



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

High-quality, genome-wide SNP genotypic data for pedigreed germplasm of the diploid outbreeding species apple, peach, and sweet cherry through a common workflow

Stijn Vanderzande, Nicholas P Howard, Lichun Cai, Cassia Da Silva Linge, Laima Antanaviciute, Marco CAM Bink, Johannes W Kruisselbrink, Nahla Bassil, Ksenija Gasic, Amy Iezzoni, Eric Van de Weg, Cameron Peace

S4 File: Hands-on guideline on how to perform data curation using the steps described in this study

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS





United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Input data for Infinium arrays

RosBREED

DISEASE RESISTANCE × HORTICULTURAL QUALITY → SUPERIOR CULTIVARS



iSCAN data

Results for 1 array
(24 samples)



Sample sheet(s)
(.csv or .xls(x)format)



Name	Date modified	Type	Size
202017910001	9/26/2018 12:33 PM	File folder	
202017910002	9/26/2018 12:33 PM	File folder	
202017910004	9/26/2018 12:33 PM	File folder	
202017910005	9/26/2018 12:33 PM	File folder	
202017910007	9/26/2018 12:33 PM	File folder	
202017910014	9/26/2018 12:33 PM	File folder	
202017910021	9/26/2018 12:33 PM	File folder	
202017910025	9/26/2018 12:34 PM	File folder	
202017910036	9/26/2018 12:34 PM	File folder	
202017910037	9/26/2018 12:34 PM	File folder	
202017910042	9/26/2018 12:34 PM	File folder	
202017910045	9/26/2018 12:34 PM	File folder	
SweetCherry 062818.csv	7/20/2018 10:25 AM	Microsoft Excel C...	16 KB
SweetCherry 062818.xlsx	7/20/2018 10:21 AM	Microsoft Excel W...	29 KB

Manifest file received separately (.bpm format; describes probe content of the array)

RosBREED_Cherry_15k_15069617X343343_A1.bpm	7/20/2018 10:02 AM	BPM File	2,487 KB
--	--------------------	----------	----------
















RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

iSCAN data

In result folder of an array

Results for
1 sample

	202017910001.sdf	6/21/2018 10:10 AM	SDF File	32 KB
	202017910001_qc.txt	6/28/2018 10:04 AM	TXT File	11 KB
	202017910001_R01C01_1_Focus_scan#1_swath#1_...	6/28/2018 9:36 AM	JPEG image	214 KB
	202017910001_R01C01_1_Green.xml	6/28/2018 9:37 AM	XML File	2 KB
	202017910001_R01C01_1_Red.xml	6/28/2018 9:37 AM	XML File	2 KB
	202017910001_R01C01_1-Swath1_Grn.jpg	6/28/2018 9:36 AM	JPEG image	2,100 KB
	202017910001_R01C01_1-Swath1_Red.jpg	6/28/2018 9:36 AM	JPEG image	1,983 KB
	202017910001_R01C01_Grn.idat	6/28/2018 9:37 AM	IDAT File	174 KB
	202017910001_R01C01_Red.idat	6/28/2018 9:37 AM	IDAT File	174 KB
	202017910001_R01C02_1_Focus_scan#1_swath#1_...	6/28/2018 9:39 AM	JPEG image	210 KB
	202017910001_R01C02_1_Green.xml	6/28/2018 9:39 AM	XML File	2 KB
	202017910001_R01C02_1_Red.xml	6/28/2018 9:39 AM	XML File	2 KB
	202017910001_R01C02_1-Swath1_Grn.jpg	6/28/2018 9:39 AM	JPEG image	2,079 KB



RosBREED

DISEASE RESISTANCE × HORTICULTURAL QUALITY

Sample sheet

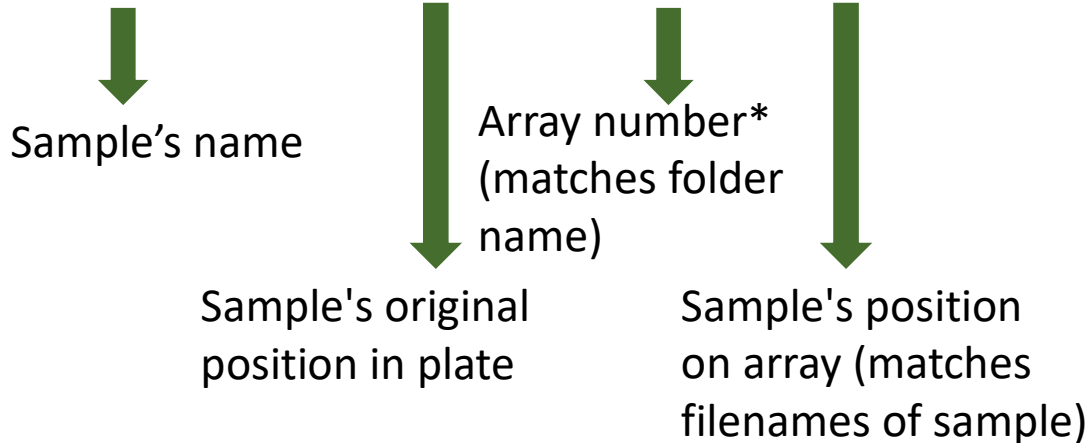
	A	B	C	D	E	F	G	H	I	J	K	L
Run info	1	[Header]										
	2	Investigator Name										
	3	Project Name										
	4	Experiment Name										
	5	Date										
Manifest info	6	**NOTE**										
	7	The following columns are required: Sample_ID, SentrrixBarcode_A, SentrrixPosition_A (and _B, _C etc.. if you have multiple manifests)										
	8	All other columns are optional; any desired additional columns may be added here or "on the fly" during analysis										
	9	Column order doesn't matter										
	10	[Manifests]										
Samples info	11	A	RosBREED_Cherry_15k_15069617X343343_A1									
	12	[Data]										
	13	Sample_ID	SampSample_Well	SentrrixBarcode_A	SentrrixPo	Gender	Sample_C	Replicate	Parent1	Parent2	Source	
	14	Rainier	A01	202017910001	R01C01		Chip1				Cherry	
	15	Sandra_Rose	A02	202017910001	R02C01		Chip1				Cherry	
	16	PC7144-11	A03	202017910001	R03C01		Chip1				Cherry	
	17	7147-9	A04	202017910001	R04C01		Chip1				Cherry	
	18	PC7309-4	A05	202017910001	R05C01		Chip1				Cherry	
	19	99F/132R4	A06	202017910001	R06C01		Chip1				Cherry	
	20	Ambrunes	A07	202017910001	R07C01		Chip1				Cherry	
21	PC8011-6	A08	202017910001	R08C01		Chip1				Cherry		
22	PC8012-5	A09	202017910001	R09C01		Chip1				Cherry		



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Sample sheet – sample data

Sample_ID	Samp	Sample_Well	SentrixBarcode_A	SentrixPosition_A	Gender	Sample_C	Replicate	Parent1	Parent2	Source
Rainier		A01	202017910001	R01C01		Chip1				Cherry
Sandra_Rose		A02	202017910001	R02C01		Chip1				Cherry
PC7144-11		A03	202017910001	R03C01		Chip1				Cherry
7147-9		A04	202017910001	R04C01		Chip1				Cherry
PC7309-4		A05	202017910001	R05C01		Chip1				Cherry



Often not filled in!

SentrixBarcode_A	S
2.01382E+11	R
2.01382E+11	R
2.01382E+11	R
2.01382E+11	R
2.01382E+11	R
2.01382E+11	R
2.01382E+11	R
2.01382E+11	R

*Sometimes saved as 'scientific number' and last digits can get lost when saved as ".csv" file (make a ".xls(x)" copy!)



Sample sheet – sample data

- Sample_ID, SatrixBarcode_A and SatrixPosition_A are required, other columns are optional and many user-defined columns can be added, e.g.:
 - Alternative names
 - Tissue source
 - Subpopulation
 - Sample quality/ploidy (when known)
- Parent names must match names used under Sample_ID (if present in data set)



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Sample sheet – sample_ID

- Check ‘Sample_ID’ of samples
 - Sometimes typing errors occurred
- Many software don’t accept spaces in names
 - Create new ‘Sample_ID’ column and save original Sample_ID as ‘Sample_ID_Original’
 - Remove spaces from new ‘Sample_ID’ column
 - Optional: Create abbreviated names for long names
- Parent and replicates names must match names used under new Sample_ID (if present in data set)



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Sample sheet – sample data

16	Sample_ID	Alternative_Name	Sample_Wt	SentrixBarcode	Sentrix	Gend	Sample	Replicate
353	Rainier		A01	202017910001	R01C01		Chip1	Rainier rep01
354	Rainier rep01		B08	201903100009	R08C01		Chip3	Rainier
355	Regina		B03	202017910007	R03C02		Chip5	
356	Salmo		B09	202017910036	R09C02		Chip 9	
357	Sandra Rose		G05	201903100003	R05C02		Chip2	Sandra Rose rep01
358	Sandra Rose rep01		A02	202017910001	R02C01		Chip1	Sandra Rose
359	Santina		H02	202017910025	R02C02		Chip8	Santina rep01
360	Santina rep01		H05	201903100003	R07C02		Chip2	Santina

Will lead to 2 versions of same report

- Replicates

- Keep one original name
 - Easy to find in files generated in next steps
- Add rep numbers for other replicates
 - Refer to original name as replicate



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Combining sample sheets

- Move all array result folders into same folder (needed to load data)
 - GenomeStudio does not allow result folders to be spread across multiple subfolders
- Open first sample sheet and save first as new “.xls(x)” file (avoid loss of original sample sheet file and loss of array number)
 - Make sure data columns match
 - Paste data rows below existing data,
 - Repeat until all data is included



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

GenomeStudio®

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS





United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Loading Data

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



GenomeStudio

Previously made projects

The screenshot displays the GenomeStudio interface. At the top, a menu bar includes File, Edit, View, Analysis, Tools, Window, and Help. Below the menu bar is a 'Start Page' section with a 'Recent Projects' table. A green box highlights this table, with an arrow pointing to the text 'Previously made projects'. The table lists several projects with columns for Project, Module, Directory, and Last Access. Below the table, a 'New Project' button is highlighted with a green box, with an arrow pointing to the text 'Start new project'. A 'File' menu is open, showing options like 'New Project', 'Open Project', 'Save Project', 'Close Project', 'Page Setup', 'Print Preview', 'Print...', 'Recent Project', and 'Exit'. A 'Directory' list is visible on the right side of the menu. At the bottom left, there is a 'Log' window showing system messages. The bottom of the slide features a decorative banner with images of various fruits (peaches, cherries, strawberries) and the text 'Ro DISEAS'.

Project	Module	Directory	Last Access
Apple BR_UMN_1ay...	Genotyping	C:\Users\stijn.vanderzande\Docu...	5/1/2018 3:01 PM
Aug18 Cherry_Goodl...	Genotyping	C:\Users\stijn.vanderzande\Docu...	8/22/2018 10:57 AM
Aug18 Cherry	Genotyping	C:\Users\stijn.vanderzande\Docu...	8/22/2018 10:57 AM
Cherry G-5K	Genotyping	C:\Users\stijn.vanderzande\Docu...	3/1/2018 3:14 PM
Nerehi 24_25	Genotyping	C:\Users\stijn.vanderzande\Docu...	3/4/2018 2:00 PM
1ay_Gax3pLextra.ind	Genotyping	C:\Users\stijn.vanderzande\Docu...	12/27/2017 11:42 AM

File Edit View Analysis Tools Window Help

- New Project
- Open Project Ctrl+O
- Save Project Ctrl+S
- Save Project Copy As... Ctrl+Shift+A
- Close Project Ctrl+Shift+C
- Page Setup Ctrl+Shift+U
- Print Preview Ctrl+Shift+V
- Print... Ctrl+P
- Recent Project
- Exit Alt+F4

Directory

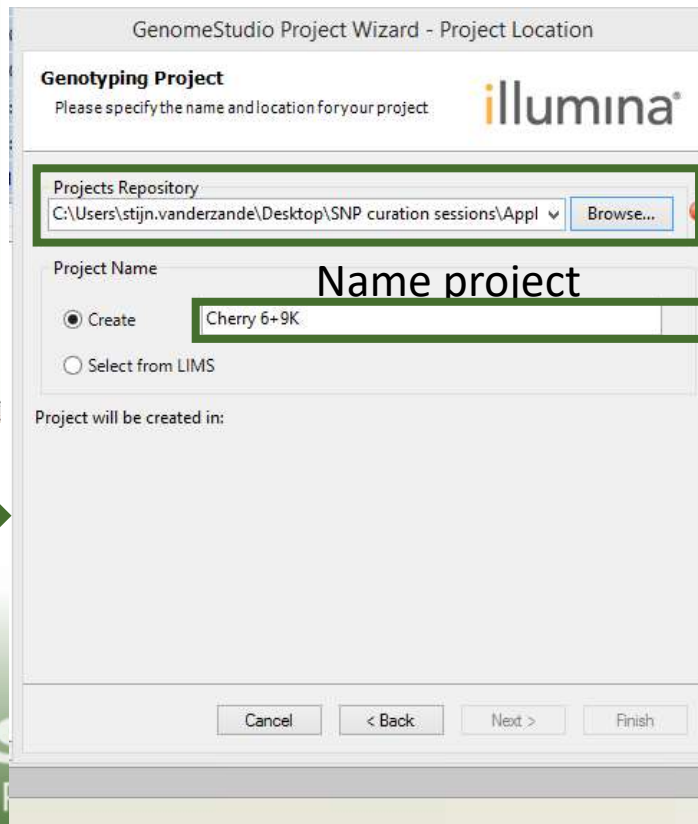
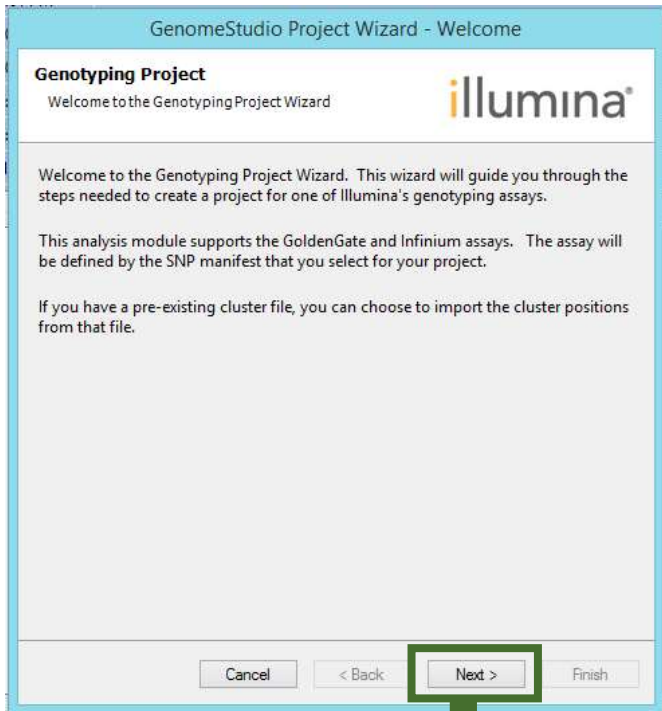
- C:\Users\stijn.vanderz
- C:\Users\stijn.vanderz
- C:\Users\stijn.vanderz
- C:\Users\stijn.vanderz
- C:\Users\stijn.vanderz
- C:\Users\stijn.vanderz

Log

Time	Severity	Message
9/27/2018 6:09:12 PM	INFO	Setting Hierarchy and Replicates
9/27/2018 6:08:12 PM	INFO	A new project has been created.

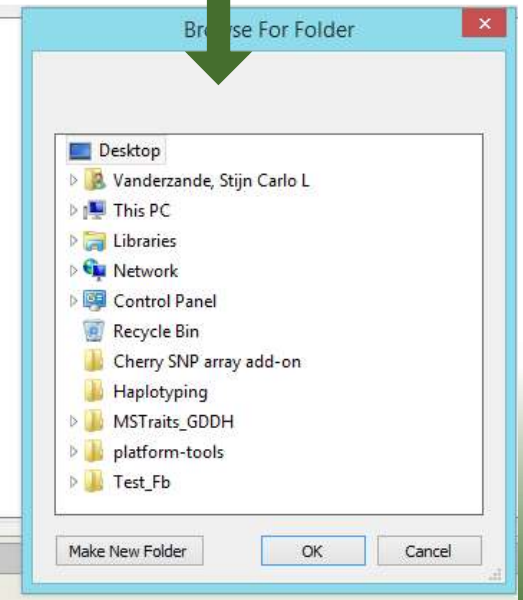
Ro DISEAS

Loading data



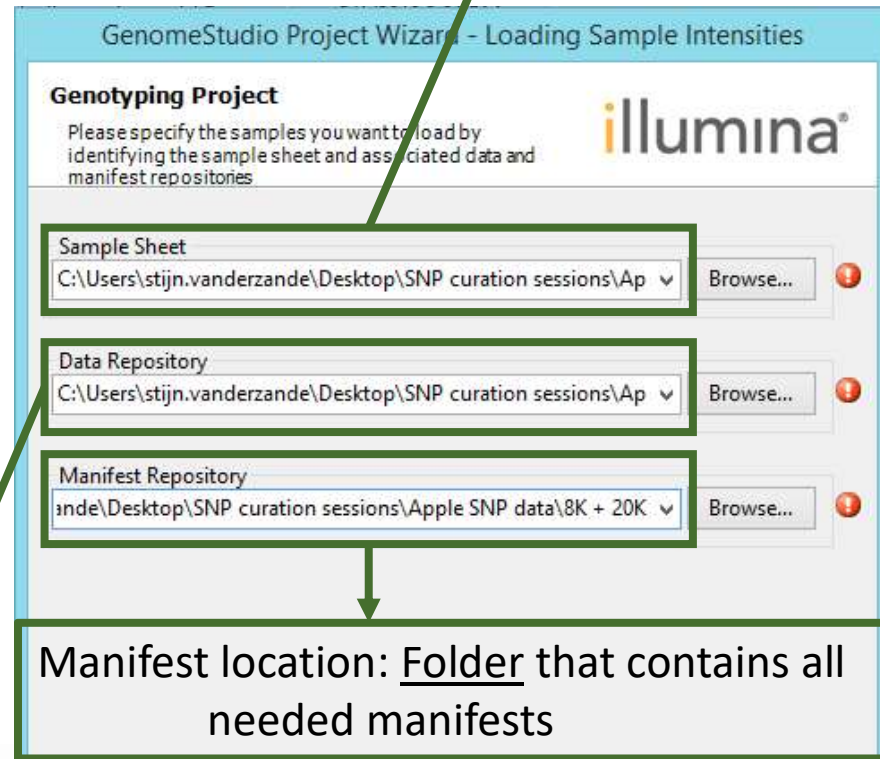
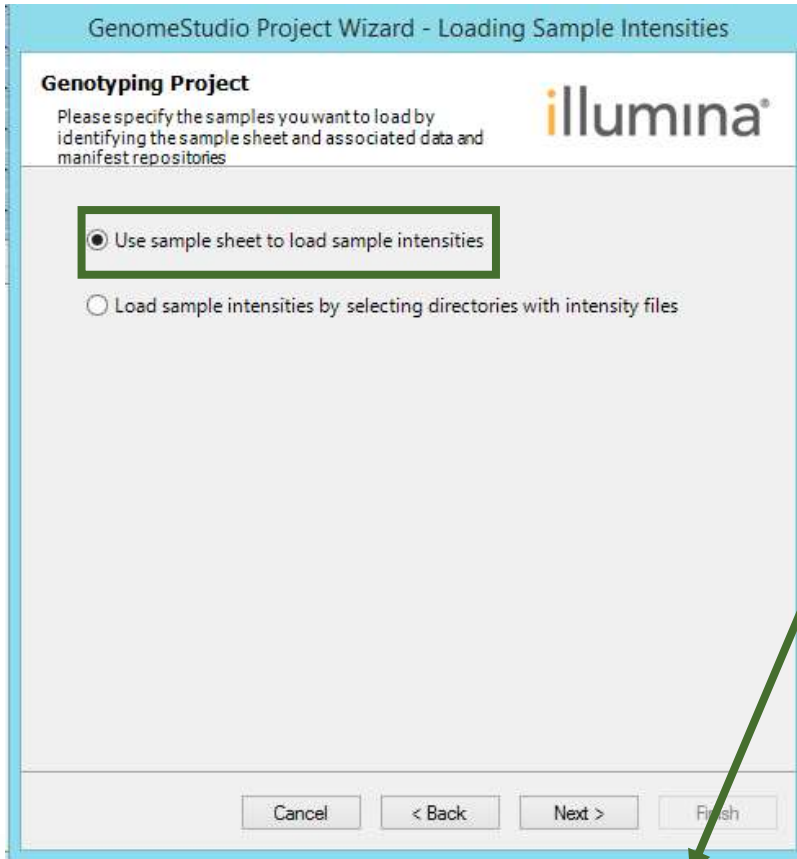
Location to save project

Name project



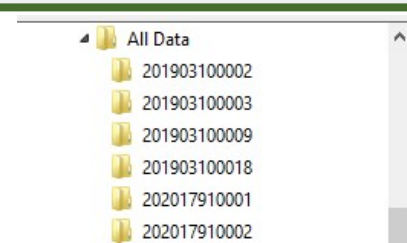
Loading data

1 sample sheet file is loaded here



Manifest location: Folder that contains all needed manifests

Data location: Folder that contains all needed SNP array folders



Loading data

GenomeStudio Project Wizard - Cluster Positions

Genotyping Project

If you have an existing cluster file that you want to import cluster positions from, enter it here. Otherwise, you can cluster the samples you've selected to determine...

Import cluster positions from a cluster file

Cluster File

Project Settings

Options

Pre-Calculate

Pre-Calculate should only be used for memory based projects.

This option will improve speed but requires 4.5x more memory.

Project Creation Actions

Cluster SNPs

Calculate Sample and SNP Statistics

Calculate Heritability

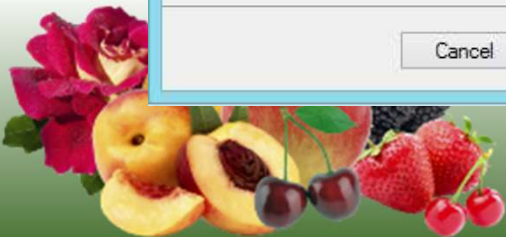
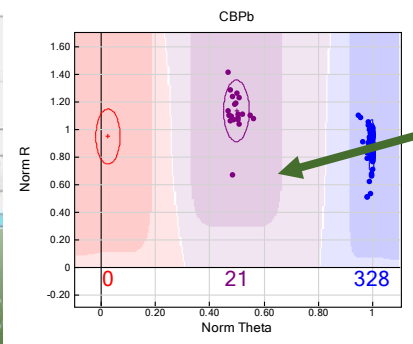
Gen Call Threshold:

Optional Script File

A file designating custom or pre-set cluster positions can be loaded here instead of the default auto-clustering

Cluster SNPs

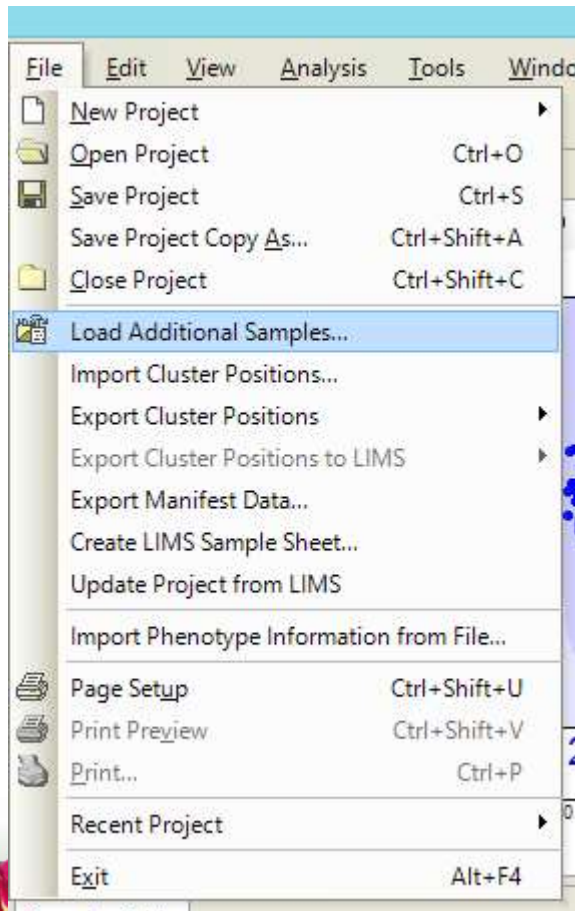
-0.15 is standard
-Defines area within which individuals will get a genotype call (outside will be 'no calls')



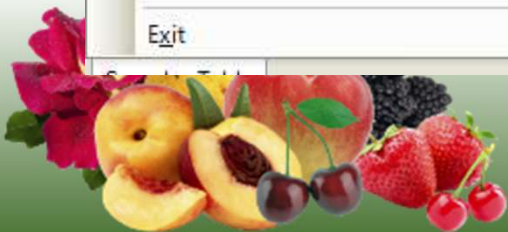
Ros
DISEASE RE

TY

Adding additional samples



- File > Load Additional Samples...
 - Same process as before



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

First look at data

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



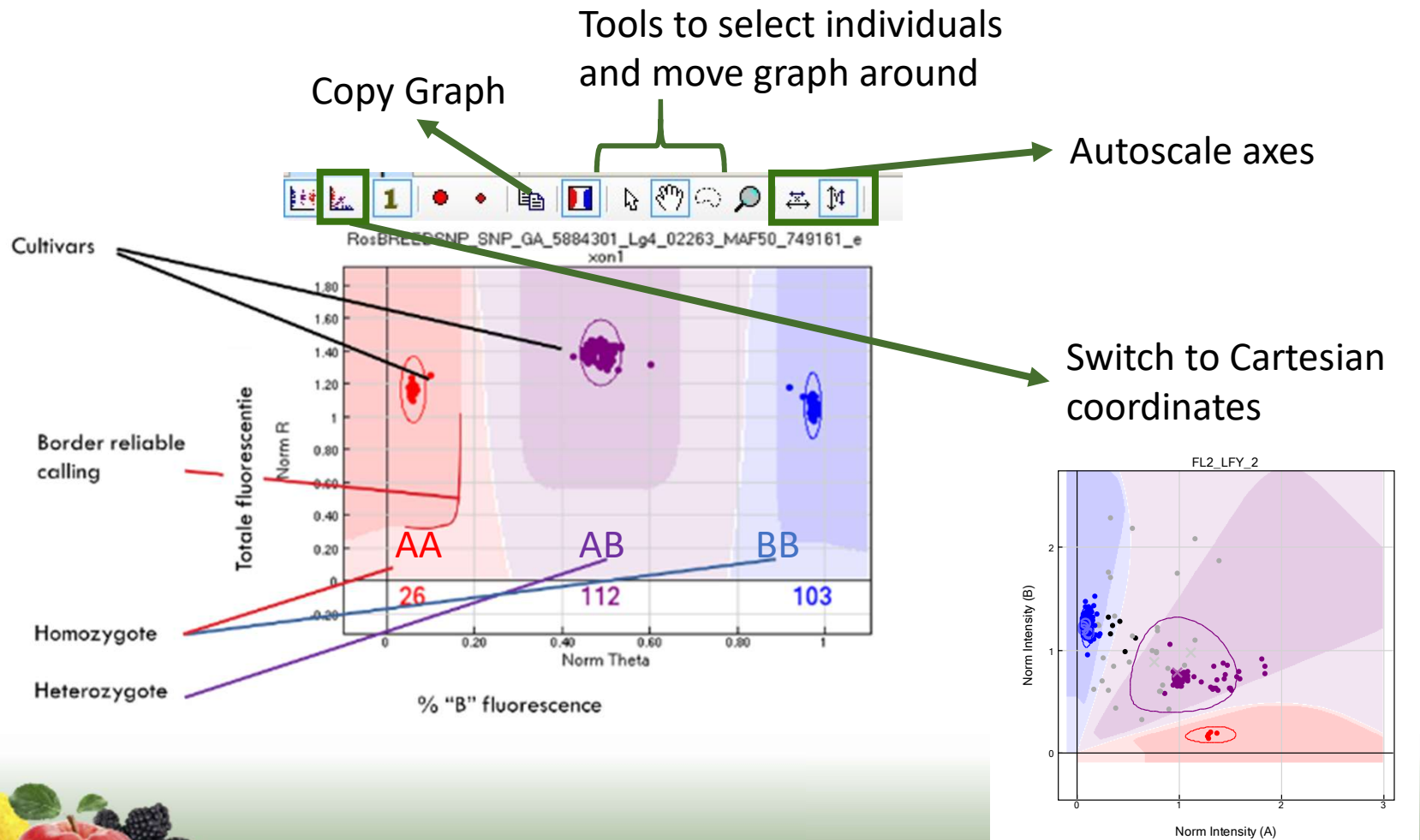
GenomeStudio® Layout

The screenshot displays the GenomeStudio 2.0 interface for a genotyping project. The main window is titled "GenomeStudio 2.0 - Genotyping - 6+9Kv5" and includes a menu bar (File, Edit, View, Analysis, Tools, Window, Help) and a toolbar. The interface is divided into several panes:

- SNP Graph:** Located in the top-left pane, it shows a plot of Norm R (y-axis, 0 to 3) versus Norm Theta (x-axis, 0 to 1). A red vertical line is positioned at approximately 0.15 on the x-axis, with a red circle around it. The plot is labeled "SNP Graph" and has values "43" and "305" at the bottom.
- Full Data Table, SNP Table, and Paired Sample Table:** The central pane displays a large table with columns for Index, Name, Address, Chr, Position, GenTrain Score, Frac A, Frac C, Frac G, Frac T, GType, Score, Theta, R, and columns for Sample 17 (LargeDeepRed), Sample 36 (UkrGriotte), and Sample 38 (YelSpanish). The table is labeled "Full Data Table, SNP Table, and Paired Sample Table".
- Samples Table:** Located in the bottom-left pane, it shows a table with columns for Sample, Original, Name, Short, Alternative Name, Germplasm Order, PoorDNA, Source, and 6K_haps available on 6+. The table is labeled "Samples Table".
- Error Table:** Located in the bottom-right pane, it shows a table with columns for Error Index and Parent1/. The table is labeled "Error Table".

