User-defined SNP calling and filtering step (from bam to vcf files)	26
7	Disclaimer	29
8	License	29
9	Acknowledgments	29

1 Introduction

Current deluge of Next Generation Sequencing (NGS) data has been accompanied by an increasing number of software and pipelines with a variety of options that are not necessarily well documented and whose actual performance or influence in the whole process is not well known. Pipeliner is a highly customisable, friendly GUI environment to compare alternative experimental designs and bioinformatics' pipelines for NGS re-sequencing studies. Pipeliner uses a simulation approach to estimate the recovery, sensitivity and false discovery rate in detecting SNPs, and the effect on population genetics' summary statistics. The current version focuses on single nucleotide polymorphisms, disregarding other sources of variation (e.g., indels or structural variants).

Pipeliner consists of two separate executables: a Graphical User Interface used for setting up analysis' parameters and plotting results; and a command-line tool used during the analysis to convert file formats, compute results, etc.

As the aim is to evaluate the performance of alternative pipelines for NGS analysis, an analysis with Pipeliner will typically require several other software packages to simulate the input datasets and analyse the NGS data.

The main steps involved in an analysis with Pipeliner are:

1. Defining the study system and experimental design. Making use of (coalescent or individual-based) simulations, the user can specify, e.g., how many individuals to sample, the evolutionary distance from the population sampled to the reference genome available, the nucleotide diversity or recombination rate of the population, patterns of selection or demographic changes, sub structuring, etc.
2. Simulating the NGS experiment, i.e. generating the NGS short reads from each individual in the population, according to a specific NGS technology and predefined average depth per individual.
3. Aligning the NGS short reads to the reference genome.
4. Performing variant calling and quality control.
5. Calculating a number of statistics (recovery, sensitivity, false discovery rate, etc) summarising the performance of the pipeline defined.
6. Estimating population genetics' summary statistics (theta, pi, neutrality tests) from the simulated NGS datasets.

All steps are easily customisable by the user, in terms of which software and which options to use. Further, Pipeliner allows the user to easily plot and compare the results obtained using different pipelines or different values for a range of parameters involved in the analysis.

If you use Pipeliner in your research results, please cite our paper:

Nevado B. and Perez-Enciso M. (*in preparation*). Pipeliner: evaluating NGS methods.

2 Obtaining and installing Pipeliner

2.1 Pipeliner

Pipeliner can be obtained from [it's website](#). The source code is available for download, and pre-compiled executables are available for some versions of Mac and Linux OS. For the command-line tool, a C++11 compiler is needed, while for the GUI the QT libraries must also be installed (www.qt-project.org).

Instructions on how to set up Pipeliner in Linux and Mac OS machines are provided together with the program.

2.2 Additional software

Pipeliner uses different software to both simulate and analyse NGS datasets. By default, Pipeliner uses the following software:

1. **ms** to simulate the population of interest and the reference using the coalescent ([website](#)).
2. **art_illumina** to simulate NGS reads ([website](#)).
3. **bwa** to align the short reads ([website](#)).
4. **samtools** for variant calling, including the utilities bcftools and vcfutils.pl ([website](#)).
5. **picard-tools**, specifically the tools AddOrReplaceReadGroups and CreateSequenceDictionary ([website](#)).
6. **mstatspop** to calculate variability and neutrality tests, optional ([website](#)).
7. **R** to plot the results ([website](#)).

These programs must be installed before doing the analysis in Pipeliner. Furthermore, any of the steps involved in the analysis can be performed with alternative software (see *User-defined steps section*).

If these programs are installed outside the computer's path, the location of the executables needs to be specified in the *Run settings* section of PipelinerGUI.

2.3 OSes

Pipeliner is available for Linux and Mac computers.

3 Pipeliner GUI

Upon starting PipelinerGUI, the user is presented with four tabs:

1. Pipeline
2. Run settings
3. Bash file
4. Plots

The *Pipeline* tab is used to set up the analysis, i.e. defining the sampling design, study system, and bioinformatics' analysis. The *Run settings* tab is used to define paths, number of replicates, and other run-time options. The *Bash file* allows previewing, editing and saving the pipeline's commands into a file. This file can then be executed, locally or remotely (e.g. on a compute cluster) and once finished, the *Plots* tab can be used for summarising and plotting of results.

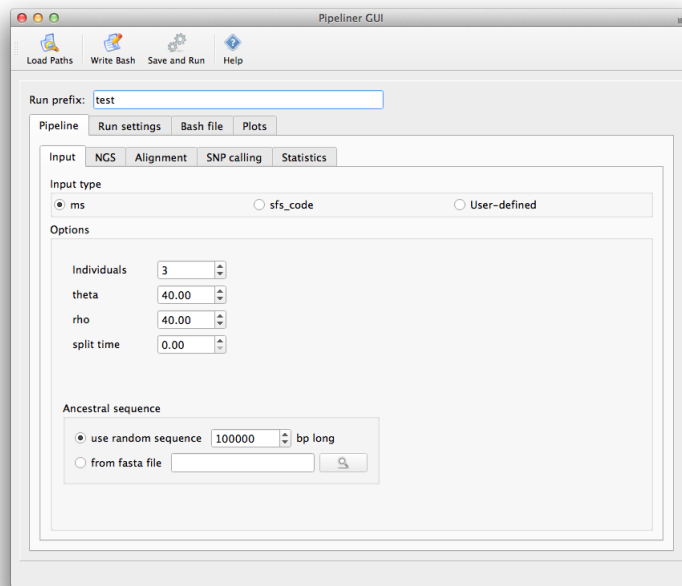
3.1 Pipeline

Used to define the both the input sequences to use in the analysis, the bioinformatics' pipeline, and which statistics to calculate.

