User Manual for



### -Trait Analysis by aSSociation, Evolution and Linkage

# Version 3

# The Buckler Lab at Cornell University

(December 22, 2011)



www.maizegenetics.net/tassel

**Disclaimer**: While the Buckler Lab at Cornell University has performed extensive testing and results are, in general, reliable, correct or appropriate results are not guaranteed for any specific set of data. It is strongly recommended that users validate TASSEL results with other software.

**Further help**: Additional help is available beyond this document. Users are welcome to report bugs, request new features through the TASSEL website. Questions are also welcome to our current team members. For more quick and precise answers, please address your questions to the most pertinent person:

Tassel User Group (recommended)	http://groups.google.com/group/tassel tassel@googlegroups.com
<b>General Information</b>	Ed Buckler (Project leader) esb33@cornell.edu
Data import, GDPC, Pipeline	Terry Casstevens tmc46@cornell.edu
Statistical analysis	Peter Bradbury pjb39@cornell.edu Zhiwu Zhang zz19@cornell.edu

Contributors: Yogesh Ramdoss, Michael E. Oak, and Karin J. Holmberg, N. Stevens, and Yang Zhang.

The TASSEL project is supported by the National Science Foundation and the USDA-ARS.



Main Web Site: <u>http://www.maizegenetics.net/tassel</u> Open source code: <u>http://sourceforge.net/projects/tassel</u> Modified version of the PAL library is used: <u>http://www.cebl.auckland.ac.nz/pal-project</u> Database access is achieved by GDPC middleware <u>http://www.maizegenetics.net/gdpc</u>

#### **Table of Contents**

INT	RODUCTION	6
<u>1</u>	GETTING STARTED	7
1.1	INSTALLATION	7
1.1.1	1 WEB START	7
1.1.2	2 STAND-ALONE	8
1.1.3	3 OPEN SOURCE CODE	8
1.2	PANELS	8
<u>2</u>		10
0.1		10
2.1	GDPC	10
2.2		11
2.2.1	1 BLOB	12
2.2.2	2 Нармар	12
2.2.3	3 Plink	12
2.2.4	4 Flapjack	13
2.2.5	5 POLYMORPHISM	13
2.2.6	6 Phylip	14
2.2.7	7 NUMERICAL DATA	14
2.2.8	8 SQUARE NUMERICAL MATRIX	16
2.2.9	9 GENETIC MAP	16
2.3	EXPORT Export	16
2.4	SITES SITES	17
2.5	SITE NAMES SITE NAMES	18
2.6		19
2.7	INDUTE SNDS RAI IMPUTE SNPS	19
2.8 2.9	TRANSFORM ?+5 Transform	20
2.9.1	1 GENOTYPE NUMERICALIZATION ?+5 Transform	20 20
2.9.2	2 TRANSFORM AND/OR STANDARDIZE DATA	21
2.9.3	3 IMPUTE PHENOTYPE	22
2.9.4	4 PCA	23
2.10	SYNONYMIZE TAXA NAMES $p \leftrightarrow q$ Synonymizer	23
2.11		25
2.12	2 INTERSECTION JOIN O Join	26
3		27

3.1		27
3.2	LINKAGE DISEQUILIBRIUM 🚺 Link. Diseq.	28
3.3		29
3.4		29
3.5	KINSHIP Kinship	30
36	GENERAL LINEAR MODEL	30
0.0		
3.7	MIXED LINEAR MODEL	32
3.8	RIDGE REGRESSION	34
<u>4</u> ]	RESULT MODE	36
4.1		36
4.2		36
4.3	2D PLOT	37
4.4		38
4.5		39
<u>5</u> ]	MENUS	41
5.1	FILE MENU	41
5.1.1	1 SAVE DATA TREE	41
5.1.2	2 OPEN DATA TREE 2 Same Data Tree Ag	41
5.1.3	J SAVE DATA TREE AS 4 Oden Data Tree	41
514	5 SAVE SELECTED AS	41
5.2	Contingency Test	42
5.3	PREFERENCES	42
<u>6</u>	TUTORIAL	43
6.1	Missing Phenotype Imputation	43
6.2	PRINCIPAL COMPONENT ANALYSIS	45
6.3	ESTIMATION OF KINSHIP USING GENETIC MARKERS	49
6.4	ASSOCIATION ANALYSIS USING GLM	50
6.5	ASSOCIATION ANALYSIS USING MLM	54
<b>0.0</b>	IMPORTING DATA FROM A DATABASE (VIA GDPC)	57
0.0.1	I CONNECTING WITH A DATABASE 2 DATA OHEDV	5/
6.6.2	2 DATA QUENT 3 IMPORTING GDPC DATA INTO TASSEI	58 61
6.6.4	4 SAVING GDPC QUERY RESULTS	63
7	APPENDIX	64

#### <u>7</u> <u>APPENDIX</u>

7.1	NUCLEOTIDE CODES (DERIVED FROM IUPAC)	64
7.2	TASSEL TUTORIAL DATA SETS	65
7.3	BIOGRAPHY OF TASSEL	66
7.4	FREQUENTLY ASKED QUESTIONS	68
1.	WHAT DO I DO IF TASSEL MISBEHAVES?	68
2.	WHERE DO I TURN FOR MORE INFORMATION?	68
3.	HOW DO I JOIN THE FUN: TASSEL ON SOURCEFORGE?	68
4.	HOW DO I CHANGE THE AMOUNT OF MEMORY USED? WHAT DO I DO WHEN THE "EXCEPTION	
	JAVA.LANG.OUTOFMEMORYERROR" APPEARS?	68
5.	WHEN I CLICK ON THE MOST CURRENT VERSION OF TASSEL WEB START, A PREVIOUS VERSION APPEARS.	
	WHAT SHOULD I DO?	69
6.	WHAT SHOULD I SUBSTITUTE FOR MISSING VALUES IN TASSEL?	69
7.	IS IT POSSIBLE TO CHANGE DATA NAMES IN THE DATA TREE?	69
8.	HOW CAN I CREATE A TASSEL ICON ON DESKTOP?	69
9.	WHY DO I GET EMPTY SQUARES IN MLM ASSOCIATION ANALYSIS?	69
10.	WHY SHOULD I EXCLUDE ONE COLUMN OF THE POPULATION STRUCTURE?	69
11.	CAN KINSHIP REPLACE POPULATION STRUCTURE?	69
12.	WHY DO TASSEL AND SPAGEDI GIVE DIFFERENT KINSHIP ESTIMATES?	70
13.	CAN I GET MARKER R SQUARE USING SAS PROC MIXED OR TASSEL MLM?	70
14.	DOES MLM FIND MORE ASSOCIATIONS THAN GLM?	70
15.	DO I NEED MULTIPLE TEST CORRECTION FOR THE P VALUE FROM TASSEL?	70
16.	CAN TASSEL HANDLE DIPLOID GENOTYPE DATA?	70
17.	HOW TO CITE TASSEL?	70
RE	FERENCES	71
INI	DEX	73

### INTRODUCTION

While TASSEL has changed considerably since its initial public release in 2001, its primary function continues to be providing tools to investigate the relationship between phenotypes and genotypes<sup>1</sup>. As indicated by its title – Trait Analysis by aSSociation, Evolution and Linkage – TASSEL has multiple functions, including association study, evaluating evolutionary relationships, analysis of linkage disequilibrium, principal component analysis, cluster analysis, missing data imputation and data visualization.

One of the design elements driving TASSEL development has been the need to analyze ever larger sets of data<sup>2</sup>. For example, the MLM (mixed linear model) function for association analysis originally used an EM (expectation-maximization) algorithm, which is a common method for solving mixed models but is relatively slow. Subsequently developers implemented the EMMA algorithm to increase computing speed<sup>3</sup>. Model compression was added to that to improve speed and statistical power for association study<sup>4</sup>. Another technique that optimizes variance components once and then uses the estimates to test markers now provides the ability to screen the large numbers of markers used in genome-wide association studies (GWAS). The method was independently described by Zhang et al. and Kang et al. in 2010. This method was named P3D by Zhang et al.<sup>4</sup> and EMMAX by Kang et al.<sup>5</sup>

TASSEL was designed for a wide range of users, including those not expert in statistics or computer science. A GWAS using the mixed linear model method to incorporate information about population structure<sup>6-8</sup> and cryptic relationships<sup>9</sup> can be performed by in a few steps by "clicking" on the proper choices using a graphic interface. All the processes necessary for the analysis are performed automatically, including importing phenotypic and genotype data, imputing missing data (phenotype or genotype), filtering markers on minor allele frequency, generating principal components and a kinship matrix to represent population structure and cryptic relationships, optimizing compression level and performing GWAS.

The command-line version of TASSEL, called the Pipeline, provides users the ability to program tasks using a script instead of the graphic user interface (GUI). This feature allows researchers to define tasks using a few lines of code and provides the ability to use TASSEL as part of an analysis pipeline or to perform simulation studies.

Due to the increasing availability of open data sources, TASSEL utilizes a data browser from the Genomic Diversity and Phenotype Connection (GDPC) project<sup>10</sup> to provide an interface to relational databases. As a result, TASSEL users can access any data source that provides a GDPC service. Using this middleware, which provides a common graphical interface, TASSEL users can avoid writing SQL queries to access data. Currently, GDPC provides connections to Panzea, Gramene, Germinate, and GRIN (USDA's Germplasm Resources Information Network).

TASSEL is written in Java, thereby enabling its use with virtually any operating system. It can be installed using Java Web Start technology by simply clicking on a link at <u>www.maizegenetics.net/tassel</u>. A stand-alone version of TASSEL can also be downloaded to use in pipeline mode or in any situation where the user wishes to start the software from a command line.

### 1 Getting Started

A quick way to get started using TASSEL is to load the tutorial data and try performing analyses. However, because some of the necessary steps may not be intuitive, we recommend that new users follow the tutorial at end of this manual. The objective of this section is to provide information necessary to install and start TASSEL software and to provide a brief overview of the interface.

Most functions are organized into three modes (Data, Analysis and Results) which correspond to the first three buttons on the TASSEL interface as shown below. Clicking one of these buttons changes the functions represented by the second row of buttons. Those three modes are described in detail in the subsequent sections of this manual. The screen shot shows TASSEL after the tutorial files have been loaded.

💰 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3	3.0.	37																	G		X
File Toole Help GDPC																					
Concerco Data Analysis Results Delete Wit	lizaro	i Mi Show Mem	ory	٦							<u> </u>		(	0%		_	ר 🛙			<b>1</b>	
											-						96		-	ركار	
🐨 CODC 🖃 Load 🏦 Export 🖓 Sites 🕨 Taxa 🍸 Traits	s	👫 Impute SNPs 🖓 ?+5 Tr	ans	forn	n (	p⇔¢	q sy	non	ymize	er (		υJo	oin	$\bigcirc$	n Jo	in	00	Sep	bara	te	
🗁 Data 🖉	~	O     Dhusical Bacilians     O     Sit	- NI	mba	**	0	0000	-		~~											
😑 🗁 Sequence			.0 140	mbe	13	0	.ocus	<u> </u>													
d8_sequence																			-		
• mdp_genotype		■ 1		616					100					1040					2464		<u>۲</u>
map_genotype	H			010					123	۷				1040					2404		
Polymorphisms			_	_	~	~	-	5	(O N			5	7	12	14	19	16	2	9	13	3
			ö	÷	8	ë	4	ω.	iii iii	i ò	ടങ്	ë	5	ė s	ni 🚽	i igi	ö	Ě	ω	6	Ω.
mdp_population_structure											1				ì				Ù		
• mdp_traits		38-11	N	N	N	N	N	N	N N		N	N	N	NP	I N	N	N	N	N	NN	
🖨 🖓 Matrix		A272	G	C	T	C	A	C		G	T	C	A		. A	A	C	A	T		-
• mdp_kinship		A441-5 0554	G	C	T	-	~	금			T	C	A		. A	A .	E	A	T		-
Tree		A554	G	c	T	<del>c</del>	Δ	È			T	c	μ μ		·	μ <u>π</u>	10		T		-
Fusions		A619	G	c	Ť	c	A	c	GA	G	Ċ	c	A	$\frac{c}{c}$	. A		È	A	T		
Synonymizer		A632	G	c	T	c	A	<del>c</del>	CA	G	T	c	A	c	A	A	c	A	T	C I	-
		B103	G	С	Т	С	A	c	C A	G	Т	С	A	C	: A	A	C	A	T	C F	τ. I
·	<ul> <li>1</li> </ul>	B104	G	С	Т	С	A	C	C A	G	Т	С	Α	C	: A	A	C	A	Т	$\subset f$	I
Number of commence of	-[	B14A	G	С	Т	С	Α	C	CA	G	Т	С	Α	C	: A	A	C	Α	Т	C F	•
Number of sequences: 91		B37	G	С	Т	С	A	C	CA	G	Т	С	Α	C	: A	A	C	A	T	C F	4
Number of sites: 2466		B68	G	С	Т	С	A	C	C A	G	T	C	A	C	: A	A	C	A	T	#	<u>+</u>
Data type: IUPACNucleotide		B73	G	C	T	C	A	<u>c</u>	CA	G	T	C	A	C	: A	A	LC.	A	T		<u>-</u>
		B84	G	C	T	C	A	C	CA	G	T	C	A	C	: A	A A	C C	A A	T		-
		697	6		1		- H				   NI		H NI		. A	H NI		-	1		-
		C105	G	ГЧ С	T	- N	0	C I			T	C	0		· ·				T		-
Table Cancel	i ŀ	CM105	N	N	N	N	N	N	N N		N	N	N	NP	. <u>-</u>		N	N	N	NP	-
	-11	CM174	G	C	T	C	A	c	GA	G	C	C	A	C	: A	A	C	A	T	CF	-
Tree Plot Cancel		CM7	G	С	Т	С	A	c	G A	G	C	С	A	CO	: A	A	C	A	T	C F	. I
		CML10	G	С	Т	С	A	C	C A	G	Т	С	Α	C	: A	A	C	A	T	CF	<u>،</u>
2D Plot Cancel		CML247	G	С	Т	С	A	C	CA	G	T	С	Α	C	: A	A	C	A	T	C F	۰ I
		CML254	G	С	Т	С	Α	C	C A	G	T	C	Α	C	: A	Α	C	Α	T	C f	•
Th HTOC		CML258	G	С	Т	С	A	C	C A	G	C	C	Α	C	: A	Α	C	А	T	C F	•
Chart Carcel		CML261	G	С	T	C	A	C	G A	G	C	C	A	C	: A	A	LC	A	T	C F	<u> </u>
	2	CML277	G	C	T	C	A	c	CA	G	T	C	A	CO	: A	A	Ļ¢	A	T		-
		CML281	G	C	T	C	A	C	CA	G	T	C	A	CIC	.   A	A	C	A	T	C   A	•
Datasets were saved to test																					

### 1.1 Installation

The graphic version of TASSEL can be installed in one of the three ways: using Java Web Start, as a stand-alone application, or using the source code

#### 1.1.1 Web start

TASSEL can be installed using Java Web Start technology, which automatically checks for the most recent version of TASSEL each time the application is executed. In addition, Java Web Start will ensure that the correct version of the Java Runtime Environment is running, thus avoiding complicated

installation and upgrade procedures. Users should use Web Start unless they have a specific reason to use one of the other installation methods.