At the bottom of the interface, there is a status bar with the text "Log" and a toolbar with icons for Select All, Copy, Save, Clear, X, Errors, Warnings, Info, and Log.

SNP Graph



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Full data table

Import (additional columns) Column Chooser

Copy Export Plotting functions filter/unfilter sample columns

Index	Name	Address	Chr	Position	GenTrain Score	Frac A	Frac C	Frac G	Frac T	Sample 1 9/47				Sample 2 5/62				
										GType	Score	Theta	R	GType	Score	Theta	R	
1	CBP	7074...	PA...	14629699	0.8098	0.235	0.265	0.255	0.245	AB	0.8304	0.5927	0.8864	BB	0.8304	0.9867	0.7032	BB
2	CBP_2	5871...	PA...	14629699	0.7834	0.235	0.265	0.255	0.245	AB	0.7893	0.5749	0.8928	BB	0.7893	0.9851	0.7371	BB
3	CBPb	2569...	PA...	37474346	0.9265	0.235	0.225	0.206	0.333	BB	0.9175	1.0000	0.8438	BB	0.9175	0.9988	0.8398	BB
4	CBPc	2877...	PA...	3108977	0.7207	0.176	0.206	0.235	0.382	NC	0.0678	0.3507	1.5965	AB	0.6285	0.4337	1.5368	BB

SNPs

Quality score

Genotype Call

Position in SNP Graph



SNP data table

Index	Name	Chr	Position	ChiTest1 00	Het Excess	AA Freq	AB Freq	BB Freq	Call Freq	Minor Freq	Aux	P-C Errors	P-P-C Errors	Rep Errors	10% GC	50% GC	SNP	# Calls	# no calls	Plus/Min us Strand	Custom Cluster #	Adi
1	CBP	PA...	14629699	0.8328	0.0211	0.0029	0.1232	0.8739	1.0000	0.0645	0	0	1	0	0.8304	0.8304	[T/C]	349	0		3	70
2	CBP_2	PA...	14629699	0.8328	0.0211	0.0029	0.1232	0.8739	1.0000	0.0645	0	0	1	0	0.7893	0.7893	[T/C]	349	0		3	58

How well did the SNP perform?

Are parents and offspring matching?
(see later)

Was the SNP polymorphic?

SNP	# Calls	# no calls	Plus/Min us Strand	Custom Cluster #	Address	GenTrain Score	Orig Score	Edited	Cluster Sep	AA T Mean	AA T Dev	AB T Mean	AB T Dev	BB T Mean	BB T Dev	AA R Mean	AA R Dev	AB R Mean	AB R Dev	BB R Mean	BB R Dev	Address2	Norm ID
[T/C]	349	0		3	7074231	0.8098	0.8098	0	0.9001	0.0706	0.0224	0.5832	0.0346	0.9914	0.0074	0.4874	0.1000	0.9331	0.1088	0.7781	0.1000	0	2
[T/C]	349	0		3	5871731	0.7834	0.7834	0	0.9221	0.0760	0.0224	0.5850	0.0328	0.9888	0.0075	0.4652	0.1000	0.9608	0.1093	0.7820	0.1000	0	2

“Array specific” data



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

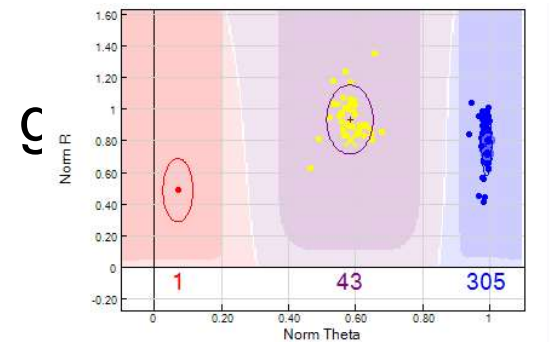
Sample Table

- All columns of sample sheet
- Sample statistics (need to be calculated)
- Linked with SNP graph
 - Selecting in table will highlight in SNP graph vice versa



Sample ID	Alternative Name	Source	Index	Call Rate	Gender
9/47		Cherry	1	0.000000	Unknown
5/62		Cherry	2	0.000000	Unknown
11/118		Cherry	3	0.000000	Unknown
13/20		Cherry	4	0.000000	Unknown
6/240		Cherry	5	0.000000	Unknown
7146-16		Cherry	6	0.000000	Unknown
7147-1		Cherry	7	0.000000	Unknown
7147-9		Cherry	8	0.000000	Unknown
8008-10		Cherry	9	0.000000	Unknown
8008-5		Cherry	10	0.000000	Unknown
8011-2		Cherry	11	0.000000	Unknown
8011-3		Cherry	12	0.000000	Unknown
8011-4		Cherry	13	0.000000	Unknown
99F/132R4		Cherry	14	0.000000	Unknown
99F/150R1A		Cherry	15	0.000000	Unknown

Rows=349 Disp=349 Sel=1 Filter=Filter is not active.



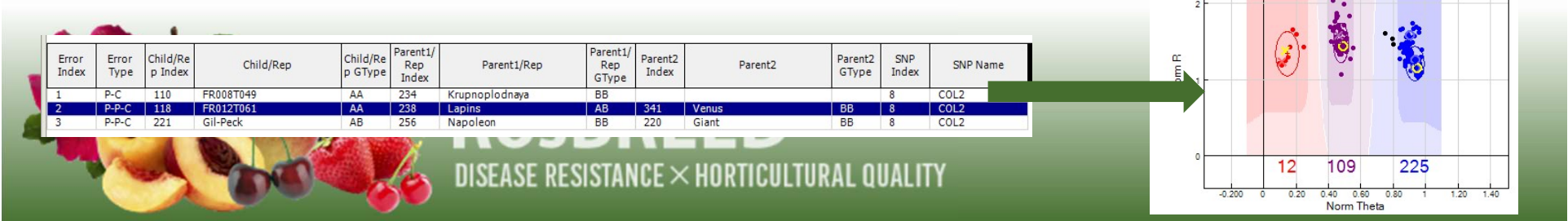
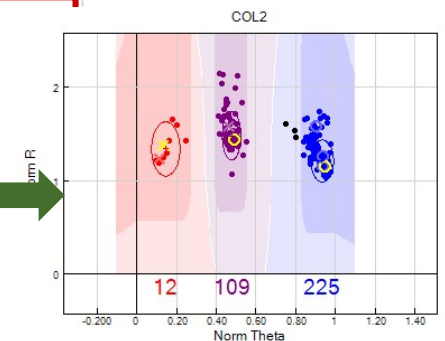
Sample ID	Alternative Name	Source	Index	Call Rate	Gender
Black Star		Cherry	27	0.8604617	Unknown
Black Tartarian		Cherry	28	0.8354598	Unknown
Black Gold		Cherry	29	0.8488089	Unknown
Black Republic		Cherry	30	0.8634118	Unknown
Burbank		Cherry	31	0.7537429	Unknown
Burbank Black		Cherry	32	0.8634118	Unknown
Cashmere		Cherry	33	0.8641493	Unknown
Cavalier		Cherry	34	0.8311822	Unknown
Celeste		Cherry	35	0.8623055	Unknown
Chelan		Cherry	36	0.8644443	Unknown
Coe		Cherry	37	0.7829486	Unknown
Compact Stella		Cherry	38	0.8640018	Unknown
Coops Special		Cherry	39	0.8204145	Unknown
Corum		Cherry	40	0.8623793	Unknown
Cowiche	DF70	Cherry	41	0.8626697	Unknown

Error Table

- Shows errors between replicates and parent(s) and offspring
 - Will be discussed further
 - Linked to SNP graph
 - Replicates: square
 - Parent(s) and offspring: circle(s) and cross, respectively

Error Index	Error Type	Child/Rep Index	Child/Rep	Child/Rep GType	Parent1/Rep Index	Parent1/Rep	Parent1/Rep GType	Parent2 Index	Parent2	Parent2 GType	SNP Index	SNP Name
1	P-C	110	FR008T049	AA	234	Krupnoplodnaya	BB				8	COL2
2	P-P-C	118	FR012T061	AA	238	Lapins	AB	341	Venus	BB	8	COL2
3	P-P-C	221	Gil-Peck	AB	256	Napoleon	BB	220	Giant	BB	8	COL2

Error Index	Error Type	Child/Rep Index	Child/Rep	Child/Rep GType	Parent1/Rep Index	Parent1/Rep	Parent1/Rep GType	Parent2 Index	Parent2	Parent2 GType	SNP Index	SNP Name
1	P-C	110	FR008T049	AA	234	Krupnoplodnaya	BB				8	COL2
2	P-P-C	118	FR012T061	AA	238	Lapins	AB	341	Venus	BB	8	COL2
3	P-P-C	221	Gil-Peck	AB	256	Napoleon	BB	220	Giant	BB	8	COL2





United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Sample quality and ploidy

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY

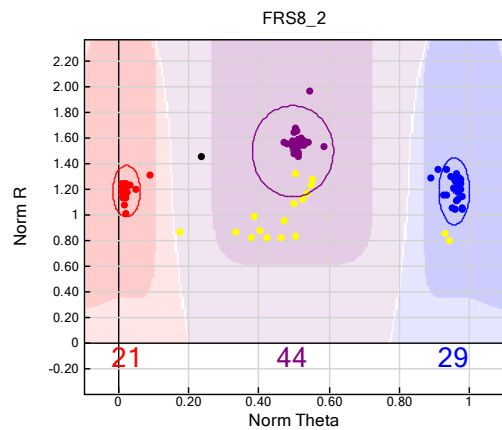


SUPERIOR
CULTIVARS

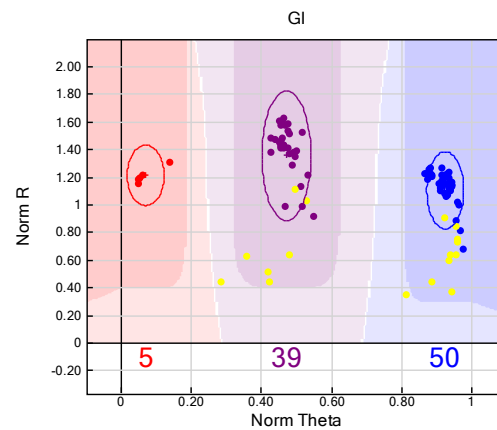


Checking Sample Quality

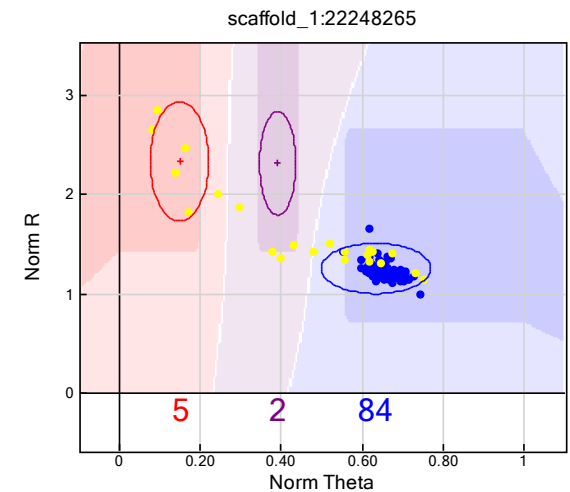
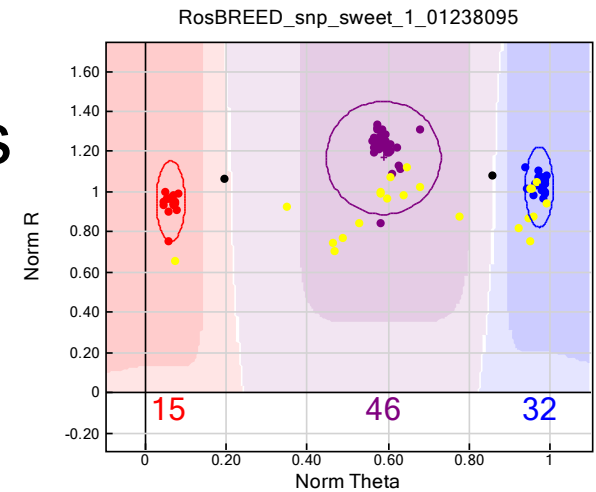
- Visual inspection of SNP graphs



Some individuals are regularly located outside main clusters



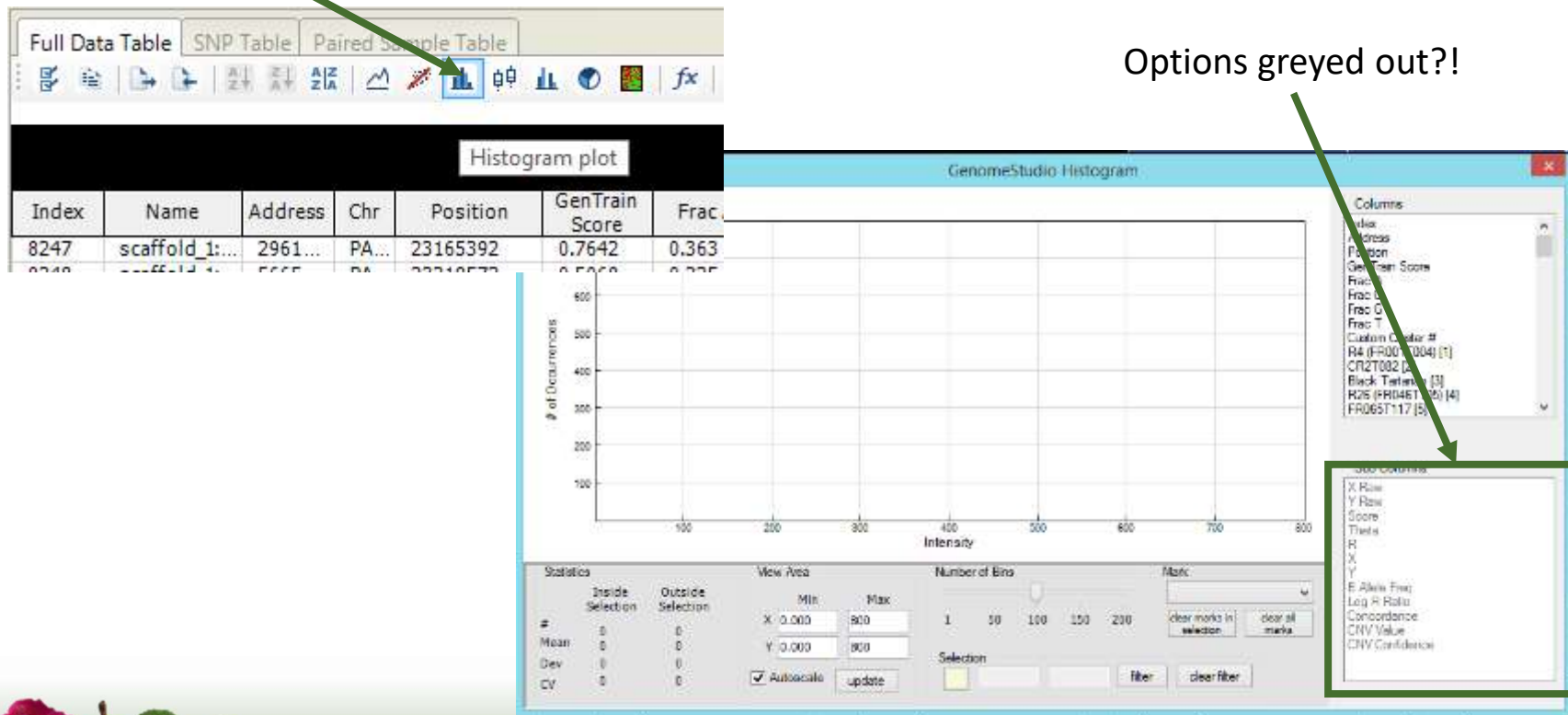
More individuals outside main clusters
*use 'ctrl' + left mouse clicks to add new individual to selection



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Checking Sample Quality

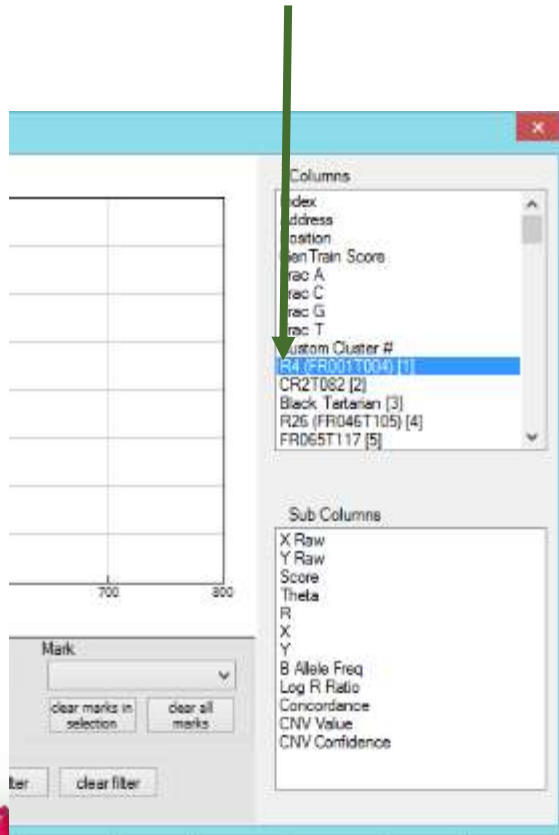
1. Histogram function of Full Data Table



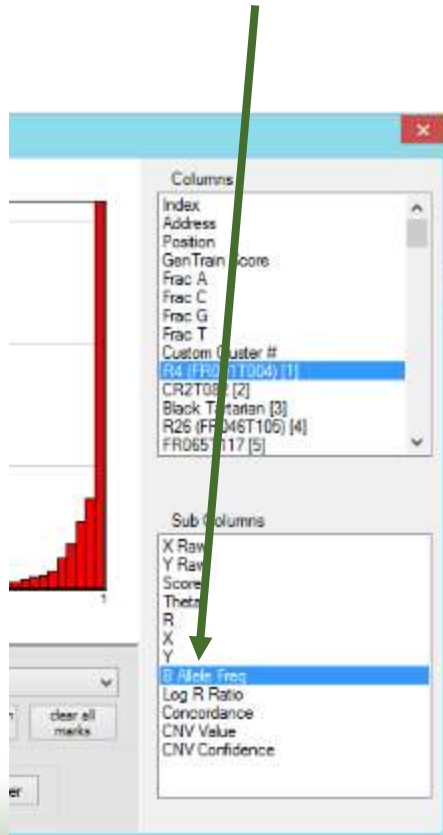
RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Checking Sample Quality

2. Select first individual



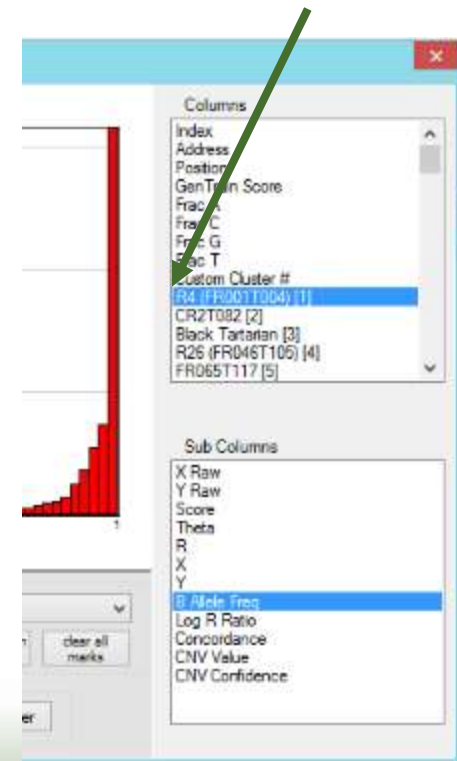
3. Choose "B Allele Freq"



4. Select individual again

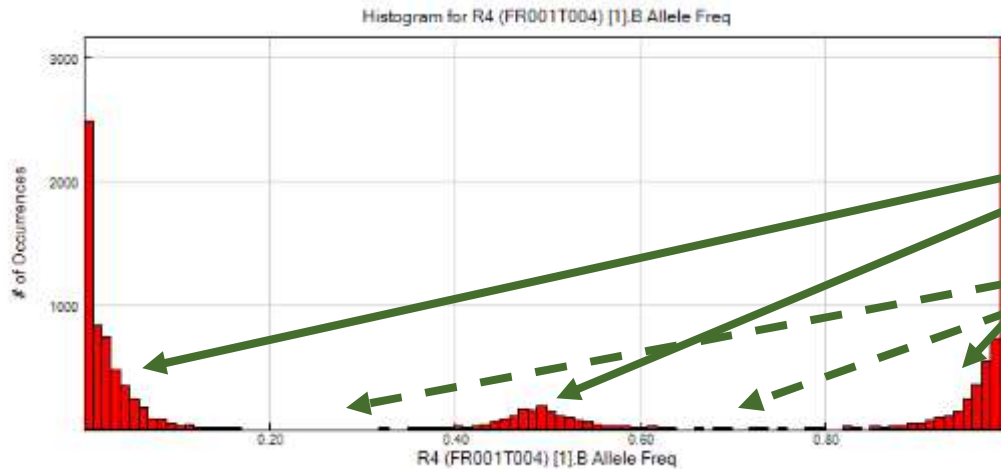
*Allows for scrolling through ind.

*Use 'up' and 'down' keys to scroll through ind.