3.1.1 Input

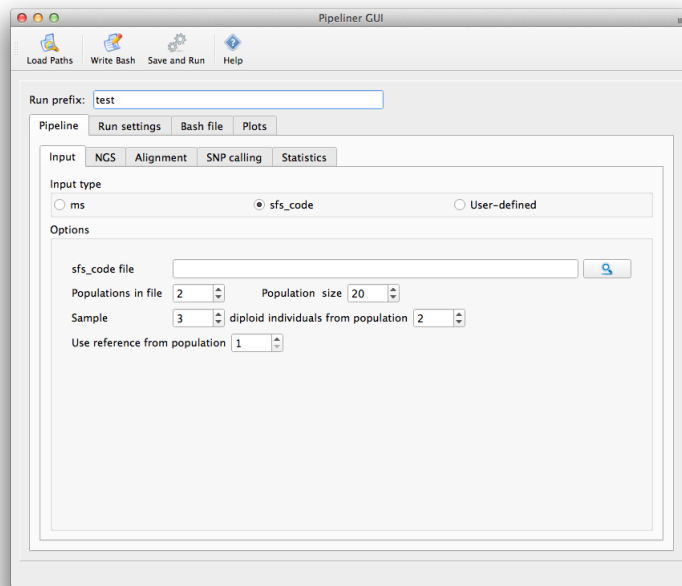
In the input step, the user defines how to obtain the DNA sequences for analysis. These include two haplotypes per diploid individual and one haplotype to be used as reference genome. Three input modes are available to achieve this: ms coalescent simulation, sfs_code forward simulation, and User-defined commands.

- ms coalescent simulation: The user chooses settings to simulate the population of interest using a subset of options from the program ms.



Options available:

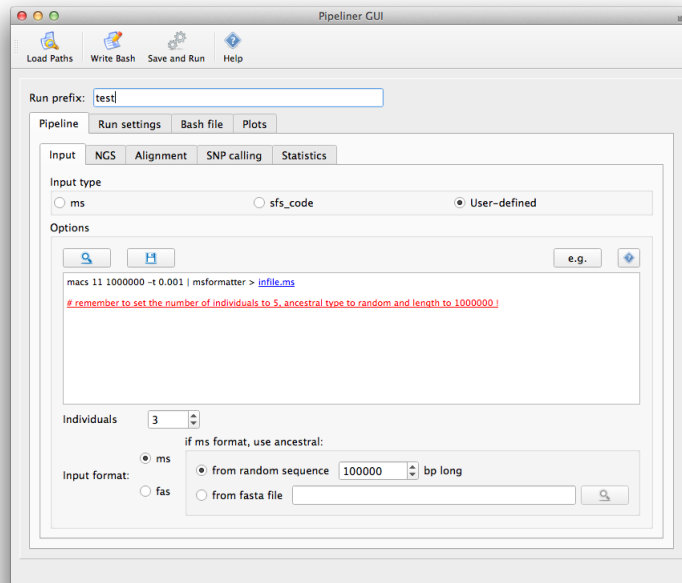
- Number of individuals sampled.
- Variability and recombination rates.
- Length of the simulated region.
- Time of divergence between the population of interest and the reference.
- Ancestral sequence upon which the SNPs obtained with coalescent simulations will be applied (provided by the user, or obtained by setting the length and choosing the random sequence option, which will generate a sequence with uniform number and distribution of the four bases a, c, g, t).
- SFS_code: The user provides an output file from the forward simulator program sfs_code ([website](#)). Note that, **if simulating more than 1 population, all populations must have the same number of individuals sampled.**



Options available:

- Location of sfs_code output file.
- Number of populations simulated.
- Number of diploids simulated per population (must be the same for every population).
- Number of diploid individuals to use in the analysis.
- Which population to analyse. The n diploid individuals will be taken from this populations, starting with the second diploid individual in the population (the first diploid individual is *reserved* for using as reference).
- Which population to use as the reference. Pipeliner will take the first sequence from this population and use it as reference genome for alignment, SNP calling etc...
- User-defined: Other software can be used to obtain the input data, by providing a command line specifying how to simulate it. This command line must be executable directly, i.e. it should work if it is copied and pasted into a terminal window. It should work with any simulation software insofar as the result is in either fasta or ms-like format. Whichever population scenario the user selects, **the first haplotype in the input file will be used as the reference**, and the remaining haplotypes represent the population, with **each 2 consecutive haplotypes treated as 1 diploid individual**. If using ms-like input, the user must also define which ancestral sequence to use - random or user provided. *See the*

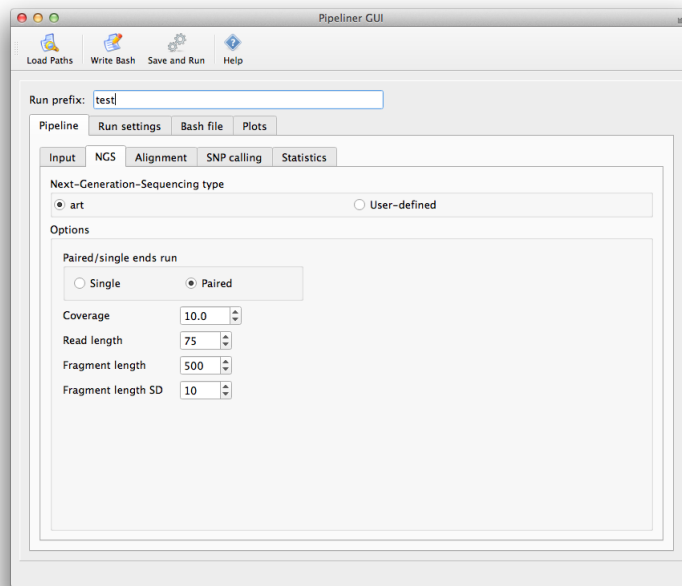
User-defined steps section for details.



3.1.2 NGS simulation

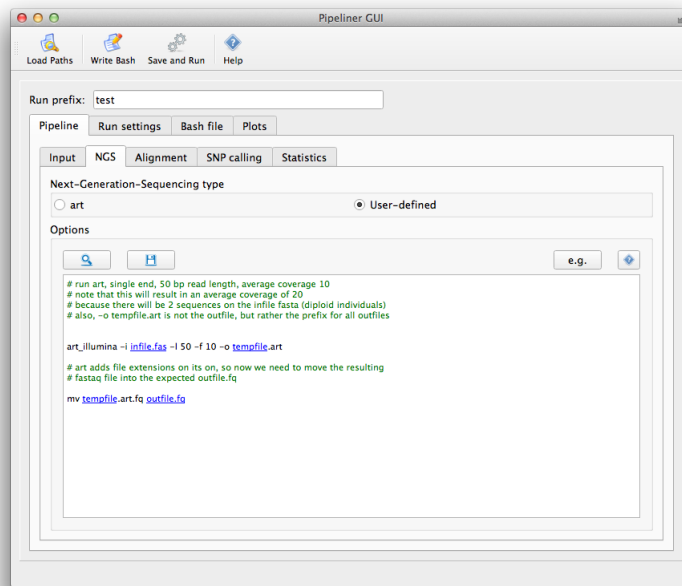
Used to obtain the NGS short reads for each diploid individual analysed. At the end of this step, there should be one (or two, if simulating paired-end sequencing) fastq files per individual. Two options are currently available: default art_illumina and User-defined.