To begin, Java Web Start (JWS) must be installed (prior to the installation of TASSEL). JWS is included as part of Java Runtime Environment (JRE) 5.0 and above. PC's and Mac's will most likely have JWS already installed. If you need to install Java, the most recent version is available at <u>http://www.java.com</u>. The easiest way to tell if it is installed on your computer is to try running TASSEL from the following link:

#### http://www.maizegenetics.net/tassel

If you will be using TASSEL frequently and would prefer to launch the application from your desktop rather than by revisiting the website, Java Web Start can be used to manually launch TASSEL each time and/or to create a shortcut. Access the Java Application Cache Viewer by going to **Start > Settings > Control Panel > Java**. From the **General** tab, click on **Settings** in the **Temporary Internet Files** section and then click on **View Applications...** and the Java Application Cache Viewer will appear. (Another way of achieving this is by going to **Start > Run** and typing in **javaws**). The TASSEL icon should now be visible and can be used to launch the application. Shortcuts can be created from the menu of the Java Application Cache Viewer: **Application > Install Shortcuts**.

#### 1.1.2 Stand-alone

Downloading a "stand-alone" version is recommended for anyone who has a slow Internet connection. While Java Web Start is a very good way of deploying software, it does not ask the user before attempting to download updates. Thus, a slow Internet connection may start a download process that requires an unreasonable amount of time to complete. If you are not interested in disabling your network connection each time before starting TASSEL, we recommend downloading the stand-alone version which does not attempt to update the program. However, given that TASSEL is a Java application, a Java Runtime Environment (version 1.6.0 or greater) is still required. To get the stand-alone version, download tassel3.0\_standalone.zip from the TASSEL web site. To run the stand-alone version, double-click on the JAR file (sTASSEL.jar). Alternatively, from a command prompt (in Windows go to **Start** > **Run** and type in "cmd" or "command"), change into the tassel3.0\_standalone directory and execute this command:

start\_tassel.bat (For Windows)
start\_tassel.pl (For UNIX)

### 1.1.3 Open source code

Open source code for the TASSEL software package is available at: <u>http://sourceforge.net/projects/tassel</u>. The package uses a number of other libraries that are included in the TASSEL distribution. These include a modified version of the PAL library (<u>http://www.cebl.auckland.ac.nz/pal-project/</u>), the COLT library (<u>http://dsd.lbl.gov/~hoschek/colt/</u>), and jFreeChart (<u>http://www.jfree.org/jfreechart/</u>). GDPC middleware (<u>http://www.maizegenetics.net/gdpc</u>) provides database access.

### 1.2 Panels

TASSEL is organized into five main panels. (1) The Control Panel at the top contains menus and buttons to control functions. (2) The Data Tree Panel is located beneath the Control Panel on the left side. This panel organizes data sets and results. Data set(s) displayed in the Data Tree Panel must first be selected before a desired function or analysis can be performed. To select multiple data sets, press the CTRL key while selecting the data sets. (3) The Report Panel is located below the Data Tree Panel. It displays

information about a selected data set from the Data Tree Panel, such as the type of data and how it was created. (4) The Progress Monitoring Panel below the Report Panel shows the progress of running tasks and has buttons that can be used to cancel tasks. (5) The Main Panel occupies the right side of the viewing area. It displays the content of a selected data set from the Data Tree Panel.

Functions in TASSEL are accessed by buttons and menus on the Control Panel.

The three buttons on the top left are the Mode Selectors (Data, Analysis and Results). The buttons below the Mode Selectors changed when a new Mode Selector is clicked. The modes are described in section 2-4. To the right of the Mode Selectors are the Progress Bar, and the Delete, Print, Save and Help buttons.

## 2 Data Mode

Data mode serves the purpose of importing and managing data. Data mode is the default mode when TASSEL starts. Click on the Data button to switch to this mode.

Tassel has two ways of importing data. One way is via GDPC to import data from databases. The other way is via flat files formatted as genotypes (e.g. hapmap, flapjack, and plink), phenotypes (trait data), population structure and kinship matrices.

The preliminary data manipulations include filtering data by site or taxa, joining data and data transformation.



Genotype and phenotype data generated from numerous genomic research projects are still valuable resources for the public, even after results are published. Some of these data have been migrated to several databases and can be accessed using Genotype Data and Phenotype Connection (GDPC). GDPC is middleware that eliminates the need for end users of data to understand various database schemas and write SQL queries to extract data. Instead, the GDPC browser provides a single, easy-to-use interface which can extract genotype and phenotype data from a variety of sources<sup>10</sup>. Currently, GDPC has connections to the following databases:

- Gramene diversity for maize, wheat and rice<sup>11-14</sup> http://www.gramene.org/db/diversity/diversity\_view
- Panzea<sup>15-17</sup> <u>http://www.panzea.org</u>
  GRIN
  - http://www.ars-grin.gov

GDPC can be used within TASSEL or as a stand-alone application. To display GDPC in TASSEL, click on the GDPC button in Data mode.

TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.32											• ×
File Tools Help GDPC		-	-								
a store construction Data 🚸 Analysis 🔛 Results Delete W	/izard		Show N	lemory	]			0%		3 🖬 🝸	
🕬 GDPC 🗐 Load 🏦 Export 🍸 Sites 🕨 Taxa	ү Trai	its 👫	Impute SN	Ps ?+5	Transform	n p⊷q	Synonymize	r 🔘 v .	Join 💽		
Data A	00	ත en Sar	CAS In	율 nport	الله Export	Get Da	ta Load	ے Deselect All	Add Conn	Load (	Get Started
<ul> <li></li></ul>		Taxa	1	Taxon P	arents		Loci		Genotype Evo	eriments	
Point     A and helts / Effered and analytics shutter / Dalat		Environme	1 nt Experimer	ts	Studies	Local	ties Geno	types	Phenotype CXP	Connection	s Loo
Polymorphisms	18 1	Merge									
In the second seco				Fea2.4	Fea	2.4	Fea2.2	Fea2.1	Fea2.1	Fea2.2	Fea2.
Elbard ada perulation structure		JSG Y LC	05-40	G:G:G:G	G G:G:	:G:G:G:	C:C:C:C:C:	C:C:C:C	C:C:C:C:C:	C:C:CT	A:A:A:
micreo_micp_population_soluciale		JSG Y LO	S-159	G:G	G:G:	:G:G:G:	C:C:C:C:C:	C:C	C:C:C:C:C:	C:C	A:A =
Matrix		JSG Y LC	DS-43	G:G:G:G	G G:G:	:G:G:G:	C:C:C:C:C:	C:C:C:C	C:C:C:G:C:	C:C:C:C	A:A:A:
		C-17-	78	G:G:G:G	G:T:	G:G:G:	C:C:C:C:C:	C:G:G:G	G:C:G:G:C:	C:C:C:T	A:A:A:
Tree		M162	w								
- + Fusions		CML3	22								
Synonymizer		TILO	1	G:G				C:C		C:C	A:A
Result		TILO	2	G:G				C:C		C:C	A:A
-      Diversity		PI5666	588	G:T	G:G:	:G:G:G:	C:C:C:C:C:	C:G	C:C:C:C:C::	C:C	A:A
- A SNP Assavs		TILO	3	G:G				C:C		C:C	A:A
Number of commences 1.61	T	TILO	4	G:G				G:G		C:C	A:A
Number of sequences: 161		PI5666	586	G:G:G:G	G:G:G:	:G:G:G:	C:C:C:C:C:	C:C:C:C	C:C:C:C:C:	C:C:C:C	A:A:A:
Number of sites: 9		TILO	5	G:G				C:G		C:C	A:A
Data type: IUPACNucleotide		TILO	6	G:G				C:C			A:A
null		MAS-	15	G:G:G:G	G G:G:	:G:G:G:	C:C:C:C:C:	C:C:C:C	C:C:C:C:C:	C:C:C:C	A:A:A: +
				•							- F
	Та	axa (workin	g list)	Genot	type Experim	nents (wor	king list)				
				Fea2.1 Fea2.1			d				
Load Cancel				Fea2.1 Fea2.1 Fea2.2							
Sites Cancel E				Fea2.2 Fea2.3							
Load Cancel				Fea2.4 Fea2.4							
Traits Cancel				Fea2.5				<del>.</del>			
Program Status											

Data is available for import once the user has defined the desired filters and data is visible in either the **Genotypes** or **Phenotypes** tab. To load data, activate either the **Genotypes** or **Phenotypes** tab (depending on the data you wish to import) and then click the **Load** button

For additional information about GDPC, please see <a href="http://www.maizegenetics.net/gdpc">http://www.maizegenetics.net/gdpc</a>

### 2.2 Load Load

This function provides options to import files for genotypes, phenotypes, populations structure, and kinship matrices. Several common sequence formats are accepted for genotype data, including BLOB, Hapmap, Plink, and Flapjack, and a general format for polymorphism data. Some file types used by TASSEL version 2 are also supported for backward compatibility. Phenotype and population structure can be imported as numerical trait data or covariates. Kinship must be loaded as a square numerical matrix.

Users can either specify the file type or use the "Guess" option to let the program determine the file type. As an example, we describe how the "Guess" function can be used to import all the files from the tutorial data set. The tutorial data can be downloaded from the TASSEL website or using this link: http://www.maizegenetics.net/tassel/docs/TASSELTutorialData3.zip.

To use the data, the zip file must be uncompressed and saved in a folder that the user specifies. To import data click the LOAD button. The File Loader dialog box will then pop up to let user choose the files and specify a format. For the files in the tutorial data set, the default (Guess) function will load all the files correctly. Multiple files can be imported simultaneously by highlighting them first (holding Shift or Control key while clicking) and then clicking the Open button.



### 2.2.1 BLOB

A Binary Large Object (BLOB) is a collection of binary data stored as a single entity. In TASSEL, BLOBs are used to compress large data sets into more manageable sizes. For sequence data, three types of BLOBs are used: SNP value BLOB, position BLOB and SNP ID BLOB. The three BLOBs are used to store individual SNP values, SNP position within the genome and the SNP identifiers respectively.

A BLOB is composed of two components, a header and a body. The header for each BLOB is 1024 bytes long, while the length of the body depends on the type of BLOB and on the amount of data being stored.

For a more detailed description on the structure and information contained within the header and body, refer to the GDPDM BLOB Specifications. (http://www.maizegenetics.net/gdpdm/docs/20100526/GDPDMBlobSpecification 20100526.pdf)

### 2.2.2 Hapmap

Hapmap is a text based file format for storing sequence data. All the information for a series of SNPs as well as the germplasm lines is stored in one file. The first row contains the header labels, and each additional row contains all the information associated with a single SNP. The first 11 columns describe attributes of the SNP, while the following columns describe the SNP value for a single germplasm line. The first 12 columns of the first row should look like this, where "Line 1" is the beginning of germplasm line names.

rs#	alleles	chrom	pos	strand	assembly#	center	protLSID	assayLSID	panelLSID	QCcode	Line 1
-----	---------	-------	-----	--------	-----------	--------	----------	-----------	-----------	--------	--------

While all 11 header columns are required, not all 11 of the columns need to be filled in for TASSEL to correctly interpret the data. The only required fields are "chrom", Chromosome name, and "pos", Position.

For TASSEL to correctly read Hapmap data, the data must be in order of chromosome and position within each chromosome, and the file should be TAB delimited. If some of the data is missing the correct number of TABs must still be present, so that TASSEL can properly assign data to columns.

#### 2.2.3 Plink

Plink is a whole genome association analysis toolset, which comes with its own text based data format. The data is stored in a set of two files, a .map file and a .ped file.

The .ped file contains all the SNP values and has six mandatory header columns for Family ID, Individual ID, Paternal ID, Maternal ID, Sex and Phenotype. TASSEL only requires that the Individual ID field be filled in. Each row of the .ped file describes a single germplasm line. Notice in Plink, an unknown character is represented with a '0'. However in TASSEL an unknown character is represented with a '0'. However in TASSEL will automatically convert between the '0' and the 'N'. Any exported Plink files will represent the heterozygous indel with a '+' (insertion) and a '-' (deletion).

The .map file describes all the SNPs in the associated .ped file, where each row provides information on one SNP. The .map file must contain exactly four columns: Chromosome, rs#, Genetic distance and Position. TASSEL does not require the Genetic distance field to be filled in.

Both files should be TAB delimited.

For a more detailed description on the data format, please visit the Plink basic usage and data formats webpage: (http://pngu.mgh.harvard.edu/~purcell/plink/data.shtml).

### 2.2.4 Flapjack

Flapjack is a software tool for graphical genotyping and haplotype visualization. The program is capable of outputting data in its own text based data format. Like Plink, the data is stored in a set of two files, a .map file and a .geno genotype file.

The genotype file contains all the SNP values. Each column in the first row contains a SNP ID, except for the first column, which is blank. The first column of the following rows contains the germplasm line names. TASSEL requires that all fields be filled out in order for data to be read correctly.

The .map file describes all the SNPs associated with the genotype file. Each row describes a single SNP. There are three columns in the .map file for Flapjack, SNP ID, Chromosome and Position, all of which are required for TASSEL to run correctly.

Both files should be TAB delimited.

For a more detailed description on the Flapjack data file format, please visit the Flapjack data import website: (http://bioinf.scri.ac.uk/flapjack/help/gui.dialog.DataImportDialog.shtml).

### 2.2.5 Polymorphism

A general format that accepts almost any type of marker data can also be used. Any alphanumeric character is allowed. Diploid data can be represented by separating alleles with a colon (":"), for example A:A, A:B, or B:B. All loci in a file **must** have the same ploidy level. The first line starts with the symbol <Marker> followed by the marker names. Subsequent lines must start with the name of the individual or taxon genotyped followed by the marker scores in the same order as the header. Comments can be inserted at the beginning of the file as long as any comment lines begin with the symbol "#". Columns are TAB delimited. Numeric values are allowed but, by default, will be treated as classification variables not as covariates in analyses.

Examp	le 1:					
<marke< td=""><td>er&gt;</td><td>m1</td><td>m2</td><td>mЗ</td><td>m4</td><td>m5</td></marke<>	er>	m1	m2	mЗ	m4	m5
33-16	A	В	В	A	А	
38-11	A	A	В	A	А	
4226	A	В	В	А	В	

In some cases, a user may wish to have marker values treated as numerical covariates. If the first line of the file is "<Numeric>", then the data will be imported as numeric data but used as marker data in GLM and MLM.

#### Example 2:

<numer< th=""><th>ric&gt;</th><th></th><th></th><th></th><th></th><th></th></numer<>	ric>					
<marke< td=""><td>er&gt;</td><td>m1</td><td>m2</td><td>m3</td><td>m4</td><td>m5</td></marke<>	er>	m1	m2	m3	m4	m5
33-16	0	1	1	0	0	
38-11	0	0	1	0.3	0	
4226	0	1	1	0.5	0	

Note to TASSEL 2.1 users: The polymorphism format specified in TASSEL v2.1 is still supported to provide backward compatibility.

#### 2.2.6 Phylip

The Phylip format used by TASSEL version 2.1 will continue to be supported. Details on Phylip format are described at the following website: <u>http://evolution.genetics.washington.edu/phylip/doc/sequence.html</u>

#### 2.2.7 Numerical data

This type of format is used for trait and covariate data such as population structure. Similar to sequence alignment genotype data, numerical data also consists of two parts: a header that defines data structure and a body containing the main data. Tabs should be used as delimiters. However, any white space character such as blank will be treated as a delimiter as well. As a result, embedded blanks in names will cause data to be imported incorrectly. We suggest representing missing values using "NA", or "NaN". However, any text value (e.g. "?") will be interpreted as missing data. There are several formats for numerical data to fit the requirement for modeling. Trait data (dependent variables) can be imported by starting the first line with "<Trait>" and following that with the trait names. Additional classifiers may also be included in subsequent header rows by starting the row with "<Header name=xxx>" followed by a name for each column of data. For instance, to define environments, start the second header row with "<Header name=env>".

