

User's Manual of QCF90 (Draft)

Yutaka Masuda

November, 2020

Contents

1	Introduction	4
1.1	What is this?	4
1.2	Features	4
1.3	Difference between QCF90 and PREGSF90	5
1.4	Condition of use	6
2	Basic usage	7
2.1	Overview of the program	7
2.1.1	Option specification	7
2.1.2	Help	8
2.1.3	Script file	8
2.1.4	Dry-run mode	9
2.2	Quality control of genomic and pedigree data	9
2.2.1	Genomic data	9
2.2.2	Pedigree data	11
2.2.3	Both genomic and pedigree data	11
2.2.4	Already renumbered files	12
2.2.5	Optional map file	12
2.3	Status file and clean files	13
2.3.1	1-pass way	13

2.3.2	With a status file	14
3	Options	14
3.1	Source file specification	15
3.1.1	SNP Marker file	15
3.1.2	Pedigree file	15
3.1.3	Cross-reference ID (XrefID) file	16
3.1.4	Map file	16
3.1.5	Pre-calculated allele frequency (MAF) file	16
3.1.6	Pre-generated status file	16
3.2	Quality-control for markers and animals	17
3.2.1	Quality control items	17
3.2.2	A criterion of call rate for markers	17
3.2.3	A criterion of call rate for animals	18
3.2.4	A criterion of allele frequency for markers	18
3.2.5	Criterion of Hardy-Weinberg statistic for markers	18
3.2.6	List for unqualified markers	18
3.2.7	List for unqualified animals	19
3.2.8	Removal of unqualified markers	19
3.2.9	Removal of unqualified animals	20
3.2.10	Precise checks for input files	20
3.3	Checks for identical (duplicated) genotypes	20
3.3.1	Genomic identity search	20
3.3.2	Threshold of identity	21
3.3.3	Use of randomly selected markers	21
3.3.4	Use of user-specified markers	21
3.3.5	Marker ID used in the check	22
3.4	Chromosome configuration	22
3.4.1	Sex chromosome	22
3.4.2	Removal of markers on particular chromosomes	22
3.4.3	Skipping markers on particular chromosomes	22
3.5	Pedigree configuration	23
3.5.1	Position of animals	23
3.5.2	Unknown parents	23

3.6	Job control and performance options	24
3.6.1	Script file	24
3.6.2	Strict check	24
3.6.3	Dry run	24
3.6.4	Fast reading the marker file	24
3.6.5	Internal format of marker-storage on memory	25
3.6.6	Number of threads	26
3.7	File output options	26
3.7.1	log file	26
3.7.2	status file	26
3.7.3	“dot” file	27
3.7.4	Create clean files	27
3.7.5	Removing markers with specific status	27
3.7.6	Removing animals with specific status	27
3.8	Miscellaneous	28
3.8.1	Version number	28
3.8.2	Help message	28
3.8.3	Comprehensive help message	28
4	Technical details	28
4.1	Algorithms	28
4.2	Workflow of the program	29
4.3	File format	31
4.3.1	Status file	31
4.3.2	SNP marker file	33
4.3.3	Allele frequencies	34

1 Introduction

1.1 What is this?

QCF90 is a software tool for quality control on genomic and pedigree data used, especially in animal breeding. The “quality-control” detects and optionally removes errors or inaccurate information in the data. The data include single nucleotide polymorphisms (SNP) as genetic markers and a pedigree list with an animal and its sire and dam. The program reports markers and genotyped animals with problems, e.g., low call-rate, Mendelian inconsistency among relatives, et cetera. The program can remove such unqualified markers and individuals from the original files.

QCF90 can efficiently handle a large number of animals genotyped with 50K SNP chips. The program stores a genotype as a 2-bit code, and multiple genotypes are packed into an array in memory. Some quality controls are performed directly for the packed genotypes using bit-wise operations. This strategy saves both computing time and memory requirements.

QCF90 is a command-line tool like Linux/Unix tools and some scientific software (e.g., `plink`). The user supplies all required information as command-line options. The program generates similar reports as PREGSF90, a BLUPF90 family program. The manual explains how the two programs differ in a later section.

It is written in Fortran 2003 and compiled with Intel Fortran Compiler. The compiled binary is available on the website of Animal Breeding and Genetics Group at the University of Georgia (<http://nce.ads.uga.edu>) for Linux, macOS, and Windows.

1.2 Features

QCF90 checks the physical format for genomic data. Once it confirms that the format is correct, QCF90 calculates the following several statistics to determine “unqualified” markers and genotyped animals. A user can compute specific statistics only.

- Call rate for each marker
- Call rate for each animal

- Statistic for Hardy-Weinberg equilibrium (χ^2 value)
- Allele frequency

The program can also check the duplicated genotypes (identical genomic data) in the data. Only the check can use a subset of markers. The following two modes are available.