- Default art_illumina: Use the program ART with Illumina technology and the built-in empirical profiles.



Options available:

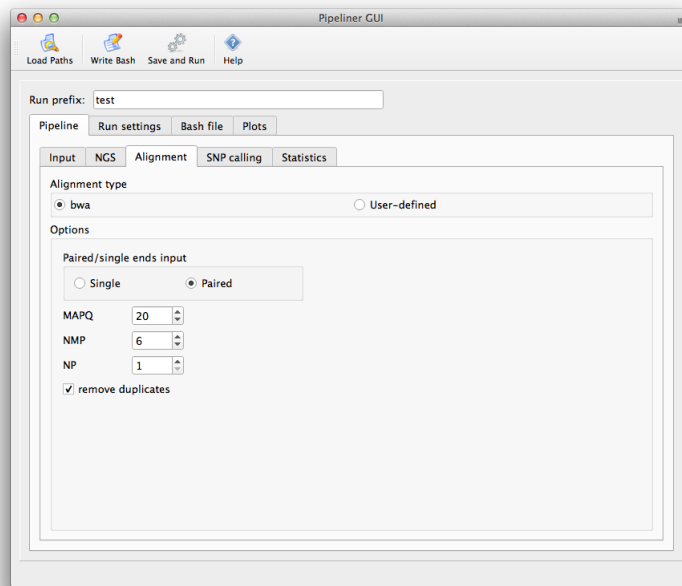
- Single-end or paired-end run.
 - Average coverage to use per diploid individual.
 - Length of the short reads.
 - Length of the fragments sequenced from both ends, and the standard deviation of this value (SD). Only for paired-end runs.
- User-defined: Allows using alternative software to obtain the NGS short reads for each individual. This can be done by using alternative software (e.g. simseq or FluxSimulator) or by using 454 or solid profiles from ART (note that the correct binaries for ART-solid and ART-454 will need to be installed on the computer). *See the User-defined steps section for details.*



3.1.3 Alignment

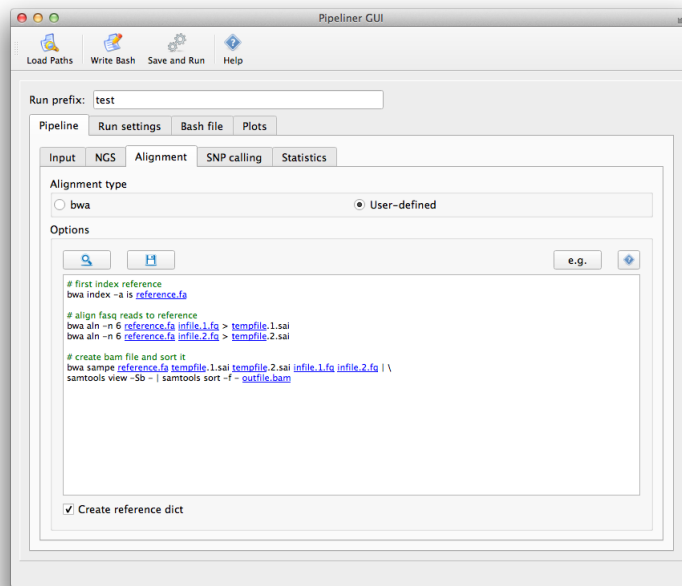
Used to specify how to align the NGS short reads from each individual to the reference sequence. This step should start from the fastq files obtained before, and finish with a single bam formatted file (bam/sam file specifications) for each individual. Two options available: default using bwa, or User-defined step.

- Default bwa: Uses the software BWA to perform the alignment of short reads.



Options available:

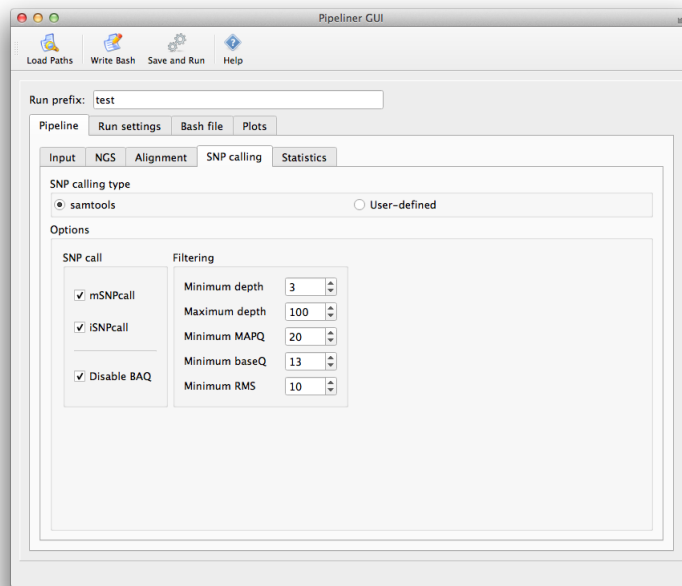
- Input type -single or paired-ends. Should match the output of the NGS simulation step.
- Mapping quality threshold - reads with mapping quality below this value will be discarded.
- Number of mismatches allowed between each read and the reference.
- How many processors to use while doing the alignment.
- User-defined: Use alternative alignment software, or use options from BWA that are not available in the default section. *See the User-defined steps section for details.*