Comment lines may be inserted at the beginning of the file as long as each comment line begins with the character "#".

#### 2.2.7.1 Trait format

This format does not require users to provide information on number of rows and columns. The file stats with key word <Trait> followed by names of columns. The column for line should not be labeled.

Example 1, simple list of trait values:

<Trait> EarHT dpoll EarDia 811 59.5 NA NA 33-16 64.75 64.5 NA 38-11 92.25 68.5 37.897 4226 65.5 59.5 32.21933 4722 81.13 71.5 32.421 A188 27.5 62 31.419 ...

Example 2, traits data collected in multiple environments:

```
<Trait>
           EarHT PlantHT
                            EarHT PlantHt
<Header name=env> Loc1 Loc1 Loc2 Loc2
811 59.5 NA
                NA
                      NA
33-16 64.75 121.5 NA
                      NA
38-11 92.25 153.8 37.897
                            83.4
4226 65.5 130.1 32.21933
                            82.1
4722 81.13 165.7 32.421
                            90.1
A188 27.5 110.2 31.419
                            79.6
```

#### 2.2.7.2 Covariate Format

Covariate data uses the same format as trait data except that the first line must be "<Covariate>". This line tells TASSEL that the variables in this file will be used as covariates not as dependent variables. This is the format to use for population structure covariates.

Example: <Covariate> <Trait> Q1 Q2 Q3 33-16 0.014 0.972 0.014 38-11 0.003 0.993 0.004 4226 0.071 0.917 0.012 4722 0.035 0.854 0.111 A188 0.013 0.982 0.005 ...

#### 2.2.7.3 TASSEL version 2.1 formats

Version 2.1 formats for numeric data will continue to be supported to provide backward compatibility. However, that format does not identify covariates as such. As a result, any covariates imported using this format will need to be properly identified using the "Trait filter" function described later in the manual.

#### 2.2.7.4 Repeated measurements

A format for repeated measurements may be implemented in the future.

#### 2.2.8 Square Numerical Matrix

Kinship can be calculated externally from pedigrees by using SAS Proc Inbreeding<sup>18</sup> or from markers by using software packages such as SPAGedi<sup>19</sup>. The following format is provided to import the resulting kinship estimates:

If n represents the number of taxa, the format for kinship files is as follows:

n				
Taxa1Name	r11	r12	•••	r1n
Taxa2Name	r21	r22	•••	r2n
 TaxanName	rn1	rn2	•••	rnn

Here rij (i, j=1,2, ..., n) is the element in the kinship matrix located at row i and column j.

Missing values are not allowed for kinship matrix.

Important note: The current format is different from the format used in TASSEL version 2.0 or lower.

#### 2.2.9 Genetic Map

A Genetic Map is a list of markers with chromosome and map position and, optionally, physical position. It can be used by GLM and MLM to provide genetic positions in the output files. It is not used as part of the analysis. The input format is:

First line: <Map> (as is, including the brackets) Following lines: marker name, chromosome name, genetic position, physical position (actual data)

Example:			
<map></map>			
marker1	<b>c</b> 1	21.3	2456873
marker2	<b>c</b> 1	52.1	52345691

There is no header line as such. Marker name, chromosome name, and genetic position are required. Physical position is optional and not used at this time. It is there because it is anticipated that information from this map may be used to convert between physical and genetic position at some time in the future.

## 2.3 Export Export

Options are provided to export sequence data: BLOB, Hapmap, Plink, Flapjack, Phylip (Sequential or Interleaved). Phenotypes and covariate data is exported as numerical trait data. Table Reports are exported as a tab delimited table.

This button has the same function as the "Save selected as" on the File menu. For numerical data, the function of Export is similar to the Table function in Results mode.

### 2.4 Sites Sites

The alignment can be filtered in several ways. Monomorphic sites can be eliminated, and regions of a sequence can be eliminated.

	Filter Alignmen	t
Minimum Count:	68	out of 91 sequences
Minimum Frequency:	0.05	
Position Type:	Position index	Physical Position (AGP)
Start Position:	0	0
End Position:	2465	2465
Remove minor SNP states		
Generate haplotypes via skiling window Generate haplotypes us skiling window Haplotype Length Step Length		

**Minimum Count** - the minimum number of taxa in which the site must have been scored to be included in the filtered data set (GAP or missing data do not count).

**Minimum Frequency** - the minimum frequency of the minority polymorphisms for the site to be included in the filtered data set.

Start Position, End Position – establishes the range of sites for filtering.

**Extract Indels** - if selected, indels are extracted from the alignment. If not selected, only point substitutions are extracted.

**Remove minor SNP states** – converts tertiary and rarer states to missing data ("?"), thereby forcing sites to have only two types of segregating sites at a locus. This may help remove sequencing errors.

Generate haplotypes via sliding window – creates haplotypes from an ordered set of SNPs.

### 2.5 Site Names Site Names

First select genotypic data from the data tree. The resulting dialog displays the site names associated with the selected data. By using either the CTRL or SHIFT key in conjunction with the mouse, the user can select or deselect site names. Once desired site names have been moved to the "Selected" window using the "Add ->" button, the "Capture Selected" or "Capture Unselected" buttons will create a new data set containing only the desired taxa.

Using the search box...

- \* is the wildcard.
- \* is always implied at end of search string.
- Search string is case sensitive. For example: use [Aa]bc to match taxa beginning with Abc or abc.
- PZ[AB] Will match anything starting with PZA or PZB.

0 0	Site	Name Filter		
Available			Selected	
PZB00859.1			PZA03613.2	
PZA01271.1			PZA03613.1	
PZA03613.2			PZA02962.13	
PZA03613.1				
PZA03614.2	(	Add ->		
PZA03614.1				
PZA00258.3				
PZA02962.13	<b>A</b>			
PZA02962.14	Ŧ			
PZ[AB]				
Capture Selected	Capture L	Inselected	Remove	Cancel

### 2.6 Taxa 🕨 Taxa

First select genotypic, phenotypic, or population structure data from the data tree. The resulting dialog displays the taxa associated with the selected data. By using either the CTRL or SHIFT key in conjunction with the mouse, the user can select or deselect taxa. Once desired taxa have been moved to the "Selected" window using the "Add ->" button, the "Capture Selected" or "Capture Unselected" buttons will create a new data set containing only the desired taxa.

Using the search box...

- \* is the wildcard.
- \* is always implied at end of search string.
- Search string is case sensitive. For example: use [Aa]bc to match taxa beginning with Abc or abc.
- A[56] Will match anything starting with A5 or A6

Available			Selected	
A554	0		A556	
A556			A6	
A6			A619	
A619			A632	
A632	$\mathbf{\Psi}$	Add ->	A634	
A634				
A635				
A641				
A654	•			
A[56]				
Capture Selected	Capto	ure Unselected	Remove	Cancel

## 2.7 Traits Traits

Clicking the "Traits" button on the "Data" toolbar launches the Trait Filter dialog. This dialog is used with numerical data sets to (1) change the trait type, (2) view, but not change whether the trait is discrete or continuous and (3) drop one or more traits from the data set. In addition, the dialog can be used to view the trait properties without changing them. If the "OK" button is clicked, a new data set is created that incorporates the changes, the original data set remains unchanged, and the dialog closes. If the "Cancel" button is clicked no data set is created, the original data set remains unchanged, and the dialog closes.

Allowable trait types are data, covariate, factor and marker. Generally, data and covariate traits will be continuous (not discrete) and factor will be discrete. Markers in a numerical data set will be continuous. Discrete valued markers are better imported as sequence or polymorphisms.

Clicking "Exclude All" unchecks the "Include" box for all traits. Clicking "Include All" checks the "Include" box for all traits. The "Exclude Selected" and "Include Selected" buttons do the same thing for traits that have been highlighted by selecting them with the mouse.

**Important:** Once a numerical data set has been joined with genotypes, it can no longer be modified using the trait filter function.

Trait	Туре	Discrete	Include	
Q1	covariate			
Q2	covariate		✓	
Q3	covariate			
	Exclude Selecte	ed Include Selec	ted	
	Exclude	All Include All		
		K Cancel		

## 2.8 Impute SNPs Impute SNPs

This function is used to impute missing genotypes. A sequence data type is required to use the function.

### 2.9 Transform ?+5 Transform

This suite of functions allows multiple data manipulation on genotype and phenotype (numerical) data. When a genotype data set is selected, the data are transformed to numbers. When a numerical data set is selected, mathematical transformation, data imputation and principal component analysis (PCA) can be performed. The Transform columns tags will be displayed in a Data dialog box with three tabs: Trans, Impute and PCA.

## 2.9.1 Genotype Numericalization ?+5 Transform

Two options are provided to transform genotype from character to numerical as shown in the following dialog box.

4	
	Numerical Genotype
	Collapse Non Major Alleles
	Separate Alleles
	Create Dataset Close

#### 2.9.1.1 Collapse Non Major Alleles

This function assigns 1 to the major allele and 0 to any other alleles. The converted genotypes are saved in a new numerical data set.

#### 2.9.1.2 Separate Alleles

This function assigns an indicator (1 for present and 0 for absent) for each allele. The converted genotypes are saved in a new numerical data set.

#### 2.9.2 Transform and/or Standardize Data

The **Trans** dialog box is the default selection, as shown below. In the **Column** list, select the column(s) you wish to transform. Then select the type of transformation you wish to execute. Selecting the **Standardize** checkbox will transform data by subtracting the column mean from the value of the trait and then dividing by the column's standard deviation. Clicking on the **Create Data set** button will result in the placement of a dataset containing only the selected columns in the Data Tree.

Column	Percent Missing Data	Trans Impute PCA
arHT.null	0.66	Raise to Power 2
poll.null		
arDia.null	12	Take Log Base 10
		E Standardize

### 2.9.3 Impute Phenotype

The k-nearest-neighbor  $algorithm^{20}$  is used to impute missing phenotype data. If data is missing for a taxon for one of the traits, the algorithm finds other taxa (neighbors) that are most like it for the non-missing traits. It uses the average of the neighbors to impute the missing data. Click on the **Impute** tab to display the following:



#### 2.9.4 PCA

Principal component analysis (PCA) can only be performed on a numerical data set without missing values. Two methods are available: correlation or covariance. This determines whether a correlation or covariance matrix will be used as the basis for the analysis. The default, correlation, is a reasonable choice for genetic data. The number of PCA axes in the output data set can be controlled by selecting either of the minimum eigen value associated with each axis, the minimum percent of the variance captured by an axis or the number of axes. The resulting axes will be sorted by the amount of variance each captures.

Т		X
Column	Percent Missing Data	Trans Impute PCA
EarHT.NA	0.00	
dpoll.NA	0.00	Method
EarDia.NA	0.00	<ul> <li>Correlation</li> </ul>
		<ul> <li>Covariance</li> </ul>
		Output
		Eigenvalue ≥ 0
		O Var Prop % ≥ 0.33333
		◯ Components = 3 🗘
	Create Dataset	Close

# 2.10 Synonymize Taxa Names P+9 Synonymizer

• This button makes taxa names uniform to permit the joining of data sets.

The join functions that generate fused data sets work by matching taxa names. Consequently, if multiple names exist for a given taxon (an added suffix, alternative spellings, different naming conventions, etc.) then the two data sets will not join correctly. To help remedy this, the Synonymizer function allows the taxa names of one data set to replace similar taxa names in the second data set. It relies on an algorithm that calculates the degree of similarity between names, using the name from the first set which is most similar to that in the second data set.

When using the Synonymizer, keep in mind that order of selection matters. Always select the data set with the names you wish to use (the "real" name) *first*, and then, while holding down the CTRL key, click on the second data set with the taxa names you wish to change (the "synonym"). Then click on the **Synonymizer** button. A synonym data set will be placed on the Data Tree panel under **Synonyms**. Each name in the data set selected second is now listed in the **TaxaSynonym** column. Next to this column is a **TaxaRealName** column listing the highest scoring match derived from the "real" name data set. The **MatchScore** column gives an indication of the amount of similarity between the two names (where 0 is no similarity and 1.0 is identity).

ATASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.36							
File Tools Help GDPC							
create Data 🚯 Analysis 🛺 Results Delete Wizard	M Sh	ow Memory				0%	
Traits 828 Ir	npute SNPs	2+5 Transform	D⇔đ Svn	onymizer	Mu Join		OO Separate
		-			-		
- Sequence	TaxaSyno	TaxaRealN	RefIDNum	MatchScore			
d8 sequence	TZI11		-1	0.0	_		*
mdp. genotype	TZI 10	TZI 10	82	1.0	_		
Polymorphisms	B79	B73	12	0.5	_		=
Numerical	AB28A		-1	0.0	_		
mdp population structure	B77	873	12	0.5	_		
mdp traits	M37W	M37W	52	1.0	_		
Matrix	B76	B73	12	0.5	_		
mdn kinshin	KI3	KI3	47	1.0	_		
• Tree	B75	B73	12	0.5	_		
Fusions	GT112	GT112	37	1.0	_		
	CI31A	KI3	47	0.33333333			
d8 sequence Synonyms	B73	B73	12	1.0	_		
Regit	VA59	A554	3	0.33333333			
( Kedure	CI44	KI44	49	0.66666666			
	CML92	CML5	29	0.57142857			
-	CML91	CML5	29	0.57142857			
	CI7	CI 187-2	16	0.25			
Table Title: Taxa Synonym Table	CML5	CML5	29	1.0			
Number of columns: 4	SC55	SC55	78	1.0			
Number of rows: 301	OH7B	OH7B	72	1.0			
Number of elements: 1204	IL677A	IL677A	43	1.0			
Taya simonime	CI3A	KI3	47	0.4			
Taxa synonyms	CM37	M37W	52	0.66666666			
Synonym lable	B68	B68	11	1.0			
301 unique matches	B115		-1	0.0			
73 unmatched:	B64	B68	11	0.5			
	H100		-1	0.0			
	DE811	KI11	45	0.28571428			
	CH9	H99	38	0.5			
Load Cancel	CO255	K55	44	0.33333333			
	NC372		-1	0.0			
	NC370		-1	0.0			
Synonymizer Cancel	U267Y	U267Y	85	1.0			
	B109		-1	0.0			
	WF9	WF9	90	1.0			-
Program Status							

**Caution!** Before the synonyms are applied, we strongly encourage the user to check the match score, especially for those taxa with low match scores. To do that, the user selects the synonym file and clicks the "Synonymizer" button. The incorrect matches, usually the ones with the lowest match scores, can be rejected at this point. Sorting on the match score column first makes this a fairly easy process.

In the event that some of the taxa are not interpreted correctly, matches can be modified manually. Select the taxa you wish to modify on the left side, and then choose a replacement taxa from the right side. Click

the arrow button **to** substitute the taxa. Taxa with no synonym can be identified by selecting then clicking "No Synonym". Click **OK** to save the changes.

🛓 Thresh	old for sync	nymizer						x
Synor	nymizer				_			
TaxaSy	TaxaRe	RefIDNum	MatchSc			TZI 18		*
TZI11		-1	0.0			12110		
TZI 10	TZI 10	82	1.0		1	1210		-
B79	B73	12	0.5			20.11		=
AB28A		-1	0.0			38-11 CT112		
B77	B73	12	0.5	1		W117UT		
M37W	M37W	52	1.0					
B76	B73	12	0.5			CI197-2		
KI3	KI3	47	1.0	1		U107-2		
B75	B73	12	0.5			WEAN		
GT112	GT112	37	1.0	1	<	W107A		
CI31A	KI3	47	0.333333			W152D		
B73	B73	12	1.0			U267V		
VA59	A554	3	0.333333			TV601		
CI44	KI44	49	0.666666			T232		
CML92	CML5	29	0.571428			5618		
CML91	CML5	29	0.571428			5010		
CI7	CI 187-2	16	0.25			SC00		
CML5	CML5	29	1.0			SA24		
SC55	SC55	78	1.0			06100		
OH7B	OH7B	72	1.0			Q0133		
IL677A	IL677A	43	1.0	Ŧ		D30		Ŧ
	No S	Synonym		1		]	Sort Alphabetically Apply threshold	
		ОК				Canc	cel	

Once it has been determined that the taxa names were matched correctly, the synonyms can be applied. With the synonyms selected, hold down the CTRL key while clicking on the second/synonym data set (the data set whose names you would like to change). Then once again click on the **Synonymizer** button to apply the new names to the data set.