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

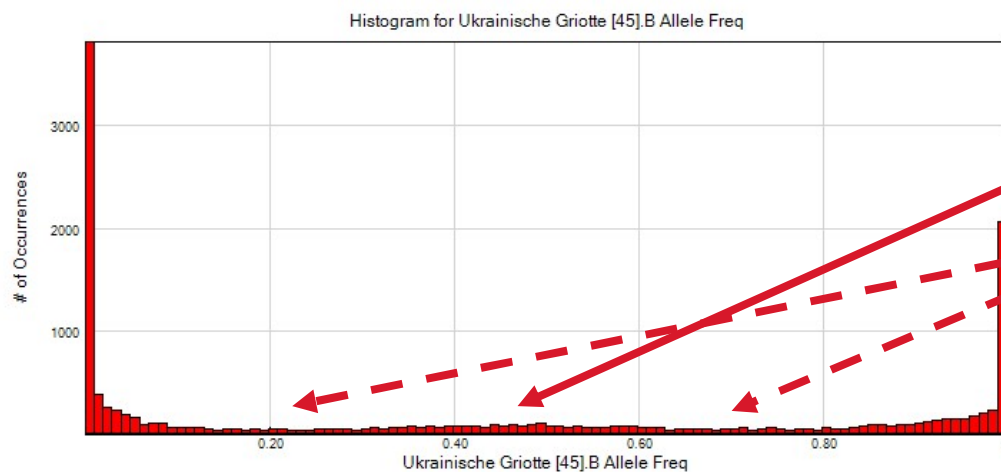
Checking Sample Quality



Sample with good quality

-Three (clear) peaks

-Barely any occurrences between three peaks

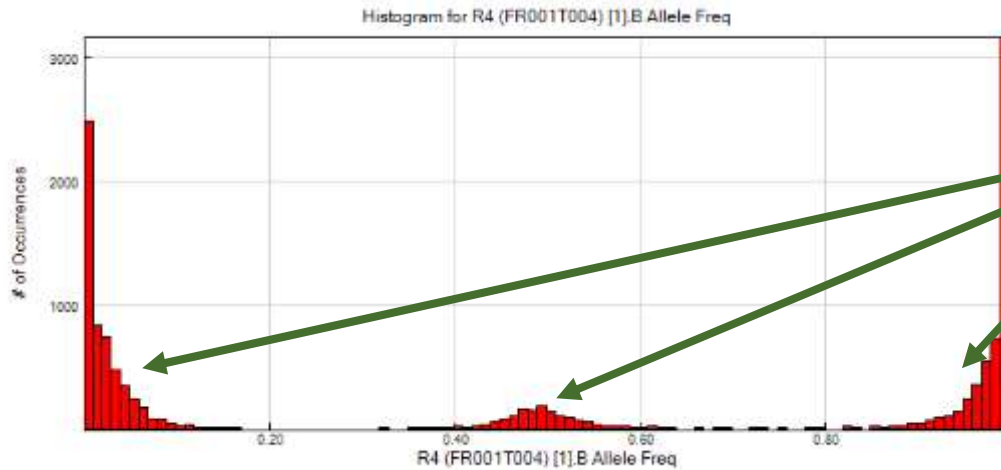


Sample with poor quality

-No clear heterozygous peak

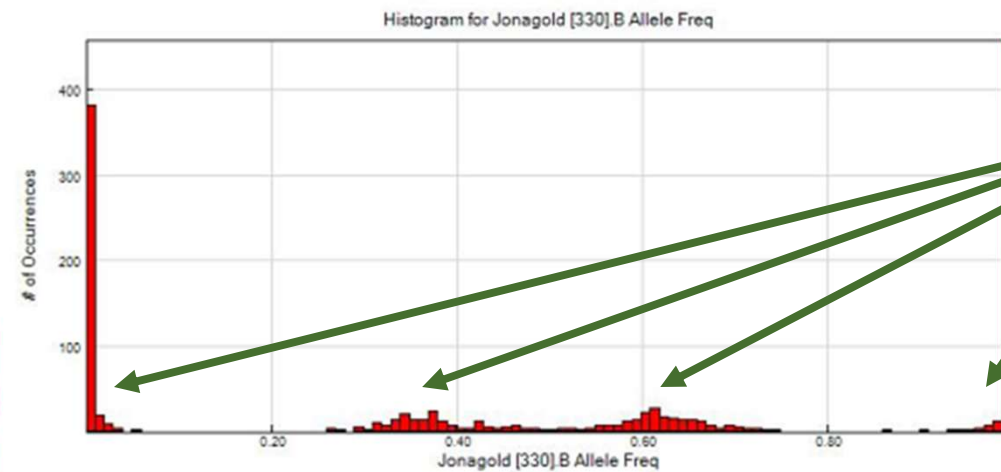
-Many occurrences between three peaks
(middle section looks ~ flat)

Checking Sample Ploidy



Diploid sample

- Three (clear) peaks
- AA, AB, and BB

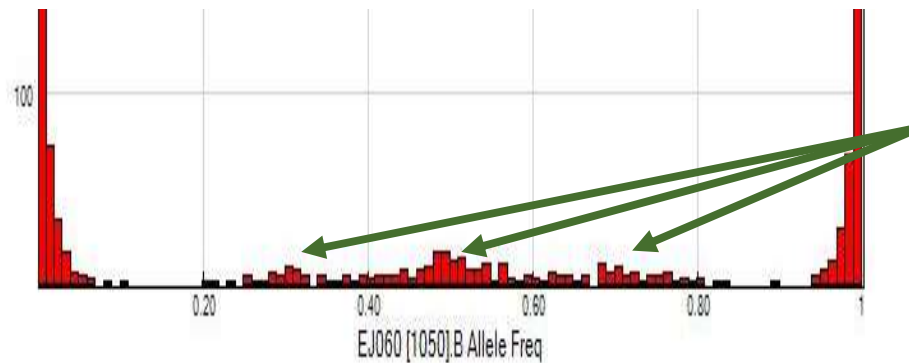
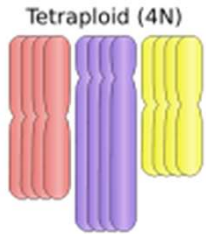


Triploid sample

- Four peaks
- AAA, AAB, ABB, and BBB

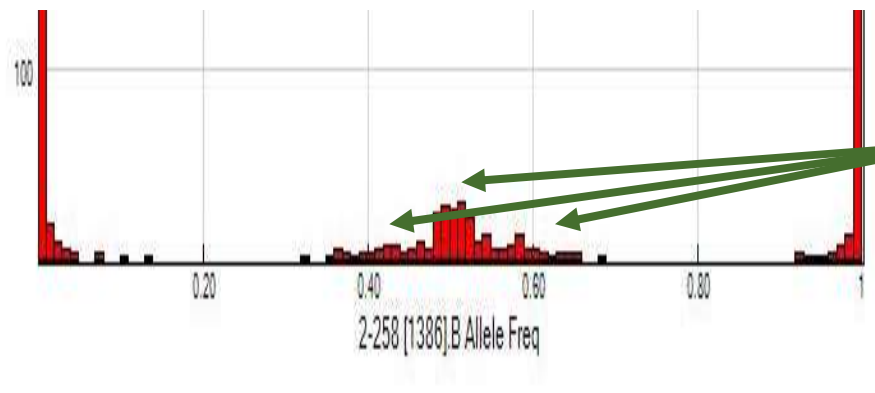
Checking Sample Ploidy

- Other ploidy



Tetraploid sample

- Five peaks?
- Different from poor quality?
- Could be a mix of 2 samples



Aneuploid sample

- 3 peaks + 2 small peaks?
- Not easy to see



Chagné et al. (2015) method

1. Choose 'Filter rows' for SNP Table in GenomeStudio®



2. Apply the following filter parameters (continued on next slide):

Call Freq	≥ 0.90	Only SNP that have good call rate
Minor Freq	>0.01	Only polymorphic SNPs
50% GC	≥ 0.40	Only SNPs for which most individuals were close to main cluster
GenTrain	≥ 0.65	Only SNPs with good overall clustering
ClusterSep	≥ 0.40	Only SNPs with well separated clusters
AA T Mean	≤ 0.125	Only SNPs with AA cluster in expected position
AA T Dev	≤ 0.028	Only SNPs with 'narrow' AA cluster (no multiple clusters)
AB T Mean	≥ 0.375 ≤ 0.625	Only SNPs with AB cluster in expected position
AB T Dev	≤ 0.056	Only SNPs with 'narrow' AB cluster (no multiple clusters)
BB T Mean	≥ 0.875	Only SNPs with BB cluster in expected position
BB T Dev	≤ 0.028	Only SNPs with 'narrow' BB cluster (no multiple clusters)



DISEASE RESISTANCE \times HORTICULTURAL QUALITY

Chagné et al. (2015) method

2. Apply the following filter parameters (continued):

2b. Choose operation

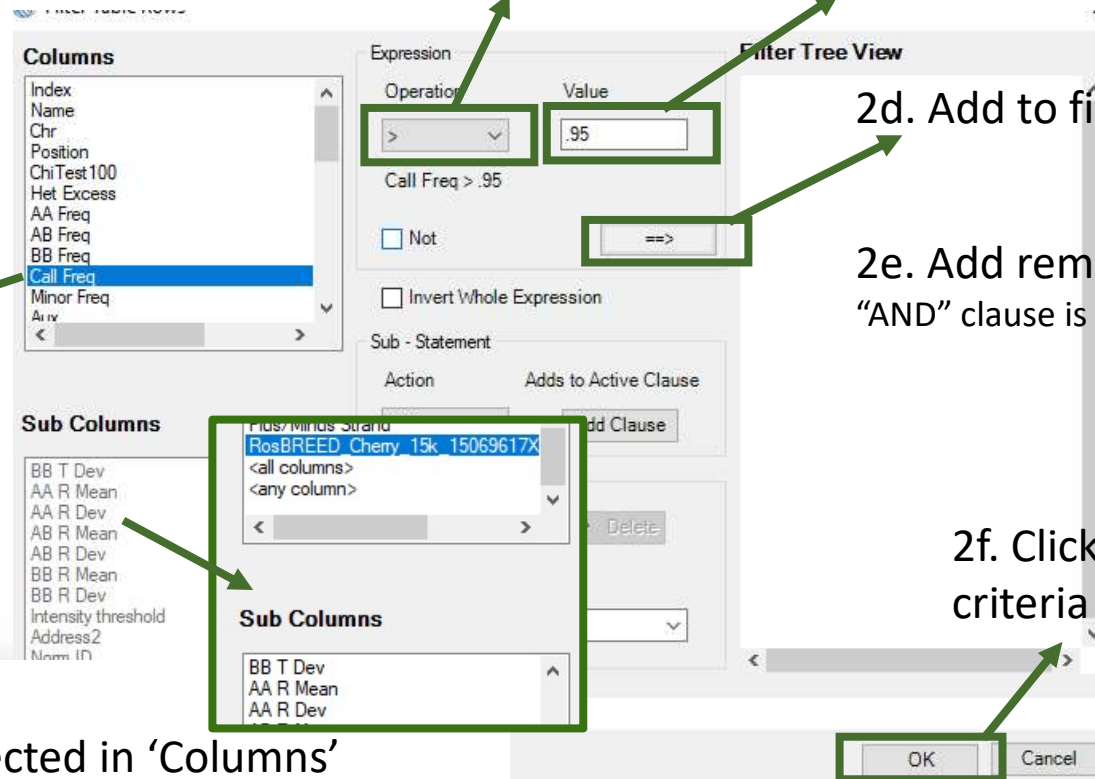
2c. Choose value

2d. Add to filter criteria

2e. Add remaining criteria
“AND” clause is added automatically

2f. Click ‘OK’ when all criteria are added

2a. Choose parameter



Note:

Manifest must be selected in ‘Columns’ section before subcolumns become available

Chagné et al. (2015) method

3. Go to Full Data Table and select 'Column Chooser'
 *Filtering of SNP Table remains for Full Data Table

4. Hide any column except Index, Chr, Position, and sample columns

Index	Name	Address	Chr	Position	GenTrain Score	Frac A	Frac C	Frac G	Frac T	GType	Score	Theta
1	CBP	7074...	PA...	14629699	0.7572	0.235	0.265	0.255	0.245	NC	0.0298	0.8421
2	CBP_2	5871...	PA...	14629699	0.7186	0.235	0.265	0.255	0.245	AB	0.2146	0.7951
3	CBPb	2569...	PA...	37474346	0.8755	0.235	0.225	0.206	0.333	BB	0.7744	0.9701
4	CBPc	2877...	PA...	3108977	0.6339	0.176	0.206	0.235	0.382	NC	0.0751	0.0881
5	CBPc_2	4875...	PA...	3108977	0.6340	0.176	0.206	0.235	0.382	NC	0.1161	0.0411
6	CLF	2680...	PA...	15025492	0.6902	0.314	0.206	0.225	0.255	NC	0.0071	0.7641
7	CLF_2	6478...	PA...	15025492	0.6882	0.314	0.206	0.225	0.255	NC	0.0192	0.7951

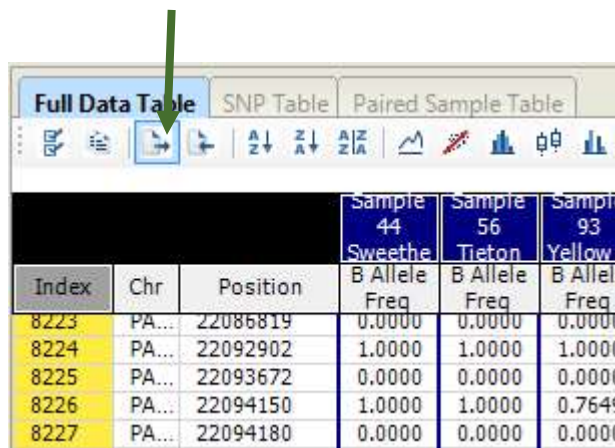
6. Show B Allele Freq subcolumn

5. Hide standard subcolumns

7. Press 'OK'

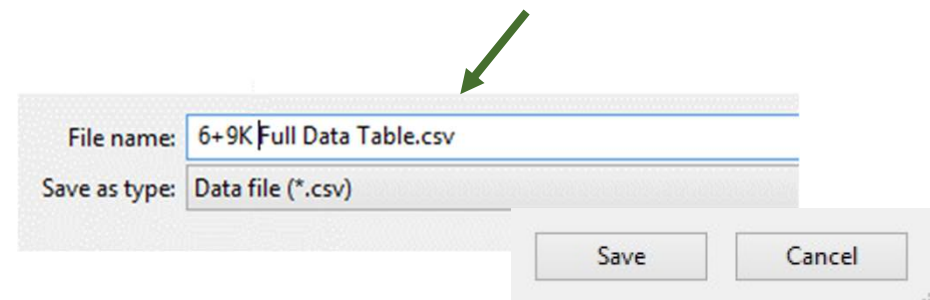
Chagné et al. (2015) method

8. Export table



Index	Chr	Position	Sample 44 Sweethe B Allele Freq	Sample 56 Tieton B Allele Freq	Sample 93 Yellow B Allele Freq
8223	PA...	22086819	0.0000	0.0000	0.0000
8224	PA...	22092902	1.0000	1.0000	1.0000
8225	PA...	22093672	0.0000	0.0000	0.0000
8226	PA...	22094150	1.0000	1.0000	0.7649
8227	PA...	22094180	0.0000	0.0000	0.0000

9. Give a meaningful name, choose save location, and save as '.csv' file

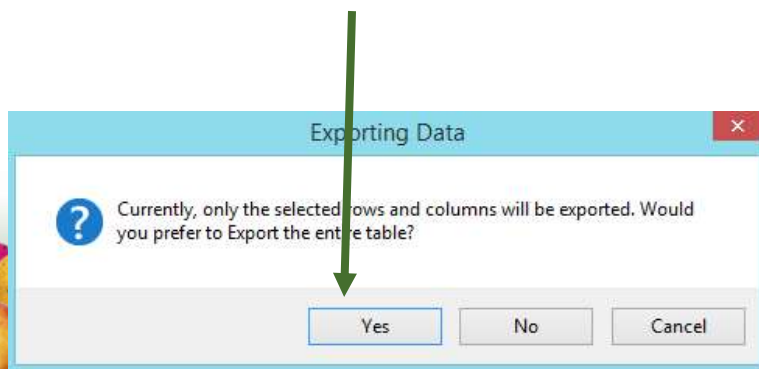


File name: 6+9K Full Data Table.csv

Save as type: Data file (*.csv)

Save Cancel

10. If asked, choose 'Yes' to export entire table

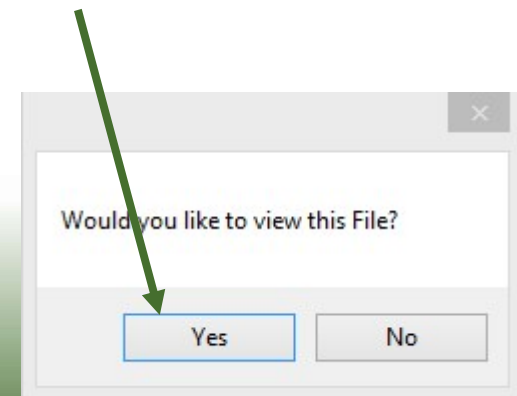


Exporting Data

Currently, only the selected rows and columns will be exported. Would you prefer to Export the entire table?

Yes No Cancel

11. Choose 'Yes' to view file in Excel

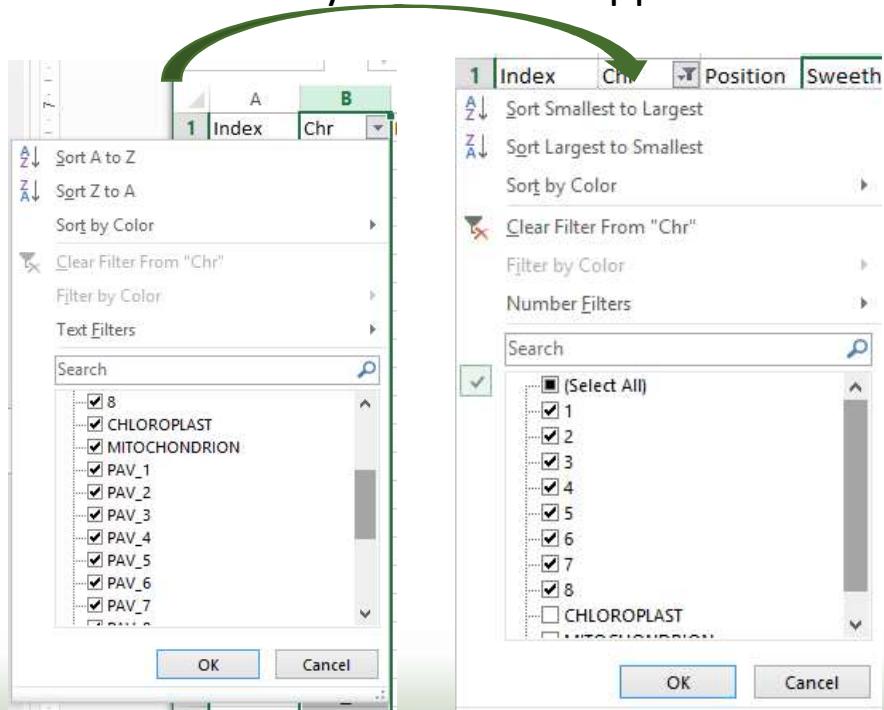


Would you like to view this File?

Yes No

Chagné et al. (2015) method

12. Convert Chromosome numbers to numeric system where applicable



13. Insert a new column in front of sample columns

The image shows a portion of an Excel spreadsheet with columns C, D, E, and F. Column C is labeled 'Position' and column D is labeled 'Sweethe'. A context menu is open over column D, with the 'Insert' option highlighted in green.

	C	D	E	F
	Position	Sweethe	Cut	
8	14629699	0.820586	Copy	
8	14629699	0.751753	Paste Options:	
1	37474346	0.97149		
2	3108977	0.007736		
2	3108977	0	Paste Special...	
7	15025492	0.826106	Insert	
7	15025492	0.855948	Delete	
3	1473284	0.373288	Clear Contents	
3	1473284	0.382824		



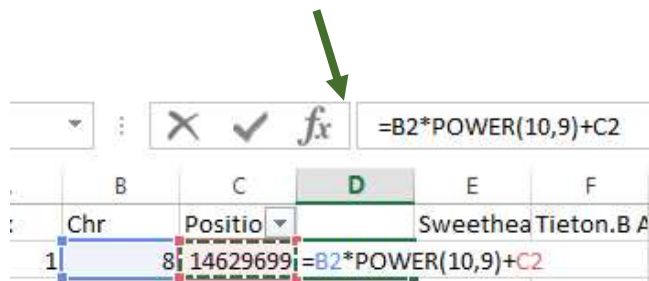
RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Chagné et al. (2015) method

14. Apply the following formula for the new column:

Cell value = 1,000,000,000 * 'Chromosome number' + 'position on chromosome'

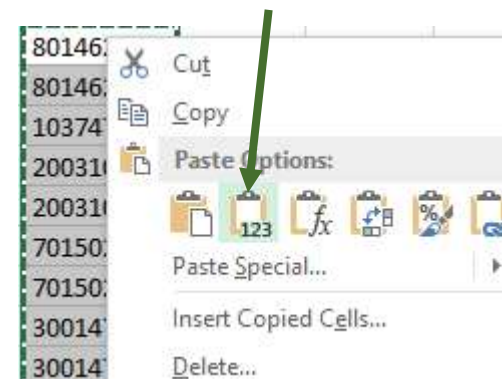
*Choose a power of 10 that is larger than any position on chromosome



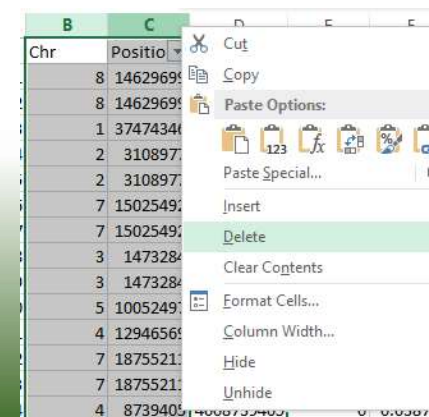
C	D	
sitio		Sw
629699	8014629699	0.
629699	8014629699	0.
474346	1037474346	1
108977	2003108977	0.
108977	2003108977	
025492	7015025492	0.
025492	7015025492	0.
473284	3001473284	0.
473284	3001473284	0.
052497	5010052497	
946569	4012946569	0.
755211	7018755211	0.

15. Replace formula by fixed values

*Copy column and paste as values



16. Delete original chromosome and position column



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

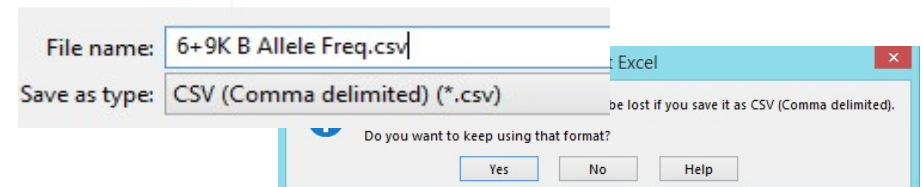
Chagné et al. (2015) method

File should look like this:

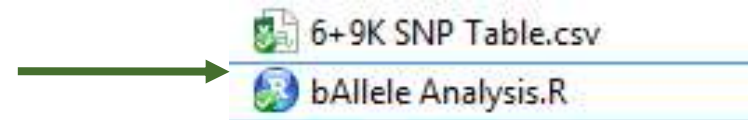
- 1st column: SNP Identifier (e.g. SNP Index)
- 2nd column: SNP position (including chromosome)
- 3rd column – end: B-allele Freq of each individual

	A	B	C	D	E	F	
1	Index	Position	Sweethea	Tieton.B	A Yellow Sp	Rainier.B	Ukr
2	1	8014629699	0.820586	0.658852	0.935553	0.926576	0.9
3	2	8014629699	0.751753	0.752404	0.970744	0.808932	0.9
4	3	1037474346	0.97149	0.817936	0.918174	0.964891	0.9
5	4	2003108977	0.007736	0.63033	1	0.637792	
6	5	2003108977	0	0.699511	1	0.723935	

17. Save as new 'csv' file:



18. Copy R-script 'bAllele Analysis.R' (Suppl. Document 2) in same folder as newly created .csv file



19. Open R-script 'bAllele Analysis.R' in RStudio



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Chagné et al. (2015) method

20. Change file-name on line 1 into correct name of newly created .csv file

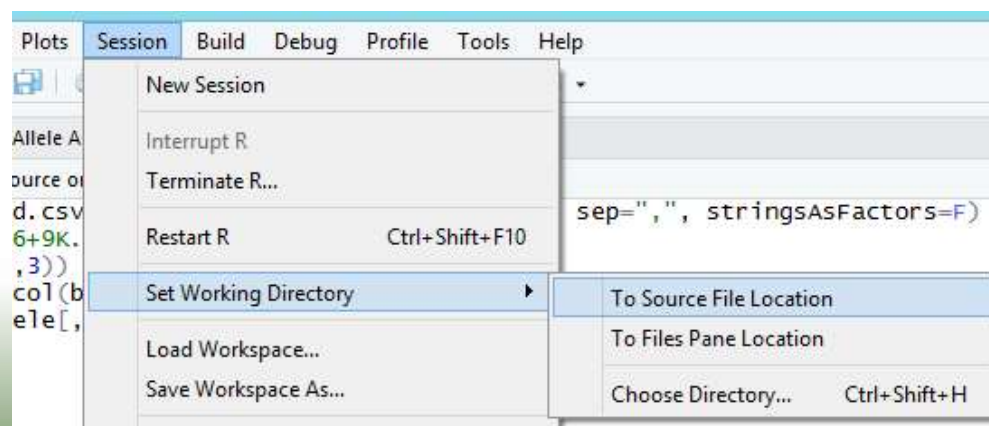
```
bAllele<- read.csv(file="SNP data_Triploids and deviations.csv", sep="," , stringsAsFactors=F)  
bAllele<- read.csv(file="6+9K B Allele Freq.csv", sep="," , stringsAsFactors=F)
```

21. Change name of pdf-output on line 2 into desired name

```
pdf("bAllele.pdf")  
pdf("bAllele 6+9K.pdf")
```

22. Change working directory into location that has .csv file and R script

*Under 'Session', choose 'Set Working Directory', then choose 'To Source File Location'

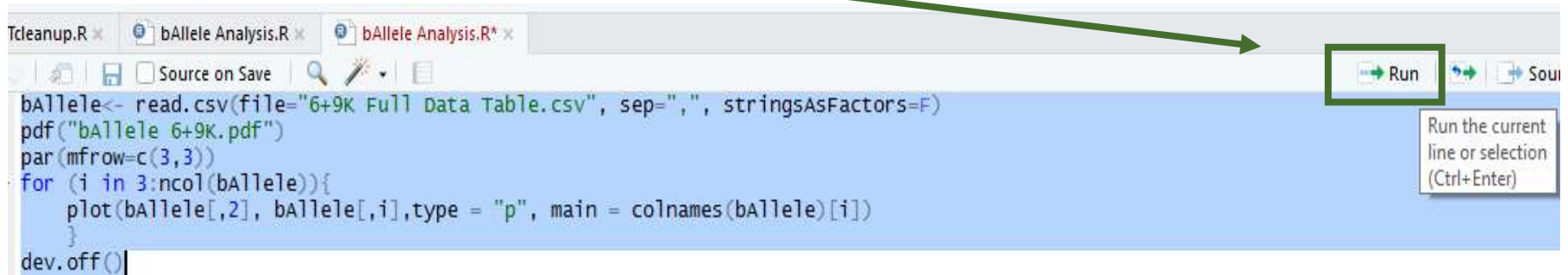


ROSBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Chagné et al. (2015) method

23. Run entire script

*Select all lines and press run

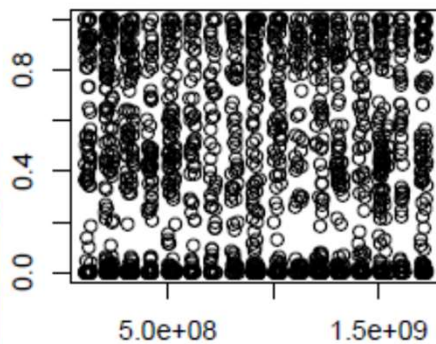


```
bAllele<- read.csv(file="6+9K Full Data Table.csv", sep=";", stringsAsFactors=F)
pdf("bAllele 6+9K.pdf")
par(mfrow=c(3,3))
for (i in 3:ncol(bAllele)){
  plot(bAllele[,2], bAllele[,i],type = "p", main = colnames(bAllele)[i])
}
dev.off()
```

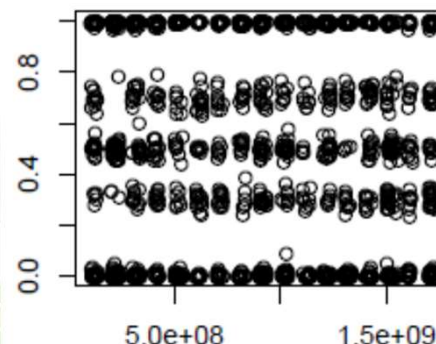
24. Open created pdf file

*It's in the same folder as the R-script and the .csv file with the B Allele Frequencies

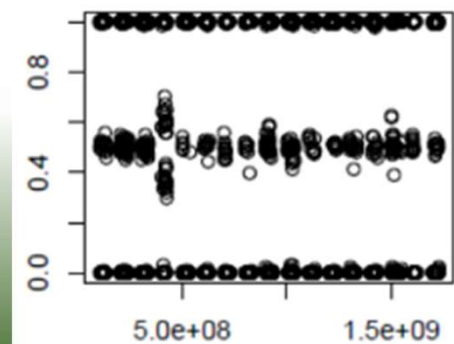
Poor quality
Monark.B_Allele_Freq



Tetraploid?