3.1.4 SNP calling

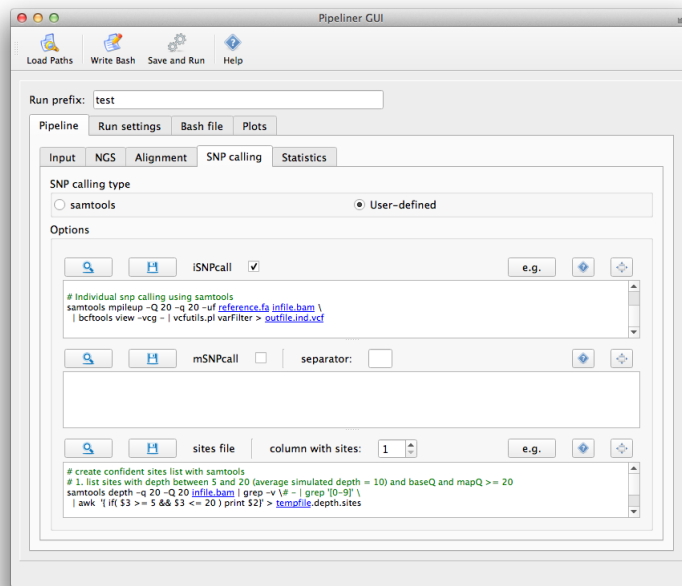
Used to define SNP calling and filtering options. This step starts from the aligned short-reads obtained before (.bam files) and finishes with a .vcf format file ([vcf file specifications](#)). SNP calling can be done separately for each individual (iSNPcall) as well as jointly for all individuals under analysis (multiple-sample SNP calling, mSNPcall). Two options available: default using samtools, or User-defined step.

- Default samtools: Use the software samtools and bcftools to perform the SNP calling and filtering.



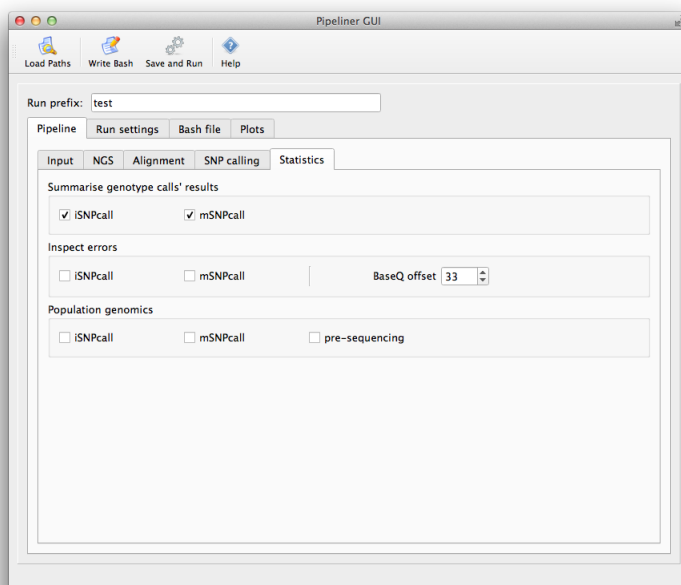
Options available:

- SNP calling using each individual's data separately, and/or using all individual's data at the same time (multi-sample SNP calling)
- Disable Base Alignment Quality (BAQ) option of samtools - leaving BAQ enabled should be useful in removing false positive SNP calls nearby indels.
- Minimum and maximum coverage to call SNPs (i.e. sites with less or more than these values will be excluded from the results).
- Minimum base quality of reads (reads below this threshold are not used for SNP calling).
- Minimum RMS mapping quality for SNPs (for filtering SNP calls using `vfcutils.pl varFilter`).
- User-defined: Define alternative SNP calling software, or perform SNP calling and filtering with extra options of samtools/bcftools not available in the default panel. *See the User-defined steps section for details.*



3.1.5 Statistics

In this section the user defines how to summarise the performance of the pipeline defined. Three main summary options are available:



- Summarise genotype calls' results. Enabling this option will report, for each individual of each replicate, a number of statistics concerning the power and sensitivity of the pipeline defined. To do so Pipeliner compares the original sequence data obtained by coalescent or forward simulations (pre-sequencing dataset), to the data obtained after the SNP calling step (post-sequencing dataset). The statistics reported include **Recovery** - percentage of original genotypes that are correctly identified in the post-sequencing datasets; **Sensitivity** - percentage of genotypes that pass the different filters, i.e. are included in the list of confident sites, and are correctly identified in the post-sequencing datasets; and **False Discovery Rate** - percentage of genotype calls performed that are incorrect (Box 1). These statistics are calculated as averages across individuals, and are reported separately for invariable positions - where the individual is homozygous for the reference allele - homozygous SNPs - where the individual is homozygous for the alternative allele - and heterozygous SNPs - where the individual carries one reference and one alternative allele. In addition, Pipeliner can report the frequency of different outcomes, e.g., the percentage of heterozygous SNPs that are incorrectly identified as homozygous for the reference allele or homozygous for the alternative allele, etc. *Enabling this option will create a summary file - or two if both iSNPcall and mSNPcall are enabled - that can be loaded into PipelinerGUI's Plot section after the analysis is finished.*

Box 1: Example of statistics calculated by Pipeliner

Original.fas

A	A	A	A	A	A	A	A	A	A	A	→ Reference
C	A	T	A	A	T	A	A	A	A	A	} Individual 1
C	A	A	A	A	T	A	A	A	A	A	
C	A	A	A	A	A	A	A	A	G	G	} Individual 2
C	A	A	A	A	A	A	A	A	A	G	

Individual.1.sites

1
2
3
4
5
8
9
10

Individual.2.sites

4
5
6
7
8
9
10

all_inds.vcf

chr1	1	.	A	C	113	.	(...)	GT:PL:GQ	1/1:146,21,0:10	0/0:10,0,0:10
chr1	3	.	A	T	113	.	(...)	GT:PL:GQ	0/1:146,0,30:10	0/0:0,0,0:11
chr1	6	.	A	T	113	.	(...)	GT:PL:GQ	0/1:10,0,10:10	0/0:0,17,10:11
chr1	9	.	A	G	113	.	(...)	GT:PL:GQ	0/0:0,15,30:10	1/1:10,17,0:11
chr1	10	.	A	G	113	.	(...)	GT:PL:GQ	0/0:0,15,30:10	0/1:10,0,10:11

post-sequencing dataset

C	A	T	A	A	N	N	A	A	A	} Individual 1
C	A	A	A	A	N	N	A	A	A	
N	N	N	A	A	A	A	A	G	A	} Individual 2
N	N	N	A	A	A	A	A	G	G	