### 2.11 Union Join Out

This button joins multiple data sets by a union of their taxa. Missing data will be inserted if taxa are missing from one data set.

Select multiple data sets using the CTRL key in conjunction with mouse clicks, and then click on the union button to join the data sets.

Because this function uses taxa names to join data sets, any variation in taxa names can prevent proper joining. Taxa names can be made uniform by using the "Synonymizer".

## 2.12 Intersection Join

This button joins multiple data sets by the intersection of their taxa. Taxa must be present in both data sets to be included.

Select multiple data sets using the CTRL key in conjunction with mouse clicks, and then click on the intersection button to join the data sets.

Because this function uses taxa names to join data sets, any variation in taxa names can prevent proper joining. Taxa names can be made uniform by using the "Synonymizer".

## 3 Analysis Mode Analysis

Analysis mode consists of the following options:

### 3.1 Diversity <sup>Toversity</sup>

This button executes a basic diversity analysis.

Average pairwise divergence ( $\pi$ ), segregating sites, and  $\theta$  estimates (4N $\mu$ ) can be calculated, as well as sliding windows of diversity.

To run a diversity analysis, click on a raw sequence alignment, and then select Analysis  $\rightarrow$  Diversity.

🛓 Diversity Surveys	×
Silent	
Noncoding	
Intron	Start Base 0
Synonymous	End Base 2465
Non-transcribed	
Noncoding Indels	
Transcribed	Sliding Window
Nonsynonymous	Step 100
Coding	Window 500
Coding Indels	
Veral	
Indels	
Run	Close

In the resulting Diversity Surveys dialog box, the various site classes available for analysis are listed on the left. If the sequence has no annotation, then only the "Overall" and "Indels" options will be active. A sliding window of diversity can also be calculated across the region. To produce a sliding window, check the box next to "Sliding Window," and then enter the desired step size and size of the sliding window.

Results can be plotted using **Results**  $\rightarrow$  **Chart** or viewed in a table via **Results**  $\rightarrow$  **Table**.

## 3.2 Linkage Disequilibrium 🔽 Link. Diseq.

This button generates a linkage disequilibrium data set from SNP data.

NOTE: It is important to use only filtered data sets (apply **Data**  $\rightarrow$  **Sites** first) when estimating linkage disequilibrium, as a raw alignment with numerous invariant bases will take a very long time and consume a large amount of memory to calculate.

🛃 Linkage Disequilibrium	x
Rapid Permutations (slightly biased p-values)	
Full Matrix LD	
Sliding Window LD	
50 LD Window Size	
Run Close	

Linkage disequilibrium between any set of polymorphisms can be estimated by clicking on a filtered set of polymorphisms and then using **Analysis**  $\rightarrow$  Link. Diseq. At this time, *D'*, *r2* and *P*-values will be estimated. The current version calculates LD between haplotypes with known phase only (unphased diploid genotypes are not supported; see PowerMarker or Arlequin for genotype support).

**D'** is the standardized disequilibrium coefficient, a useful statistic for determining whether recombination or homoplasy has occurred between a pair of alleles.

 $r^2$  represents the correlation between alleles at two loci, which is informative for evaluating the resolution of association approaches.

D' and r2 can be calculated<sup>21</sup> when only two alleles are present. If multiple alleles are present, a weighted average of D' or r2 is calculated between the two loci<sup>22</sup>. This weighted average is determined by calculating D' or r2 for all possible combinations of alleles, and then weighting them according to the allele's frequency. *Note: It is not entirely certain that this procedure fully accounts for allele number effects.* 

**P-values** are determined by two methods. If only two alleles are present at both loci, then a two-sided Fisher's Exact test is calculated. *Note: Previous editions of TASSEL used a one-sided test, but TASSEL version 1.0.8 and later use a two-sided test.* 

If more than two alleles are present, permutations are used to calculate the proportion of permuted gamete distributions that are less probable then the observed gamete distribution under the null hypothesis of independence<sup>21</sup>.

When calculating linkage disequilibrium, users have the option of employing "**Rapid Permutations.**" If this option is selected, the algorithm will compute either a fixed number of permutations or run until 10 permutations are found that are more significant than the observed P-value. While this slightly reduces P-values, it also saves a large amount of computational time. If an unbiased p-value is desired, then the user must unselect the "**Rapid Permutations**" check box.

"Full Matrix LD" calculates LD for every combination of sites in the alignment. "Sliding Window LD" calculates LD for sites within a window of sites surrounding the current site. The LD Window Size determines the width of the window on one side of the current site.

Linkage disequilibrium results can be plotted using **Results**  $\rightarrow$  **LD Plot** or viewed in a table via (**Results**  $\rightarrow$  **Table**).

## 3.3 Cladogram VCladogram

This function generates a tree or cladogram data set.

TASSEL produces neighbor-joining trees using only simple parsimony substitution models.

To retrieve cladogram data, first select genotypic data from the Data Tree panel and then click on the "Analysis" button, followed by the "Cladogram" button. The resulting tree data and the corresponding matrix will appear as separate data sets on the Data Tree panel.

Results can be plotted using **Results**  $\rightarrow$  **Tree Plot**.

## 3.4 SNP Extract SNP Extract

"SNP Extract" extracts SNPs from a raw sequence alignment into a useful format for export. Additionally, this function provides information for designing genotyping assays.

Below is a detailed explanation of the SNP Extractor Dialog:

SNP Extractor Dialog	x
Minimum Site Frequency:	0.85
Minimum SNP Frequency:	0.05
Minimum Surrounding Bases:	150
Minimum Good SBE Bases:	18
☑ Filter SNPs to Biallelic	
ОК	Cancel

Minimum Site Frequency: the minimum frequency for which the site must have a good base

**Minimum SNP Frequency**: the minimum frequency of the minority polymorphisms for the site to be included in the resulting data set

Minimum Surrounding Bases: the minimum number of good bases on at least one side of the SNP

Minimum Good SBE Bases: the minimum number of good bases on at least one side of SNP

Filter SNPs to Biallelic: converts tertiary and rarer states to missing data ("?"), thereby forcing sites to have only two types of segregating sites at any particular locus. This helps to remove bad sequence effects.

Results are displayed on the Data Tree panel and include SNPs along with their context. Additional information is also provided, including: the location of the nearest polymorphisms on either side, polymorphism information content ("**PIC**") and "**Haplotype PIC**". "**Overall score**" is essentially an estimate of the ability to design a single-base pair extension reaction in the region.

These results can be exported by using a table (**Results**  $\rightarrow$  **Table**).

## 3.5 Kinship Kinship

The function generates a kinship matrix from a set of random SNPs. To do so, first highlight SNP data then click on the "**Analysis**" button, followed by the "**Kinship**" button. The resulting kinship data will be added as a data set on the Data Tree panel.

When a genotype file is selected, the kinship matrix is generated by first using the TASSEL Cladogram function to calculate a distance matrix. Each element  $d_{ij}$  of the distance matrix is equal to the proportion of the SNPs which are different between taxon i and taxon j. The distance matrix is converted to a similarity matrix by subtracting all values from 2 then scaling so that the minimum value in the matrix is 0 and the maximum value is 2. Kinship can be derived from a set of random SNP data (a minimum of several hundred SNPs spread over the whole genome is recommended).

Warning: This method currently works correctly only for homozygous inbred lines. The method will be modified in the near future to work with heterozygous taxa. At that point, this warning will be removed.

Users may also load their own kinship data using **Data**  $\rightarrow$  **Load**. Kinship matrices can be calculated using the SPAGeDi software package (<u>http://www.ulb.ac.be/sciences/ecoevol/spagedi.html</u>). Comparisons of methods for calculating kinship can be found in the literature (*e.g.* Stich et al. 2008).

## 3.6 General Linear Model 🏝 💷

This function performs association analysis using a least squares fixed effects linear model.

TASSEL utilizes a fixed effects linear model to test for association between segregating sites and phenotypes. The analysis optionally accounts for population structure using covariates that indicate degree of membership in underlying populations. A main effects only model is automatically built using all variable in the input data. A separate model is built and solved for each trait and marker combination. Any factors, covariates, reps or locations are included in every model as main effects. How the data is

used must be defined either in the input data files or using the **Trait Filter** after the data has been imported but before it has been joined with a genotype.

General Linear Model (GLM) can be run using a numeric data set only, numeric data joined to genotype data. If only numeric data is selected, best linear unbiased estimates (BLUEs or least square means) will be generated for the taxa for each trait. [Note: only factors and covariates intended to control field variation should be included at this stage. Population structure covariates which are intended to control for marker effects should only be included when markers are also in the analysis.] If numeric data with genotypes are analyzed, each trait by marker combination will be tested and two reports will be produced, one containing trait by marker F-tests and the other containing allele estimates.

To run GLM, select a data set and then click the GLM button. A dialog box will pop-up to allow the user to indicate that a permutation test should be run and to allow the number of permutations to be changed. The permutation test will be run using the method suggested by Anderson and Ter Braak (2003), which calculates the predicted and residual values of the reduced model (contained all terms except markers) then permutes the residuals and adds them to the predicted values. When the GLM options dialog is closed, the user is presented with a dialog allowing the output to be saved to a file rather than stored in memory and displayed by TASSEL. This option is useful when the output is expected to be very large and risks exceeding available RAM.

The following table shows an example of the Marker Test output as viewed with Results/Table:

T Marker Test			R. J. Barlinson	•		2. LA	No. New York				(and a 🖓 )	×
Trait	Marker	Locus	Locus_pos	marker_F	marker_p	markerR2	markerDF	markerMS	errorDF	errorMS	modelDF	modelMS
EarDia	PZB00859.1	1	157104	1.663	0.199	0.007	1	23.522	223	14.144	3	31.921 🔺
EarDia	PZA01271.1	1	1947984	0.005	0.942	0	1	0.076	222	14.081	3	22.88 📟
EarDia	PZA03613.2	1	2914066	0.144	0.705	0.001	1	2.126	227	14.791	3	27.224
EarDia	PZA03613.1	1	2914171	0.014	0.905	0	1	0.208	228	14.66	3	24.487
EarDia	PZA03614.2	1	2915078	0.018	0.894	0	1	0.261	215	14.742	3	35.401
EarDia	PZA03614.1	1	2915242	2.663	0.104	0.012	1	39.907	213	14.984	3	38.967
EarDia	PZA00258.3	1	2973508	0.42	0.517	0.002	1	6.089	209	14.487	3	25.438
EarDia	PZA02962.13	1	3205252	0.374	0.541	0.002	1	5.374	217	14.357	3	25.183 🔻
•												P.
	Print Export (CSV) Export (Tab)											

In addition to displaying the F-statistics and p-values for the requested F-tests, the table also contains markerR2, mean squares (MS) and degrees of freedom (DF) for the marker effect, for the model (corrected for the mean), and for error. If taxa are replicated (across reps or environments), then the markers are tested using the taxa within marker mean square. If taxa are unreplicated, then the residual mean square is used. MarkerR2 is the marginal R2 for the marker calculated as SS Marker (after fitting all other model terms) / SS Total, where SS stands for sum of squares. The following table shows an example of the Allele Estimates output as viewed with Results/Table:

		000	Locus	Locus_pos	Allele	Estimate
EarDia	PZB00859.1	181	1	157104	С	0.804
EarDia	PZB00859.1	46	1	157104	A	0
EarDia	PZA01271.1	117	1	1947984	G	0.039
EarDia	PZA01271.1	109	1	1947984	С	0
FarDia	D7A03613 0	60	1	2014066	c	0.213

For each marker and trait combination, each marker allele is listed along with the number of observations for taxa carrying that allele (Obs), the locus (usually chromosome) and locus position of that marker, the allele, and the estimate of the effect of that allele. Because of the way that GLM codes alleles, the last allele estimate for a marker is always zero and the other allele estimates are relative to that.

## 3.7 Mixed Linear Model

This conducts association analysis via a mixed linear model (MLM).

A mixed model is one which includes both fixed and random effects. Including random effects gives MLM the ability to incorporate information about relationships among individuals. When a genetic marker based kinship matrix (K) is used jointly with population structure (Q), the "Q+K" approach improves statistical power compared to "Q" only<sup>9</sup>. MLM can be described in Henderson's matrix notation<sup>23</sup> as follows:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

where y is the vector of observations;  $\boldsymbol{\beta}$  is an unknown vector containing fixed effects, including genetic marker and population structure (Q); u is an unknown vector of random additive genetic effects from multiple background QTL for individuals/lines; X and Z are the known design matrices; and e is the unobserved vector of random residual. The u and e vectors are assumed to be normally distributed with null mean and variance of

$$\operatorname{Var}\begin{pmatrix}\mathbf{u}\\\mathbf{e}\end{pmatrix} = \begin{pmatrix}\mathbf{G} & \mathbf{0}\\\mathbf{0} & \mathbf{R}\end{pmatrix}$$

where  $\mathbf{G} = \mathbf{K}$  with as the additive genetic variance and  $\mathbf{K}$  as the kinship matrix. Homogeneous variance is assumed for the residual effect which means  $\mathbf{R}=\mathbf{I}$ , where is the residual variance. The proportion of genetic variance over the total variance is defined as heritability (h<sup>2</sup>).

$$h^2 = \frac{\sigma_a^2}{\sigma_a^2 + \sigma_e^2}$$

When K is derived from pedigrees, the elements of K equal 2\*Probability(IBD), where IBD means that two alleles drawn at random are identical by descent. Generally, K calculated from markers is an IBS matrix. The resulting multiplier is then not  $\sigma_a^2$  but some unknown constant times  $\sigma_a^2$ . Some methods for calculating K, such as those implemented in SPaGEDI, actually use markers to develop an estimate of the IBD relationship matrix. For those values of K, the resulting variance estimate can be considered an estimate of  $\sigma_a^2$  as long as the assumptions of the method used to derive K are not violated for the population being analyzed. One implication is that two different K matrices may give very different estimates of  $\sigma_a$  and heritability yet produce the same model fit and test of marker association.

TASSEL implements several methods to improve statistical power and reduce computing time. The Restricted Maximum Likelihood (REML) estimates of  $\delta_a^2$  and  $\delta_e^2$  are obtained through the Efficient Mixed-Model Association (EMMA) algorithm<sup>24</sup> which is much faster than the expectation and maximization (EM) algorithm<sup>25</sup>.

TASSEL also implements a method called compression which reduces the dimensionality of the kinship matrix to reduce computational time and improve model fitting. When MLM is used without compression (compression = 1), each taxon belongs to its own group. At the other extreme, GLM can be interpreted as maximum compression (compression = n) with all taxa in a single group. In that case, it is not possible to estimate the random effect independently of error and  $\delta_a^2$  is absorbed into  $\delta_e^2$ . Between these two extremes, taxa can be grouped using cluster analysis based on kinship. When n individuals are compressed into s clusters (groups), the kinship among individuals is replaced with the kinship among groups. At some grouping levels, dependent on the trait and population being analyzed, this compressed MLM has improved statistical power compared to the regular MLM<sup>4</sup>. The optimum grouping with the best model fit for MLM without fitting genetic markers has the best statistical power for an association test of markers<sup>4</sup>. TASSEL allows users to specify the compression level (average number of individuals per group), or to have the program determine the optimum grouping.