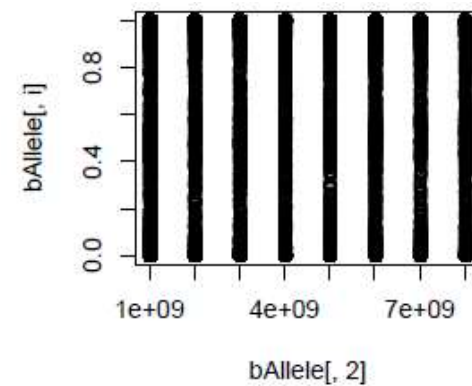
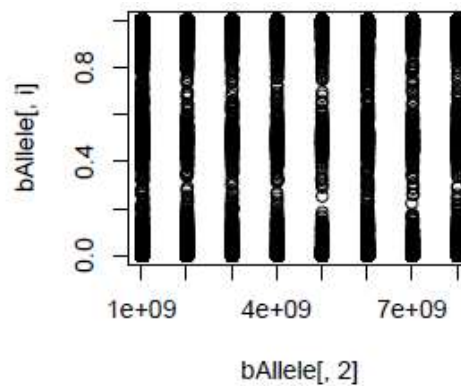
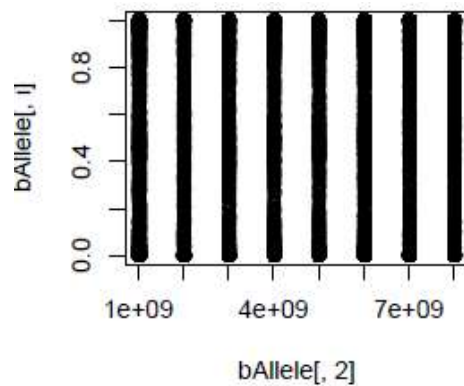


Aneuploid
X2.258.B_Allele_Freq

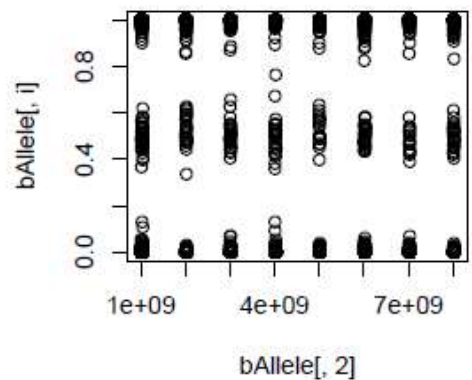
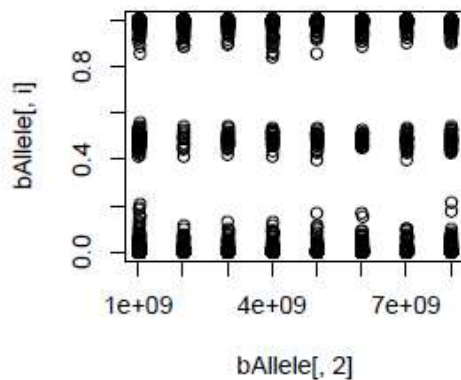
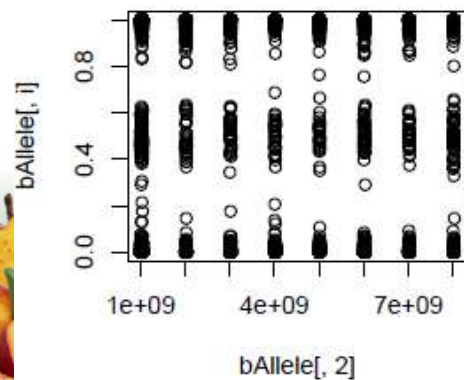


Chagné et al. (2015) method

Without initial filtering in SNP Table



With initial filtering in SNP Table



Segmental aneuploids

- Lack a large segment of a chromosome
 - Large extra segment?
- Cannot be identified with methods described above
- Will lead to many errors further on for one chromosomal segment



RosBREED

DISEASE RESISTANCE × HORTICULTURAL QUALITY



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

SNP subset

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS

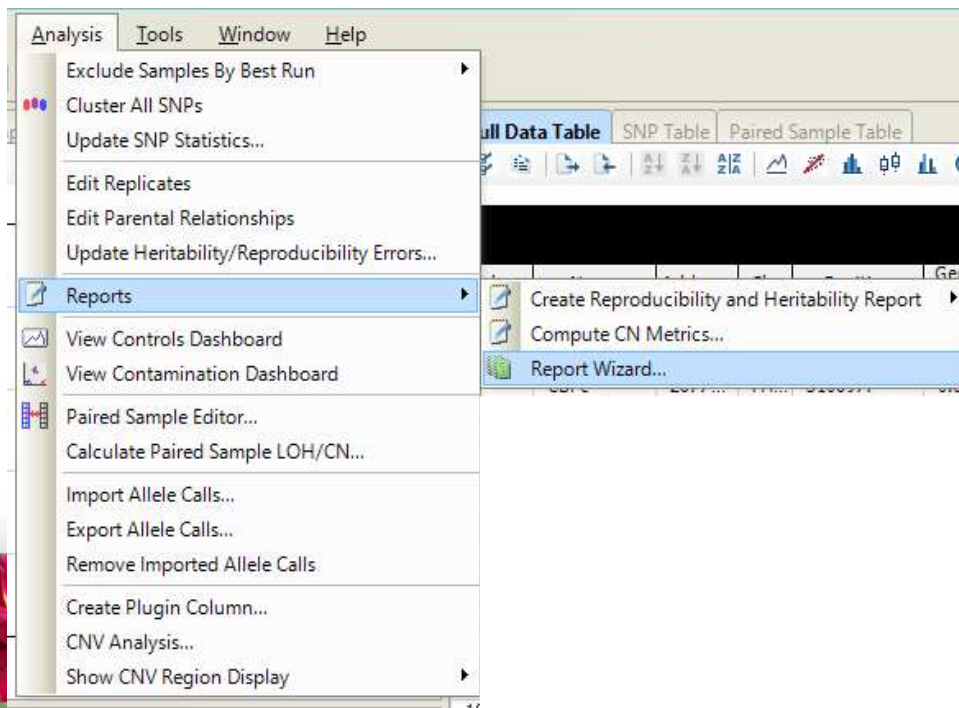


Creating input files ASSIsT

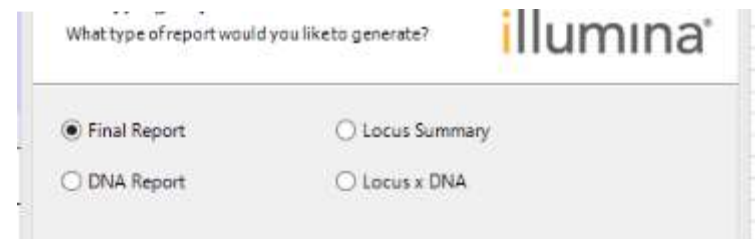
- Manual can be found under ‘...\\ASSIsT_Win_v1.01\\docs’

Creating 1st input file (Final report)

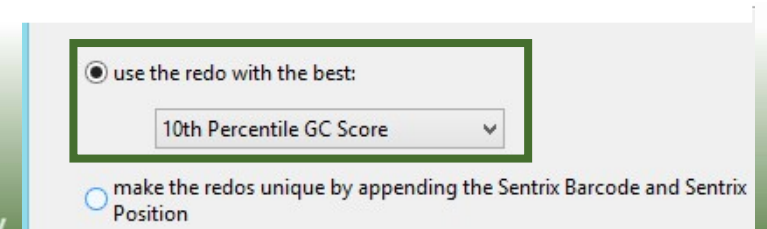
1. Open Report Wizard in GenomeStudio



2. Choose Final Report



3. Choose 'redo with the best' and '10th Percentile GC score'



Creating input files ASSIsT

4. Select all samples/arrays



Chip1
 Chip2
 Chip3
 Chip4

5.-Set format to Standard

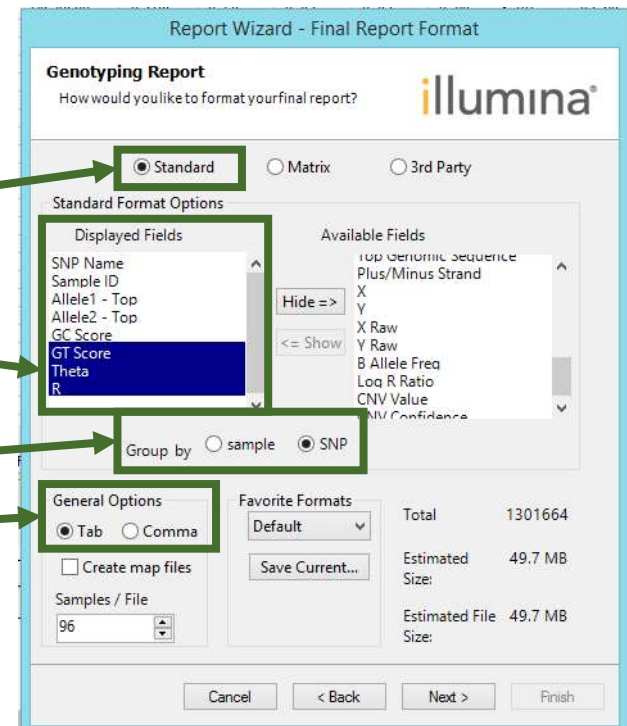
-Include the following columns:

GT Score, Theta, R


**Order should be the same as figure*

-Select 'Group by SNP'

-Select 'Tab' as delimiter



Report Wizard - Final Report Format

Genotyping Report
How would you like to format your final report? 

Standard Matrix 3rd Party

Standard Format Options

Displayed Fields

- SNP Name
- Sample ID
- Allele1 - Top
- Allele2 - Top
- GC Score
- GT Score
- Theta
- R

Available Fields

- Top Genomic Sequence
- Plus/Minus Strand
- X
- Y
- X Raw
- Y Raw
- B Allele Freq
- Log R Ratio
- CNV Value
- CNV Confidence

Hide => <=< Show

Group by sample SNP

General Options

Tab Comma

Create map files

Samples / File: 96

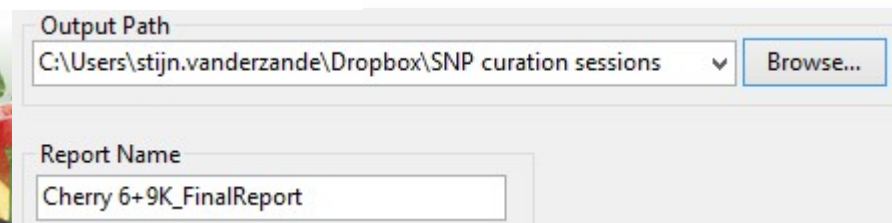
Favorite Formats: Default

Save Current...

Total: 1301664
Estimated Size: 49.7 MB
Estimated File Size: 49.7 MB

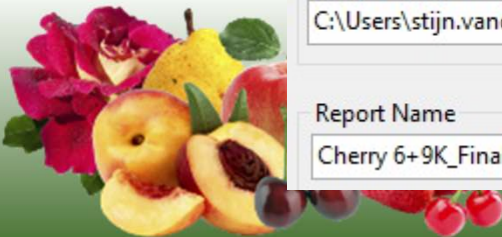
Cancel < Back Next > Finish

6. Save Report



Output Path
C:\Users\stijn.vanderzande\Dropbox\SNP curation sessions

Report Name
Cherry 6+9K_FinalReport

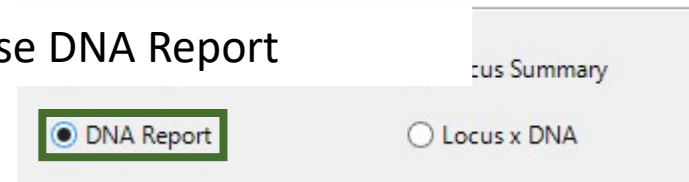


DISEASE RESISTANCE * HORTICULTURAL QUALITY

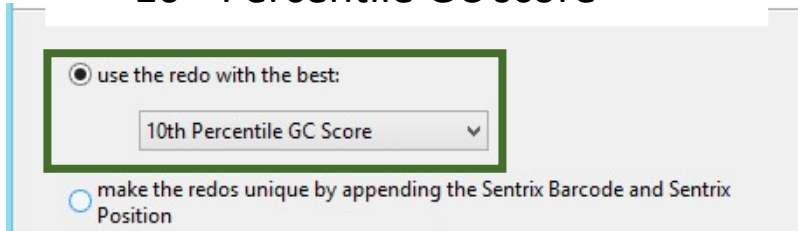
Creating input files ASSIST

Creating 2nd input file (DNA report)

7. Open Report Wizard in GenomeStudio and choose DNA Report



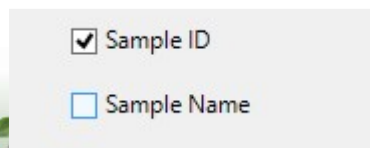
8. Choose 'redo with the best' and '10th Percentile GC score'



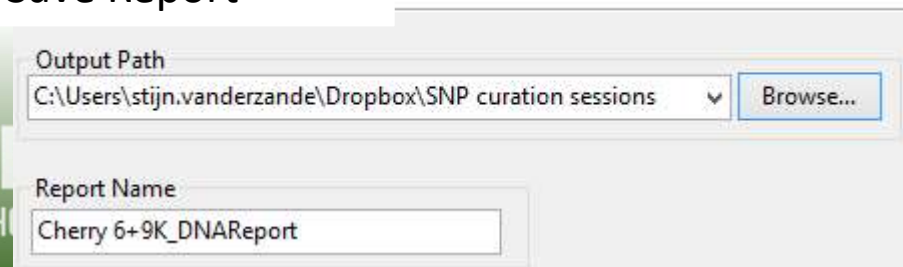
9. Select all samples/arrays



10. Choose Sample ID



11. Save Report



Input files ASSIsT

Creating 3rd input file (Pedigree file)

12. Create a three-column file:

-1st column has the individual's name

-2nd column has the female parent

-3rd column has the male parent

-Header row should be “//SampleID [tab] Mother [tab] Father”

*Copy SampleID and parent columns from sample sheet into Notepad++ and add the header line

*Parents do not have to be included in sample list

	A	B	C
1	//SampleID	Mother	Father
2	Abundance	Napoleon	
3	Corum		
4	Cuvelier		
5	Persian		
6	Benton	Stella	Moreau
7	Bing	BlackRepublican	Napoleon
8	BlackRepublican	Napoleon	BlackTartarian
9	BlackTartarian		
10	Lambert	Napoleon	Blackheart

```
//SampleID Mother Father
Abundance Napoleon
Corum
Cuvelier
Persian
Benton Stella Moreau
Bing BlackRepublican Napoleon
BlackRepublican Napoleon BlackTartarian
BlackTartarian
Lambert Napoleon Blackheart
```



Input files ASSIsT

Optional - Creating 4th input file (Map file)

13. Create a three-column file:

-1st column has the SNP's name/SNP's ID

-2nd column has the chromosome (numerical)

-3rd column has the position (physical (bp, Mbp) or genetic (cM))

-Header row should be “//SNPid [tab] Chromosome [tab] Position”

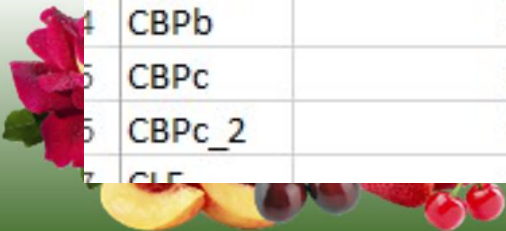
*Copy Name, Chr and Position from exported Full Data Table and convert chromosome to numeric if needed

*Use genetic position if available

*Give fictional chromosome number to chloroplast/mitochondrial/... SNPs


	A	B	C	D	E
1	//SNPid	Chromosome	Position		
2	CBP	8	14629699		
3	CBP_2	8	14629699		
4	CBPb	1	37474346		
5	CBPc	2	3108977		
6	CBPc_2	2	3108977		
7	CLF	7	15025402		

1	//SNPid	Chromosome	Position
2	CBP	8	14629699
3	CBP_2	8	14629699
4	CBPb	1	37474346
5	CBPc	2	3108977
6	CBPc_2	2	3108977
7	CLF	7	15025402

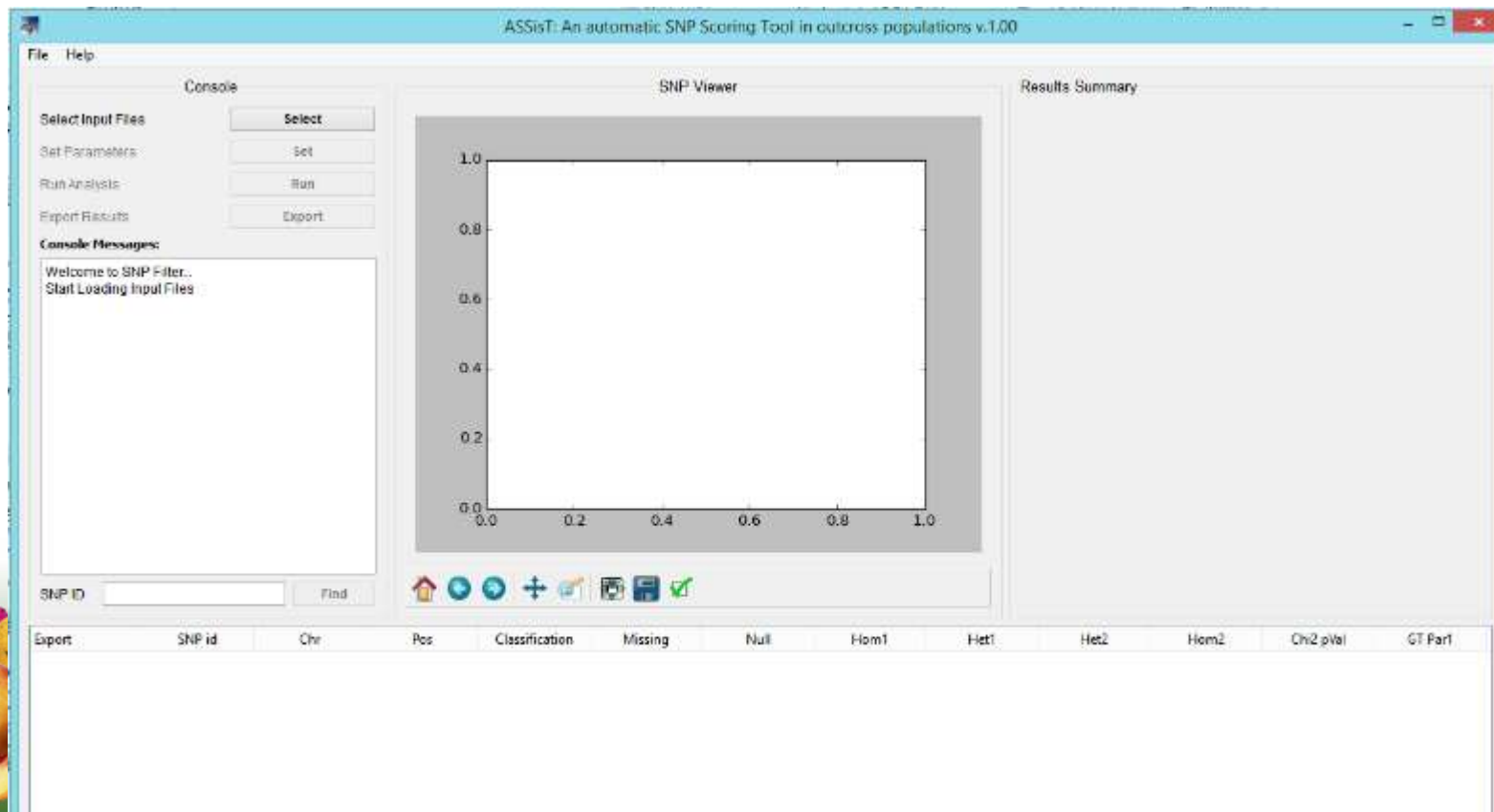


Running ASSIsT

1. Open ASSIsT



_ssl.pyd	5/5/2016 6:33 PM	PYD File	879 KB
_tkinter.pyd	5/5/2016 6:33 PM	PYD File	34 KB
_win32sysloader.pyd	5/5/2016 6:33 PM	PYD File	8 KB
assist.ico	5/5/2016 6:33 PM	Icon	222 KB
assist.png	5/5/2016 6:33 PM	PNG image	209 KB
ASSIsT_Win_v1.01.exe	5/5/2016 6:33 PM	Application	5,025 KB
ASSIsT_Win_v1.01.exe.manifest	5/5/2016 6:33 PM	MANIFEST File	1 KB
bz2.pyd	5/5/2016 6:33 PM	PYD File	67 KB
LIBEAY32.dll	5/5/2016 6:33 PM	Application extens...	1,075 KB



ASSIsT: An automatic SNP Scoring Tool in outcross populations v.1.00

File Help

Console

Select Input Files

Set Parameters

Run Analysis

Export Results

Console Messages:

Welcome to SNP Filter..
Start Loading Input Files

SNP ID

SNP Viewer

Results Summary

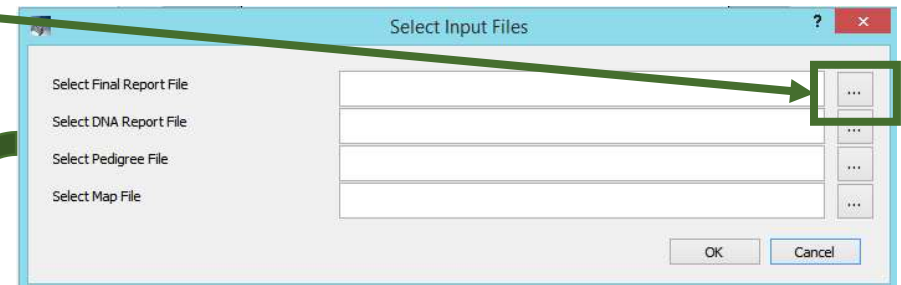
Export	SNP id	Chr	Pos	Classification	Missing	Null	Hom1	Het1	Het2	Hom2	Ch2 pval	GT Part
--------	--------	-----	-----	----------------	---------	------	------	------	------	------	----------	---------

Running ASSIsT

2. In console section, click on 'Select' next to 'Select Input Files'

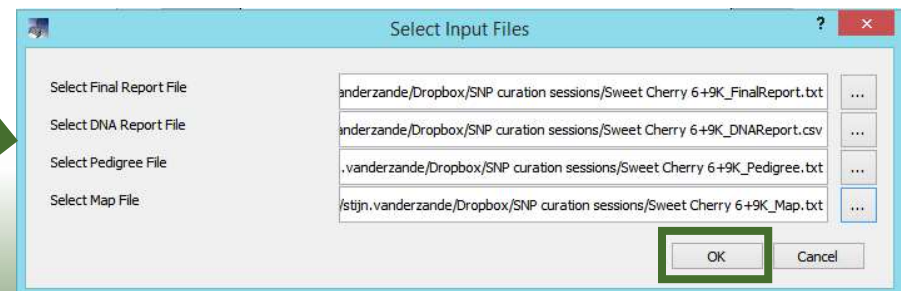


3. Click on '...' boxes to choose each file created earlier. Once each file is chosen, click 'OK'.



Console Messages:

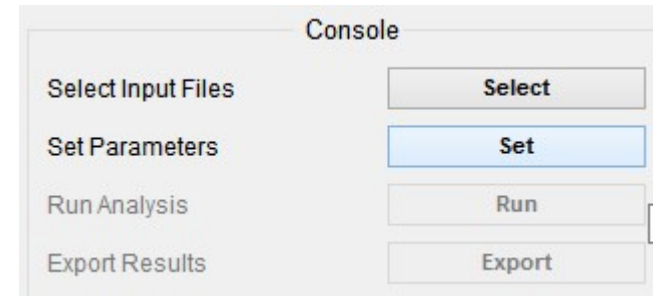
Welcome to SNP Filter...
Start Loading Input Files
Editing Input files... done.
Files looking OK.
Can now proceed to setting parameters.



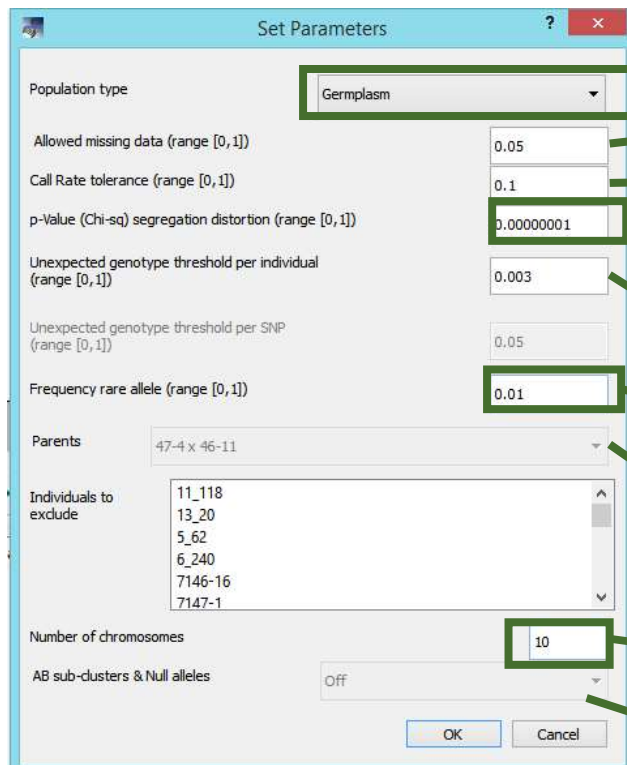
RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Running ASSIsT

4. In console section, click on 'Set' next to 'Set Parameters'



5. Set Parameters as follows:



-CP (F1) for (large) F1 population
-Germplasm in most other cases for us

Proportion of allowed missing data

How much can an individual's call rate differ from population mean

How much segregation distortion is allowed?
*lower means more distortion is allowed

Proportion of allowed inconsistencies between child and parents per SNP

Maximum frequency to define an allele as rare
*Only when Germplasm is chosen

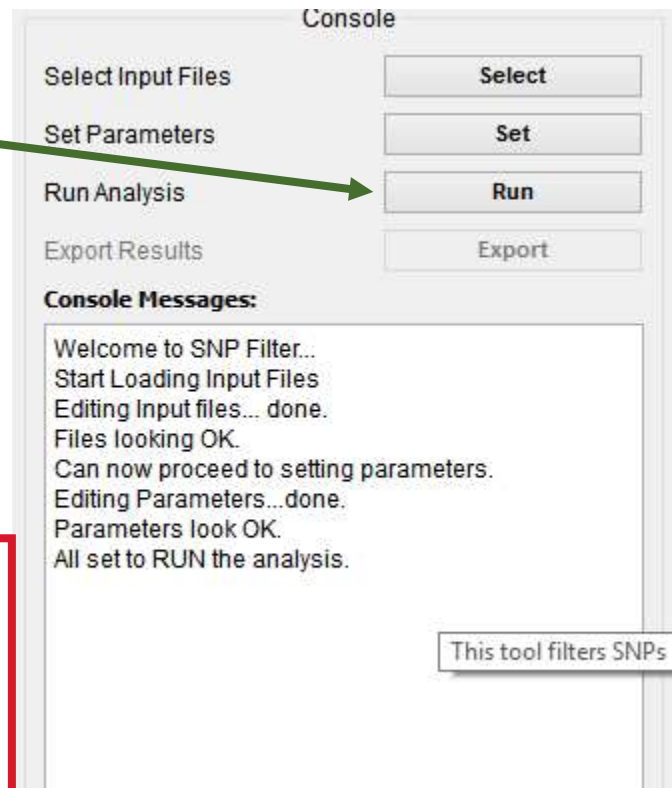
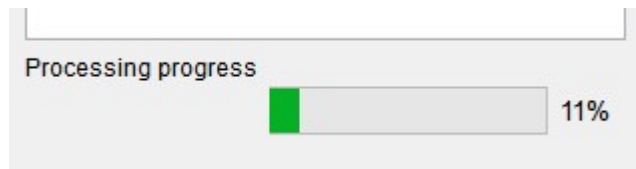
Parents (CP (F1), BCx) or grandparents (F2) of the analyzed experimental population.