Statistics

Original number of homozygous SNP genotypes = 4

Number of homozygous SNP genotypes covered = 2

Number of correctly identified homozygous SNP genotypes = 1

Recovery homozygous SNP genotypes = $1/4 = 25\%$

Sensitivity homozygous SNP genotypes = $1/2 = 50\%$

Original number of heterozygous SNP genotypes = 2

Number of heterozygous SNP genotypes covered = 2

Number of correctly identified heterozygous SNP genotypes = 1

Recovery heterozygous SNP genotypes = $1/2 = 50\%$

Sensitivity heterozygous SNP genotypes = $1/2 = 50\%$

Number of homozygous SNP genotypes in post-sequencing dataset = 2

Number of incorrect homozygous SNP genotypes in post-sequencing dataset = 1

FDR homozygous SNP genotypes = $1/2 = 50\%$

Number of heterozygous SNP genotypes in post-sequencing dataset = 2

Number of incorrect heterozygous SNP genotypes in post-sequencing dataset = 1

FDR heterozygous SNP genotypes = $1/2 = 50\%$

- Inspect errors. If enabled, Pipeliner will report, for each site that is incorrectly identified in the post-sequencing datasets, the error type - homozygous SNP called heterozygous, heterozygous SNP called invariable etc.. - the original configuration i.e. the alleles present in reference and individual, and the mpileup line. The

mpileup line shows the reference allele, the alternative allele, the bases seen covering the site and their base qualities, and can be used to pinpoint specific errors incurred by the pipeline defined. *This option is only available when running the analysis in a Linux machine, and the results must be checked manually using a text editor.* The first two lines of the example *errors* file from Box 1 are provided below:

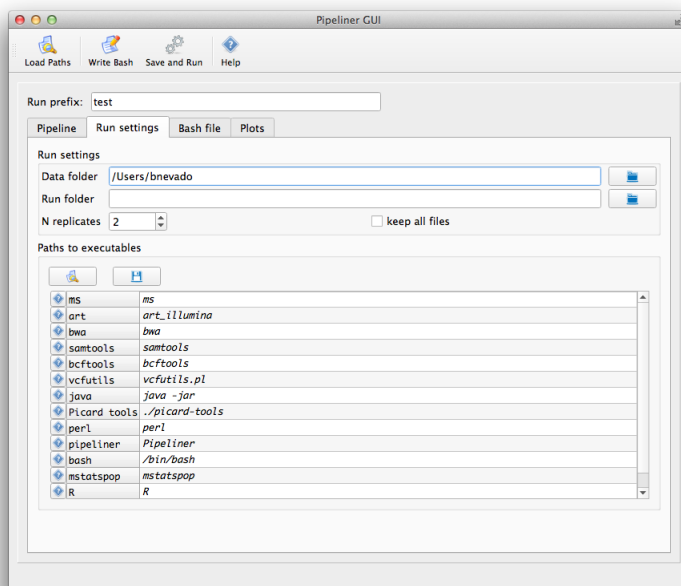
```
REP IND SITE ERR_TYPE CONFIG : MPILEUP LINE (BASE QUALS)
1 2 9 ra_aa aga_gg : chr1 9 A 5 G$GgGt -8=8* (12 23 28 23 9)
```

The second line shows that in replicate 1, individual 2 contained a heterozygous genotype in site 9. This site was called homozygous for alternative allele (ra_aa), with reference allele being an A and alternative allele a G (aga_gg). Following the colon ":" comes the line in the mpileup file where the error occurred, showing the chromosome (chr1), site (9), reference allele (A), depth (5) followed by the bases seen and their base-qualities. The numbers within parenthesis at the end are the base-qualities.

- Population genomics. If selected, Pipeliner will calculate the following population genetics' statistics: population variability (theta), nucleotide variability (pi), Tajima's D, Fu and Li's D and Fay and Wu's H. Neutrality tests are calculated using the reference genome as outgroup. Statistics are calculated using the program mstatspop, and can be estimated for each SNP calling method defined in the pipeline - iSNPcall and mSNPcall - as well as with the original, pre-sequencing datasets i.e. without errors nor missing data. *Results will be written into a summary file that can be loaded into the Plot section of PipelinerGUI, together with the results of genotype calls.*

3.2 Run settings

The Run settings tab is used to specify several run-time options.



Options available:

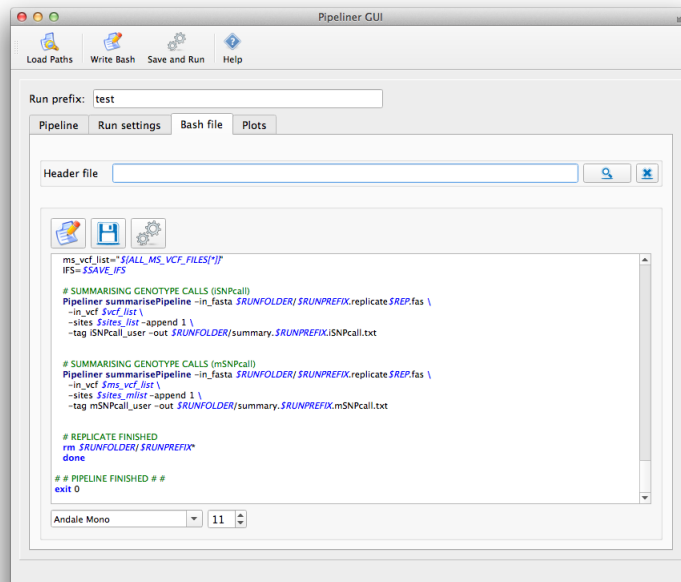
- Data folder: where to store the result files. **The folder defined here must exist prior to running the analysis, and the absolute path must be given.**
- Run folder: where to run the analysis. Default: same as Data folder. Under some computer configurations, it might be convenient to run the analysis on a different folder (e.g. on the computer's temporary folder) and copy back the results to Data folder once the analysis finishes. If Run folder does not exist, it will be created by Pipeliner, and deleted once the analysis finishes.
- Run prefix: the prefix used for all files created during the analysis, except summary files. This is useful for cleaning up once the run finishes, or to manually inspect the resulting files from the analysis. For summary files, the run prefix is included in the middle of the file's name.
- Number of replicates: Number of times to run the entire pipeline before summarising the results.
- Keep all files. The default behaviour is to delete all files except summary and log files at the end of the run. If this option is checked, all files (fasta, fastq, bam files, vcf files, as well as all temporary files) will be available in the Data folder once the run finishes.