- Duplication checks within a genomic file
- Duplication checks between two files

QCF90 optionally accepts a pedigree file. The program checks the physical format of the file and finds logical errors (e.g., loops) in pedigree. If the user provides both genomic and pedigree files, the program performs the following checks.

- Mendelian conflicts between a parent and its progeny

The supplied pedigree file is not necessarily pre-processed by RENUMF90. The pedigree file is allowed to have alphabets and symbols as well as integers.

After the checks, the program makes a report (a status file) on the marker quality and consistency. The program creates “clean” files by removing all “unqualified” markers and animals from the original files. The clean file is available in the same run on quality control. The program can also read the status file and create clean files without computations for quality control in a separate run. In this way, the user can check the details and remove particular markers and animals with specific failures.

1.3 Difference between QCF90 and PREGSF90

All BLUPF90 family programs can check the quality of genomic data, remove “unqualified” markers, and calculate relationship matrices. PREGSF90 is an independent tool only to performs genomic-related checks and computations. This program has the same interface and design as BLUPF90, so it is useful when all files have already prepared for the BLUPF90 programs. The program assumes that all the required files have been processed with RENUMF90.

QCF90 has the same functionality and gives (almost) the same results as PREGSF90. If the file format is correct, the user does not have to use RENUMF90 to prepare the files. QCF90 is newly designed and specialized to quality check on large-scale genomic data.

See the following descriptions for details.

- Interface. PREGSF90 needs a parameter file to describe the instruction for quality-control. QCF90 directly reads options from the command line.
- Files. PREGSF90 always needs at least 4 files, i.e., phenotypic, pedigree, genomic, and cross-reference (XrefID) files; no matter what, you do not need such files. QCF90 needs a genomic file, or a pedigree file, or both files.
- Format. PREGSF90 requires files prepared by RENUMF90 (a “renumbering” process). QCF90 does not need the “renumbering” process. It can also accept the “renumbered” files.
- Checks. PREGSF90 shows minimal messages on erroneous genomic data and none on abnormal pedigree. QCF90 shows explicit messages both on genomic and pedigree errors.
- Efficiency. PREGSF90 needs a large amount of memory if the data is significant. QCF90 needs less memory and runs faster in many cases.
- Purposes. PREGSF90 performs more than quality control. QCF90 focuses on quality control.

1.4 Condition of use

This program is allowed to use for non-commercial, academic, or personal purposes only. For the other use, please contact Ignacy Misztal (ignacy@uga.edu) for details, unless you have already had a license to use BLUPF90 programs in your purpose.

The software is provided “as is,” without warranty of any kind, express or implied fitness for a particular purpose, and non-infringement. In no event shall the authors be liable for any claim, damages or other liability, whether in an action of contract, tort or otherwise, arising from, out of or in connection with the software or the use or other dealings in the software.

2 Basic usage

2.1 Overview of the program

2.1.1 Option specification

QCF90 is a command-line tool like Linux/Unix commands. The user should give all required information as command-line options in your terminal (or a command window). Here we introduce a general usage of the program. The next section shows the options related to actual quality control.

The program can be invoked by typing the command name `qcf90`.

```
qcf90
```

By default, the program shows a list of options and stops immediately. The program performs quality control by putting the option characters, followed by `qcf90`. For example, the following command performs the quality control for genomic data stored in a file `snpdata.txt` (here, we do not explain what this command does; see the later sections for details).

```
qcf90 --snpfile snpdata.txt
```

The sequence of characters `--snpfile` is called *option*. The option always starts with two letters `--`. For this particular option, an argument (`snpdata.txt`) provides additional information. Some options do not need an argument, and some need two or more arguments. The user should put one or more white spaces among the command name (`qcf90`), the options, and the arguments.

The following example specifies two options in the same command line.

```
qcf90 --snpfile snpdata.txt --pedfile pedigree.txt
```

The user should put white spaces to separate between the option and arguments. This program assumes the marker and the pedigree file share the same ID for animals.

The program accepts the files generated with `renumf90`. You should specify the cross-reference (XrefID) file when you combine the marker and the pedigree files. See the later sections for details.

2.1.2 Help

The option `--help` shows selected options on the screen, and the program stops.

```
qcf90 --help
```

With this option, the program ignores all the other options. The Complete option-list is available with `--long-help`.

```
qcf90 --long-help
```

2.1.3 Script file

QCF90 can read options from a file as well as the command line. The file is a *script file* as used in the `plink` software. In a script file, each line has to have one option with its arguments. For example, assume a script file `make_clean_files.txt` with the following lines.

```
--snpfile snpdata.txt  
--pedfile pedigree.txt  
--save-clean
```

Use `--script` to specify the script file.

```
qcf90 --script make_clean_files.txt
```

It is equivalent to the following command.

```
qcf90 --snpfile snpdata.txt --pedfile pedigree.txt --save-clean
```

QCF90 reads the options from a script file the first, then from the command line the next. If the script file and the command line have the same option with different arguments, the program primarily uses the command-line option. In other words, the command line options overwrite the options in the script file.