Chromosome number
*Take "extra" chromosomes into account

Check for Null alleles and AB sub-clusters
*Only for CP (F1) and BC

Running ASSIsT

6. Run the analysis



Console

Select Input Files

Set Parameters

Run Analysis

Export Results

Console Messages:

Welcome to SNP Filter..
Start Loading Input Files
Editing Input files... done.
Files looking OK.
Can now proceed to setting parameters.
Editing Parameters...done.
Parameters look OK.
All set to RUN the analysis.

This tool filters SNPs

Note:

ASSIsT may stop responding when switching to other programs during the analysis. ASSIsT is still working and after a few minutes, it will show the results of the analysis

ASSIsT: An automatic SNP Scoring Tool in outcross populations v.1.00 (Not Responding)



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Results ASSIsT

SNP Graph similar to GenomeStudio
*non-interactive

The screenshot shows the ASSIsT software interface. On the left is a control panel with buttons for 'Select', 'Set', 'Run', and 'Export', and a 'Console Messages' window. The center features a 'SNP Viewer' window displaying a scatter plot of r vs Theta for the SNP 'RosBREED_snp_sweet_5_04486238'. The plot shows data points for genotypes AA (red), BB (blue), AB (purple), and NC/OO (grey). On the right is a 'Results Summary' window with a table of SNP classification for 383 individuals. At the bottom is a table of results for each SNP, ordered by classification and then by chromosome and position.

Results Summary

Number of SNPs: 13559
 Number of offsprings (Outcross): 383 (0)
 Total number of individual in FinalReport: 383
 Number of excluded samples (Poor quality DNA): 0

SNP Classification (383 individuals):

	#	%
Approved	6159	45.4
- Robust	1479	10.9
- OneHomozygRare_HWE	1677	12.4
- OneHomozygRare_NotHWE	1507	11.1
- DistortedAndUnexpSegreg	1496	11.0
Discarded	7400	54.6
- Monomorphic	632	4.7
- Failed	2716	20.0
- ShiftedHomo	4001	29.5
- NullAllele-Failed	51	0.4

Export	SNP id	Chr	Pos	Classification	Missing	Null	Hom1	Het1	Het2	Hom2	Chi2 pVal	GT Par1
<input checked="" type="checkbox"/>	RosBREED_snp_...	5	4.41373e+06	Robust	15	0	56	168	0	142	0.588	--
<input checked="" type="checkbox"/>	RosBREED_snp_...	5	4.48624e+06	Robust	13	0	63	159	0	146	0.086	--
<input checked="" type="checkbox"/>	RosBREED_snp_...	5	4.71063e+06	Robust	17	0	54	171	0	143	0.804	--
<input checked="" type="checkbox"/>	scaffold_5:4729...	5	4.73025e+06	Robust	22	0	42	153	0	164	0.491	--

Results for each SNP
Ordered by Classification, then by Chromosome and Position

Results ASSIsT

Results Summary

Number of SNPs: 13559

Number of offsprings (Outcross): 383 (0)

Total number of individual in FinalReport: 383

Number of excluded samples (Poor quality DNA): 0

SNP Classification (383 individuals):

	#	%
Approved	6159	45.4
- Robust	1479	10.9
- OneHomozygRare_HWE	1677	12.4
- OneHomozygRare_NoHWE	1507	11.1
- DistortedAndUnexpSegreg	1496	11.0
Discarded	7400	54.6
- Monomorphic	632	4.7
- Failed	2716	20.0
- ShiftedHomo	4001	29.5
- NullAllele-Failed	51	0.4

of excluded samples

SNPs used for further curation

Discarded SNPs*

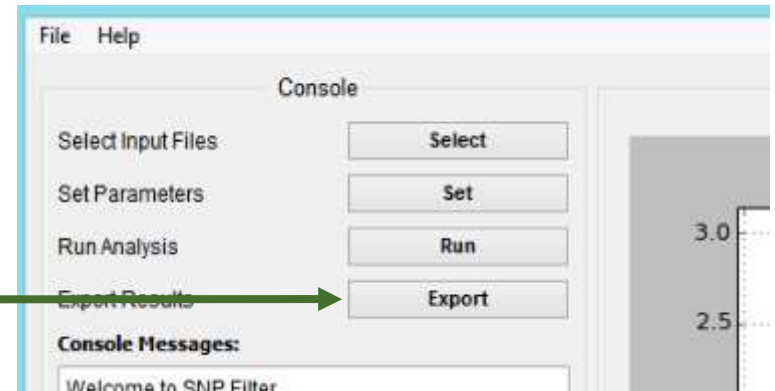


RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

*Note:
Shifted-Homo SNPs were added back in but many required manual adjustment of clustering

Export Results

1. Click on "Export" button in console section

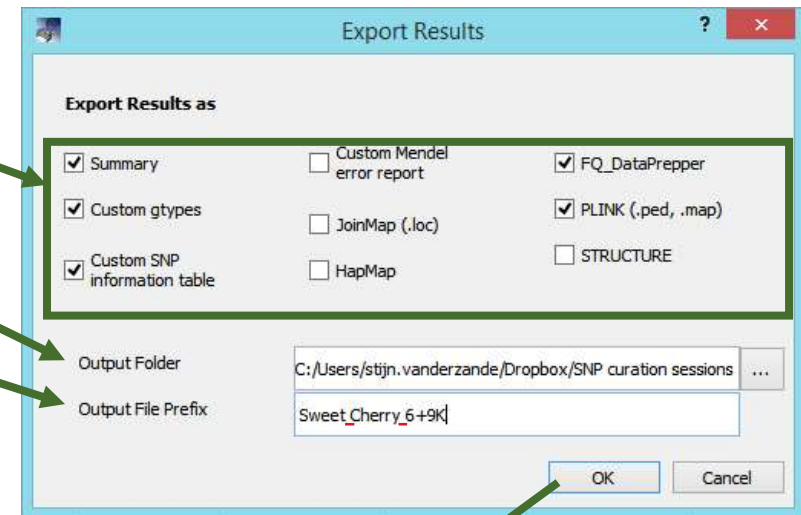


2.-Choose the files you want (next slide)

-Choose folder to save files in

-Choose name for files

***Avoid spaces in file names**



<input type="checkbox"/>	Name	Date modified	Type	Size
<input checked="" type="checkbox"/>	Sweet Cherry 6+9K_FQ_DataPrepper.txt	10/2/2018 6:43 PM	TXT File	10,322 KB
<input checked="" type="checkbox"/>	Sweet Cherry 6+9K_gtypes.csv	10/2/2018 6:43 PM	Microsoft Excel C...	7,414 KB
<input checked="" type="checkbox"/>	Sweet Cherry 6+9K_plink_in.map	10/2/2018 6:43 PM	MAP File	227 KB
<input checked="" type="checkbox"/>	Sweet Cherry 6+9K_plink_in.ped	10/2/2018 6:43 PM	PED File	9,223 KB
<input checked="" type="checkbox"/>	Sweet Cherry 6+9K_snp_info_table.csv	10/2/2018 6:43 PM	Microsoft Excel C...	1,274 KB
<input checked="" type="checkbox"/>	Sweet Cherry 6+9K_summary.txt	10/2/2018 6:43 PM	TXT File	2 KB
<input type="checkbox"/>	Sweet Cherry 6+9K_Map.txt	10/2/2018 5:23 PM	TXT File	453 KB

Export Results

Summary

Summarizes Parameter settings and Results (as shown on the right in ASSIsT)

```

Parameter Set
Population type:   Germplasm
Allowed missing data: 0.1
Call Rate tolerance: 0.1
p-Value (Chi-sq) segregation distortion: 1e-08
Unexpected genotype threshold: 0.009
Frequency rare allele: 0.01
Number of chromosomes: 10
AB sub-clusters & Null alleles: Off
    
```

Approved	6159	45.4
Robust	1479	10.9
OneHomozygRare_HWE	1677	12.4
OneHomozygRare_NotHWE	1507	11.1
DistortedAndUnexpSegreg	1496	11.0
Discarded	7400	54.6
Monomorphic	632	4.7
Failed	2716	20.0
ShiftedHomo	4001	29.5
NullAllele-Failed	51	0.4

Custom gtypes

Gives genotype for each individual and each **approved** SNP

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	
1	SNP id	Chr	Pos	Classificat	Missing	Null	Hom1	Het1	Het2	Hom2	Chi-Squar	1_118	13_20	5_62	6_240	714
2	RosBREE	1	43011	Distorted	25	0	14	265	0	77	0			AG	AG	AG
3	scaffold	1	43481	OneHomc	15	0	2	288	0	76	0			AC	AC	AC
4	RosBREE	1	51332	OneHomc	13	0	357	9	0	2	0			AA	AA	AA
5	S1_54448	1	54446	OneHomc	4	0	358	18	0	1	0			AA	AA	AA
5	scaffold	1	56261	OneHomc	12	0	1	8	0	360	0			GG	GG	GG
7	RocRFF	1	94695	OneHomc	11	0	356	17	0	2	0			AA	AA	AA

SNPs are ordered according to their position

Genotype score and individuals missaligned?!

Genotypes calls

Custom SNP information table

Gives results of analysis for each SNP (as shown on the bottom in ASSIsT)

PLINK (.ped, .map)

Two files (.ped and .map) needed to run PLINK and identify unknown duplicates (see further)

FQ_DataPrepper

Helps create FlexQTL files needed for further data curation (see next sessions)



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

duplicates

RosBREED

DISEASE
RESISTANCE



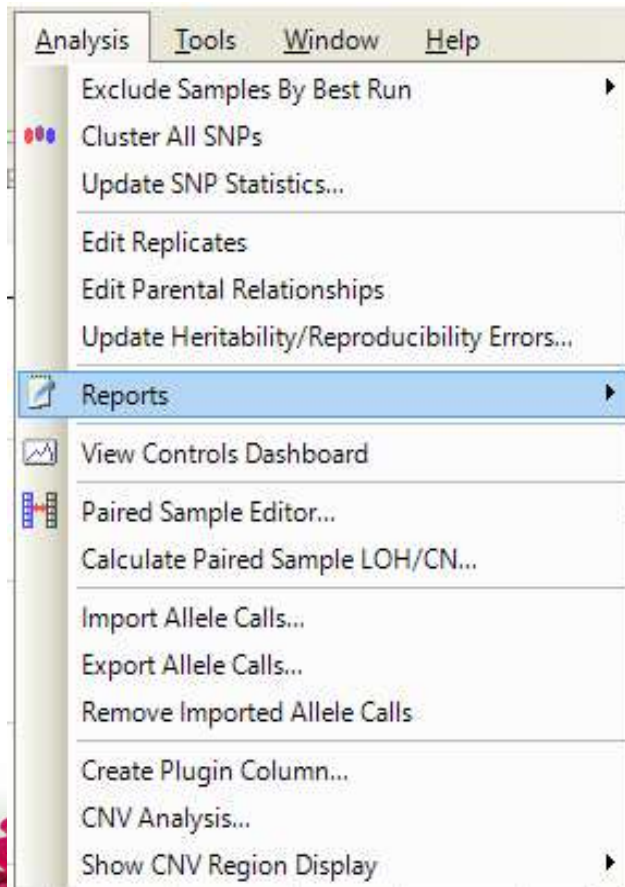
HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



Known duplicates



1. Create Reproducibility and Heritability report in GenomeStudio



2. Choose name and save file

Sweet Cherry 6+9K Reproducibility and Heritability Report.csv

CSV Files (*.csv)



Known duplicates

1. Open Reproducibility and Heritability report
Duplicate comparison found under 'Duplicate Reproducibility'

Rep1_DNA_Name	Rep2_DNA_Name	# Correct	# Errors	Total	Repro_Freq
Bing	Bing_rep01	8829	647	9476	0.965258
Bing_rep01	Bing	8829	647	9476	0.965258
FR070T105	FR070T105_rep01	11139	2	11141	0.99991
FR070T105_rep01	FR070T105	11139	2	11141	0.99991
FR072T074	FR072T074_rep01	11150	0	11150	1

Name Individual

Name Duplicate

Number consistent
genotype calls

Number inconsistent
genotype calls

Proportion inconsistent
genotype calls



RosBREED

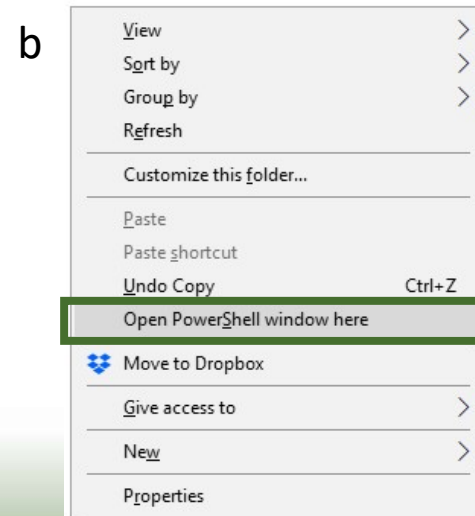
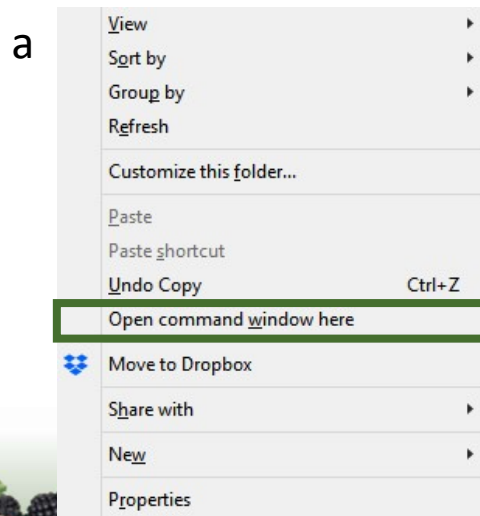
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Finding unknown duplicates

1. Copy .ped and .map file into PLINK folder
ASSIST output

	Date modified	Type	Size
Sweet Cherry 6+9K_plink_in.map	10/2/2018 6:43 PM	MAP File	
Sweet Cherry 6+9K_plink_in.ped	10/2/2018 6:43 PM	PED File	

2. Hold 'shift' and right click in the PLINK folder to get
 - a. the "Open command window here" option
 - b. the "Open PowerShell window here" option



3. (left) click on "Open command window here"/"Open PowerShell window here" to open the respective window

Using Plink – a. command window

4. Start with “plink.exe” then add additional commands with space between each. Press enter when all commands are given

*change “filename” to name of .ped and .map file

```
C:\Users\stijn.vanderzande\Documents\Software\PLINK 1.90>plink.exe --file Sweet_Cherry_6+9K --missing-genotype - --genome full
```

```
© 2005–2016 Shaun Purcell, Christopher Chang GNU General Public License v3
Logging to plink.log.
Options in effect:
--file Sweet_Cherry_6+9K
--genome full
--missing-genotype -

12248 MB RAM detected; reserving 6124 MB for main workspace.
.ped scan complete (for binary autoconversion).
Performing single-pass .bed write (6159 variants, 383 people).
--file: plink-temporary.bed + plink-temporary.bin + plink-temporary.fam
written.
6159 variants loaded from .bin file.
383 people (0 males, 0 females, 383 ambiguous) loaded from .fam.
Ambiguous sex IDs written to plink.nosex.
Using up to 8 threads (change this with --threads).
Before main variant filters, 383 founders and 0 nonfounders present.
Calculating allele frequencies... done.
Total genotyping rate is 0.968226.
6159 variants and 383 people pass filters and QC.
Note: No phenotypes present.
IBD calculations complete.
Finished writing plink.genome .
```

Expects same name for .ped and .map file

*e.g. Sweet_Cherry_6+9K.ped and Sweet_Cherry_6+9K.map

*NO blank spaces in name

Input file characterization

--file *filename*

Define files, change *filename* to name of ped and map file

--no-fid

when FID column is missing

--no-sex

when sex column is missing

--allow-no-sex

when sex column is missing, needed for some analyses

--no-pheno

when phenotype column is missing

--missing-genotype *N*

Define missing genotype (*N*) when different from 0

Analysis

--genome full

calculate IBD (IBS) for all individuals

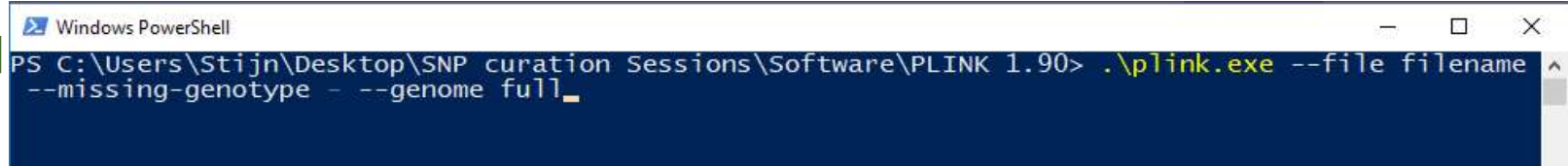


ROSE
DISEASE

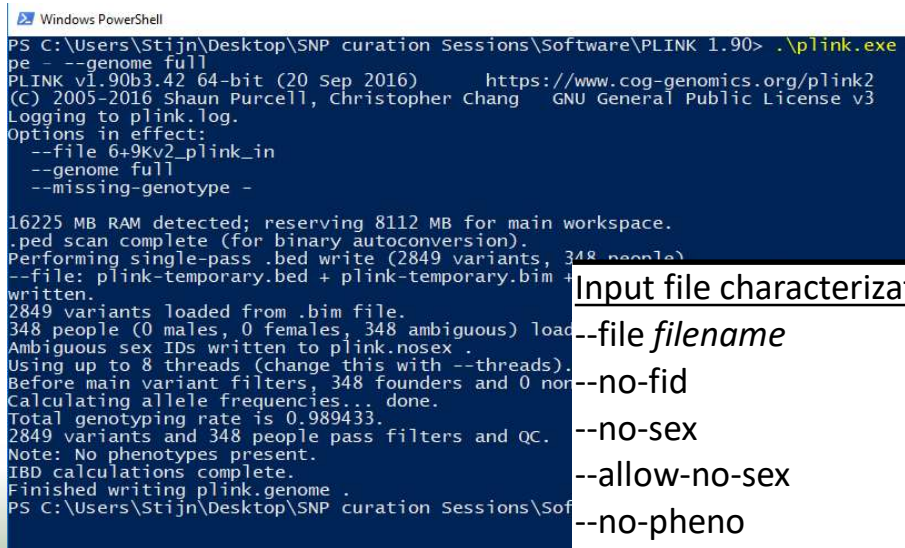
Using Plink – b. PowerShell window

4. Start with “.\plink.exe” then add additional commands with space between each. Press enter when all commands are given

*change “filename” to name of .ped and .map file



```
Windows PowerShell
PS C:\Users\Stijn\Desktop\SNP curation Sessions\Software\PLINK 1.90> .\plink.exe --file filename
--missing-genotype - --genome full_
```



```
Windows PowerShell
PS C:\Users\Stijn\Desktop\SNP curation Sessions\Software\PLINK 1.90> .\plink.exe
--genome full
PLINK v1.90b3.42 64-bit (20 Sep 2016)      https://www.cog-genomics.org/plink2
(C) 2005-2016 Shaun Purcell, Christopher Chang  GNU General Public License v3
Logging to plink.log.
Options in effect:
  --file 6+9Kv2_plink_in
  --genome full
  --missing-genotype -

16225 MB RAM detected; reserving 8112 MB for main workspace.
.ped scan complete (for binary autoconversion).
Performing single-pass .bed write (2849 variants, 348 people)
--file: plink-temporary.bed + plink-temporary.bim +
written.
2849 variants loaded from .bim file.
348 people (0 males, 0 females, 348 ambiguous) load
Ambiguous sex IDs written to plink.nosex .
Using up to 8 threads (change this with --threads).
Before main variant filters, 348 founders and 0 nor
Calculating allele frequencies... done.
Total genotyping rate is 0.989433.
2849 variants and 348 people pass filters and QC.
Note: No phenotypes present.
IBD calculations complete.
Finished writing plink.genome .
PS C:\Users\Stijn\Desktop\SNP curation Sessions\So
```

Expects same name for .ped and .map file
*e.g. Sweet_Cherry_6+9K.ped and Sweet_Cherry_6+9K.map
*NO blank spaces in name

Input file characterization

- file *filename* Define files, change *filename* to name of ped and map file when FID column is missing
- no-fid when FID column is missing
- no-sex when sex column is missing
- allow-no-sex when sex column is missing, needed for some analyses
- no-pheno when phenotype column is missing
- missing-genotype *N* Define missing genotype (*N*) when different from 0

Analysis

- genome full calculate IBD (IBS) for all individuals



Using Plink

5. Open plink.genome with Excel

✓ plink.genome	10/3/2018 10:07 AM	GENOME File	12,151 KB
plink.log	10/3/2018 10:07 AM	LOG File	2 KB

6. Select first column, use “Text to Columns” to create separate columns

The screenshot shows the Excel 'Text to Columns' wizard. The 'Original data type' section has 'Delimited' selected. The 'Delimiters' section has 'Space' checked. A green arrow points from the 'Text to Columns' button in the ribbon to the 'Delimited' radio button. Another green arrow points from the 'Space' checkbox to the 'Delimiters' section header.

Original data type
Choose the file type that best describes your data:
 Delimited - Characters such as commas or tabs separate each field.
 Fixed width - Fields are aligned in columns with spaces between each field.

Delimiters
 Tab
 Semicolon
 Comma
 Space
 Other:

7. Sort according to the “PI_HAT” column, “Largest to Smallest”

The screenshot shows the Excel 'Sort' dialog box. The 'Sort by' dropdown is set to 'PI_HAT' and the 'Order' dropdown is set to 'Largest to Smallest'. A green arrow points from the 'Sort' button in the ribbon to the dialog box.

Column	Sort On	Order
Sort by PI_HAT	Values	Largest to Smallest

Plink Results

Individual 1		Individual 2		IBS	HE	DST	PPC	RATIO	IBS0	IBS1	IBS2	HOMH	HETHET
FID1	IID1	FID2	IID2	PI_HAT									
	MBigarreau_dup02		MBigarreau_dup01	1	-1	1	1	NA	0	0	2805	0	232
	Lambert		KootenayLambert	1	-1	1	1	NA	0	0	2849	0	218
	99F/132R4		99F131_RJ	1	-1	1	1	NA	0	0	2849	0	260
	MertonBigarreau		MBigarreau_dup02	0.9995	-1	0.999824	1	NA	0	1	2842	0	234
	MertonBigarreau		MBigarreau_dup01	0.9995	-1	0.999822	1	NA	0	1	2806	0	232
	Skeena		Santina_dup01	0.9985	-1	0.999473	1	NA	0	3	2842	0	197
	Santina_dup01		Santina	0.9985	-1	0.999473	1	NA	0	3	2842	0	197
	MertonHeart		MertonHeart_dup01	0.9984	-1	0.999422	1	NA	0	3	2590	0	249
	Hedelfingen		Hedelfingendup01	0.998	-1	0.999295	1	NA	0	4	2833	0	261
	StarBlush		StarBlush_dup01	0.997	-1	0.99894	1	NA	0	6	2824	0	231
	SandraRose_dup01		SandraRose	0.9892	-1	0.99618	1	NA	0	21	2728	0	222
	Vega		Cv1	0.976	-1	0.991929	1	NA	1	40	2561	0	213
	FR063T023		FR055T064	0.8279	-1	0.939648	1	97	1	300	2201	1	97

0.97 as cutoff for duplicates

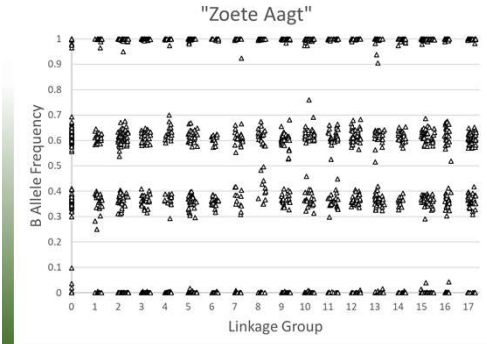


RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

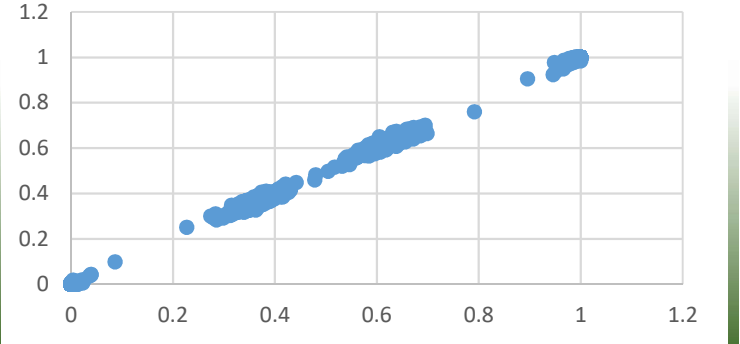
Duplicate triploids

- Calls should be the same between duplicates
 - Are AB calls identical?
 - Could be AAB or ABB
 - B allele frequency
 - Should both be 0.33 or 0.66
 - Correlation between identically called individuals

IID1	FII IID2	PI_HAT
1 "PommeHervi"	1 "StreepingAlken"	1
1 "ZoeteAagt"	1 "BelledeFumes"	1



“Zoete Aagt” vs. “BelledeFumes”





United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Pedigree check

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



Checking parent-child relationships

1. Open 'Reproducibility and Heritability' report created earlier for known duplicates
 - A. Parent-child errors found under 'P-C Heritability'

P-C Heritability	Child_DNA_Name	Parent_DNA_Name	# Correct	# Errors	Total	P-C Heritability Freq
	Abundance	Napoleon	0	0	0	N/A
	Burbank	EarlyPGuigne	0	0	0	N/A
	Coe	OxHeart	0	0	0	N/A
	Summit	Van	0	0	0	N/A
	Kordia	Schneiders	11593	2	11595	0.9998275
	Sunburst	Summit	11601	8	11609	0.9993109

Only when parent in dataset

Name child

Name Parent

Number genotype calls without PC error

Number genotype calls with PC error

Proportion genotype calls with PC error



ROSBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Checking parent-parent-child relationships

1. Open 'Reproducibility and Heritability' report created earlier for known duplicates
 - B. Parent-child errors found under 'P-P-C Heritability'

Child_DNA_Name	Parent1_DNA_Name	Parent2_DNA_Name	# Correct	# Errors	Total	P-P-C Heritability Freq
Chelan	Stella	Moreau	11554	12	11566	0.999
Index	Stella	Bing	11551	12	11563	0.999
Benton	Stella	Moreau	11541	15	11556	0.9987
Vic	Bing	Schmidt	11542	20	11562	0.9983
Venus	Hedelfingen	Windsor	11475	31	11506	0.9973
Santina	Stella	Summit	11205	339	11544	0.9706
Santina_dup01	Stella	Summit	11102	350	11452	0.9694