Similar to GLM, MLM performs an association test for each combination of traits and markers. TASSEL provides users several options: 1) to estimate genetic and residual variance for each combination; 2) to get these estimates once for each trait without fitting genetic markers and then to use those estimates to test markers; 3) to use a prior heritability estimate provided by the user. The second option, named P3D (population parameters previously determined), has the same statistical power as the first option<sup>4</sup>. Using the P3D method or using a prior heritability can be much faster than calculating heritability for each marker.

Using MLM is very similar to using GLM. The difference is that in addition to choosing the joint data set (or numerical data set), kinship data must also be highlighted before clicking the MLM button to show the MLM option dialog. The option of "No Compression" is the regular MLM which is equivalent to "Custom level=1". For data sets with large numbers of taxa, the optimal compression option may be considerably slower than no compression or user supplied compression. This is because the algorithm solves the model once for each of a series of compression levels in order to determine the optimal one.

All MLM analyses create two output tables, model statistics and model effects. If compression is used, the analysis creates three tables.

MLM_statistics_for_Filtered_mdp_traits + Filtered_mdp_population_structure + mdp_genotype_chr1_157104-3706018									×		
Trait	Marker	Locus	Site	df	F	p	errordf	markerR2	Genetic Var	Residual Var	-2LnLikelihood
dpoll	None			0			257		8.068	14.585	1,477.183
dpoll	PZB00859.1	1	157104	1	0.001	0.979	250	0	8.068	14.585	1,477.183
dpoll	PZA01271.1	1	1947984	1	4.339	0.038	248	0.015	8.068	14.585	1,477.183
dpoll	PZA03613.2	1	2914066	1	0.132	0.716	255	0	8.068	14.585	1,477.183
dpoll	PZA03613.1	1	2914171	1	2.829	0.094	256	0.01	8.068	14.585	1,477.183
dpoll	PZA03614.2	1	2915078	1	0.044	0.834	243	0	8.068	14.585	1,477.183
dpoll	PZA03614.1	1	2915242	1	0.788	0.375	241	0.003	8.068	14.585	1,477.183
dpoll	PZA00258.3	1	2973508	1	0.732	0.393	240	0.003	8.068	14.585	1,477.183
dpoll	PZA02962.13	1	3205252	1	0.967	0.326	244	0.004	8.068	14.585	1,477.183
dnoll	P7402962 14	1	3205262	1	0.026	0.873	239	n	8 068	14 585	1 477 183
•											•
Print Export (CSV) Export (Tab)											

The statistics table shows the results of the tests for each trait. The first line is for the model with no markers. Following that is a single line for each marker tested. The columns labeled "df", "F", and "p" are the degrees of freedom, F, and p-value from the F distribution for the test of the marker. The column "errordf" is the degrees of freedom used for the denominator of the F-test. The column labeled "markerR2" is the R2 for the marker calculated based on a formula for R2 for a generalized least squares GLS) model as shown here.

$$R^{2} = \frac{\left(\hat{Y}_{full} - \hat{Y}_{reduced}\right)^{T} V^{-1} \left(\hat{Y}_{full} - \hat{Y}_{reduced}\right)}{\left(Y - \overline{Y}\right)^{T} V^{-1} \left(Y - \overline{Y}\right)}$$

The columns "Genetic Var", Residual Var", and "-2LnLikelihood" list  $\sigma 2a$ ,  $\sigma 2e$ , and minus two times the model likelihood, respectively. When the P3D option is used, all of the values are the same for a given trait because they are only calculated once. A second table lists the estimated effects of each allele for each marker similar to the output for GLM. The compression results table shown below shows the likelihood, genetic variance, and error variance for each compression level tested during the optimization process. The meaning of groups and compression is discussed above in the description of the compression method. The compression level with the lowest value of -2LnLk is used for testing markers.

MLM_compression_for_Filtered_mdp_traits + Filtered_mdp_population_structure + mdp								
Trait	# groups	Compression	-2LnLk	Var_genetic	Var_error			
dpoll	259	1	1,480.402	7.362	8.146			
dpoll	248	1.044	1,479.47	7.635	7.81			
dpoll	243	1.066	1,479.505	7.656	7.855			
dpoll	238	1.088	1,481.049	7.338	8.416			
dpoll	234	1.107	1,482.935	6.957	9.069			
dpoll	229	1.131	1,483.301	6.904	9.261			
dpoll	224	1.156	1,482.597	6.866	9.394			
dpoll	220	1.177	1,486.718	6.172	10.576			
dpoll	215	1.205	1,485.526	6.407	10.342			
dpoll	211	1.227	1,486.045	6.21	10.709			
dpoll	207	1.251	1,488.214	5.897	11.345	Ŧ		
4					Þ			
Print Export (CSV) Export (Tab)								

# 3.8 Ridge Regression 篷 GS

This function performs ridge regression to predict phenotypes from genotypes. It is one of the methods used for genomic selection (GS).

The input dataset must contain one or more phenotypes and numeric marker data. Optionally, it may also contain factors and covariates. The analysis is run by selecting the input dataset then clicking the "GS" button. Because no additional user input is needed, the analysis will run immediately after the button is clicked. All traits will be analyzed separately using all of the genotypes, factors, and covariates in the dataset. The output will consist of two new datasets for each trait. One of the datasets will contain genomic estimated breeding values (GEBVs) for each taxon and the other will contain BLUPs for each marker in the genotype file. The output datasets will appear in the "Numerical" folder, which holds the input data as well. The output datasets can in turn be used for subsequent analysis. For example, it could be joined with the input data so that the predicted values could be graphed against the original values.

Understanding the input data requirements is important to ensure that the results of the analysis will be correct and useful. Genotypes must be numeric with one column for each marker. It is expected that the markers are bi-allelic, with the homozygotes coded as 1 and -1 and the heterozygotes coded as 0. However, any reasonable coding scheme will work. For instance, missing data could be replaced by a probability resulting from imputation. If any genotype data is missing, it will be imputed as the average of

the marker scores across all taxa for that marker. If a user prefers to use a different method of imputation, then the missing genotypes must be imputed before importing the data into TASSEL.

GEBVs will be calculated for all taxa in the dataset, including any lines that have missing phenotype data. A typical use of genomic selection is to predict GEBVs for a set of unphenotyped lines based on the performance of a training set. To do that a dataset containing both the genotypes to be predicted and the genotypes of the training set can be joined with a dataset containing the phenotypes of the training set using a union join. All taxa in the phenotype set should have genotypes. If an individual without genotype data is included, all the marker data for that individual will be imputed, which is not a generally useful thing to do.

## 4 Result Mode

Results mode consists of the functions to present data as table or graphics.

### 4.1 Table Table

Allows data to be displayed in a spreadsheet view and exported into a flat file.

To create a table, select a data set from the Data Tree panel, then click on the "Results" button followed by the "Table" button (**Results**  $\rightarrow$  **Table**). Shown below is an example in which diversity estimates are displayed.

T Diversity estimates									
Site_Type	StartSite	EndSite	Site	SiteCount	SegSites	Pi	Theta	Hapl	
All	0	499	249	500.0	2	4.1471571906352743E-4	7.853054255286058	NotAvail	~
All	100	599	349	500.0	2	4.7396082178688333E-4	7.853054255286058	NotAvail	
All	200	699	449	499.9347826086956	3	0.0012539335391229384	0.001178111805050	NotAvail	
All	300	799	549	499.9347826086956	4	0.0012974146839950498	0.001570815740066	NotAvail	
All	400	899	649	499.9347826086956	7	0.0020208359569745516	0.002748927545116	NotAvail	
All	500	999	749	499.9347826086956	7	0.0020208359569745516	0.002748927545116	NotAvail	
All	600	1099	849	499.9347826086956	6	0.002464276341277462	0.002356223610100	NotAvail	
All	700	1199	949	500.0	8	0.0018146201624462499	0.003141221702114	NotAvail	
All	800	1299	1049	500.0	9	0.0018580984233158102	0.003533874414878	NotAvail	
All	900	1399	1149	500.0	12	0.0014381270903010516	0.004711832553171	NotAvail	
All	1000	1499	1249	498.7601528905877	17	0.0026352213180236305	0.006691689460662	NotAvail	
All	1100	1599	1349	498.7601528905877	18	0.0021560977308205304	0.007085318252466	NotAvail	
<								>	
Print Export (csv) Export (tab)									

Data can be sorted by clicking on the column header of interest. A secondary sort can be done by holding down the CTRL key and clicking on a second column.

Data can be exported to flat files that are either comma-separated (Comma Separated Values = CSV) or tab-delimited. Both these formats can then be imported into a spreadsheet program such as Excel. Tables can also be printed.

## 4.2 Tree Plot Tree Plot

Displays the results of cladogram analysis.

After running Analysis  $\rightarrow$  Cladogram, select the desired data set and then click Tree Plot in the Results mode (Results  $\rightarrow$  Tree Plot). Trees can be visualized in either a Normal or Circular layout.


These images can be printed, saved in JPEG format, or saved as a Scalable Vector Graphics (SVG) file.



Displays 2D plots and determines color thresholds.

This function is useful for plotting associations in multiple environments.

First, select the desired result set. Using the drop down boxes provided, populate rows with "Environment," columns with "Site," and value with "PermuteP." The cutoff value for coloring can be chosen either by inputting a value in the text box or by using the slider tool to the right of the text box. Users can "mouse over" any box to view the value associated with that box, as shown here:

T 2-D chart												
🔀 📇 🎦 🏠 Cell size 🛛 14 🔽 🗋 only upper triangle 📄 P-Value 🏢												
	Row: 🔽	Environme	ent 🔽 Colum	n: 🗹 Site	~	Value: Perm	uteP 💌					
Cutoff 0.001							ուղուղուղ		Save SVG			
PermuteP	18	511	625	880	1000	1459	1490	1570	1618			
CLAYTON.ID15												
HOMESTEAD.ID1												
		0450300.00						CTEAD ID1.	1400/0_0			
Btatistics - Min: U.U. Max: (	J.98 Mean: U	.2658788 SD	0.32200772				HOME:	STEAD.IDT:	1490(0.0)			

If P-value coloring is desired, simply check the P-value box as shown below:

T 2-D chart									X		
🔀 📇 🗛 Cell size 🛛 14 🔽 🗋 only upper triangle 🔽 P-Value ?											
	Row: 🔽	Environme	nt 🔽 Columr	n: 🗹 Site	~	Value: Permu	ıteP 🔽				
Cutoff 80,9								nhunhunhun	Save SVG		
PermuteP	18	511	625	880	1000	1459	1490	1570	1618		
CLAYTON.ID15 HOMESTEAD.ID1						K					
Statistics - Min: 0.0 Max: 0	).98 Mean: 0	.2658788 SD:	0.32200772			HOM	ESTEAD.ID1:	1459(0.04)			

By checking the P-value box, Cutoff selection tools will be disabled and fields will instead be colored according to the following grayscale:

Dialog 🛛 🔀
< 0.001
0.001 - 0.01
0.01 - 0.05
> 0.05
OK Cancel

This key can be shown by clicking on the "?" icon next to the P-value check box.

# 4.4 LD Plot LD Plot

Displays the results from the linkage disequilibrium analysis.

After selecting the desired result from the Data Tree panel, click on the "LD Plot" button while in "Results" mode (**Results**  $\rightarrow$  LD Plot).

The graph that is generated displays LD between all possible pairs of sites. The black diagonal represents LD between each site and itself. The default setting graphs  $r^2$  in the upper right and *p*-values in the lower left. This default can be modified by clicking on the radio buttons in the lower left. The left side of the

graph contains a text description of the gene (or chromosome) and the site within the gene (or genetic position within the chromosome). At the bottom of the graph is a display of the position of each site along the gene or chromosome. This display can be hidden by deselecting the "Schematic" checkbox. Legends describing the color scheme appear on the right hand side of the graph.



LD plots can be printed, saved in JPEG format, or saved as a Scalable Vector Graphics (SVG) file. An SVG file is useful for creating publication quality graphics which can be easily sized using an editor such as Adobe Illustrator, Corel Draw, or OpenOffice.org Draw 2.0+.

# 4.5 Chart Chart

Chart provides a variety of graphs for visualizing numeric data.

This feature can be used to display histograms, XY plots, bar charts and/or pie charts. Any numeric table data can be charted, including LD results, phenotypic data, diversity results, and association results.

Histograms: Use the graph type combo box to select the desired graph type (Histogram) from the list of options. Up to two different series of data can be plotted together. Users may specify the number of bins to be used in the histogram.



Scatter plots: Use the graph type combo box to select the desired graph type (XY Plot) from the list of options. Select data to be plotted in X and Y axes using the appropriate drop down boxes. If two data series are plotted simultaneously on the Y axis, the "2 Y Axes" checkbox will provide an axis for each.



## 5 Menus

The menus in TASSEL include File, Tools, GDPC, and Help menus. The File menu is mainly used to save the entire data tree which includes the data loaded into TASSEL and the data created within TASSEL. A previously saved data tree can be loaded to TASSEL. This function provides the users the capability to save their intermediate results. The tools menu contains contingency test and option to set preference.

GDPC (Genomic Diversity and Phenotype Connection) is a software package to retrieve data from open database sources such as SNPs and phenotypic data. It can also be started using the "GDPC" button in data mode. Its use is described earlier in the manual.

## 5.1 File Menu

Individual data sets on the data tree and the entire data tree can be saved. An individual data set is saved in the genotype format for sequence data or numerical format for phenotype, covariate, and kinship. The data tree is saved in a binary format.

## 5.1.1 Save Data Tree

This feature allows you to save the entire contents of the Data Tree panel to a default location. This is helpful when the user does not wish to recreate a Data Tree panel that is already well populated with information the next time they initializes the program. To save a Data Tree, select **File > Save Data Tree**.

## 5.1.2 Open Data Tree

To restore a Data Tree that was saved previously saved, select **File > Open Data Tree**.

## 5.1.3 Save Data Tree As...

To save the contents of a Data Tree to a specific location or to give it a specific name, select File > Save Data Tree As....

## 5.1.4 Open Data Tree...

To restore a Data Tree from a specific location, select **File > Open Data Tree...** 

NOTE: The information outlined above for saving a Data Tree is applicable to files that are, in general, version specific. When a new version of TASSEL is released, a data tree saved with a previous version might not load to the version. For longer term storage, the best practice is to save individual data sets rather than the entire data tree.

## 5.1.5 Save Selected As...

To export data to one of the supported file types, select File > Save Selected As...

## 5.2 Contingency Test

Contigency/F	isher Exact Te v Or Fishe	st er Exact Te	est	L
Rows	Columns	F	leps	
2	<b>∨</b> 2	~		
Enter Data & I	lit Calculate			
	0		0	
	0		0	
	Left	-Tailed P:		
	Right	-Tailed P:		
	Two	-Tailed P:		
	Cal	culata Ca	ocel	

This utility calculates a chi-square contingency test or Fisher exact test (when using only the  $2 \times 2$  table of observations) using the same algorithm as is used in determining linkage disequilibrium.

#### 5.3 Preferences

The **Quality Score Colors** tab, found in the **Preferences** dialog box, allows the user to set cutoff values for visualizing quality score values on a sequence alignment or a set of called SNPs.