2.1.4 Dry-run mode

QCF90 has many options to control the computations. It is useful for the user to know what is happening with specified options without any actual computations nor modifications. QCF90 has an option `--dry-run` to do this. The user can specify regular options in addition to the dry-run option.

```
qcf90 --snpfile snpdata.txt --pedfile pedigree.txt --dry-run
```

QCF90 kindly shows what it calculates and which file it creates. The program shows all messages on the screen, and it does not create any files. This option is recommended to all users before running strict quality controls.

2.2 Quality control of genomic and pedigree data

2.2.1 Genomic data

Here we use an example marker file `snpdata.txt`. The first several lines in this file are as follows.

```
A57309 122211022220002111112212111011012011011
A57310 122211022220002111112212111021012011011
A57311 112200110211102200110122121000100020202
A57312 122211022220002111112212111011012011011
A57313 012211112221111100111122021000022111101
```

The file should have the same format for the BLUPF90 family programs (See the technical details chapter for details). Briefly, the data consists of at least two fields per row: 1) ID for genotyped animal and 2) the animal's genotypes. The marker data must start with the same position (for example, the above genotypes starts at the 9th column) in all rows. The file can have any fields between ID and genotypes, but the user inserts spaces between the fields.

The option `--snpfile` specifies a marker file. The following example reads a marker file `snpdata.txt` and checks the format, applies the default quality control to the genomic data, and writes a report.

```
qcf90 --snpfile snpdata.txt
```

By default, the program performs through the following menu.

1. Check for format for genomic data.
2. Calculate call rate, allele frequency, and Hardy-Weinberg statistics for all available markers.
3. Remove markers with a low call-rate.
4. Remove markers with low allele frequency.
5. Remove markers with 0 or 1 allele frequency (monomorphic markers).
6. Calculate call rate for animals using remained markers.
7. Remove animals with a low call-rate.
8. Recalculate allele frequency with remaining markers and animals.
9. Make a report.

An option defines criteria to remove markers and animals. Once markers or animals get unqualified, the program excludes such objects in the subsequent calculation. QCF90 stops if there is a severe problem (e.g., format error) in the file. In this case, the program will show a detailed explanation of what is wrong and how to fix the problem.

During the quality control, the program creates the following five files by default.

- `freqdata.count` = allele frequency for all available markers (saved in step 2 above)
- `Gen_call_rate` = call rate for animals (saved in step 6)
- `freqdata.count.after.clean` = allele frequency after quality control (saved in step 8)
- `qcf90.status` = status for markers and animals (saved in step 9)
- `qcf90.log` = a copy of the message displayed on the screen

The first three files are compatible with the files created with PREGSF90. The status file is readable and contains the test-result for each marker and animal.

2.2.2 Pedigree data

The program can check pedigree information without genomic data. The option `--pedfile` specifies a pedigree file with at least three fields for animal, sire, and dam identification codes. Hereafter, we use `pedigree.txt` as a pedigree file with the following lines.

```
A70269 A58649 A63290
A44516 A36831 A29636
A7769 A8 A999
A19781 A14427 A8539
A52523 A32283 A25587
```

By default, the 1st field is for animal ID, the 2nd field for sire ID, and the 3rd field for dam ID (you can change the order using an option). The ID can contain any alphabets, symbols, and numbers (but tested only with ASCII characters). The following command entirely checks the pedigree information.

```
qcf90 --pedfile pedigree.txt
```

The program checks three flaws in pedigree. The first one is “missing entries,” i.e., an animal that appears as a parent in any entries but does not have its entry. The second is sire-dam inconsistency, i.e., an animal that appears as a sire in some entries but as a dam in other entries. The third is pedigree loops, i.e., an animal that appears in its ancestor. The program shows warning with case 1 and 2; it stops in case 3.

QCF90 writes a file `qc_pedigree_loop.dot` if there are pedigree loops. An external software, GraphViz (<http://www.graphviz.org>), can read this file and draw a graph that shows the loops. It helps the user to find a problematic animal in a pedigree loop.

2.2.3 Both genomic and pedigree data

The user can specify both options in the same command line. The program checks Mendelian inconsistency in addition to the standard genomic and pedigree checks.

```
qcf90 --snpfile snpdata.txt --pedfile pedigree.txt
```

The program checks Mendelian inconsistency (conflicts) among pairs of parent-progeny in addition to separate quality control for genomic and pedigree data. The program performs this Mendelian check just after removing animals with low call-rate (after step 7 in the genomic quality control shown above).

The program assumes the pedigree is correct and compares genotypes between a parent and its progeny. The algorithm is described in Hayes (2011) and is equivalent to the implementation in PREGSF90. The following file is available as a report.

- `Gen_conflicts` = detailed information for Mendelian checks

Note that the animal ID must be common both in genomic and pedigree files. It is not a case if the pedigree file is prepared with RENUMF90. If using such a file, the user should supply the XrefID file (See the next section).