Name child

Name Parent 1

Name Parent 2

Number genotype calls without PC error

Proportion genotype calls with PPC error

Number genotype calls with PPC error



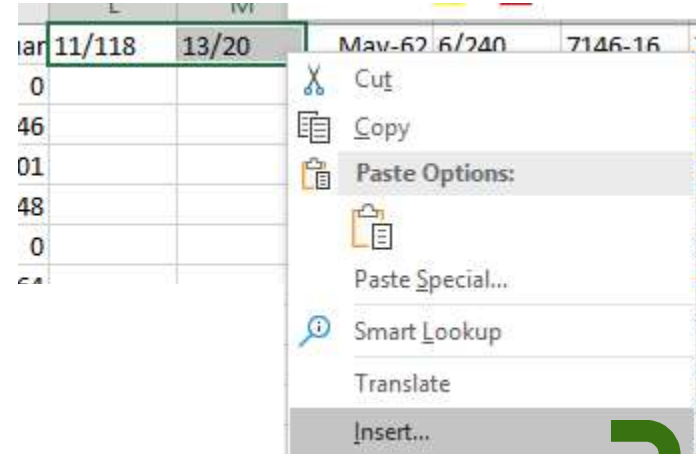
RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Finding new P(P)C relationships

1. Open ASSIsT's ".gtypes" output in Excel

6+9K v2_gtypes.csv - Excel

2. Ensure genotypic data aligns with individuals' ID



3. Copy and transpose data so that individuals are in rows and markers in columns

	B	C	D	E	F	G
	RosBREED scaffold_1	scaffold_1	scaffold_1	scaffold_1	scaffold_1	scaffold_1
13/20	AA	AA	AA	AA	GG	AA
May-62	AA	AA	AA	AA	GG	AA
6/240	AG	AA	AA	AA	GG	AA
7146-16	AA	AG	AG	AG	AG	AC
7147-1	GG	AG	AG	AG	AG	AA

Insert

Insert

- Shift cells right
- Shift cells down
- Entire row
- Entire column

OK Cancel

Finding new P(P)C relationships

4. Bring in pedigree information for all individuals (parents do not have to be included in data set)

5. Set missing parents to “-”

6. Use “Ind”, “Parent1” and “Parent2” as column headers

Ind	Parent1	Parent2	RosE
11/118			AA
13/20			AA

7. Save file as “.csv” in same folder as “Parent Check” R script (Suppl. Document 2)



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Finding new P(P)C relationships

8. Load libraries in R

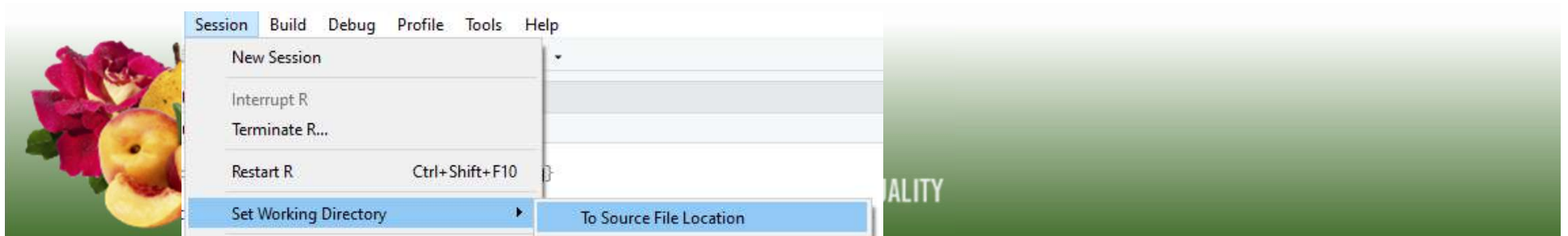
```
library(stringr)  
library(svmisc)
```

9. Define functions in R by running them

```
##Functions  
AdjustMV <- function(GTData){  
  
CheckParErr <- function(GenotypeIndPar){  
  
CheckParParErr <- function(GenotypeIndParPar){  
  
CheckPar <- function(GenotypesIndParPar){  
  
FindParGT <- function(IndName, Genotypes){  
  
CheckParAll <- function(IndToCheck, ChrReport=NULL){  
  
FindPosPar <- function(IndName, GTData, treshold){  
  
FindPosParComb <- function(GTData, tresholdPE=0, tresholdPPE=0){
```

Run

10. Set working directory to source file location



Finding new P(P)C relationships

11. Define:

- Possible Alleles under AlleleList
 - \$ or \$\$ does not work for null alleles
- Characters used for missing genotypes under “MissGT”
- Characters used for missing alleleles under “MissAllele”

```
AlleleList <- c("A","B","C","null", "G", "T")  
MissGT <- c("NC","00", "--")  
MissAllele <-c("N","0","-")
```

12. Provide filename and load into R

```
## Load data files  
GenotypeData <- as.data.frame(read.csv(file="sweet_cherry_6+9K.csv", head=T, sep="," , stringsAsFactors=F))
```

13. Run “FindPosParComb”

Loaded data file

Threshold to accept PC relationship

Threshold to accept PPC relationship

```
> FindPosParComb(GenotypeData, tresholdPE = 70, tresholdPPE = 100)  
[1] "Starting step 1 of 3"  
[1] "Starting step 2 of 3"  
Progress: 1 on 40
```

3 steps:

- Identify possible mothers for ind. with missing mothers
- Identify possible fathers for ind. with missing fathers
- Identify parents for ind. with both parents missing



Checking and finding grandparent-grandchild relationships

- Excel - Suppl file 1 from van de Weg et al. (2018)
 - If offspring is 'AB' and
 - 1 known parent is 'AA' OR
 - Two known grandparents through single parent are both 'AA'
 - Then: one unknown grandparent cannot be 'AA'
 - Minimize 'AA' calls gives indication of possible grandparent
 - For any 'AA' in putative grandparent, second unknown grandparent cannot be 'AA'
 - Also true for 'AB' in individual and 'BB' in known parents/grandparents



van de Weg et al. (2018) method

1. Filter individual for AB

2. Filter parent 1 for AA

3. Parent 2 should have 0 AA

The image shows a multi-step data filtering process in a software interface. The interface includes a table of individuals and their parents, a filter menu, and a genotype counts table.

Table 1: Individual and Parent Information

Individual	Enterprise	Individual	Enterprise
1661-2	McIntosh	1661-2	McIntosh
Coop-7	PRI1661-2	Coop-7	PRI1612-1

Table 2: Genotype Counts

Genotype	Count
A B+B B	481
AA	0
AB	306
B B	175

Table 3: Individual and Parent Information (Detailed)

Individual	Enterprise	Mother	Father
Coop-7	PRI1661-2	PRI1018-9	PRI1036

The interface also shows a filter menu with the following options:

- Sort A to Z
- Sort Z to A
- Sort by Color
- Clear Filter From "(Column I)"
- Filter by Color
- Text Filters
- Search
- (Select All)
- AA
- AB
- BB



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Curating remaining Mendelian-inconsistent errors

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS





United States
Department of
Agriculture

National Institute
of Food and
Agriculture

FlexQTL Data Prepper

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



FlexQTL Data Prepper

- FlexQTL Data Prepper
 - Requires 3 input files
 - Marker file
 - Adjusted ASSIsT output file
 - Pedigree file
 - Adjusted ASSIsT input file
 - Data file
 - Generated by ASSIsT



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL Data Prepper

Creating 1st input file (Marker file)

1. Create a three-column file:

- 1st column has the SNP's name/SNP's ID
- 2nd column has the chromosome (numerical)
- 3rd column has the genetic position (cM))
- Header row should be "MarkerId", "Group", "Position"
- Save as ".csv" file

*Start from "..._plink_in.map" ASSIsT output file

*Change header row

*Change physical position to genetic position if needed

6+9Kv2_plink_in.map

6+9Kv5_map.csv

	A	B	C
MarkerId		Group	Position
RosBREED_snp_tart_1_0021		1	0
scaffold_1:425480		1	3.62
scaffold_1:427005		1	3.64
scaffold_1:427260		1	3.65
scaffold_1:429080		1	3.68
scaffold_1:433354		1	3.75
RosCOS1201-071_snp_swee		1	4.04
scaffold_1:459954		1	4.2



Ros
DISEASE R

FlexQTL Data Prepper

Creating 2nd input file (Pedigree file)

2. Create a three-column file:

- 1st column has the individual's name
- 2nd column has the female parent
- 3rd column has the male parent
- Header row should be "Name", "Parent1", "Parent2"
- Save as ".csv" file

*Start from ASSIST pedigree input file

*Change header row

A	B	C
Name	Parent1	Parent2
Abundance	Napoleon	
BlackRepublican	Napoleon	BlackTartarian
EmperorFrancis		
Gil-Peck	Napoleon	Giant
Hedelfingen		
Kordia	Schneiders	
Lambert	Napoleon	Blackheart
Summit	Van	
Sunburst	Summit	
Van	EmpressEugenie	BlackRepublican



FlexQTL Data Prepper

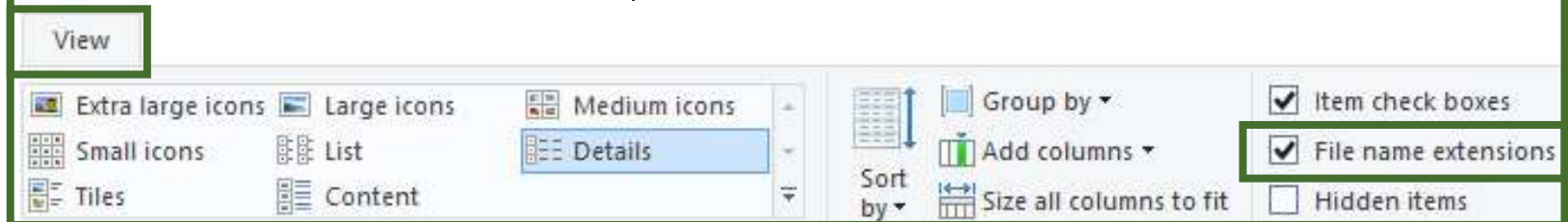
Creating 3rd input file (Data file)

3. Adjust "..._FQ_DataPrepper.txt" file format from ".txt" to ".csv"



Note:

If file format extension is not visible, enable it under the "View Tab" in the file's folder



4. Open "..._FQ_DataPrepper.txt" file in Excel



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL Data Prepper

5. Adjust header row so it matches genotypic data

SNPId	-	--	Abundanc	Ambrunes	Benton	Bing	Walpurgis	WhiteGol	Windsor	YelGlass
scaffold_1	AA	AA	AA	AA	AA	AA	AA	AA		
RosBREED	AA	AA	AA	AA	AA	AA	AA	AA		
S1_54448	AA	AA	AA	AA	AA	AA	AA	AA		
s7_888812	GG	GG	GG	GG	GG	GG	GG	GG		
scaffold_1	GG	GG	GG	GG	GG	GG	GG	GG		

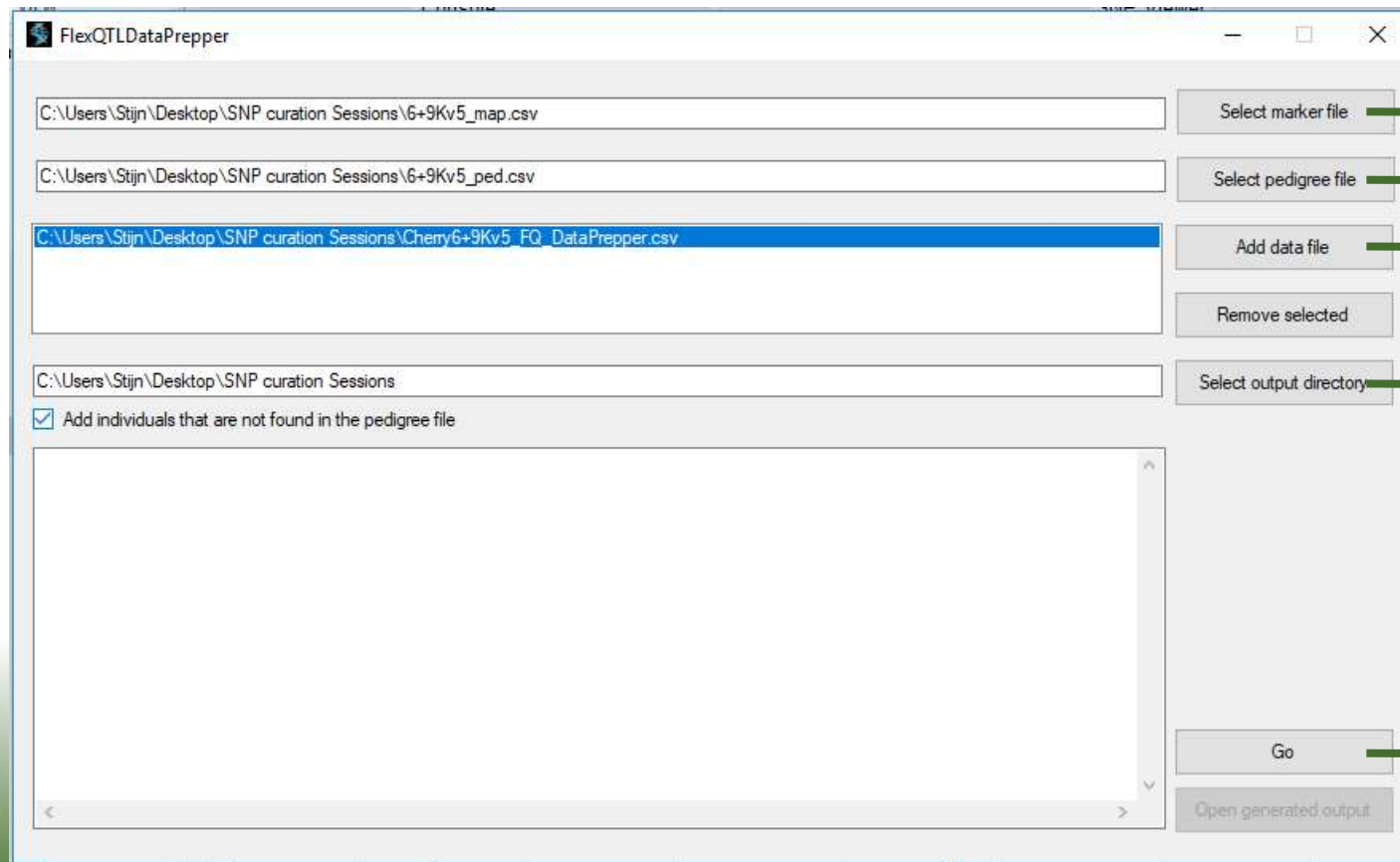
6. Change cell A1 to "ID"

7. Save file

ID	Abundanc	Ambrunes	Benton	Bing
scaffold_1	AA	AA	AA	AA
RosBREED	AA	AA	AA	AA
S1_54448	AA	AA	AA	AA
s7_888812	GG	GG	GG	GG
scaffold_1	GG	GG	GG	GG

FlexQTL Data Prepper

8. Open FlexQTLDataPrepper.exe



- 9. Load marker file
- 10. Load pedigree file
- 11. Load data file
- 12. Choose output directory
- 13. Select this option
- 14. Press 'Go'



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

FlexQTL™ Input files

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



FlexQTL™ input files

- Three files
 - Data file (.dat)
 - Map file (.map)
 - Parameter file (.par)
- Text-based files
- Everything after “;” will be ignored
 - Easy excluding of rows in file



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

- Data file (.dat file)
 - No header needed (if added, start with “;”)
 - Column 1: Population
 - Column 2: Individual
 - Column 3: Parent 1
 - Column 4: Parent 2
 - Both parents needed (or none)
 - Use dummy individuals when only one parent is available
 - Column 5 – X: nuisance variables
 - Columns X+1 – Y: Phenotypic variables
 - Columns Y+1 – End: Genotypic data
 - Two columns per marker
 - One column per marker with space(s) between alleles



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

- Map file (.map)
 - Column 1: marker name
 - Column 2: genetic position within LG
- Each LG starts with “group X”
 - X is LG number
- Identical number of markers needed as in data file



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

- Parameter file (.par)
 - Always called “flexqtl.par”
 - Defines parameters for FlexQTL
 - Some are adjusted automatically through Visual FlexQTL (see further)
 - Some can be adjusted within FlexQTL
 - Need to be “correct” when running FlexQTL independently from Visual FlexQTL
 - E.g. on Linux



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

- FlexQTL Data Prepper output needs additional adjustments
 - Parents without genotype data need to be added
 - Every individual in pedigree needs a data row for FlexQTL
 - Individuals with only one known parent need a second dummy parent
 - FlexQTL only accepts no or 2 parents known

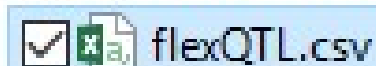


RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

Adjusting FlexQTLDataPrepper file to create 1st FlexQTL input file (data file)

1. Open “.csv” file generated by FlexQTLDataPrepper



2. Find individuals with only one parent known/given and add a second ‘dummy’ parent

-Use “M_[Individual’s name]” for unknown mother

-Use “F_[Individual’s name]” for unknown father

-Use “UP_[Individual’s name]” if unclear whether known parent is mother or father

	B	C	D	
Individual		Parent1	Parent2	Ro
1 Abundance		Napoleon	F_Abundance	A C
1 Summit		Van	F_Summit	A A
1 Sunburst		Summit	F_Sunburst	A A



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

3. Open “.csv” file generated by FlexQTLDataPrepper



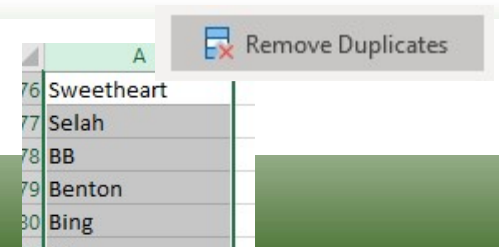
4. Copy both parental columns (columns C and D) to a new sheet

	A	B	C	D	
1	;BaseCoh	Individual	Parent1	Parent2	R
2		1 Abundance	Napoleon	F_Abundance	A
3		1 Ambrunes		0	0 G
4		1 BlackRepublican	Napoleon	BlackTartarian	A
5		1 EmperorFrancis		0	0 A
5		1 Gil-Peck	Napoleon	Giant	A

5. Copy “Parent 2” column below “Parent 1” column

	A
45	EE
46	Rainier
47	Rainier
48	F_Abundance
49	0
50	BlackTartarian

6. Remove duplicates* from combined column
*found under “data” tab in Excel



CROSBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

7. Count how often each parent is in datafile:

use “COUNTIF” function

Range is “Individual” column (column B) of original sheet

Criteria is parent to be checked (from created column in new sheet)

A	B	C	D
Ambrunes	=COUNTIF(flexQTL!B:B,Sheet1!A1)		
BB	COUNTIF(range, criteria)		
BedfordProl	0		
Benton	1		
Bertiolle	0		
Bing	1		
BlackBeauty	1		

8. Filter for individuals not in the data file (count = 0)

Ambrunes		▼
BedfordProl	0	
Bertiolle	0	
Cristobalina	0	
EmpressEugenie	0	



R
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

9. Copy parents not in data file at the bottom of the data file

	A	B	C	D	E	
00	1	Ulster	Napoleon	Gil-Peck	A A	A G
01		BedfordProl				
02		Bertiolle				
03		Cristobalina				
04		EmpressEugenie				

A	B
Ambrunes	
BedfordProl	0
Bertiolle	0
Cristobalina	0
EmpressEugenie	0



10. Fill in “Population” column (Column A), “Pedigree info” (Column C and D), and Genotypic data columns (Column E onwards) of added parents

-”0” for unknown pedigree

-Make sure newly-added pedigree info is also in data set

-“- -” for genotypic data

	A	B	C	D	E	F	G	H	I
00	1	Ulster	Napoleon	Gil-Peck	A A	A G	A G	A G	A G
	1	BedfordProl	0	0	--	--	--	--	--
	1	Bertiolle	0	0	--	--	--	--	--
	1	Cristobalina	0	0	--	--	--	--	--
	1	EmpressEugenie	0	0	--	--	--	--	--



RUSDREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

11. Add 2 columns between parentage information (Column C+D) and genotypic data (Column E)
12. Fill 1st new column (Column E) with “1”s
*Used as dummy nuisance column
13. Fill 2nd new column (Column F) with random numeric values
*Used as dummy phenotype column

D	E	F	G
parent2	DummyNuis	DummyPheno	RosBREED sc
Abundance	1	=RANDBETWEEN(0,100)	G
0	1	RANDBETWEEN(bottom, top)	;
ackTartarian	1	87	A G G

14. Copy full data sheet into text editor (e.g. Notepad++)
15. Save file in text format
16. Change new file's format extension from “.txt” to “.dat”



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

FlexQTL™ input files

2nd FlexQTL input file (map file)

17. Use flexQTL.map file generated by FlexQTL Data Prepper

3rd FlexQTL input file (parameter file)

18. Save Suppl. Table 4c as text file

* “datafile”, “mapfile”, “indiC” and “nmrkrC” parameter settings will be updated automatically when loading files into Visual FlexQTL

19. Change file extension from “.txt” to “.par”

20. Change file name to “flexqtl”

FlexQTL does not accept any other name!



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Installing (Visual) FlexQTL™

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



Installation

- Install Visual FlexQTL
 - Will also install FlexQTL
- Install following packages in R
 - Data.tables
 - Plyr
 - Lattice

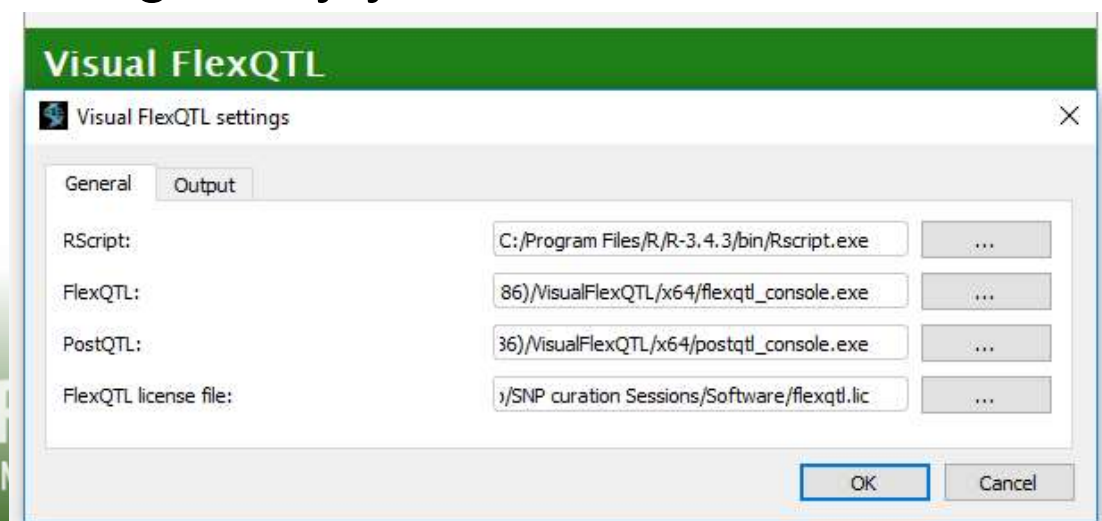


RosBREED

DISEASE RESISTANCE × HORTICULTURAL QUALITY

Installation

- Tools > Settings
 - General Tab
 - Where is Rscript.exe on pc
 - Where is flexqtl_console.exe on pc
 - Where is postqtl_console.exe on pc
 - Where can license be founds
 - Needs updating every year



Installation – Win8 or later

- Install Visual C++ Redistributable for Visual Studio 2012
 - <https://www.microsoft.com/en-us/download/details.aspx?id=30679>

Visual C++ Redistributable for Visual Studio 2012 Update 4

Download

Choose the download you want

File Name

VSU_4\vc redistrib_x64.exe

VSU_4\vc redistrib_x86.exe

VSU4\vc redistrib_arm.exe

Most new pc
need this one



RosB
DISEASE RESISTA



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Finding Mendelian- inconsistent errors

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



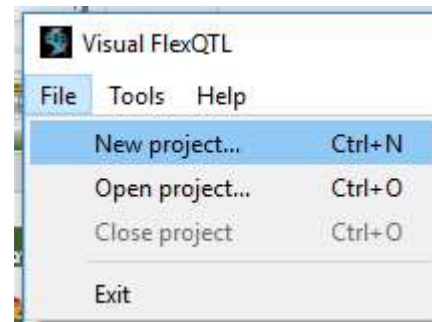
SUPERIOR
CULTIVARS



Creating project in Visual FlexQTL™

1. Open Visual FlexQTL™

2. Under “File”, choose “New Project”



Create new project

Create new project

To create a project, select a project name and specify the location of the project in your file system. If you want to create a project for an existing map and/or data file, then you can set these files to be imported into your project.