T Preferences		
Quality Score Colors		
Low Range Cutoff:		
	0	Set Color
-1 9 19 29 39 49 59 69 79 89 99		
Middle Range Cutoff:	0 - 20	Set Color
High Range Cutoff:		
	20	Set Color
-1 9 19 29 39 49 59 69 79 89 99		
ОК С.	ancel	

To set a desired threshold, simply adjust the slider on the left side of the dialog. Ns, "-" (dashes), and alignments without any quality score information have a default value of -1 (minus one).

## 6 Tutorial

This tutorial reviews several common scenarios for using TASSEL in order to help the user better understand its capabilities for data manipulation and association analyses. The TASSEL software package includes a tutorial data set that can be downloaded from the TASSEL website (please unzip all files to a directory of your choice). This tutorial data set contains data for phenotype, genotype, population structure, and kinship.

## 6.1 Missing Phenotype Imputation

The phenotype file **mdp\_traits** will be used to demonstrate the process of imputing missing data. Note that the data set below contains missing values (NaN).

ACC STATEST (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.37						x
File Tools Help <u>G</u> DPC						
Constant D Analysis Results Delete Wizard	M	Show Memo	ry		0%	
GDPC 🗐 Load 🏦 Export 🍸 Sites 🕨 Taxa 🍸 Traits 🎇	Impu	te SNPs	+5 Transform	p⇔q sync	onymizer 🔘 u	Join
d8_sequence_chr_6-2404		Таха	EarHT	dpoll	EarDia	
Porymorphisms		811	59.5	NaN	NaN	
	-	33-16	64.75	64.5	NaN	
mdp_pppalation_pd dcddre	-	38-11	92.25	68.5	37.897	
Matrix		4226	65.5	59.5	32.21933	
• mdp kinship	-	4722	81.13	71.5	32.421	
	_	A188	27.5	62	31.419	
Table Title: Phenotypes	<b>A</b>	A214N	65	69	32.006	
Number of columns: 4	=	A239	47.88	61	36.064	
Number of rows: 301		A272	35.63	70	NaN	
Number of elements: 1204	-	A441-5	53.5	67.5	35.008	
		A554	38.5	66	33.41775	
Load Cancel	1	A556	28	65	31.929	
Program Status		142	1100 5	100 5	104 5475	

To impute missing data, first select the  $mdp\_traits$  data set in the Data Tree Panel and then click the **Transform** button (**Data**  $\rightarrow$  **Transform**). The "Transform Column Data" window will open. Click on the **Impute** tab in this window. Finally, click on the **Create Data set** button to create the new data set with missing values imputed.

Note that missing values are now filled.

ATASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.37			A R server of	a conserve		X
File Tools Help <u>G</u> DPC						
Constant D Analysis Results Delete Wizard	M	Show Mem	ory		0%	
GDPC IL Load 🖆 Export 🝸 Sites 🕨 Taxa 🍸 Traits 🐯	Impu	te SNPs	?+5 Transform	p⇔q s	ynonymizer 🔘 U	Join
d8_sequence_chr_6-2404     Belymershime	•	Таха	EarHT	dpoll	EarDia	
		38-11	92.25	68.5	37.897	
	=	4226	65.5	59.5	32.21933	
mdp_population_strate		4722	81.13	71.5	32.421	
mdp traits 3 imputed		A 188	27.5	62	31.419	
Matrix	-	A214N	65	69	32.006	
	_	A239	47.88	61	36.064	
Imputed Phenotypic Values.	*	A441-5	53.5	67.5	35.008	
Taxa with insufficient data: 35		A554	38.5	66	33.41775	
K = 30.8 cutoff):	=	A556	28	65	31.929	
	-	A6	109.5	80.5	31.5175	
		A619	36	61	40.63	
Load		A632	60	61	35.953	-
Program Status		4634		Ino		

## 6.2 Principal Component Analysis

Principal component analysis (PCA) is a statistical tool that transforms a set of correlated variables into a smaller number of uncorrelated variables called principal components (PCs). The first PC captures as much of the variation as possible, and the succeeding PCs account for a decreasing fraction of the remaining variance. Another application of PCA is to use PCs derived from genetic markers to represent population structure<sup>8</sup>. This method requires much less computing time than maximum likelihood estimation. As most marker data are characters, numericalization must be performed first. A common approach for converting character marker scores is to set one of the homozygotes to 0, the other homozygote to 2, and the heterozygote to 1. For haploids, the conversion can be simply performed by coding one allele as 0 and the other as 1. The TRANSFORM function in TASSEL converts the major allele to 0. All the other alleles are collapsed to a single class and coded as 1. PCA requires that all variables should have variation and should not have missing values. As a result, filtering genotype to eliminate monomorphic markers and imputing missing values may be necessary. Imputing missing values can be done before or after numericalization. Here we demonstrate how to generate PCs from the genotype file in the tutorial data.

- 1. Remove monomorphic sites: Make sure TASSEL is in **Data** mode. Highlight the genotype and click **Site**. Set the minimum frequency to 0.05 and have "Remove minor SNP status" checked. Click **Filter**.
- 2. Numericalization: Highlight the filtered genotype and click **Transform**. Use the default option of "Collapse non major alleles." Click **Create data set**.
- 3. Imputation of missing values: Highlight the numerical genotype and click **Transform** and then click **Impute** Tab. Use the default options. Click **Create data set**.
- 4. PCA: Highlight the imputed numerical genotype, click **Transform**, and then click **PCA** Tab. Change the default option to "Components=3" by choosing **Components** and type 3 in the text box. Click **Create data set**.

🕌 Filter Alignment			
	Filter Alignment		<u>ه</u>
Minimum Count:	210	out of 281 sequences	
Minimum Frequency:	0.05		
Position Type:	Position index	Physical Position (AGP)	I Numerical Genotype I
Start Position:	0	157104	
End Position:	2560	148907116	
Extract Indels  Remove minor SNP states Generate haplotypes via sliding window Haplotype Length Step Length	of 2561 sites		<ul> <li>Collapse Non Major Alleles</li> <li>Separate Alleles</li> </ul>
	Filter Cancel		Create Dataset Close

📓 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.39																											
File Tools Help <u>G</u> DPC																											
Constance Data Analysis 🚂 Results Delete Wizard M Show Memory 0%																											
📰 GDPC 🗐 Load 🏦 Export 🍸 Sites 🕨 Taxa 🍸 Traits 🗱 Impute SNPs 🏞 Transform 🖓 + 🦞 Synonymizer 🌆 U Join 🔘 o Join																											
🔁 Data 🔺 🔨	<ul> <li>Physical Positions</li> </ul>	🚫 Site N	umbe	rs 🤇	Locu	ıs (	) Alle	eles						(Er	nter p	hysic	al po	sition	)					Se	arch		
<ul> <li>mdp_genotype</li> <li>mdp_genotype</li> <li>mdp_genotype</li> </ul>	157104		2494	8772		4	1 19740	0440			, 74532	108			9932	 :3776			12	4115	144			1489	<b>7</b> 0711	2	►
mdp_genotype_ctr1_157104-14890711     Polymorphisms     Dumerical     mdp_population_structure		36356797	36357534	39668467	40524105	42821031	43853993	44466196 44466243	44466246	44400414	43421000 48153258	48153805	48154058	50820416 50830416	50830673	50830782	50837246	50837488 55565777	55576390	55818939	56252241	56252478	57104591	57263770 57640200	57640764	57640944	57646430
map_traits     map_araits     map_inship     map_kinship		870-1	871:1	872:1	874:1	875:1	876:1	877:1 878:1	879:1	000	882:1	883: 1	884:1	885:1 886:1	887:1	888: 1	889:1	890:1	892:1	893:1	894:1	895: 1	896: 1	897:1	899:1	900: 1	901:1
Fusions	33-16	T	T	AC	: т	A	C N	I G	GO	:   c	A	G	G	:   т	c	G	c   i	C A	G	c	c	G	G	N A	A	G	Α 🔼
Synonymizer	38-11	G	T	G	: Т	G	NG	i G	G	: т	Α	Α	G	: т	С	Α	c i	T C	C	С	С	G	G	A G	G	G	A
🚞 Result 🧊	4226	T	T	A N	I T	A	CG	i G	CC	: т	Α	Α	G	: T	C	G	A	c c	C	T	С	А	G	A G	G	C	A
	4722	T	T	A C	T	N		A G	CI	·   c	A	G	GN	I T	C	G	C		G	N	N	N	G	WN	A	G	A
	A188	T	T	GO	C C	N		G	GO		G	A	G	T	C	G	C	T A	G	C	C	G	G	AA	A	G	A
Number of sequences: 281 🔷	A214N			AC	.	A		9 6	GU	.	A	A	GU	. C	C	G			C		C	N C	G	W G	<u> </u> G	G	A
Number of sites: 2561 🧮	A239		1 T		+-	G		4 6		+	A .	G		- T		A .			10		-	6	0	A G		문	H NI
Data type: IUPACNucleotide 💦 🚃	A441-5	T		A (	+-	N				-+-	<u> </u>	G		- <del>  </del> _		0			+	T	c	6	6	A G	+	G	
× · · · · · · · · · · · · · · · · · · ·	A554	T	$\left  \frac{1}{T} \right $	GC	T	G		A G	CT	-   ' -   T		G	GC	T	È	G	<u> </u>		10	T	è	Ğ	G	AG	G		A
	A556	T	N	GY	Ċ	A	N A	G	C T	· T	A	G	G	T	G	G	A	C A	Ē	T	T	N	G	WG	N	c	A 🗸
Program Status																1									í.		

📓 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.39												
File Tools Help <u>G</u> DPC												
cross cross Data 🚸 Analysis 👔	Results	Delete	Wizard	M Sł	iow Memory	]		0	%		?	
🔛 GDPC 🗐 Load 🏦 Export	Sites	Þ Taxa	👕 Traits	Impute SNPs	?+5 Transf	orm p⇔q sy	nonymizer	🔘 u Join	🛈 n Join	🛇 🗘 Separate	•	
🕽 Data 🔨	Taxa	S0	51	52	53	54	S5	S6	57	58	59	
P C Sequence	33-16	0	1	1	0	0	0	1	0	0	1	
<ul> <li>mdp_genotype</li> </ul>	38-11	0	0	1	0	0	0	0	0	0	0	
	4226	0	1	1	0	0	0	0	0	0	1	
mdp_genotype	4722	0	0	1	0	0	0	NaN	0	0	0	
mdp_genotype_chr1_157104-148907116	A188	1	1	1	0	0	0	0	0	0	0	
Polymorphisms	A214N	0	1	0	1	0	1	0	0	0	0	
P 🗁 Numerical	A239	1	1	0	0	1	1	0	0	0	0	
mdp_population_structure	A272	1	1	0	0	1	1	1	0	0	0	
mdp_traits	A441-5	0	1	1	0	0	0	0	0	0	0	
mdp_genotype_chr1_157104-148907116_	A554	0	0	0	0	1	0	0	0	0	0	
Matrix	A556	0	1	1	0	0	0	NaN	NaN	NaN	0	
• mdp_kinship	A6	1	1	0	0	1	1	0	0	0	0	
Tree	A619	0	0	1	0	0	0	0	1	1	0	
Fusions	A632	0	1	0	1	0	1	0	0	0	0	
Synonymizer	A634	0	1	0	1	0	1	0	0	0	0	
A Decult	A635	0	1	0	1	0	1	0	0	0	0	
Wumper of Columns: (55)	A641	1	1	0	0	1	0	0	0	0	0	
Number of course 201	A654	NaN	0	0	0	1	0	0	0	0	0	
Number of rows: 281	A659	1	0	0	0	1	0	0	0	0	0	
Number of elements: 719922	A661	0	0	0	0	1	NaN	0	0	0	0	
<u> </u>	A679	0	1	0	1	1	0	1	0	0	NaN	
	A680	in in	1	In	1	11	in	1	In	In	in 💌	
											>	
Program Status												

Column	Percent Missing Data		Trans Impute PCA
S0.null	2.1		Manhatten Distance
			Euclid Distance
			Unweighted Average
			Weighted Average
			O Weighted Average
		N	lumber of Neighbors (K): 3
		M	lin. Freq. of Row Data: 0.80
		-	

📓 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.39												
File Tools Help <u>G</u> DPC												
state state state state and the state of the											] 🕐	
🕬 GDPC 🗐 Load 🏦 Export	<b>Sites</b>	Þ Taxa	🍸 Traits 🛞	👌 Impute SN	Ps ?+5 Tra	nsform <b>p</b> ⇔q	) Synonymizer	🔘 u Join	🛈 n Join	🛛 🖓 Separ	ate	
۵ د	Taxa	50	S1	52	53	54	55	S6	57	58	59	
Sequence	33-16	0	1	1	0	0	0	1	0	0	1	~
mdp_genotype	38-11	0	0	1	0	0	0	0	0	0	0	
mdp_genotype	4226	0	1	1	0	0	0	0	0	0	1	
mdp_genotype	4722	0	0	1	0	0	0	0	0	0	0	
mdp_genotype_chr1_157104-148907116	A188	1	1	1	0	0	0	0	0	0	0	
Polymorphisms	A214N	0	1	0	1	0	1	0	0	0	0	
Numerical	A239	1	1	0	0	1	1	0	0	0	0	
mdp_population_structure	A272	1	1	0	0	1	1	1	0	0	0	
mdp_traits	A441-5	0	1	1	0	0	0	0	0	0	0	
mdp_genotype_chr1_157104-148907116_Collas	A554	0	0	0	0	1	0	0	0	0	0	
mdp_genotype_chr1_157104-148907116_Collas	A556	0	1	1	0	0	0	0	0	0	0	
Matrix	A6	1	1	0	0	1	1	0	0	0	0	
mdp_kinship	A619	0	0	1	0	0	0	0	1	1	0	
Tree	A632	0	1	0	1	0	1	0	0	0	0	
Fusions	A634	0	1	0	1	0	1	0	0	0	0	
Cuponuminor	A635	0	1	0	1	0	1	0	0	0	0	
	A641	1	1	0	0	1	0	0	0	0	0	
Imputed Phenotypic values.	A654	0	0	0	0	1	0	0	0	0	0	
Taxa with insufficient data: 7	A659	1	0	0	0	1	0	0	0	0	0	
K = 30.8% cutoff):	A661	0	0	0	0	1	0.33333	0	0	0	0	_
×	A679	0	1	0	1	1	0	1	0	0	0	
	A680	n	1	In	1	1	n	1	n	n	n	<u> </u>
		<										>
Program Status												



Three items will be added to the data tree after running PCA. The first are the PCs. The second are the eigenvalues. And, the last are the eigenvectors. Here we use the Chart Function in the Result mode to graph the first three PCs, the individual eigenvalue contributions (sometimes called a skree plot) and the cumulative eigenvalue contributions. The eigenvalues are of interest because they equal the variance explained by each of the PCs.





## 6.3 Estimation of Kinship using genetic markers

While PCs can be used to capture major population subdivisions, kinship can be used to capture more subtle relationships. This section shows how to create a kinship matrix based on the same SNP data used to calculate PC's.

- 1. Remove monomorphic sites: Highlight the genotype and click **Site** in **Data** mode. Set the threshold on MAF to 0.05, check "Remove minor SNP status," then click **Filter**.
- 2. Estimate kinship: Highlight the filtered genotype and click **Kinship** in **Data** mode. A kinship matrix will be added to the data tree under Matrix category.