2.2.4 Already renumbered files

QCF90 accepts files prepared (“renumbered”) with RENUMF90. QCF90 recognizes this situation by specifying an XrefID file using the option `--xrefid`. In such a case, the program cares for both the original animal ID (in genomic data) and the renumbered ID (in pedigree data). The following example uses the XrefID file.

```
qcf90 --snpfile snpdata.txt --pedfile renadd02.ped --xrefid snpdata.txt_XrefID
```

The file `renadd02.ped` is supposed to be generated with RENUMF90.

2.2.5 Optional map file

QCF90 can read a map file for markers. The map file is useful when some markers are on the sex chromosomes that should be excluded from the Mendelian checks. The map file has at least three fields with white-space separators: the first field as the sequential integer number of the marker, the second field, the chromosome number (integer), and the third field, the physical location on the chromosome.

The map file can have the optional fourth field for the marker name. Note that qcf90 uses only the chromosome code (the 2nd column), and you can put any data in the other columns. The following is an example of a map file; see the BLUPF90 manual for details.

```
1 1 125
2 1 1028
3 1 3372
(snip)
5720 7 63701
5721 7 65417
```

The option `--mapfile` specifies the map file and `--sex-chr` determines the sex-chromosome numbers. See the following example.

```
qcf90 --snpfile snpdata.txt --mapfile map.txt --sex-chr 30
```

In this example, QCF90 relates a marker and a chromosome and knows the chromosome number 30 is sex-related. The program excludes those markers on the sex chromosomes temporarily from the Mendelian checks. The option `--sex-chr` accepts multiple numbers separated by comma (or space) like `--sex-chr 30,31` or `--sex-chr 30 31`.

2.3 Status file and clean files

QCF90 puts a status code (so-called “flag”) for “unqualified” markers and animals and removes them from the subsequent quality controls. By default, QCF90 does not change the original files, but it reports the flags in a status file `qcf90.status` when the program stops successfully. There are two ways to remove the “unqualified” markers and animals from the original files and generate new files to have the qualified subjects.

2.3.1 1-pass way

The first way is to do it in 1-pass. If you put the option `--save-clean`, the program creates “clean” files containing only qualified markers and animals. See

the following example.

```
qcf90 --snpfile snpdata.txt --pedfile pedigree.txt --save-clean
```

The program creates a new marker, XrefID, and map files when the user specifies the corresponding input option. The name of a clean file is the original name plus a suffix `_clean`. In the above case, the clean genomic file is `snpdata.txt_clean`. The user can change the suffix with an additional argument with `--save-clean`. For example, the following command changes the suffix to `.new`.

```
qcf90 --snpfile snpdata.txt --pedfile pedigree.txt --save-clean .new
```

2.3.2 With a status file

The second option is to use the status file. When the user specifies a status file with `--statfile`, QCF90 reads all flags from the file without performing any quality-control computations. In this case, the user can specify particular flags of markers and animals, and the program excludes the “unqualified” objects from the clean files. The following command creates clean files by removing markers with flags 1, 2, and 5 and animals with flag 1 (the following command-line is wrapped, and the user should type it in one line).

```
qcf90 --snpfile snpdata.txt --pedfile pedigree.txt --save-clean  
      --statfile qcf90.status  
      --cleanup-marker-flags 1,2,5 --cleanup-animal-flags 1
```

This operation is useful for the selective removal of markers and animals.

3 Options

Brief descriptions for frequently-used options are available with the `--help` option on the command line. Option `--long-help` shows all available options.

3.1 Source file specification

3.1.1 SNP Marker file

`--snpfile markerfile`

The argument `markerfile` is the name of a marker file. The file format should be the same as the one used in the BLUPF90 programs. See the BLUPF90 manual or Technical Details in this manual for details.

The program can perform the quality control on markers and animals: physical format error, low call rate, low allele frequency, and Hardy-Weinberg equilibrium test. Tests for Mendelian inconsistency is available only when a pedigree file is specified.

3.1.2 Pedigree file

`--pedfile pedigreefile`

The argument `markerfile` is the name of a pedigree file. The file must have at least three fields for an animal's ID, its sire's ID, and its dam's ID. The user can specify the position of each ID using `--ped-position`. A missing parent must have a single character 0 or negative integer (- at the beginning of ID and followed by an integer code); the program does not accept any other characters as missing ID by default. For example, Any of 00, NA, or . are never recognized as a missing animal. The code for missing (unknown) parents can be changed with an option `--upg` (see later sections).

The pedigree does not have to be prepared ("renumbered") with RENUMF90. Any characters will be used as ID (except 0 for missing animals). If a genomic file is also supplied, the ID must be shared in both files. It is not a case when the "renumbered" pedigree is given. In this case, put `--xrefid` with an XrefID file.

The program can perform the following checks for pedigree: physical format error, confusion of sire and dam, missing entry of parents, and pedigree loops. The program saves the structure of pedigree loops to a file `qc_pedigree_loops.dot` that Graphviz (<http://www.graphviz.org>) can read and draw the pedigree chart. Tests for Mendelian inconsistency is available only when a genotype file is specified.