Project name and location

Project name: 6+9K

Project location: C:/Users/Stijn/Desktop/SNP curation Sessions/FlexQTL

Project folder: sers/Stijn/Desktop/SNP curation Sessions/FlexQTL\6+9K

Import data files

Parameter file: :sktop/SNP curation Sessions/FlexQTL/6+9kv5/flexqtl.par

Map file: :top/SNP curation Sessions/FlexQTL/6+9kv5/6+9Kv5.map

Data file: :ktop/SNP curation Sessions/FlexQTL/6+9kv5/6+9Kv5.dat

OK

Cancel

3. Choose name for project

4. Choose project directory

Filled in automatically

4. Load parameter file

5. Load map file

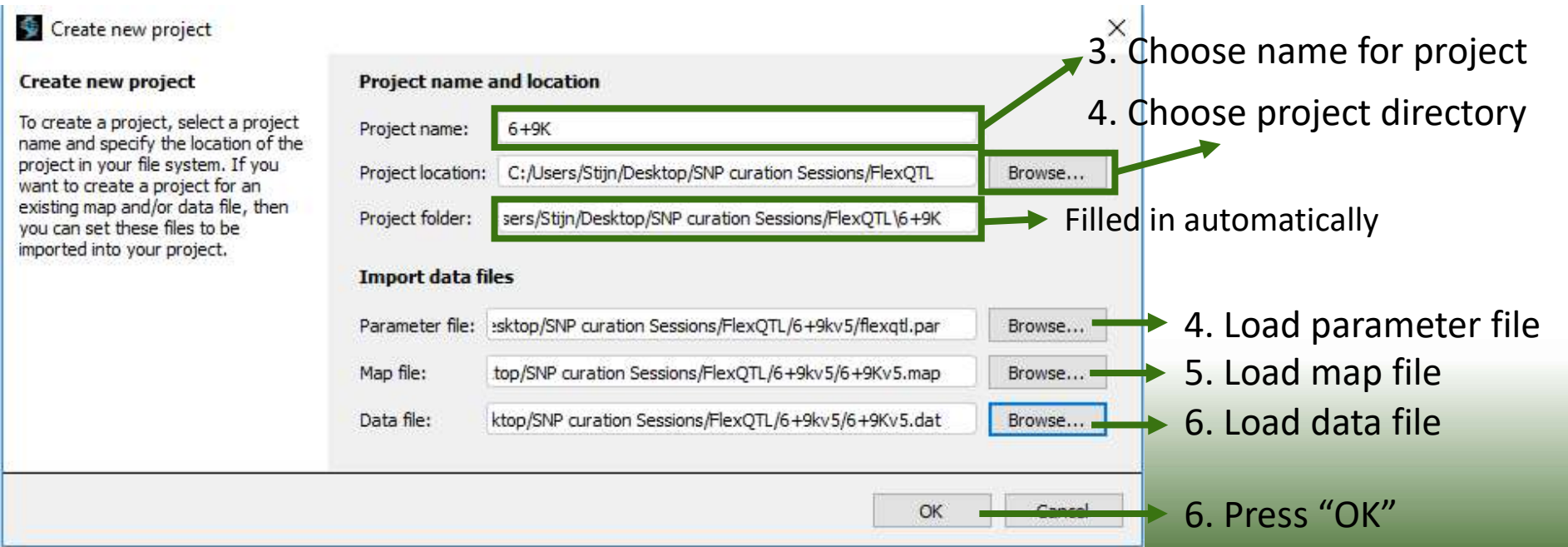
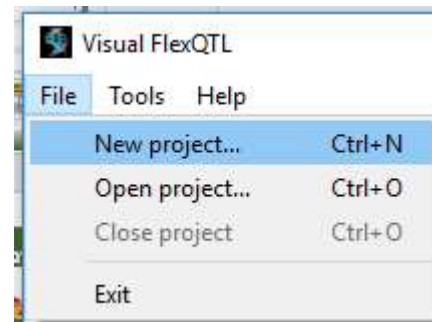
6. Load data file

6. Press “OK”

Creating project in Visual FlexQTL™

1. Open Visual FlexQTL™

2. Under “File”, choose “New Project”

A screenshot of the 'Create new project' dialog box in Visual FlexQTL. The dialog is titled 'Create new project' and contains the following fields and buttons:

- Project name and location:**
 - Project name: 6+9K (highlighted with a green box, with an arrow pointing to it from the annotation '3. Choose name for project')
 - Project location: C:/Users/Stijn/Desktop/SNP curation Sessions/FlexQTL (with a 'Browse...' button next to it, and an arrow pointing to it from the annotation '4. Choose project directory')
 - Project folder: sers/Stijn/Desktop/SNP curation Sessions/FlexQTL\6+9K (highlighted with a green box, with an arrow pointing to it from the annotation 'Filled in automatically')
- Import data files:**
 - Parameter file: :sktop/SNP curation Sessions/FlexQTL/6+9kv5/flexqtl.par (with a 'Browse...' button next to it, and an arrow pointing to it from the annotation '4. Load parameter file')
 - Map file: :top/SNP curation Sessions/FlexQTL/6+9kv5/6+9Kv5.map (with a 'Browse...' button next to it, and an arrow pointing to it from the annotation '5. Load map file')
 - Data file: :ktop/SNP curation Sessions/FlexQTL/6+9kv5/6+9Kv5.dat (with a 'Browse...' button next to it, and an arrow pointing to it from the annotation '6. Load data file')
- Buttons: 'OK' and 'Cancel' (with an arrow pointing to the 'OK' button from the annotation '6. Press "OK"')

3. Choose name for project

4. Choose project directory

Filled in automatically

4. Load parameter file

5. Load map file

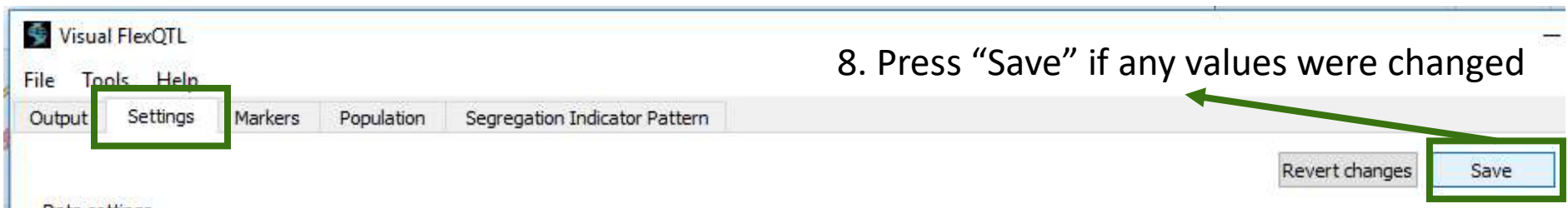
6. Load data file

6. Press "OK"

Checking parameter settings

7. Check parameter settings under the 'Settings' tab

8. Press "Save" if any values were changed



Allow marker loci segregation distortion:

Checked to allow for segregation distortion

Delete double recombinants:

Unchecked to keep genotypic data for single marker double recombinations

Create pedimap visualisation files:

"2" for an early stop; generates files for error checking but does not do full FlexQTL™ analysis (can take days)

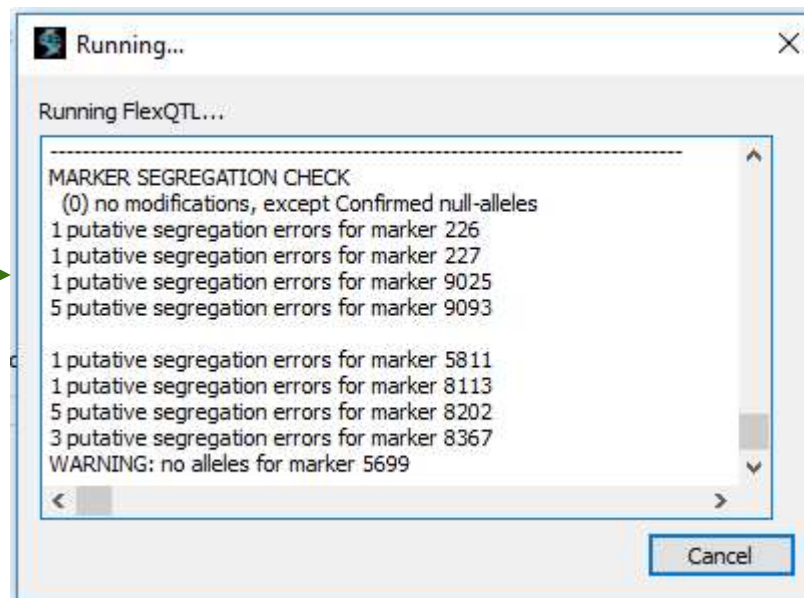
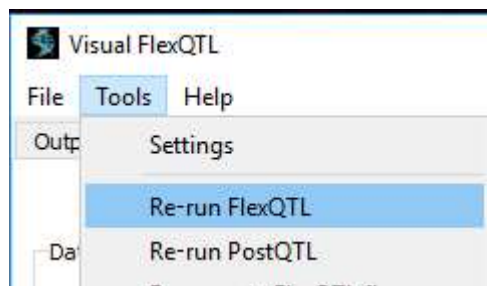


RosBREED

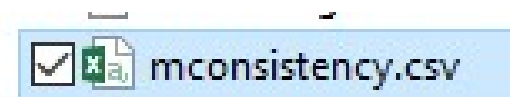
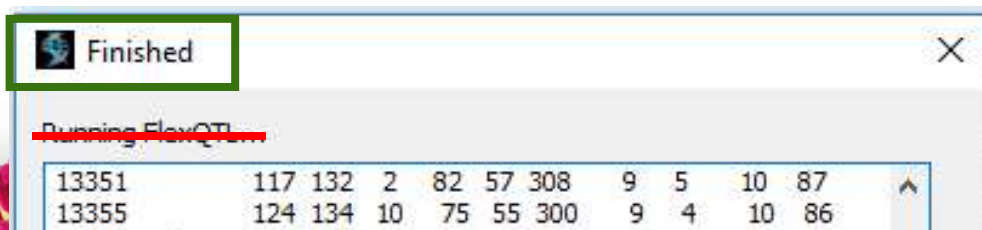
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Running FlexQTL™

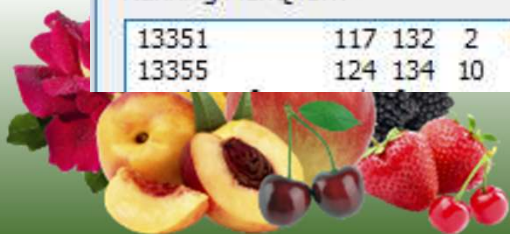
8. Run FlexQTL (Tools>Re-run FlexQTL)



9. Once run, go to project directory and open “mconsistency.csv”



10. Save file in “.xls(x)” format



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Mconsistency.csv

Information found in mconsistency file (1)

Code_Pec	Symbol	ShortNam	Description
0	.	OC	observed & consistent
1	&	MA	missing & augmented
2	}	OEA	observed & deleted error & augmented
3]	pnuA	putative null-allele & augmented
4	~	pnuM	putative null-allele & missing
5	-	IM	incomplete & missing
6	/	OpE	observed & putative error
7	{	OE	observed & deleted error
8	[pnuE	putative null-allele & error
9	#	ME	missing & error

Consistencies and errors reported by the mconsistency file

marker position	SUM({[#)	{	[#	}	/	SUM(}/)	~
SUM({[#)	1180							
{		351						
[0					
#				829				
}					484			
/						15		
SUM(}/)							499	
~								0

Summary of errors found in data set

Mconsistency.csv

Information found in mconsistency file (2)

position	SUM({[#]	{	[#	}	/	SUM({/)	~
Napoleon	0	0	0	0	0	0	0	0
F_Abunda	0	0	0	0	0	0	0	0
Abundanc	0	0	0	0	0	0	0	0
Ambrunes	2	2	0	0	6	0	6	0

Summary of errors found for each individual

Summary of errors found for each marker

marker	4071	8853	8854
position	100	100.0362	100.0364
SUM({[#]	1	1	0
{	0	0	0
[0	0	0
#	1	1	0
}	0	0	0
/	0	0	0
SUM({/)	0	0	0
~	0	0	0

Observation (consistent genotype or erroneous genotype) for each individual and marker combination

marker	4071	8853	8854	8855
position	100	100.0362	100.0364	100.0365
Napoleon	-	&	&	&
F_Abunda	-	-	-	-
Abundanc	.	-	-	-
Ambrunes	.	.	-	.
Angela
Blackhear	-	-	-	-
Lambert
J12420	-	-	-	-
Stella
Morreau

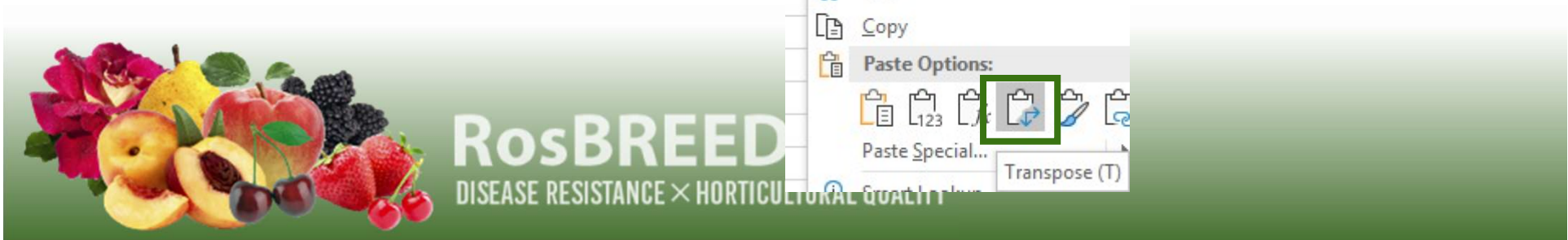
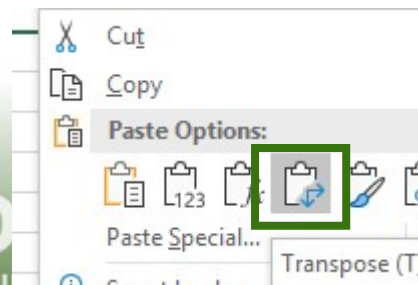


Mconsistency.csv

11. Copy data from row 15 onward and all columns

13	;	9 #	ME	missing & error							
14											
15	marker								4071	8853	
16	position								100	100.0362	
17		SUM({[#	{	[#	}	/	SUM(}/)	~		
18	SUM({[#	1180								1	1
19	{		351							0	0
20	[0						0	0
21	#				829					1	1
22	}					484				0	0
23	/						15			0	0
24	SUM(}/)							499		0	0
25	~								0	0	0
26	11_118	0	0	0	0	0	0	0	0	.	.
27	13_20	0	0	0	0	0	0	0	0	.	.
28	5_62	0	0	0	0	0	0	0	0	.	.
29	6_240	0	0	0	0	0	0	0	0	.	.
30	99F_131RJ	0	0	0	0	0	0	0	0	.	.

12. Transpose copied data in new sheet



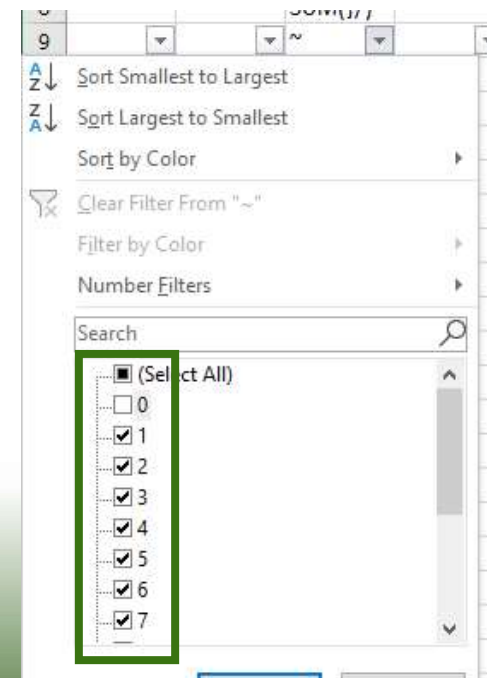
Mconsistency.csv

13. Sum all errors for each marker in column C

A	B	C	D	E	F	G	H	I	J	K	
marker	position		SUM({[#	{	[#	}	/	SUM(}/)	~	11
		~								0	
4071	100	=D10+J10+K10		0	0	1	0	0	0	0	0
8853	100.0362	1	1	0	0	1	0	0	0	0	0

14. Add filter and filter for markers with at least 1 error
 *optional: sort full data according to number of errors from high to low

15. Copy all data and transpose into new sheet
 *original mconsistency format but only markers with errors are kept



Mconsistency.csv

16. Freeze panes so that summary for each marker and individual is locked

17. Add filter to last summary row

A	B	C	D	E	F	G	H	I	J	K	
marker									4071	8853	
position									100	100.0362	100
	SUM({[#] {	[# }	/	SUM(}/)					1	1	
SUM({[#]	1180								1	1	
{		351							0	0	
[0						0	0	
#				829					1	1	
}					484				0	0	
/						15			0	0	
SUM(}/)							199		0	0	
1	~										
2	11_118	0	0	0	0	0	0	0	0	0	
3	13_20	0	0	0	0	0	0	0	0	0	
4	5_62	0	0	0	0	0	0	0	0	0	



Mconsistency.csv

18. For each marker, filter for characters that are associated with an error to find which individuals are causing an error in column A

	A	B	C	D	E	F	G	H	I	AEJ
1	marker									2977
2	position									600.4256
3		SUM({{#} {	[# }	/	SUM({/)	~				11
4	SUM({{#}	1180								3
5	{		351							3
6	[0						0
7	#				829					0
8	}					484				7
9	/						15			1
10	SUM({/)							499		8
11	%									
12										
13	BlackRepu	1	1	0	0	0	0	0	0	{
14	Lapins	3	3	0	0	7	0	7	0	}
15	Rainier	17	17	0	0	15	0	15	0	}
16	Regina	0	0	0	0	12	0	12	0	}
17	Selah	4	4	0	0	8	0	8	0	{
18	PMR-1	0	0	0	0	6	1	7	0	/
19	Cowiche	32	32	0	0	56	0	56	0	{
20	FR010T000	0	0	0	0	1	0	1	0	}

Error-associated characters

cc	Symbol	Shortname	Description
2	}	OEA	observed & deleted error & augmented
6	/	OpE	observed & putative error
7	{	OE	observed & deleted error
8	[pnuE	putative null-allele & error
9	#	ME	missing & error

Filtered for error-associated characters

Individuals who caused an error for this marker

- *Could also be their parents (or ancestors if parents have missing data)
- *Could also be their offspring (very often the case)

Resolving errors

19. Check genotypic data to find what is causing the issue

*Note: sometimes imputed missing data is causing the issue

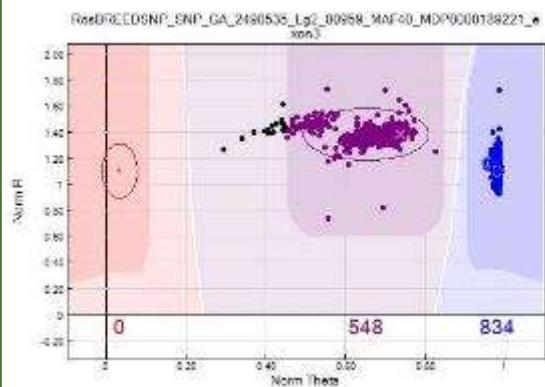
e.g. parental genotypic data is missing but grandparents are both “AA” -> parent is imputed as “AA” however individual is “BB” -> error reported

*Note: In case of multiple parents/ancestors with missing data that prevented earlier pedigree checking, pedigree could still be incorrect

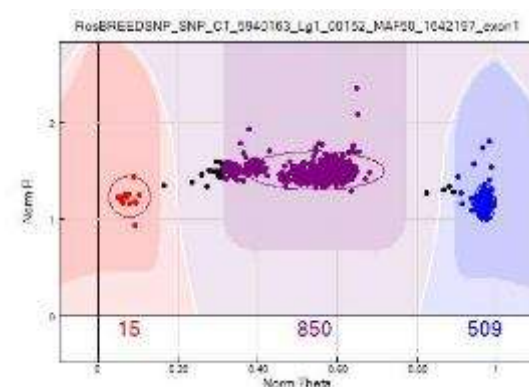
20. Check GenomeStudio® to:

- Ensure genotypic data matches
- Genotype clustering is correct (e.g. additional clusters or null alleles observed)
- Single individuals are assigned the correct genotype

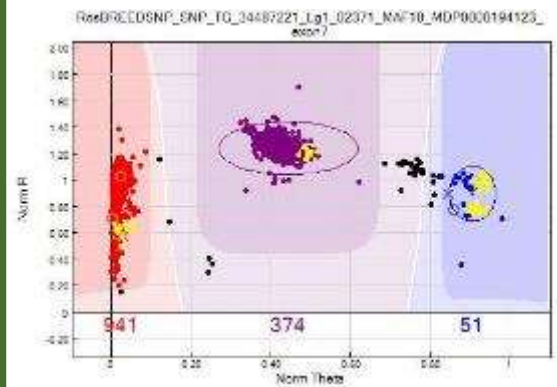
Incorrect clustering



Incorrect clustering due to additional clusters



Incorrect clustering due to null alleles



Resolving errors

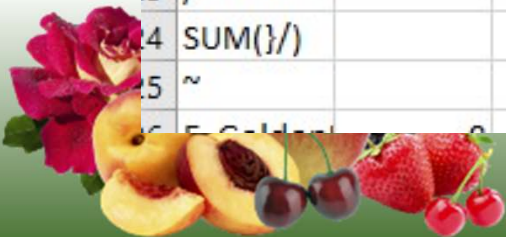
21. Update data file (pedigree or genotypic data)

*Update map file if markers were removed from data set

22. Re-run FlexQTL™ and check mconsistency file to ensure errors are resolved

22. Repeat until no errors are reported

6	position									100
7		SUM({[#	{	[#	}	/	SUM(}/)	~	
8	SUM({[#	0								0
9	{		0							0
10	[0						0
11	#				0					0
12	}					0				0
13	/						0			0
14	SUM(}/)							0		0
15	~								0	0





United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Finding Mendelian- consistent errors

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY

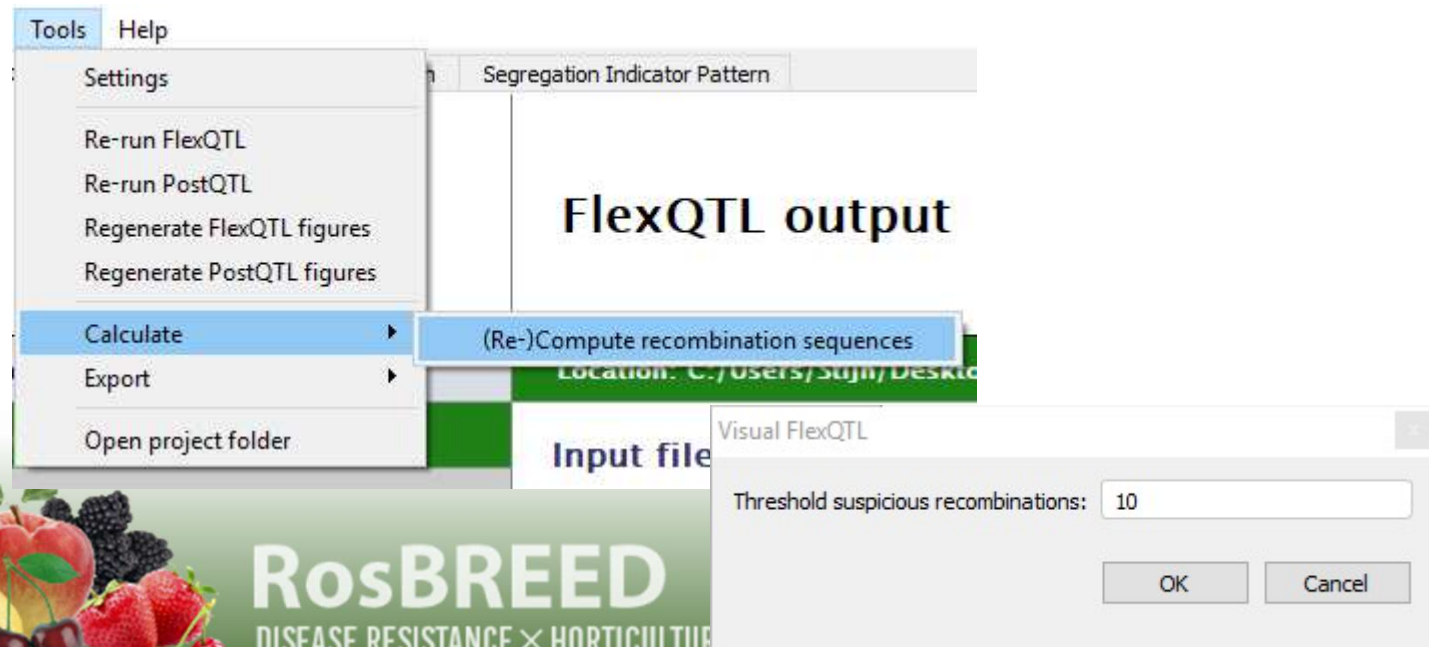


SUPERIOR
CULTIVARS



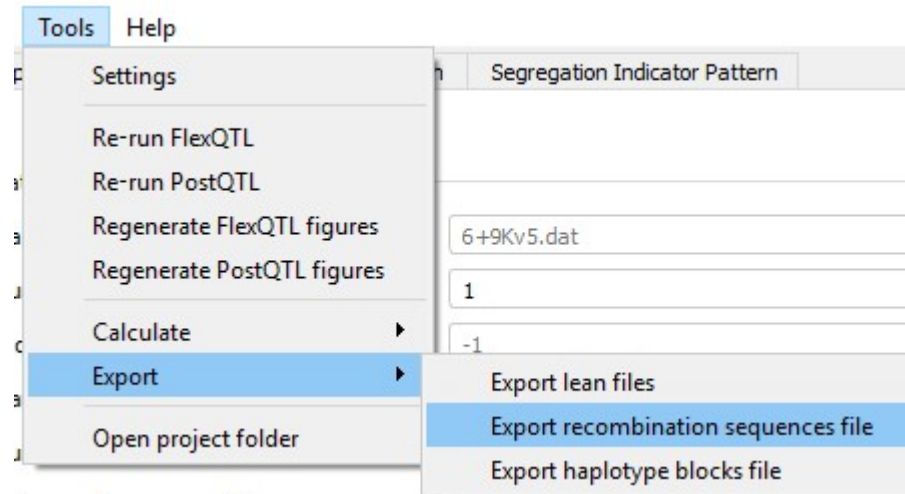
Finding double recombinations

1. Run FlexQTL™ with latest map and data file that did not lead to any reported errors
 - *Parameter settings remain the same
2. Set interval for which double recombinations (DR) (Tools>Calculate>(Re-)Compute recombination sequences)
 - *Default is all DR within 10 cM

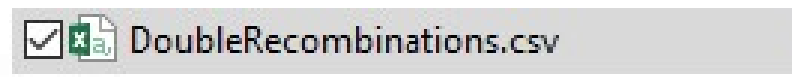


Finding DR

3. Export all double recombinations within defined interval (Tools>Export>Export recombination sequences file; then select directory to save file in)



4. Open generated file called “DoubleRecombinations.csv”



5. Save as “.xls(x)” format



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Finding DR

Information found in DoubleRecombinations file

Individual in which DR is observed

How suspicious/unlikely is this DR

Chromosome on which DR is observed

Length of DR interval

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Individual	Chromosome	Homolog	M1	M2	M3	M4	P1	P2	P3	P4	R	Suspicion
2	Abundance	1	0	8112	8202	8202	453	47.82	53.41	53.41	56.17	8.35	0.073106
3	Abundance	1	0	8202	453	453	8255	53.41	56.17	56.17	57.31	3.9	1
4	Abundance	1	0	509	8343	8359	8367	76.07	78.47	80.13	80.55	4.48	0.999835

Homolog on which DR is observed

0: Homolog coming from mother

1: Homolog coming from father

-M/P1: Last marker name(M)/position(P)
before 1st recombination

-M/P2: First marker name(M)/position(P)
after 1st recombination

-M/P3: Last marker name(M)/position(P)
before 2nd recombination

-M/P4: First marker name(M)/position(P)
after 2nd recombination

Finding DR

- Sort file according to chromosome and position of DR
- Look first for regions with many DR
- Use 'Segregation Indicator Pattern' tab of Visual FlexQTL™ to visualize DR
Use drop-down menu on top-right to choose "Recombination sequences (based on FlexQTL SIP)"

Visual FlexQTL

File Tools Help

Output Settings Markers Population Segregation Indicator Pattern

Threshold suspicious recombinations: 10

Marker scores

- Marker scores
- SIP (computed by FlexQTL)
- Recombination sequences (based on FlexQTL SIP)
- Parent indicator (computed by Visual FlexQTL)
- SIP (computed by Visual FlexQTL)

#	Individual	Parent	Offspring	1897	5444	4177	1964
				0	0	0	0



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Investigating DR

Information found under “Segregation Indicator Pattern” tab

Indicator of grandparental origin:

-0: Allele originates from mother of parent (grandmother; yellow)

-1: Allele originates from father of parent (grandfather; blue)

Recombination region (0-to-1 or 1-to-0 switch; orange)

Parent’s name

Individual’s name

DR region (0-to1-to-0 or 1-to-0-to-1 switches; red)