📓 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.39											
File Tools Help <u>G</u> DPC											
Gaxa Gaxa Coxa Coxa Data	Results	Delete	Wizard	M Sh	ow Memory		l	0'	%		?
		🏢 Table	W Tree Ple	ot 🔛 2D P	lot 🐮 LD P	lot 🚹 Chai	rt				
map_genocype	Taxa	33-16	38-11	4226	4722	A188	A214N	A239	A272	A441-5	A554
map_genotype	33-16	2.0	0.41059635	0.41382671	0.47004667	0.35458766	0.30197215	0.33163612	0.31519530	0.45907538	0.41320502
mdp_genotype mdp_genotype_cbr1_157104-14890711	38-11	0.41059635	2.0	0.36100872	0.43054512	0.26257421	0.29979028	0.47599188	0.34147128	0.29127360	0.34688548
Polymorphisms	4226	0.41382671	0.36100872	2.0	0.38863963	0.29315161	0.25695761	0.37291940	0.16620914	0.31066169	0.50956082
	4722	0.47004667	0.43054512	0.38863963	2.0	0.33695692	0.29274166	0.42737888	0.36800713	0.46969122	0.52179566
mdp population structure	A188	0.35458766	0.26257421	0.29315161	0.33695692	2.0	0.25301853	0.34102312	0.29219921	0.33751064	0.34114748
mdp traits	A214N	0.30197215	0.29979028	0.25695761	0.29274166	0.25301853	2.0	0.37141276	0.22059220	0.14934564	0.24796287
mdp genotype chr1 157104-14890711	A239	0.33163612	0.47599188	0.37291940	0.42737888	0.34102312	0.37141276	2.0	0.39160082	0.33260847	0.42418659
mdp_genotype_chr1_157104-14890711 =	A272	0.31519530	0.34147128	0.16620914	0.36800713	0.29219921	0.22059220	0.39160082	2.0	0.44690342	0.40563040
🚊 🗁 Matrix	A441-5	0.45907538	0.29127360	0.31066169	0.46969122	0.33751064	0.14934564	0.33260847	0.44690342	2.0	0.33335121
mdp_kinship	A554	0.41320502	0.34688548	0.50956082	0.52179566	0.34114748	0.24796287	0.42418659	0.40563040	0.33335121	2.0
kin_mdp_genotype_chr1_157104-14890	A556	0.37545628	0.38747139	0.32322427	0.37214558	0.33681686	0.23745026	0.43953514	0.31439560	0.28893715	0.37408670
Tree	A6	0.30856833	0.28090708	0.28621390	0.42040214	0.22954367	0.17005219	0.35647018	0.39689910	0.38189907	0.30641707
Fusions	A619	0.34902247	0.28599197	0.30876406	0.45524754	0.30112220	0.19064762	0.34940884	0.20266868	0.31075000	0.28263598
Synonymizer	A632	0.23011784	0.34705504	0.21609578	0.29009216	0.26324743	0.95100562	0.25506456	0.21524838	0.17867383	0.34738935
🛅 Result 🛛 🕑	A634	0.29463120	0.35781792	0.28973072	0.29546773	0.30267760	1.10550823	0.33799043	0.22008600	0.21976495	0.35153913
	A635	0.29911961	0.31287794	0.24363542	0.30643638	0.23921689	0.98623945	0.33400231	0.21855159	0.22835054	0.36080213
	A641	0.39726813	0.36262502	0.26407750	0.35695987	0.28045903	0.79088096	0.32380022	0.19282299	0.29113453	0.47129657
e	A654	0.4144/822	0.34265472	0.51796310	0.3/4/5564	0.38022672	0.27452505	0.47240984	0.37837519	0.41451092	0.63930566
genotype chr1 157104-148907116 🚽	A659	0.41781542	0.37743951	0.35619768	0.42522581	0.36078896	0.36676025	0.64187328	0.25444722	0.27058073	0.46082075
	A601	0.34618751	0.34212471	0.31164700	0.50161412	0.20331735	0.2231/8/3	0.41240909	0.31558391	0.35093751	0.36040054
	A6/9	0.05732332	0.14708046	0.07293583	0.07423771	0.00035550	0.240/61/3	0.12458721	0.07765475	0.10017728	0.19000655
	ADOLI	<	10.0900/097	11.1149/01144	0.02379902		u	0.091216.5.5	III.II49090.10		
Program Status											

## 6.4 Association analysis using GLM

We use three files from the tutorial data set to perform association analysis using the **GLM**. The first file is the dwarf8 gene sequence with 2466 sites on 91 maize inbred lines. The second one is the population structure of 282 maize inbred lines. The last one is phenotypes for three traits, for 282 maize inbred lines. The statistical model is:

Flowering time = Population structure + Marker effect + residual

- 1. Remove monomorphic sites: Highlight the genotype and click **Site** in **Data** mode. Set the threshold on MAF as 0.05, then click **Filter**.
- 2. Trait selection: Highlight the phenotype and click **Trait** in **Data** mode. Uncheck all the traits except flowering time (DPOLL). Make sure that the Type is set to Data. Click **OK** to create a filtered phenotype.
- 3. Covariate selection: The population structure is presented as the proportion of each population. There are three populations represented as Q1, Q2, and Q3. They sum to 100%. This creates linear dependency if we use all of them as covariates. We can eliminate the dependency by removing one of them. In this demonstration, we exclude the last one. Highlight the filtered phenotype and click **Trait** in **Data** mode. Uncheck the last population (Q3). Make sure that the Type is set to Covariate. Then click **OK** to create a filtered population structure data.
- 4. Joining data: Highlight the three filtered data sets by holding the Control key while selecting the individual data. Then click Intersection ( $\cap$ ) Join on Data mode to create a combined data set.
- 5. Association analysis: Highlight the joint data set then click **GLM** in **Analysis** mode to perform association analysis. Two reports will be added to the data tree.



🛓 Filter Traits / Modify	Trait Properties			🍰 Filter Tra	aits / Modify Tr	ait Properties			
Trait EarHT	Type	Discrete	Include		Trait Q1	Type covariate	Discrete	Include	
dpoll EarDia	data data				Q2 03	covariate covariate			
	Exclude Selected Exclude All OK	Include Selecter Include All Cancel	1			Exclude Selected Exclude All OK	Include Select Include All Cancel	ed	

🖆 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.38																					
File Tools Help <u>G</u> DPC																					
GENER Data	Delete	Wizard		N	<b>S</b> h	ow N	Memo	ory								0%	6			H	?
GDPC Load & Export Sites	Taxa 🍸 Tr	aits 👫 Im	pute	SN	Ps	?+5	Tra	nsf	orm	p.	•q s	ynon	ymiz	zer	0	) U	Join		۵ſ	n Jo	in
equence																					
d8_sequence     mdp_genotype					-0														_		
mdp_genotype     mdp_genotype	6			605					1204	+	_	-	1	803	-	-	-		2402		
da_sequence_chr_6-2404     Filtered_mdp_population_struct			0:6	1:9	2:24	3: 28	4:40	5: 61	<u>6: 66</u>	7:81	o. 229 9: 452	10: 754	11:804	12:1297	3:1411	14: 1666	15:1786	l 6: 2245	17: 2276	8: 2356	19: 2404
Numerical																	Ξ.		۲.		- I
mdp_population_structure	3	8-11	N	N	C	-	C	G	C	A T	C	T	G	G	A	A	A	C	A	G	
mdp_traits	A	2/2	C	+	G	÷	G	A	+			+	G	G	A	A	A		A	G	
Filtered_mdp_population_structure	A-	554	C	÷	G	÷	G	A	÷			÷	G	G	Δ	A	G	÷	A		
Filtered_mdp_traits	<u> </u>	A6	c	Ť	G	÷	G	Â	Ť			τ I	G	G	Â	Â	A	c	Â	G	
Matrix	A	619	G	c	G	c	G	A	c	AC	T	-	G	G	G	A	A	G	G	c	c
mdp_kinship	A	632	С	т	G	т	G	A	т	GC	c	т	G	G	G	G	G	c	Α	G	c
Tree Training	В	103	С	т	G	т	G	A	т	GC	c	т	G	G	A	A	G	с	A	G	c
	В	104	С	Т	G	т	G	A	Т	GC	C	Т	G	G	G	G	G	C	A	G	c
	В	14A	С	Т	G	Т	G	A	Т	GC	C	Т	G	G	G	G	G	C	A	G	C
	E	337	С	Т	G	т	G	Α	Т	GC	C	Т	G	G	G	G	G	С	A	G	С
	E	368	С	Т	G	Т	G	Α	Т	GC	C	Т	G	G	G	G	G	C	Α	G	С
	E	373	С	Т	G	Т	G	Α	Т	GC	C	Т	G	G	G	G	G	С	Α	G	C 🛨
ogram Status																					

One of the reports added to data tree is labeled "GLM\_Marker\_Test\_" followed by the name of the joint data. In addition to the information for traits and markers, the data set contains the following statistics:

marker\_F: F value from the F test on marker; marker\_p: P value from the F test on marker; markerR2: R<sup>2</sup> for the marker after fitting other model terms (population structure);

markerDF: Degree freedom of marker; markerMS: Mean square of marker; errorDF: Degree freedom of residual error; errorMS: Mean square of residual error; modelDF: Degree freedom of model; modelMS: Mean square of model.

A TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.40												
File Tools Help GDPC												
Image: Second Data     Image: Second												
mdp_population_structure	Trait	Marker	Locus	Locus_pos	marker_F	marker_p	markerR2	markerDF	markerMS	errorDF	errorMS	modelDF
mdp_traits	dpoll	6	1		6 16.85981	1.1021E-4	0.12252	1	365.98041	68	21.70727	3
Filtered_mdp_population_structure	dpoll	9			9 1.11011	0.29579	9.906E-3	1	29.58913	68	26.6542	3
Filtered_mdp_traits	dpoll	24		2	4 4.13003	4.5925E-2	3.4081E-2	1	104.92808	70	25.40613	3
Matrix	dpoll	28		2	8 19.64751	3.4242E-5	0.1345	1	414.06591	69	21.07473	3
Tree	dpoll	40		4	0 5.05807	2.7487E-2	3.99E-2	1	126.6331	74	25.03586	3
+ Fusions	dpoll	61		6	1 5.43052	2.2443E-2	4.0985E-2	1	132.53059	76	24.40477	3
Synonymizer	dpoll	66		6	5 1.00808	0.31855	8.045E-3	1	26.01491	76	25.80629	3
🕌 Result 😑	dpoll	81		8	1 1.11701	0.29391	8.9017E-3	1	28.7851	76	25.76984	3 =
Diversity	dpoll	229		22	9 1.47735	0.22781	1.1545E-2	1	38.48458	79	26.04973	3
SNP Assays	dpoll	452		45	2 13.32677	4.6818E-4	9.0775E-2	1	302.60367	79	22.70645	3
• LD	dpoll	754		75	4 1.41055	0.23898	1.3107E-2	1	34.93908	70	24.7699	3
Association	dpoll	804		80	40.55818	0.45721	4.4122E-3	1	14.70836	79	26.35069	3
GLM_marker_test_d8_sequence_chr_6-2404 +	dpoll	1297		129	7 0.75478	0.38763	6.0117E-3	1	19.9997	78	26.49751	3
GLM allele estimates for d8_sequence_chr_6-2· +	dpoll	1411		141	1 7.79846	6.5566E-3	5.6502E-2	1	188.35345	79	24.15265	3
< III +	dpoll	1666		166	6 0.29725	0.58715	2.3574E-3	1	7.85844	79	26.4374	3
	dpoll	1786		178	6 1.33392	0.25159	1.0442E-2	1	34.81022	79	26.09624	3
۰ TTT	dpoll	2245		224	5 9.6922	2.577E-3	6.8723E-2	1	229.09406	79	23.63695	3
	dnoll	2276		227	5 10 51689	1 734E-3	7 3884F-2	1	246 29701	79	23 41919	3
	•					111						+
Program Status												

Clicking "marker\_p" will sort the table by P value. The smallest P value is  $1.1021 \times 10^{-4}$  for SNP at position 6. The threshold is  $5 \times 10^{-4}$  at a significance level of 1% after Bonferroni multiple test correction (0.01/20). The denominator in the Bonferroni correction is the total number of SNPs tested. The association was significant.

The other data added to the data tree is labeled "GLM\_Allele\_Estimates\_" followed by the name of the joint data. For the most significant SNP at position 6, there were two genotypes (CC and GG). There are 62 lines with genotype CC and 10 lines with allele GG. For the trait dpoll (days to pollination), the difference between the two homozygotes was 6.63755 days.

📓 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.39										
File Tools Help <u>G</u> DPC										
CARANG COCANG Data 🚸 Analysis 📴 Resu	ults Delet	te Wizar	rd 🛛	M Show Mem	nory		0%			
	I Table	Tree Plot	2D Plo	t 🚺 🚺 ED Plo	t 📊 Chart					
SNP Assays	Trait	Marker	Obs	Locus	Locus_pos	Allele	Estimate			
LD	dpoll	6	62		6	с	6.63755	<u>^</u>		
🗐 🗁 Association	dpoll	6	10		6	G	0			
GLM_marker_test_Filtered_mdp_	dpoll	9	49		9	Т	1.44879	1		
GLM allele estimates for Filtered	dpoll	9	23		9	С	0			
Variances	dpoll	24	15		24	С	3.10862			
Stepwise 🗠	dpoll	24	59		24	G	0			
	dpoll	28	1		28	-	9.23009	=		
	dpoll	28	62		28	Т	6.76409			
	dpoll	28	11		28	c	0			
	dpoll	40	17		40	c	3.27583			
road	dpoll	40	61		40	G	0			
	dpoll	61	18		61	G	3.25231			
Traits	dpoll	61	62		61	A	0			
	dpoll	66	30		66	C	-1.2317E0			
Sites 🛄	dpoll	66	50		66	Т	0			
	dpoll	81	29		81	A	-1.3196E0			
Traits	dpoll	81	51		81	G	0			
	dpoll	229	6		229	Т	2.84783			
∩ Join 🚺	dpoll	229	77		229	C	0			
	dpoll	452	73		452	C	6.16275			
GLM	dpoll	452	10		452	Т	0	_		
	dpoll	754	69		754	Т	-2.8769E0			
Program Status										

## 6.5 Association analysis using MLM

Running MLM in tassel is similar to running GLM. The difference is that in addition to the joint data (or numerical data), MLM requires kinship data to define the relationship between individuals. The kinship matrix times a parameter equals the covariance matrix between individuals. Here we use kinship file from the tutorial data set to fit the following statistical model.

Flowering time = Population structure + Marker effect + Individuals + residual

Individuals and the residual are fit as random effects. The other terms are treated as fixed effects.

With respect to the marker effect, we will demonstrate the analysis using two sets of markers. One is the dwarf8 gene sequence used in the GLM tutorial. The other is a set of 3093 SNPs spread across the maize genome.

For the dwarf8 gene sequence, use the joint data set created by following the tutorial for GLM. Solve the mixed linear model by highlighting the joint data set and the kinship data then clicking the **MLM** button in **Analysis** mode.