3.1.3 Cross-reference ID (XrefID) file

`--xrefid xrefidfile`

The argument `xrefidfile` is the name of a cross-reference ID (XrefID) file generated with RENUMF90. This option is required when using the “renumbered” pedigree created with RENUMF90.

3.1.4 Map file

`--mapfile mapfile`

The argument `mapfile` is the name of a map file. The program requires the marker file when markers on sex chromosomes should be excluded in Mendelian checks.

The file format should be the same as the one used in the BLUPF90 programs. The file should contain at least three fields: 1) a sequential number, 2) a chromosome number, and 3) a physical location on the chromosome. All numbers should be integer values. The additional fourth field can have any characters for the marker name. The number of rows in this file should be the same as the number of markers in the marker file.

3.1.5 Pre-calculated allele frequency (MAF) file

`--maffile maffile`

The argument `maffile` is the name of a file containing pre-calculated allele frequency for each marker. The file contains at least one field, and the first field should be allele frequency, a real number. The number of rows in this file should be the same as the number of markers in the marker file.

3.1.6 Pre-generated status file

`--statfile statusfile`

The argument `statfile` is the name of a status file generated with `qcf90` in the previous run. `qcf90` reads status (qualifications) on markers and animals instead of calculating quality-statistics from the data. Then `qcf90` writes “clean” files without unqualified markers and animals. Therefore, this option needs the marker file specified with `--snpfile`.

With this option, `qcf90` does not create a new status file.

3.2 Quality-control for markers and animals

3.2.1 Quality control items

`--qc items`

The arguments `items` specify which quality control the program applies to the genomic data. Multiple arguments can be used; the arguments should be separated by a comma (,) or white space (.). The following items are acceptable.

Item	Default	Criterion	Description
<code>crm</code>	Yes	0.10	low call-rate for markers
<code>maf</code>	Yes	0.05	low allele frequency (MAF) for markers
<code>mono</code>	Yes	0.00	monomorphic markers (0 or 1 MAF)
<code>hwe</code>	No	15.00	Hardy-Weinberg test for markers
<code>cra</code>	Yes	0.10	low call-rate for animals
<code>par</code>	Yes	0.01	Mendelian inconsistency for markers and animals

The other options can change the qualifying criteria (see below for details). Detailed explanations for each quality control is available in Technical Details.

3.2.2 A criterion of call rate for markers

`--crm n`

The argument `n` specifies the criterion of “unqualified” call-rate for markers. The criterion should be a real number between 0 and 1; the default value is 0.10.

If the call-rate is lower than this value, the marker is recognized as “unqualified” and removed from the subsequent quality control.

3.2.3 A criterion of call rate for animals

`--cra n`

The argument `n` specifies the criterion of “unqualified” call-rate for animals. The criterion should be a real number between 0 and 1; the default value is 0.10. If the call-rate is lower than this value, the animal is recognized as “unqualified” and removed from the subsequent quality control.

3.2.4 A criterion of allele frequency for markers

`--maf n`

The argument `n` specifies the criterion of allele frequency for markers. The criterion should be a real number between 0 and 1; the default value is 0.10. If the AF is lower (or $1 - AF$ is higher) than this value, the marker is recognized as “unqualified” and removed from the subsequent quality control.

3.2.5 Criterion of Hardy-Weinberg statistic for markers

`--hwe n`

The argument `n` specifies the criterion of unqualified Hardy-Weinberg statistic for markers. The criterion should be a real number between 0 and 1; the default value is 0.10. If the statistic is bigger than this value, the marker is recognized as “unqualified” and removed from the subsequent quality control.

3.2.6 List for unqualified markers

`--exclude-marker-list markerlist`

This option forces specific markers supplied by `markerlist` to be unqualified. In other words, the markers on the list are unqualified and removed from the entire quality-control. This option is useful, primarily when you have known some markers should be removed from the marker file. Note that the file format depends on which other options are specified.

- If the user gives a map file, `markerlist` must be the same format as the map file.
- Otherwise, `markerlist` contains integer numbers corresponding to the positions of markers in the genomic file.

3.2.7 List for unqualified animals

`--exclude-animal-list animallist`

This option forces specific animals supplied by `animallist` to be unqualified. In other words, the animals on the list are unqualified and removed from the entire quality-control. This option is useful primarily when you have known some animals should be removed from the marker file. Note that the file format depends on which other options are specified.

- If the user gives a cross-reference ID (XrefID) file, `animallist` must be the same format to the XrefID file.
- Otherwise, `animallist` contains integer numbers corresponding to the positions of the animals in the genomic file (i.e., in which rows the animals locate).

3.2.8 Removal of unqualified markers

`--remove-markers`

`--no-remove-markers`

By default, once the program recognizes a set of markers as unqualified, it excludes the markers from the subsequent quality-control. The option (`--remove-markers`) supports the default behavior. If the other option (`--no-remove-markers`) is

specified, the unqualified markers are NOT removed from the subsequent quality control. In other words, the unqualified markers are still involved in all quality control, although such a marker has a status, “unqualified”.