2 lines for each individual, one for each parent’s homolog



Individual	Parent	Offspring	7093	5084	6503	5093	0100	0222	2213	5216	5307	8376	494	6723	2276	5407	5413	7239
0415-0023	Enterprise	0	s1	-	-	-	-	s1	s1	-	s1	-	-	s1	1	1	1	s1
0415-0024	Honeycrisp	0	0	-	-	-	-	0	0	-	0	-	-	0	-	-	-	s0
0415-0024	Enterprise	0	0	-	-	-	-	0	0	-	0	-	-	0	1	1	1	s1
0415-0025	Honeycrisp	0	s1	-	-	-	-	s1	s1	-	s1	-	-	s1	-	-	-	s1
0415-0025	Enterprise	0	s0	-	-	-	-	s0	s0	-	s0	-	-	s0	0	0	0	s0
0416-0008	Honeycrisp	0	0	-	-	-	-	0	0	-	0	-	-	0	-	-	-	0
0416-0008	Fuji	0	0	-	-	0	1	-	0	-	0	-	-	0	0	0	0	0
0416-0009	Honeycrisp	0	s1	-	-	-	-	1	s1	-	s1	-	-	s1	-	-	-	s1
0416-0009	Fuji	0	s0	-	-	0	1	-	s0	-	s0	-	-	s0	0	0	0	s0

Resolving DR

8. Check genotype calls with GenomeStudio®

-e.g. incorrect clustering of “AA” as “AB” ->

“AB” (actually “AA”) x “AB” -> 100% “AB” (actually 50% “AA” and 50% “AB”)

*does not lead to PPC errors

*often characterized by many DR in a single or few families

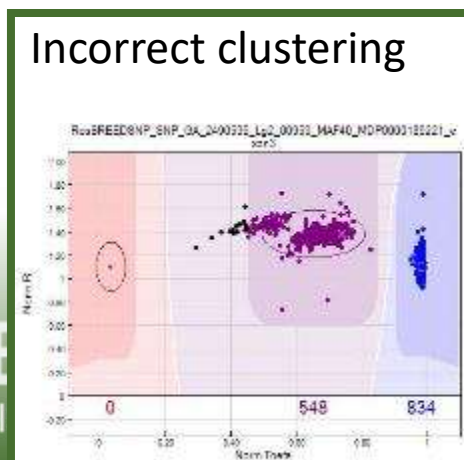
*often characterized by a single marker involved in DR

-single genotyping error in individual

*does not lead to PPC errors

*often characterized by many DR in offspring of this individual

*often characterized by a single marker involved in DR



Rose
DISEASE RESI

Resolving DR

9.a Check for errors in map order using the “Segregation Indicator Pattern” tab

Moving marker(s) a few positions resolves DR

- *Characterized by DR in a few individuals who are not necessarily related
- *Characterized by DR spanning one or multiple markers
- *Often characterized by a third recombination nearby
moving markers to this recombination would solve DR

4177	1964	201	4581	7850	737	6127	2565	4112	4209	7461	3225	3252	545	5058
0	0	0	0	0	0	0	0	0	0	0	0.877	1.912	1.912	1.912
0	0	-	-	s0	-	-	-	0	-	-	-	-	-	-
1	-	1	1	s1	1	-	1	-	-	-	1	1	1	-
1	1	-	-	s1	-	-	-	1	-	-	-	-	-	-
1	-	0	0	s1	1	1	1	-	-	-	0	0	0	-
1	1	-	-	s1	-	-	-	-	-	-	-	-	-	-

DR

3rd recombination

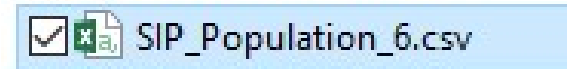
Moving markers “201” and “4581” behind marker “2565” would resolve these DR

*This cannot create new DR for other individuals (next slide)

Resolving DR

9.b Adjust marker position and check if no new DR are created

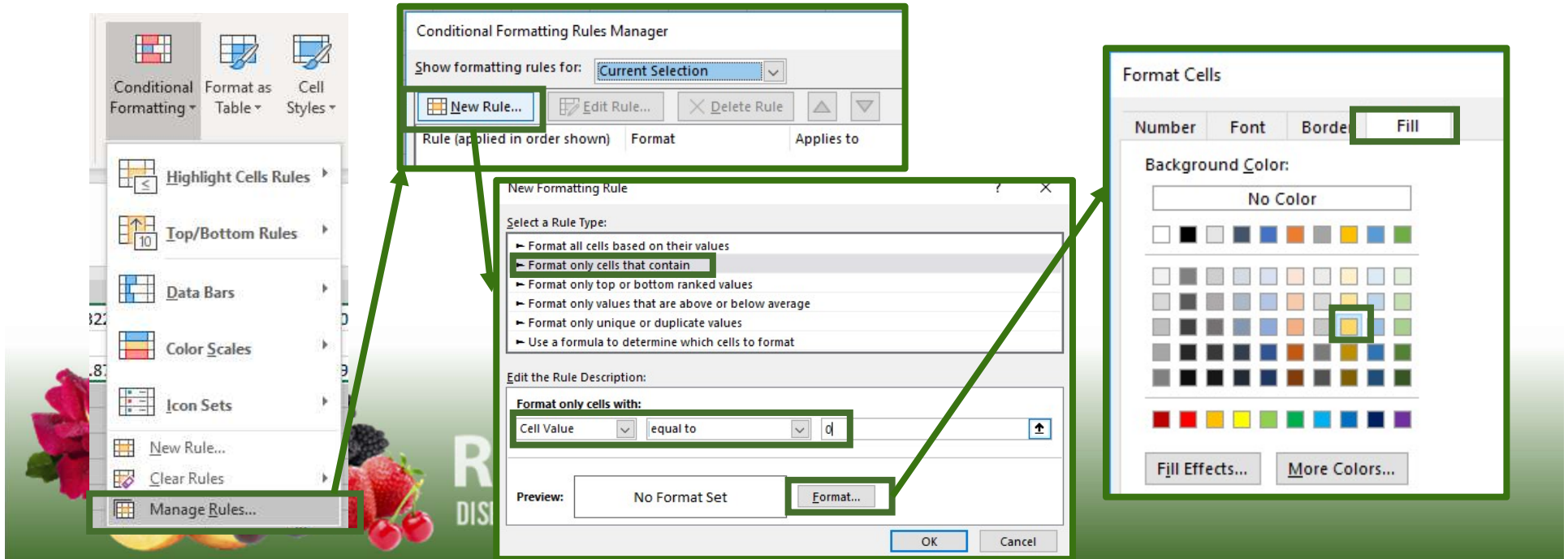
9.b.1 Open “SIP_Population_X.csv” in FlexQTL project folder
*“X” is highest number found in folder



*Interpretation the same as “Segregation Indicator Pattern” tab of Visual FlexQTL™

9.b.2 Use conditional formatting to color “0”/“s0” yellow and “1”/“s1” blue

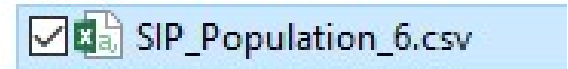
*same coloring as “Segregation Indicator Pattern” tab of Visual FlexQTL™

A composite screenshot illustrating the steps to create a conditional formatting rule in Excel. The 'Conditional Formatting' ribbon is shown on the left. The 'Conditional Formatting Rules Manager' dialog box is open, with 'New Rule...' highlighted. The 'New Formatting Rule' dialog box is also open, showing 'Format only cells that contain' selected under 'Select a Rule Type'. The 'Edit the Rule Description' section shows 'Cell Value' selected, 'equal to' in the operator dropdown, and '0' in the text box. The 'Format...' button is highlighted. To the right, the 'Format Cells' dialog box is open, with the 'Fill' tab selected and a yellow color chosen from the 'Background Color' palette. Green arrows indicate the flow from the ribbon to the Rules Manager, then to the New Rule dialog, then to the Format Cells dialog, and finally to the 'Format...' button in the New Rule dialog.

Resolving DR

9.b Adjust marker position and check if no new DR are created

9.b.1 Open “SIP_Population_X.csv” in FlexQTL project folder
*“X” is highest number found in folder



*Interpretation the same as “Segregation Indicator Pattern” tab of Visual FlexQTL™

9.b.2 Use conditional formatting to color “0”/“s0” yellow and “1”/“s1” blue

*same coloring as “Segregation Indicator Pattern” tab of Visual FlexQTL™

The image illustrates the process of creating a conditional formatting rule in Excel. It shows the 'Conditional Formatting' ribbon, the 'Conditional Formatting Rules Manager' dialog box, the 'New Formatting Rule' dialog box, and the 'Format Cells' dialog box. Green boxes and arrows highlight the 'New Rule...' button, the 'Format only cells that contain' rule type, the 'Cell Value equal to d' configuration, and the 'Format...' button in the preview section.

Resolving DR

9.b.3 Cut columns of markers to move and insert them into new position

	D	H	I	J	K	L	M	N	O	P	
marker		1964		201	4581	7850	737	6127	2565	4112	4209
group		1		1	1	1	1	1	1	1	1
position		0		0	0	0	0	0	0	0	0
48		1			1				1		
48				0	0	1		1			
48		1			1				1		
48				0	0	0		0			
48		1			1				1		
48				1	1	1		1			

	K	L	M	N	O	P	
er		6127	2565	201	4581	4112	4209
group		1	1	1	1	1	1
position		0	0	0	0	0	0
48					1		
48			1	0	0		
48						1	
48			0	0	0		
48						1	
48			1	1	1		

9.b.4 Check rest of the file (rows) to make sure this move did not create additional DR

9.b.5 Update “.dat” and “.map” file to reflect the change in marker positions

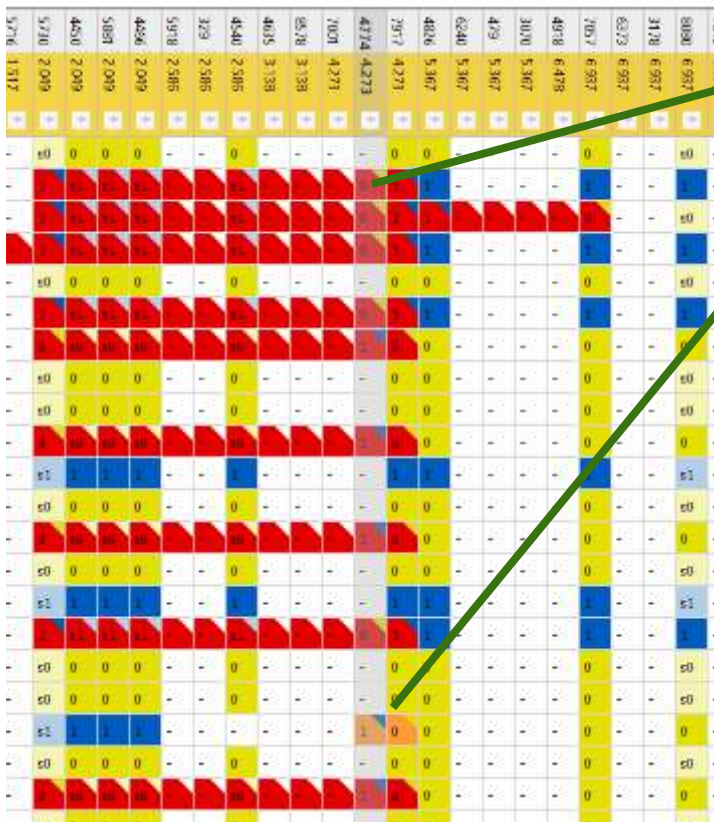


RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Resolving DR

10. Check for phasing issues

- * Often characterized by DR in almost all offspring of a single individual
- * Individuals is often a founder (parents not genotyped)
- * One or more offspring have a single recombination for same marker(s)
- * Incorrect phasing of this individual causes errors for sibs



Marker always involved in DR

One individual with single recombination

* Make genotypic data of this individual missing (gets imputed again in later stages)

Probable cause: minimizing recombination interval

Incorrect phasing – short recombination interval



Correct phasing – long recombination interval



Resolving DR

11. Re-run FlexQTL™ and re-generate “DoubleRecombinations.csv” file with Visual FlexQTL™ after resolving major DR

*Resolving DR may resolve other nearby DR

12. Repeat these steps until remaining DR are very likely true

*Unlike Mendelian-inconsistent errors, some DR can remain



RosBREED

DISEASE RESISTANCE × HORTICULTURAL QUALITY



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Haploblocking and haplotyping

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS





United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Generating haploblocks

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



Adjusting input files

- Accurate phasing requires:
 - Removal of intermediate ungenotyped progenitors
 - Well-represented ancestors (>3 offspring) whose genotype can be imputed can remain
- Haploblock determination based on historical recombination requires:
 - All individuals that need to be taken into account to have offspring in data set



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Adjusting input files

1. Remove any parentage information of ungenotyped intermediate progenitors from FlexQTL™ data input file

- *Ungenotyped founders (no known parents) can stay in data set

- *Ungenotyped progenitors with more than 3 genotyped offspring can stay in data set

2. Add dummy offspring to FlexQTL™ data input file for any individual whose recombination are important for haploblock determination (e.g. breeding selections) that don't have any offspring in the data set

- *Set genotypic data of dummy individuals to missing

- *Set parents of dummy individuals to individuals who need offspring in data set



RosBREED

DISEASE RESISTANCE × HORTICULTURAL QUALITY

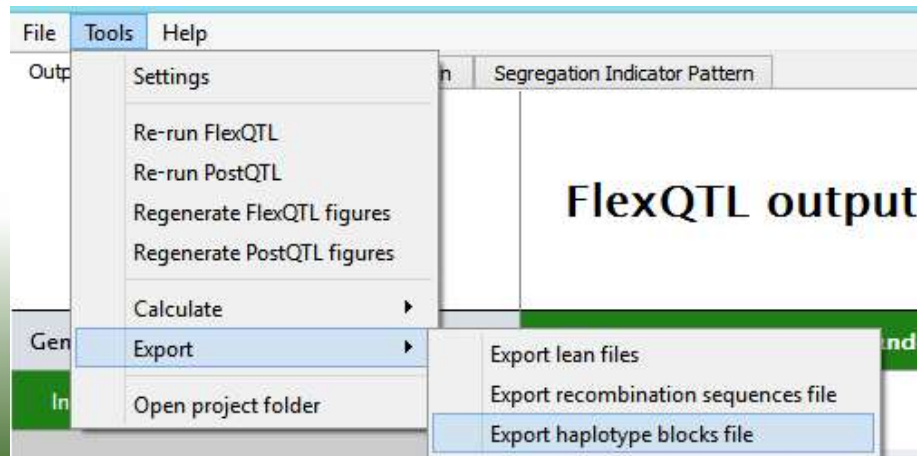
Creating haploblocks

3. Create a new FlexQTL™ project and use following settings

- skipSampleMarkers: 0
- REDprint: 0
- Markerblock: 5
- MSegDelta: 1
- DeleteDR: 1

4. Run FlexQTL™

5. Generate haploblock delimitation with Visual FlexQTL™ (Tools>Export>Export haplotype blocks file then select directory to save file in)



Creating haploblocks

6. Open generated HaploBlock.map file



7. Adjust Haploblock designations as wanted
e.g. to ensure maximum size of haploblock is 1 cM

Marker	Chromosome	Position	Haploblock
1897	1	0	HB-01-1
5444	1	0	HB-01-1
4177	1	0	HB-01-2
1964	1	0	HB-01-2
201	1	0	HB-01-2
4581	1	0	HB-01-2

Assigned haploblock
*name can be
changed freely



ROSBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Determining haplotypes

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



Creating input files

1. Remove dummy offspring from FlexQTL™ data input file used for haploblock determination
2. Create a new FlexQTL™ project and use following settings
 - skipSampleMarkers: 0
 - REDprint: 0
 - Markerblock: 5
 - MSegDelta: 1
 - DeleteDR: 1
3. Run FlexQTL™



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Creating input files

4. Copy needed input files in a single folder

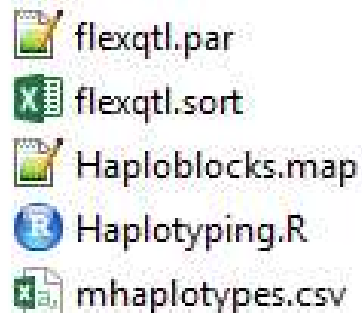
*From latest FlexQTL™ run:

- mhaplotypes.csv
- flexqtl.par
- flexqtl.sort

*From FlexQTL™ run to determine haploblocks:

- Generated (and adjusted) haploblocks.map

5. Create a new R project in the same folder to run PediHaplotyper



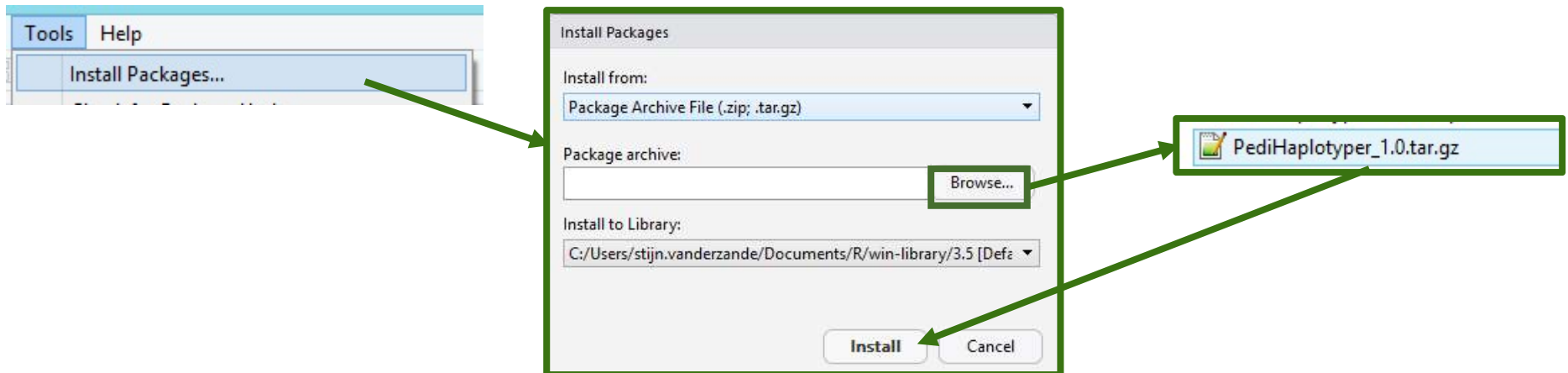
RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Running PediHaplotyper

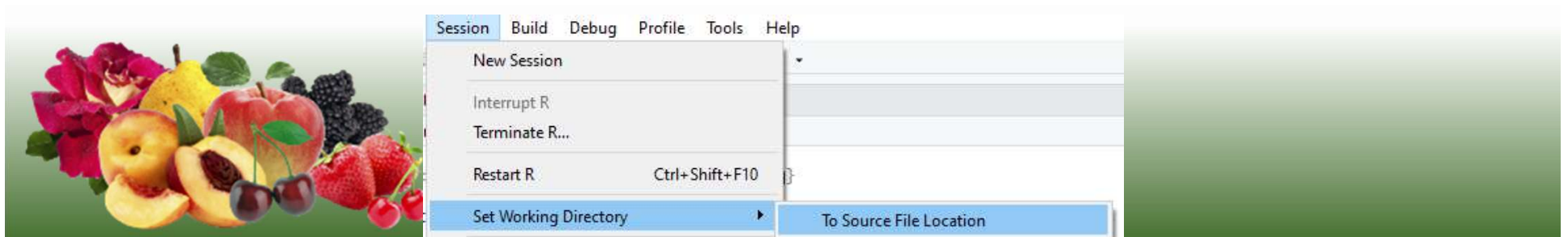
6. In RStudio, install PediHaplotyper as a package (only needed once)

*Tools>Install Packages...

*Install from: "Package Archive" and Browse to find "PediHaplotyper_1.tar.gz"




7. Set working directory to "Source File Location"



Running PediHaplotyper

8. Run the following code

```
1 library(PediHaplotyper)
2
3 fq_haplotyping_session(sessionID="20180613_HTDet_All" mapfile="Haploblocks.map")
4
```



Prefix for all output files

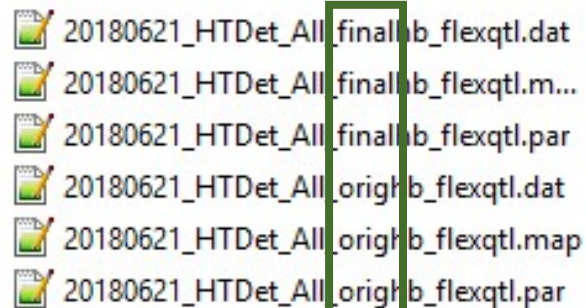
Name of haploblock definition file

```
> fq_haplotyping_session(sessionID="20180613_HTDet_All", mapfile="Haploblocks.map")
reading datafiles ...
hb 1 = HB-01-1: 27 initial alleles
hb 2 = HB-01-1a: 3 initial alleles
hb 3 = HB-01-2: 25 initial alleles
-----
writing initial HS family haploblock alleles file ...
calculation haploblock alleles ...
hb 1 = HB-01-1: convergence = yes in 4 cycles
hb 2 = HB-01-1a: convergence = yes in 2 cycles
hb 3 = HB-01-2: convergence = yes in 3 cycles
hb 4 = HB-01-3: convergence = yes in 3 cycles
-----
hb 503 = HB-17-52: convergence = yes in 2 cycles
writing all requested output files ...
haplotyping session 20180613_HTDet_All finished.
>
```

Output files

“orig”-flag vs “final”-flag

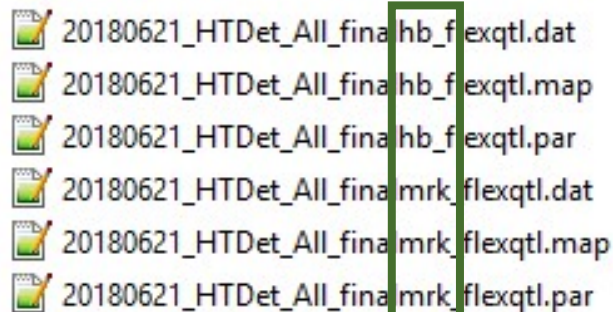
- “orig”: results after initial haplotype determination
- “final”: results add end of haplotype determination (use these files)



20180621_HTDet_All_finalhb_flexqtl.dat
20180621_HTDet_All_finalhb_flexqtl.m...
20180621_HTDet_All_finalhb_flexqtl.par
20180621_HTDet_All_orighb_flexqtl.dat
20180621_HTDet_All_orighb_flexqtl.map
20180621_HTDet_All_orighb_flexqtl.par

“mrk”-flag vs “hb”-flag

- “mrk”: results at the single marker level
- “hb”: results at haplotype level (use these files)



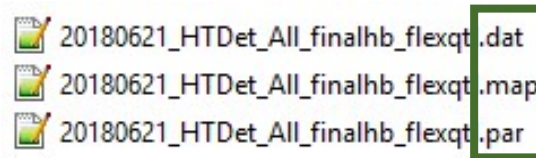
20180621_HTDet_All_finalhb_flexqtl.dat
20180621_HTDet_All_finalhb_flexqtl.map
20180621_HTDet_All_finalhb_flexqtl.par
20180621_HTDet_All_finalmrk_flexqtl.dat
20180621_HTDet_All_finalmrk_flexqtl.map
20180621_HTDet_All_finalmrk_flexqtl.par



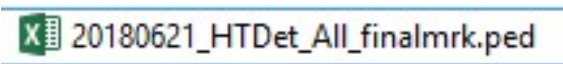
Ro
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Output files

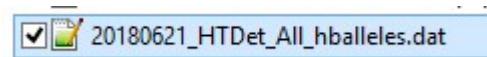
-All three FlexQTL™ input files (".dat", ".map", and ".par") created for each combination



-Pedimap input file (".ped") created for each combination



- "[...] _hballeles.dat" gives definition of each haplotype for each haploblock



Haploblock number and name

Marker names within haploblock

Number of markers in haploblock

Marker alleles that create haplotype

hbnr	haploblock	markercount	hballeleID	hballelenr	nacount	freq	4177	1897	1964	7850
1	HB-01-1	4	1	1	4	214	-	-	-	-
1	HB-01-1	4	2	2	1	1	-	A	A	A
1	HB-01-1	4	3	3	1	6	-	B	B	A
1	HB-01-1	4	4	4	0	794	B	A	A	A

Haplotype name and number

Missing values within haplotype

How often haplotype occurs in data set



United States
Department of
Agriculture

National Institute
of Food and
Agriculture

Checking for errors in haplotypes

RosBREED

DISEASE
RESISTANCE



HORTICULTURAL
QUALITY



SUPERIOR
CULTIVARS



Mendelian-inconsistent errors

1. Use “[...]_finalhb” FlexQTL™ input files generated by PediHaplotyper to create a new FlexQTL™ project
 - *parameter file should be renamed to “flexqtl.par”
2. Use “mconsistency.csv” file to identify issues
3. Use “[...]_hballeles.dat” file to investigate haplotypes that cause P(P)C errors
 - 3.a If missing data within haplotype causes inconsistency: assign correct haplotype (if both parental haplotypes are possible, use information on flanking markers and minimize recombinations)

GoldenDel										
mchr			9	9	9	9	9	9		
mpos			4.538	4.821	6.61	7.602	7.602	9.2		
ID	VAR		HB-09-7	HB-09-8	HB-09-9	HB-09-10a	HB-09-10	HB-09-11		
GoldenDel	Grimes_Gc hap1	1	3	3	2	2	4			
GoldenDel	F_GoldenL hap2	1	4	4	2	3	5			
MN1627	GoldenDel hap2	1	-	-	2	-	-			
Splendour	GoldenDel hap1	1	3	3	2	2	4			
Pinova	GoldenDel hap2	1	4	4	2	3	6			
PRI14-152	GoldenDel hap1	1	4	4	2	3	5			

Haplotype 6 equals haplotype 5 except for a missing value

	HB-09-10	439	6324	3144	6325	3146
4	A	B	A	B	A	
5	B	B	A	A	B	
6	B	B	A	A	-	



RosBR
DISEASE RESISTANCE

Mendelian-inconsistent errors

3.b If a recombination occurs within haploblock:

3.b.1 if recombination occurs within selected material: adjust haploblock borders in “haploblock.map” file and re-determine haplotypes for adjusted haploblocks

3.b.2 if recombination occurs in seedling: don't adjust haploblock borders and do not adjust haplotype (or make it missing)

	A	B	C	D	T	Z	MM	MB	ML	MD
1	Splendour									
2	mchr				1	1	1	1	1	
3	mpos				22.387	23.021	23.554	24.678	24.89	2
4	ID	VAR			HB-01-19	HB-01-19a	HB-01-20	HB-01-21	HB-01-21a	HB-01-22
65	Splendour	GoldenDel	hap1	1	4	4	4	4	3	
66	Splendour	Delicious	hap2	1	10	7	5	5	2	
3286	Sciros	Splendour	hap2	1	4	4	4	6	2	

HB-01-21	1970	4277	131	2566	6846
4	A	B	A	A	A
5	A	B	B	A	B
6	A	B	A	A	B

Recombination within HB

Mendelian-inconsistent errors

- 3.c If a genotype calling error occurred (check with GenomeStudio[®]): adjust haplotype to resemble correct genotype call
- 3.d If a marker order is incorrect: adjust marker order in “haploblock.map” file, and latest data and map file for single marker level. Re-run FlexQTL[™] to create new PediHaplotyper files and rerun pedihaplotyper for affected haploblocks



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY

Mendelian-consistent errors

4. Once Mendelian-inconsistent errors are resolved, check for Mendelian-consistent errors as was done for single marker level
5. Resolve Mendelian-consistent errors as was done for Mendelian-inconsistent errors with haplotypes
6. Enjoy your high-quality genotypic data set and make lots of great discoveries!!



RosBREED
DISEASE RESISTANCE × HORTICULTURAL QUALITY