In addition, the Run settings tab is also used to specify the location of the different executables needed to simulate and analyse the NGS datasets, including the simulation tools

selected, alignment and SNP calling software, the mstatspop software and the command-line Pipeliner tool, etc. **If these programs are not installed on the computer's path, the paths must be defined in full.** To do so, double-click the default software name, and enter the full path to the executable. Once defined, the location of the different executables can be saved into a file, for quick setting up in subsequent analysis.

3.3 Bash file

Once the Pipeline and Run settings are defined, pressing the *Write Bash* button will write the commands into the Bash file section, where they can be previewed and edited. Executing the pipeline can be done locally by choosing *Save and Run*, but the Bash commands can also be saved to a file to run on a different computer, or directly from the terminal. For runs using a computer cluster, there is also the option to load a header file, which can contain e.g. which nodes to use, or how much memory to allocate to the job. The options will depend on the cluster where the analysis will be performed, its job-handling software, etc.



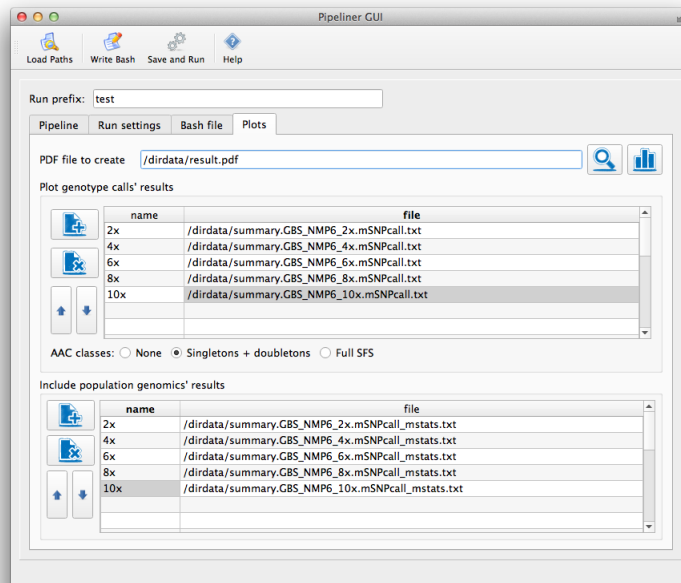
3.4 Plots

The Plots tab is used to plot and compare the output of runs with different programs/settings (though it can also be used to summarise the results of a single run).

Each analysis with Pipeliner will output one or several summary result files named `summary.runprefix(...).txt`. The files ending in `"...runprefix.iSNPcall.txt"` and `"...run-`

prefix.mSNPcall.txt" contain the raw results of genotype calls, and can be loaded into the upper section of the Plots tab. The files ending in "...runprefix.iSNPcall_mstats.txt", "...runprefix.mSNPcall_mstats.txt" and "...runprefix.preseq_mstats.txt" contain the results of the population genetics' estimates, and can be loaded into the lower section of the Plots tab. For both cases, the name column, to the left of the files' names, can be used to label the resulting plots. If no names are entered, the results are labelled with numbers, according to the order in these sections.

Results concerning the Sensitivity of genotype calls can also be summarised according to Alternative Allele Count (AAC), i.e. the number of alternative bases, for a given SNP, present in the population. AAC will affect Recovery and Sensitivity values even in the case of individual sequencing, if the SNP calling step is performed jointly for all individuals. No difference should be found when doing SNP calling individually for each diploid individual. The three options for AAC plots are: *none*, i.e. do not create plots with AAC results; *singletons + doubletons*, will create plots with 3 classes - singletons of the alternative allele, doubletons for the alternative allele, and all remaining mutations; and *full Site-Frequency-Spectrum*, which will use one class for each possible AAC (singletons, doubletons ... up to fixed for alternative allele). Plots of AAC are only possible for analysis with the same number of diploid individuals sampled - otherwise an error will be raised.



4 Pipeliner command line tools

Pipeliner's command line tools is a collection of commands to change file formats, summarise results, etc. To see a complete list of available commands and options, call Pipeliner's command line tool with the command "help".

5 A guided run

This short tutorial shows how to perform an example analysis in Pipeliner.

1. Download and install Pipeliner and default software (see section 2). Start PipelinerGUI by double-clicking the executable.
2. In the **Pipeline** > **Input** tab: Set *Input type* to *ms*, then choose 5 Individuals, and 100 *theta* and *rho*. Set the ancestral sequence type to *random*, and 100 000 bp long. Choose *Split time* of 0.
3. In the **Pipeline** > **NGS** tab: Set the run to *Paired-ends*, *coverage* to 6x and *read length* to 75 bp.
4. In the **Pipeline** > **Alignment** tab: Set the input to *Paired-ends*, *mapping quality* to 20 and number of mismatches (*NMP*) to 6.
5. In the **Pipeline** > **SNP calling** tab: Select both *iSNPcall* and *mSNPcall* modes and check the *disable BAQ* option. In the filtering section choose *minimum depth* of 3 and *maximum depth* of 100, *minimum base quality* of 13 and *minimum RMS* of 10.
6. In the **Pipeline** > **statistics** tab: check both *iSNPcall* and *mSNPcall* in the *Summarise genotype calls' results* section, and all three boxes in the *Population genomics* section.
7. In the **Run settings** tab: Select a folder where the analysis should be performed (*Data folder*). Leave *Run folder* blank, then choose a prefix for this test run, and set number of replicates to 10. Uncheck the *keep all files* option. In the section *Paths to executables*, set the paths to the different executables, and remember to save them for future analysis (use the *Save paths to file* button)
8. In the **Bash file** tab: press the *Write bash* button. The bash commands should now appear in the text editor. If any settings are incorrect, an error will be raised with information on which error and how to correct it. Otherwise, the pipeline should be correctly defined, and the analysis can be started. Press the *Save and Run* button, save the commands into a file (preferably in the Data folder defined before). The pipeline defined in the PipelinerGUI will now be executed.
9. Wait until the run finishes, and then check the log file for errors or warning.
10. If no errors were encountered during the analysis, there should be 5 files in the Data folder named *summary.runprefix(...).txt*, and we can proceed to summarise and plot the results. Get back to PipelinerGUI and choose the **Plots** tab. In the upper section of this tab, press the *Add file* button, navigate to the Data folder, and select the two appropriate summary files. In the *name* column, enter

either iSNPcall or mSNPcall, according to which summary file is present on the *file* column.