3.2.9 Removal of unqualified animals

```
--remove-animals  
--no-remove-animals
```

By default, once the program recognizes a set of animals as unqualified, it excludes the animals from the subsequent quality control. The option (`--remove-animals`) supports the default behavior. If the other option (`--no-remove-animals`) is specified, the unqualified animals are NOT removed from the subsequent quality control. In other words, the unqualified animals are still involved in all quality control, although such an animal has a status, “unqualified.”

3.2.10 Precise checks for input files

```
--check-format  
--no-check-format
```

By default, the program precisely checks the physical format for all specified files. The option (`--check-format`) supports this default behavior. If the other option `--no-check-format` turns off this feature and just quickly checks the format. It may save running time if the user knows the file format is correct.

3.3 Checks for identical (duplicated) genotypes

3.3.1 Genomic identity search

```
--check-identity  
--check-identity secondary_markerfile
```

This option checks identical genotypes among genotyped animals. Just specified `--check-identity` without an argument, the program checks all possible

pairs of animals in the original marker file supplied with `--snpfile`. For example, when the SNP file has n animals, the program compares $n(n - 1)/2$ pairs of animals.

If the argument `secondary_markerfile` is specified, the program compares the genotypes between the primary and the secondary files. In this case, the program does not compare any pairs within the same file. For example, when the primary marker file has n animals and the secondary file has m animals, the number of comparisons is mn .

3.3.2 Threshold of identity

`--threshold-identity n`

This option specifies a threshold of the identity statistic to define that two genotypes are identical. The program calculates the identity statistic as the number of the identical genotypes divided by the number of compared genotypes excluding missing. When the statistic exceeds the threshold, two genotypes are considered as duplicated. The default value is 0.99.

3.3.3 Use of randomly selected markers

`--random-markers-identity n`

The identity check takes a long time if all the markers are used in the check. This option defines the number of markers used in operation. The program randomly chooses the markers from all possible loci.

3.3.4 Use of user-specified markers

`--fixed-markers-identity markerlist`

This option specifies a set of markers used in the identity check. The user should supply the name of the marker list. The format of the list file is the same as explained in the option `--exclude-marker-list` (see this option for details).

3.3.5 Marker ID used in the check

`--save-markers-identity`

This option saves the marker ID used in the identity check. It may be useful when you use random markers.

3.4 Chromosome configuration

3.4.1 Sex chromosome

`--sex-chr n . .`

The arguments `n . .` specify the numbers corresponding to the sex-chromosome code defined in a map file. With this option, the user must provide a map file. Multiple arguments can be used; the arguments are separated by a comma (,) or white space ().

The markers on the sex chromosomes are temporarily excluded from the Mendelian checks.

3.4.2 Removal of markers on particular chromosomes

`--exclude-chr n . .`

The arguments `n . .` specify the numbers corresponding to the chromosome code defined in a map file. With this option, the user must give a map file. Multiple arguments can be used; the arguments are separated by a comma (,) or white space ().

The program removes the markers permanently. The markers are marked as *removed-by-user*.

3.4.3 Skipping markers on particular chromosomes

`--skip-chr n . .`

The arguments `n . .` specify the numbers corresponding to the chromosome code defined in a map file. With this option, the user must give a map file. Multiple arguments can be used; the arguments are separated by a comma (,) or white space ().

This option prevents QCF90 to remove the markers on the chromosomes even though those markers are unqualified. Note that this option simply passes through the markers to output; the program does include the markers in all quality control, and the markers may impact the resulting statistics. Currently, there is no way not to apply quality control to particular markers because of the performance issue.

3.5 Pedigree configuration

3.5.1 Position of animals

`--ped-position a,s,d`

The arguments `a` specifies the position of an animal's ID in a pedigree file and `s` for the sire's ID, and `d` for the dam's ID. All three positions should be different and must exist in the pedigree file.

3.5.2 Unknown parents

`--upg zero`

`--upg negative`

`--upg zero+negative`

The argument specifies what kind of ID should be treated as missing (unknown) parents.

Argument	Missing parents
<code>zero</code>	a single 0
<code>negative</code>	the negative sign - and integer code in the subsequent part
<code>zero+negative</code>	a combination of the above two (default)

3.6 Job control and performance options

3.6.1 Script file

`--script file`

QCF90 reads options from the script `file` specified here. Each row in the script has one option and its argument(s).

Even with this option, QCF90 still reads options from the command line. If the same option is specified in both sources, the command line option is used.

3.6.2 Strict check

`--halt-on-warning`

QCF90 gives you warnings and continues the process if the issue is not crucial. This option makes the program stop when any small problems are detected.

3.6.3 Dry run

`--dry-run`

QCF90 shows what is happening with specified options without any actual computations nor modifications. This option is useful to make sure the options work correctly as the user intends to do.

3.6.4 Fast reading the marker file

`--fastread`

It accelerates to read a marker file when using the packed format. This option is always recommended for any users unless there is any trouble.

3.6.5 Internal format of marker-storage on memory

`--marker-storage type`

The argument `type` specifies which marker-storage strategy the program uses. The following types are available.

Type	Default	Description
real	No	Use a double-precision (8-byte) variable per marker
int	No	Use a 1-byte integer variable per marker
pack	Yes	Use a 4-byte integer variable per 16 markers

The user does not have to care much about this option. See the Technical Details section for details.

3.6.6 Number of threads

`--num-threads n`

This option defines the number of threads used in this program if a functionality supports multi-threading. Currently, only the identity (duplication) check supports the parallel computation.

3.7 File output options

3.7.1 log file

`--save-log file`
`--no-save-log`

QCF90 saves all the log to a file `qcf90.log` by default. This option can change the name of the log file. You can also suppress to create the log with `--no-save-log`.

3.7.2 status file

`--save-status file`
`--no-save-status`

QCF90 creates a status file `qcf90.status` by default. This option can change the name of the status file. You can also suppress to create the log with `--no-save-log`.

3.7.3 “dot” file

```
--save-dot file  
--no-save-dot
```

QCF90 creates a “dot” file by default if the pedigree information has a loop. The dot file contains a graph written in “dot” language, which is interpreted by Graphviz (<http://www.graphviz.org>). The file name is the original pedigree file + dot. This option can change the name of the dot file. You can also suppress to create the log with `--no-save-dot`.

3.7.4 Create clean files

```
--save-clean  
--save-clean suffix
```

QCF90 saves “clean” files after removing all unqualified markers and animals from the original genomic, map, and the cross-reference-ID (XrefID) files. The name of a new file is the same as the original file with a suffix `_clean`. If an optional argument `suffix` is supplied, the program uses those characters as the suffix.

3.7.5 Removing markers with specific status

```
--cleanup-marker-flags n..
```

When saving clean files, the program removes only particular markers with the status specified here. The arguments `n..` are integer status codes (flags) to describe the qualification of a marker. Multiple arguments can be used; the arguments are separated by a comma (,) or white space (.). See the Technical Details for the status code.

3.7.6 Removing animals with specific status

```
--cleanup-animal-flags n..
```

When saving clean files, the program removes only particular animals with the status specified here. The arguments `n..` are integer status codes (flags) to describe the qualification of an animal. Multiple arguments can be used; the arguments are separated by a comma (,) or white space (). See the Technical Details for the status code.

3.8 Miscellaneous

3.8.1 Version number

`--version`

QCF90 shows the current version of the program and immediately stops.

3.8.2 Help message

`--help`

QCF90 shows a list of frequently-used options and stops.

3.8.3 Comprehensive help message

`--long-help`

QCF90 shows a list of all available options and stops.

4 Technical details

4.1 Algorithms

The program uses bit-wise operations on the packed markers. The computational detail is explained in Masuda et al. (in review).

4.2 Workflow of the program

The following is the complete description of the workflow. The program skips some steps unless appropriate options are specified (e.g., reading a pedigree file only with the `--ped` option). When the program defines a set of markers and animals as “unqualified,” they are excluded from the subsequent quality-control by default. If serious problems or errors occur, the program immediately stops with a diagnosis.

- Initialization
 - Read a script file (`--script`)
 - Read the command line
 - Check inconsistency or error in options
- Quick checks for files and some preparations
 - Read the first line of the pedigree file and check the format
 - Read the first line of the marker file and check the format
 - Read through the genomic file and count the number of lines
 - Prepares a hash table to store ID for genotyped animals
- Precise checks for genomic data
 - Read the first line of the marker file and check the format, again
 - Read through the genomic data and precisely check the format
- Load genomic data in memory
 - Allocate memory for genomic data
 - Read genomic data from the file and store it in memory
- Precise checks for pedigree
 - Read the entire pedigree from the file and store it in memory
 - Check missing entries for parents
 - Check sire-dam inconsistency
 - Check loops in the pedigree

- Precise checks for combined genomic and pedigree data
 - Check genotyped animals not found in the pedigree
- Checks for identical genotypes
 - Within/between file(s)
- Computations of statistics for markers
 - Allele frequency
 - Call rate
 - Hardy-Weinberg statistic
- Removal of specific markers specified by the user
 - Individual marker
 - Markers on particular chromosomes
- Quality control for markers
 - Low call rate (criterion<0.10; changed by --crm)
 - Low allele frequency (criterion<0.10; changed by --maf)
 - monomorphic markers
 - Hardy-Weinberg test
- Quality control for animals
 - Low call rate
- Quality control for Mendelian conflicts
 - Calculate Mendelian consistency for markers
 - Low-consistency markers
 - Calculate Mendelian consistency for animals
 - Low-consistency animals
- Post-quality-control process
 - Calculate allele frequency with qualified markers and animals
 - Save status for markers and animals to a file
 - Save clean files

4.3 File format

4.3.1 Status file

The program creates a status file as a report of quality control. It has which markers and animals are defined as unqualified and the reason why they are. It is a human-readable text file consisting of several sections. You can easily extract specific information out of it, but do not edit the file.