In the *AAC classes*, choose the *singletons + doubletons* option.

In the lower section of the Plots tab add the summary files (...)preseq_mstats.txt, (...)iSNPcall_mstats.txt and (...)mSNPcall_mstats.txt, and label then accordingly in the *name* column.

Now select the file to save the PDF results to, and click *Plot results*. After a brief moment, the resulting pdf file should open automatically with the default PDF viewer of the computer. This PDF will contain a graphical overview of the performance of both individual and multiple-sample SNP calling at the coverage and settings defined earlier, and the population genetics' statistics estimated with the pre-sequencing datasets, and each of the SNP calling results. In the same folder as this PDF, there will also be a file with the commands used to plot these results, which can be loaded into the R package for further refinement of the plots.

11. Once this is done, it is possible to re-run the analysis using different parameter values, either in the input options, or any of the other tabs. For instance, using different average depth (NGS simulation tab), to see how the recovery, FDR etc change with different coverages. Or choosing a different number of individuals, more stringent options in SNP calling, etc... Then, rerun the analysis, and load the new summary files, together with the previous ones, into the Plot results tab. Plotting the new results will provide a quick graphical overview of the performance of the different options/pipelines.

6 User-defined steps

For every step of the analysis performed by Pipeliner (generating the population and reference, obtaining the NGS short reads, aligning them, and calling variants) there is the option to use User-defined steps. **These User-defined steps allow using any software, insofar as standard data formats are used (fasta, fastaq, bam, vcf).** This functionality is at the core of Pipeliner, and provides the required flexibility to remain up-to-date with the fast-pace with which new software and methods are devised to analyse NGS data.

All User-defined steps are implemented in a similar way (although each expects different input/output files). Overall, the user must provide Pipeliner with the command-line options (which must be executable from the terminal) required to use the software chosen for the analysis, by entering these steps in the appropriate fields of each User-defined option. The user can specify sequential commands, by writing them on different lines within the User-defined fields. Linux redirection and pipe operators can also be used.

Within these command-lines, the user must include **'tags'** which represent the **input** and **output** files expected by Pipeliner and needed for the corresponding step. In addition to input and output files, tags can be used to denote the **reference genome** (available in some User-defined steps). Finally, the user might want to use **temporary files**, e.g. when the analysis includes more than 1 command, and the output from the first command is needed to run the second (see the examples below). The tags are entered simply by writing their names in the text fields, and should be automatically recognised and highlighted by PipelinerGUI. Specific tags to each step are described below.

With all User-defined steps, there is the possibility to save into or load from a text file the commands to be used, and comments can be included by preceding any line with a **#**. The folder `/examples` within Pipeliner's distribution contains a few examples in this format.

Note that **any command written into the User-defined steps will be executed directly on the terminal using the current user's privileges**. Thus care should be taken to avoid using commands that could damage your computer or files - such as removing unwanted files, formatting disks, etc...

6.1 User-defined Input step

In the input tab, the user can specify different software to obtain the population and the reference genome for the analysis. The software must provide a ms-like or a fasta-formatted file, in which the first line will be used as the reference genome, and each 2 consecutive lines thereafter as diploid individuals. Tags available for this step are:

1. **infile.ms** : the input file in ms-like format. Pipeliner will convert it to fasta format using the ancestral sequence provided by the user.
2. **infile.fas** : the input file in fasta format.
3. **tempfile** : any temporary file needed.

In addition, the user will need to define (in the appropriate fields) the number of diploid individuals which will be simulated and, if using ms-like input, which sequence to use as ancestral.

Some examples (these examples are provided, as plain-text formatted files, in the folder `/examples` of Pipeliner's distribution):

- User-defined coalescent simulation with the msms program. Assuming the executable for msms is installed in `/home/jsmith/applications`, the following line (entered in the User-defined section of the input tab) would simulate a population of 5 diploid individuals (plus reference sequence) which have experienced constant

population size and weak selection (example 5.1.1 in msms user's manual):

```
/home/jsmith/applications/msms -N 100000 -ms 11 1 -t 50 -SAA 200 -SaA 100  
-SF 1e-2 > infile.ms
```

- User defined coalescent simulation with the macs program. Assumes macs is installed and in your path, will simulate a 1 Mb genomic region for 5 diploid individuals.

```
macs 11 1000000 -t 0.001 | msformatter > infile.ms
```

6.2 User-defined NGS simulation step (from fasta to fastq files)

In the NGS simulation step, the user can define alternative commands to obtain the NGS short-reads. These commands will be executed once for each individual under analysis. The allowed tags on this step are:

1. **infile.fas**: the fasta-format file containing 2 sequences for the current individual.
2. **outfile.fq**: the fastq output file containing the short-reads for the current individual. Using this tag will make Pipeliner expect a single-end NGS experiment.
3. **outfile.1.fq** and **outfile.2.fq**: for paired-end NGS experiments, these two files should contain the forward and reverse direction short reads for the current individual. If used, both tags must be specified. Also, these tags and outfile.fq are mutually exclusive (and will raise an error if present on the same step).
4. **tempfile** : any temporary file needed.
5. **ind_number** : current individual number.

Some examples:

- User defined NGS simulation step using art. This example shows how to deal with programs that output data into non-customisable files. For instance, using art_illumina the user only specifies the prefix for the output files, and the program adds extension to this prefix to obtain the full name for output files. We assume art_illumina is the name of the executable, which is installed in the computer's path. The first line below runs art with average 10x depth per individual, 75 bp reads etc, and resulting files are named tempfile.1.fq and tempfile.2.fq... In the second and third lines, these files are then renamed so that Pipeliner can proceed with them for the analysis.

```
art_illumina -i infile.fas -p -m 500 -s 10 -l 75 -f 10 -o tempfile  
mv tempfile1.fq outfile.1.fq  
mv tempfile2.fq outfile.2.fq
```

6.3 User-defined alignment step (from fastq to bam files)

In the User-defined alignment step, the user must provide commands to align the NGS short reads into a single bam file for each individual. The allowed tags are:

1. **infile.fq** or **infile.1.fq** + **infile.2.fq**: the fastq files containing the short reads for the current individual. Only the first or the two last should be present, depending on whether a single- or paired-end NGS experiment is being simulated.
2. **outfile.bam**: the bam-format file containing the aligned reads for the current individual.
3. **reference.fa**: the fasta-formatted, interleaved reference genome to use to align the short reads to.
4. **tempfile** : any temporary file needed.
5. **ind_number** : current individual number.

e.g. using bwa:

Some examples:

- Paired-end alignment with bwa (plain-text example provided in folder /examples of Pipeliner's distribution):

```
# first index reference
bwa index -a is reference.fa

# align fastq reads to reference
bwa aln -n 6 reference.fa infile.1.fq > tempfile.1.sai
bwa aln -n 6 reference.fa infile.2.fq > tempfile.2.sai

# create bam file and sort it
bwa sampe reference.fa tempfile.1.sai tempfile.2.sai infile.1.fq infile.2.fq | \
samtools view -Suh - | samtools sort -f - outfile.bam
```

6.4 User-defined SNP calling and filtering step (from bam to vcf files)

In this step, the user must define the commands to call variants from the aligned short-reads files in bam format. There are three sub-steps, including the individual SNP calling, the multiple-sample SNP calling (with all individuals together) and the 'sites' step (see below). The first two steps are optional (at least one must be present), while the 'sites' step is always required.

1. **Individual SNP calling**: commands to perform SNP calling and filtering for each individual separately. Should start from the bam file for the current individual, and finish with a single vcf file for each individual. The allowed tags here are:

- a) **infile.bam**: the current individual's bam file.
- b) **outfile.ind.vcf**: resulting variant-calls for the current individual.
- c) **reference.fa**: the fasta-formatted, interleaved reference genome to use to align the short reads to.
- d) **tempfile** : any temporary file needed.
- e) **ind_number** : current individual number.

Some examples:

- Individual SNP calling with GATK (plain-text example provided in folder /examples of Pipeliner's distribution):

```
# index bam file
samtools index infile.bam
# Individual snp call with GATK. Adapt the path to GenomeAnalysisTK.jar accordingly!
java -jar /apps/GATK/GenomeAnalysisTK.jar -T UnifiedGenotyper -R reference.fa \
-I infile.bam -o outfile.ind.vcf
```

2. **Multiple-sample SNP calling**: commands to perform SNP calling with all individuals at the same time. The outcome of this step should be a single vcf file, containing the variant calls for all individuals. These variant calls will be filtered with the .sites file obtained in the next step (see below). Allowed tags are:
 - a) **all.infiles.bam**: list of bam files for all individuals. By default, the list of files will be printed separated by a single space, but the field *separator* can be used to specify otherwise. For instance, GATK expects a "-" in between each bam file name, while samtools expects bam files' names to be separated by spaces. With default options the tag *all.infiles.bam* will be replaced with the string *infile.1.bam infile.2.bam ... infile.n.bam*; with a "-" in the *separator* entry the string will be *infile.1.bam -I infile.2.bam ... -I infile.n.bam*.
 - b) **outfile.allinds.vcf**: vcf file containing variant calls for all individuals.
 - c) **reference.fa**: the fasta-formatted, interleaved reference genome to use to align the short reads to.
 - d) **tempfile** : any temporary file needed.
3. **Obtain .sites file**: in order to obtain accurate estimates of Recovery and Sensitivity, and of population genetics' statistics, it is imperative to obtain a list of "confident sites". This list should contain all sites that pass the quality control filters, irrespective of the presence or absence of variants. Sites not present in this list are coded as missing, whereas those present are coded according to their genotype call or, if no genotype call is present, as homozygous for the reference allele. This step is performed once for each individual under analysis. Allowed tags:

- a) **infile.bam**: the current individual's bam file.
- b) **outfile.sites**: list of sites passing all filters.
- c) **reference.fa**: the fasta-formatted, interleaved reference genome to use to align the short reads to.
- d) **tempfile** : any temporary file needed.
- e) **ind_number** : current individual number.

The user must also specify (in the field to the right) in which column of the output file the site numbers are stored, e.g. in a standard vcf file the sites' numbers are on column 2.

Some examples:

- Obtaining the .sites file with GATK (plain-text example provided in folder /examples of Pipeliner's distribution):

```
# obtain list of confident sites with GATK's EMIT_CONFIDENT_SITES option

java -jar ./GenomeAnalysisTK.jar -T UnifiedGenotyper -R reference.fa \
  -I infile.bam -o tempfile.ind_number.vcf --output_mode EMIT_ALL_CONFIDENT_SITES

# get confident homozygous-reference sites, with min dep of 3 and max dep of 20
cat tempfile.ind_number.vcf | grep -v "\./\." - | grep -v PL \
  | awk -F: '{ if ($NF >=3 && $NF <= 20 ) print; }' \
  | awk '{print $2}' - > tempfile.ind_number.confident_invar.sites

# get snps from vcf file (ignore indels)
java -jar ./GenomeAnalysisTK.jar -T SelectVariants -R reference.fa \
  -V tempfile.ind_number.vcf -selectType SNP -o tempfile.ind_number.iGATK.unflt_var.vcf

# filter variant calls
java -jar ./GenomeAnalysisTK.jar -T VariantFiltration -R reference.fa \
  -V tempfile.ind_number.iGATK.unflt_var.vcf --filterExpression \
  "QD < 2.0 || FS > 60.0 || MQ < 40.0 || HaplotypeScore > 13.0 || \
  MappingQualityRankSumTest < -12.5 || ReadPosRankSum < -8.0" \
  --filterName "GATK_BP" -o tempfile.ind_number.iGATK.flt_var.vcf

# get sites with confident vars, and add sites with confident homozygous-reference calls
grep PASS tempfile.ind_number.iGATK.flt_var.vcf | awk '{print $2}' > outfile.sites
cat  tempfile.ind_number.confident_invar.sites >> outfile.sites

# file outfile.sites now contains all sites that are confidently genotyped by GATK,
# irrespective of the presence/absence of variants
```

7 Disclaimer

Without any warranty of any kind etc...

8 License

Pipeliner is distributed under the GNU General Public License version 2.

This program is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; version 2 dated June, 1991 or at your option any later version. This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

9 Acknowledgments

Work funded by Consolider CSD2007-00036 "Centre for Research in Agrigenomics" and AGL2010-14822 grants (Spain) to MPE.