Comment section A *comment* is a line that the program ignores. The comment line starts with the character #, and the entire line should be ignored. At the beginning of the file, all human-readable information is in a comment section. The following example illustrates how the information is placed in the block.

```
# Quality-control status
# Oct 6, 2017 22:22:14
#
# marker file    = newsnp
#
# marker quality:
#   flag 0 = OK           :           83 markers (checked)
#   flag 1 = Call Rate    :           0 markers (checked)
#   flag 2 = MAF          :          17 markers (checked)
#   flag 3 = Monomorphic  :           0 markers (checked)
#   flag 4 = Excluded by user :         0 markers (not checked)
#   flag 5 = Mendelian error :         0 markers (not checked)
#   flag 6 = HWE          :         0 markers (not checked)
#   flag 7 = High Correlation :         0 markers (not checked)
#   flag 8 = Same Position :         0 markers (not checked)
# animal quality:
#   flag 0 = OK           :          15 animals (checked)
#   flag 1 = Call Rate    :           0 animals (checked)
#   flag 2 = Parent-progeny confl :         0 animals (not checked)
#   flag 3 = Diagonal outlier :         0 animals (not checked)
#   flag 4 = Excluded by user :         0 animals (not checked)
```

#

In this example, each marker took three tests (call rate, MAF, and monomorphic tests marked as checked), and each animal took only 1 test (call rate). The marker file has 100 SNP markers and 15 animals, of which 17 markers are unqualified, but no animals are unqualified. The comment also tells the user how to remove such unqualified objects from the marker file using this status file.

Variable section A structure of the marker file is saved just after the comment block. The number of markers and the number of animals are saved in the current version.

```
NUMBER_OF_MARKERS 100
NUMBER_OF_GENOTYPES 15
```

Marker status section The marker status section shows which marker is unqualified and why. It has three fields: 1) the sequential number of the marker, 2) the flag, and 3) the meaning of the flag. See the following example.

```
MARKER_STATUS
  1  0 OK
  2  2 MAF
  3  0 OK
  (snip..)
 99  2 MAF
100  2 MAF
```

Animal status section Like the marker status, the animal status section shows which animal is unqualified and why. It also has three fields: 1) the sequential number of the animal, 2) the flag, 3) the name of the animal in the genomic data, and 4) the meaning of the flag. See the following example.

```
ANIMAL_STATUS
  1  0 1 OK
```



```
2 0 2 OK
3 0 3 OK
(snip...)
14 0 14 OK
15 0 15 OK
```

4.3.2 SNP marker file

The SNP-marker file is a text file with at least two fields. The first field is for individual ID, and the last field is for genotypes. The user should put one or more spaces to separate the fields. The following is a simple example of a marker file.

```
1 120120120120
2 120120120120
3 120120120120
```

The field of genotypes in each row should start from the same column (i.e., fixed-width format). For example, in the above example, the genotypes start from the 6th column in all rows. In contrast to genotypes, the individual ID is not necessarily aligned. The following file is acceptable for the program (but, again, genotypes should be aligned).

```
1 120120120120
2 120120120120
3 120120120120
```

You can insert extra fields between ID and genotypes as long as genotypes are aligned to a particular column (in the above, 6th column). The following example is also valid.

```
1 A001 M 120120120120
2 B001 M 120120120120
3 C001 M 120120120120
```

There are two possible formats of genotypes.

- Integer genotypes: Each genotype has three possible values (0,1, and 2) and a missing code (5). One genotype occupies one letter in the field. No spaces between genotypes are allowed.
- Real genotypes: Each genotype can be a real number with a fixed width like 1.00 and 0.95. Keeping the compatibility with BLUPF90, the width should be 4, and the number of digits under dot should be 2. No spaces between genotypes are allowed.

The following is an example of real genotypes.

```
1 1.002.000.001.002.000.00
2 1.002.000.001.002.000.00
3 1.002.000.001.002.000.00
```

You must not mix two formats (integer and real) in the same file. The number of genotypes should be the same for all individuals.

The integer genotypes are efficiently stored in memory, but real genotypes are not. The bit-wise algorithms apply only to the integer genotypes.

4.3.3 Allele frequencies

The program computes the allele frequencies, using the original genotypes (saved to `freqdata.count`) and using the clean genotypes (`freqdata.count.after.clean`). The file has two columns: 1) the marker number and 2) its allele frequency. The marker number is a sequential integer-code, which is ordered the same as the genomic file. See the following example.

```
1 0.362348
2 0.587222
3 0.610498
4 0.783390
5 0.730564
6 0.096242
7 0.862672
8 0.717877
```

9 0.387042
10 0.174393