STASSEL (Trait Analysis by aSSociation, Evolution, and	d Linkage) 3.0.3	38	_	-								
File Tools Help <u>G</u> DPC												
GRANG Data Analysis	Delete	Wizard	M Sho	w Memory		0%						
Arrow Diversity Link. Diseq. W Cladogram Kinship Arrow GLM [] MLM												
Data 🔺	Таха	dpoll	Q1	Q2	Haplotype							
	38-11	68.5	3.0E-3	0.993	NNC-CGCAT	A						
<pre>mdp_genotype</pre>	A272	70	1.9E-2	0.122	CTGTGATGC							
mdp_genotype	A441-5 A554	66	1.9E-2	0.979	CTGTGATGC	=						
mdp_genotype	A6	80.5	3.0E-3	3.0E-2	CTGTGATGC							
as_sequence_cnr_6-2404 Eiltered_mdp_traits_+ Eiltered_mdp_populatio	A619	61	9.0E-3	0.99	GCGCGACA							
Polymorphisms	A632	61	0.993	4.0E-3	CTGTGATGC							
🗐 🖳 Numerical	B103	57.5	0.163	0.829	CTGTGATGC							
mdp_population_structure	B14A	63.5	0.998	1.0E-3	CTGTGATGC							
mdp_traits	B37	65.5	0.997	2.0E-3	CTGTGATGC							
<ul> <li>Filtered_mdp_population_structure</li> <li>Filtered_mdp_traits</li> </ul>	B68	71.5	0.998	1.0E-3	CTGTGATGC							
Matrix	B73	70	0.999	1.0E-3	CTGTGATGC							
mdp_kinship	B84	67.5	0.996	3.0E-3	CTGTGATGC							
A Tree	897	60.5	1.6E-2	0.981	CTGTGATGC							

🛃 MLM Options 📃 🔀									
Compression Level									
Optimum Level									
Custom Level:									
No Compression									
Variance Component Estimation									
P3D (estimate once)									
Re-estimate after each marker									
Run Cancel Help Me Choose									

An MLM option dialog will pop up as shown above. Choose the default options, which use P3D and compression at the optimum compression level. After the Run button is clicked, the progress bar will start moving. The time required will depend on sample size, number of traits, number of markers, and the options chosen in the MLM option dialog. After the progress bar is reset to zero, indicating completion of MLM, three reports will be added to the data tree. The first two are similar to the reports created by GLM. The most significant SNP is still the same, however the strength of association is weaker, with a P value of  $7.199 \times 10^{-4}$  (vs.  $1.1021 \times 10^{-4}$  from GLM) which does not pass the Bonferroni multiple test threshold (5x10<sup>-4</sup>).

The third report contains the MLM specific statistics, including -2 Log Likelihood, genetic variance and residual variance components under different level of compression. These statistics are illustrated by the Chart function on the Result mode as follows.



In the example, 79 are included in the final analysis. When they are clustered into 44 groups, the -2 Log Likelihood reaches a minimum, which indicates the best model fit. The screening of SNPs was performed at this optimum compression level.

**Note**: When two or more individuals are clustered into one group, the variance component for the random effect is not equivalent to the one without compression. Consequently, the heritability derived should not be interpreted as the individual based heritability.

To perform a Genome-Wide Association Study (GWAS) on the 3093 SNPs, we need to create a new joint data set containing the filtered phenotype, population structure, and the genome–wide genotype. Highlight the new joint file and the kinship data and click the **MLM** button. Choose the default options on the MLM option dialog. The analysis will take a minute or two. The output report labeled "MLM\_compression" indicates that 259 lines were used in the analysis. With 74 groups, the statistics from the best are as graphed below.



The strongest associated SNP is at 193565357 bp on chromosome 3. The P value is  $1.3027 \times 10^{-4}$ . The threshold is  $3.2331 \times 10^{-5}$  at significant level of 1% after Bonferroni multiple test correction (0.01/3093). The association was not significant. As illustrated below, the output labeled "GLM\_Allele\_Estimates" shows the marker effects assigned to genotypes for each SNP (The GLM is also the same). For example, the first SNP at 157104 bp on chromosome 1 had three genotypes (AA, CC and AC) coded as A, C, and M based on the IUPAC code, see Appendix (Nucleotide Codes).

📓 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 3.0.39									
File Tools Help <u>G</u> DPC									
CARANG Data	Results	Delete	Wizard		Show Memory			0% 🖹 📙 🝸	
		🎹 Table 🛛 📉 Tr	ee Plot	2D Plot	LD Plot	Chart			
Tree •	Trait	Marker	Locus	Site	Allele	Effect	Obs		
Fusions	dpoll	PZB00859.1	1	157104	C	3.64912	197		
Synonymizer	dpoll	PZB00859.1	1	157104	A	3.60484	53		
Nesult	dpoll	PZB00859.1	1	157104	M	0	3		
Diversity	dpoll	PZA01271.1	1	1947984	С	-1.2325E0	121		
SNP Assays	dpoll	PZA01271.1	1	1947984	G	0	127		
• D	dpoll	PZA03613.2	1	2914066	G	0.22634	75		
Association	dpoll	PZA03613.2	1	2914066	т	0	180		
MLM_statistics_for_Filtered_n	dpoll	PZA03613.1	1	2914171	т	5.20917	195		
MLM_compression_for_Filtere	dpoll	PZA03613.1	1	2914171	A	6.46891	61		
MLM_statistics_for_Filtered_n =	dpoll	PZA03613.1	1	2914171	w	0	2		
MLM_compression_for_Filtere	dpoll	PZA03614.2	1	2915078	G	-1.2702E-1	125		
MLM_statistics_for_Filtered_n	dpoll	PZA03614.2	1	2915078	A	0	118		
MLM_effects_for_Filtered_md	dpoll	PZA03614.1	1	2915242	т	0.55196	130		
MLM compression for Filtere	dpoll	PZA03614.1	1	2915242	A	0	111		
•	dpoll	PZA00258.3	1	2973508	G	-2.7856E0	65		
Table Title: MLM effects ^	dpoll	PZA00258.3	1	2973508	С	-3.3585E0	175		
Number of columns: 7	dpoll	PZA00258.3	1	2973508	S	0	2		
Number of columns. /	dpoll	PZA02962.13	1	3205252	т	-4.1007E0	218		
۲ ( III ) ۲	dpoll	PZA02962.13	1	3205252	A	-3.2237E0	26		
	dpoll	PZA02962.13	1	3205252	w	0	3		
Load 📻 🚔	dpoll	PZA02962.14	1	3205262	C	-2.0366E-1	224		
	dpoll	PZA02962.14	1	3205262	G	0	15		
Export	dpoll	PZA00599.25	1	3206090	С	0.61187	26		
	dpoll	PZA00599.25	1	3206090	т	-1.2877E-1	231		
Transform -	dpoll	PZA00599.25	1	3206090	Y	0	1		
	dpoll	PZA02129.1	1	3706018	т	0.59304	124		
4 III >	dpoll	PZA02129.1	1	3706018	С	0	131	-	
Program Status									

## 6.6 Importing Data from a Database (via GDPC)

GDPC, middleware that is integrated into TASSEL, allows the user to import data from a database. To display GDPC in TASSEL, click on the **GDPC** button in **Data** mode. General rules for working with databases include: 1) Establish a connection with the database; 2) Define a query; 3) once the desired data is in GDPC, load the data from GDPC into TASSEL.

#### 6.6.1 Connecting with a Database

To establish a connection with a database, click the **Add Conn** button followed by the button of the database you wish to add. Then click **Ok**. In the example below, we chose Panzea.



To connect to more than one database, simply repeat the process outlined above.

In the figures of following sections, only the GDPC area will be displayed if other areas are deemed irrelevant.

## 6.6.2 Data Query

GDPC is equipped with several tabs to query data, namely **Taxa**, **Taxon Parents**, **Loci**, **Genotype Experiments**, **Environment Experiment**, and **Localities**. Within each tab, any retrieved data will be displayed in the "Filtered List." Choose attributes by checking the desired boxes (located beneath the Filtered List). After an attribute is selected, values of that attribute from the database are displayed.

Here, using the Taxa tab, choose **Germplasm type** (field) and then select. After clicking the **Get Data** button, the subset of taxa from the database that meets these criteria will appear in both the **Filtered List** and the **Working List**.

-							_
	🔄 🖫 🏦 Open Save As Export	🏀 Get Data Lo	තා දී oad Desel	ect All Add Cor	nn Load	Get Star	ted
	Taxa Loci Genotype Experiments En	vironment Exp	periments Loo	alities Genotypes	Phenotypes	Connections L	Log
	Filtered List (5315) 🛛 🛱		Wor	king List (5315)	100 500	<b>1</b> 12	Pr
	<ul> <li>→-Accession</li> <li>→-33-16</li> <li>→-38-11</li> <li>→-4226</li> <li>→-4722</li> <li>→-81-1</li> <li>→-A158</li> <li>→-A158</li> <li>→-A188</li> <li>→-A214N</li> <li>→-A239</li> <li>→-A272</li> <li>→-A441-5</li> <li>→-A554</li> <li>→-A554</li> <li>→-A556</li> <li>→-A619</li> <li>→-A632</li> <li>→-A634</li> </ul>			Coession 33-16 33-16 4226 4722 4722 4722 4758 4158 4158 4158 4214N 4239 4239 4272 4441-5 4556 4556 46 4619 4632 4634			<u>Aco</u>
	<ul> <li>Genus</li> <li>Species</li> <li><u>Subspecies</u></li> <li>✓ Germplasm Type</li> <li>Accession</li> <li>Accession Number</li> <li>Source</li> </ul>	Locality (w	vorking list)	Germplasm T BC hybrid Hybrid Inbred Open pollinated Teosinte Teosinte inbred Tripsacum	ype i		

Items listed in the **Working List** can be modified by the user. To do so, first break the link between the Filtered List and the Working List by clicking on the Link/Unlink button . The button will now appear as . This activates the Add selected items , Add all items , Remove . Remove

and **Remove all** items buttons. Remove all items from the working list, then select items with a name starting with the letter D. Click on the Add selected items button to move them to the Working List. The resulting Working List is shown as follows:

Taxa Taxon Parents Loci Genoty	pe Experiments Environme	ent Experiments Studies Localities Gen
Filtered List (10531) 🛛  🕅		Working List (86) 🛛 🕸 👘
⊕-DE2(1)     ⊕-DE3(1)     ⊕-DE3(1)     ⊕-DE41(11)     ⊕-DE_3(7)     ⊕-DF_3(7)     ⊕-DF(3)(7)     ⊕-DK78002A(1)     ⊕-DK78002A(1)     ⊕-DK78002A(1)     ⊕-DK78004(1)		ame (25) = D02 (1) = D146 (1) = D940Y (27) = DE1 (11) = DE2 (1) = DE3 (1) = DE3 (7)
Accession	Locality (working list)	Germplasm Type
Genus		Hybrid
Species		Inbred developed from wild
Subspecies		Natural Segregating Population
🗹 Germplasm Type		Traditional cultivar/landrace
Accession Number		Wild
Source		
Locality		
Collector		

To filter data by polymorphism type, first click on the **Genotype Experiments** tab, check the **Polymorphism Type** and **Producer** checkbox (field), and then select **SNP** and **Jim**. Finally, click the **Get Data** button to reveal the subset of data that meets these criteria. Results for this example are shown below:

Environment Experiments	Studies Localities	Genotypes	Phenotypes	Connections [	_og
Таха	Taxon Parents		Loci	ľ	Genotype
Filtered List (693)		Worki	ing List (693)		Proper Producer
<ul> <li>Name (693)</li> <li>PHM10225.15 (1)</li> <li>PHM10321.11 (1)</li> <li>PHM10404.8 (1)</li> <li>PHM10525.11 (1)</li> <li>PHM10525.9 (1)</li> <li>PHM10750.26 (1)</li> <li>PHM10750.21 (1)</li> </ul>		Produce	er (1) Holland (693)	* *	
1 Debeweenbiews Trace				<b>&gt;</b>	
<ul> <li>Polymorphism Type</li> <li>Name</li> <li>Locus</li> </ul>		st) Poly Cat Len SNF	rmorphism Type egorical Igth	Genaissand Goodman	:e 🔺
String Representation		Sec	uence	Jim Holland	
Source Experiment				McMullen	
Producer				Pioneer	=
Align Program				Sequenom	
Primer List				Ware	▼
Reference Sequence					

Genotype data can be extracted from the database by clicking on the **Genotypes** tab, followed by the **Get Data** button. After a moment, genotype data will be displayed as follows:

Environment Experir	nents Stud	ies Localities	Genotypes	Phenoty	pes Con	nections	Log				
Таха		Taxon Parents		Loc	i		Genotyp	e Experiments			
Merge											
	PZA0309	7.4 PZA03092.7	PZA03090.31	PZA03083.7	PZA03067.	17 PZA030	064.6 PZA03063.	18 PZA03062.7	7 PZA03047.1	2 PZA(	
DJ7	T:C	T:T	T:T	T:T	C:C	A:A	T:T	C:C	A:A	G:G	٠
DK2MA22	C:C	T:T	T:T	T:T	T:T	G:G	T:T	T:T	G:G	G:G	_
DK4676A	T:T	T:C	T:G	T:T	T:C	G:G	T:C	C:C	A:G	A:G	
DK78002A	T:C	T:T	T:T	T:T	C:C	G:G	T:T	T:T	A:A	G:G	
DK78004	T:T	T:C	T:G	T:T	T:C	-:-	C:C	T:C	-:-	-:-	
DK78010	T:C	T:T	T:T	T:T	C:C	G:G	T:T	T:T	A:A	G:G	
DK78371A	T:T	0:0	T:G	-:-	T:T	G:G	T:T	0:0	G:G	A:A	•
	<u> </u>										
Taxa (working list	)	Genotype Experi	ments (workin	ıg list)							
D02	A PH	IM10225.15			<b>A</b>						
D146	= PF	IM10321.11									
D940Y	PH	M10404.8									
D940Y	DI	IM10525 11									
D040V	DL	JM10525.0									
D9401	PI	IM 10323.9									
D940Y	Pr	1M10750.20									
D940Y	Ph	IM11000.21									
D940Y	PH	M11114.10									
D940Y	Pł	IM11114.7									
D940Y	👻 Pł	IM112.8			-						

Users can either save this genotype data in several formats or upload it to TASSEL. However, before outlining these procedures, let us finish the query by exploring phenotypes. To get data from experiments conducted in 2000, first select the **Environment Experiments** tab, followed by the **Repetition** checkbox.

Select the desired repetitions in 2000 as the values to be used for filtering, then click the **Get Data** button. The subset of data that meets these criteria is returned as follows:

Environment Experime	ents Studies	Localities	Genotypes	Phenotypes	Connections	L
Таха	Ta	xon Parents		Loci		
Filtered List (369	) 🕼 🎼		Work	ing List (369)		N
- 38-11 Row_2000_W - 38-11 Row_2000_St - 38-11 Row_2000_St - 38-11 Row_2000_St - A272 Row_2000_St - A272 Row_2000_Wit - A272 Row_2000_St - A272 Row_2000_St - A441-5 Row_200_St - A441-5 Row_2000_St - A441-5 Row_2000_St - A441-5 Row_2000_St - A441-5 Row_2000_St - A441-5 Row_2000_St - A441-5 Row_2000_St - A441-5 Row_200_St - A441-5 Row_200_St - A441-5 Row_20	inter ummer_B10 ummer_C2 mmer nter mmer_C2 summer_B10		P- Locality	(3) 'arms (99) ton (186) st Lafayette (84)		<b>A</b>
4			4		1	
Locality	Loc	ality (working l	ist) Rep	etition		
Evaluation Site			2	5_vviitei	<b>^</b>	
Repetition			200	0_Summer		
Block	=		200	0_Summer_4F8	3 -	
Plot			200	0_Summer_B1	0 =	
Plant Date			200	0_Winter		
Coordinate X			200	1_Summer_F4/	A_Re	
Coordinate Y			200	1_Summer_F4/	A_R€ ∧ R¢▼	
String Representation	on 🗸 🚛		•	III		

Now extract phenotype data by clicking on the **Phenotypes** tab. Traits can only be extracted one at a time. Choose **Days to Silk** from the **Ontology** field. Make sure no Taxa are selected and all Environment Experiments are selected that were retrieved in the previous step. Click the **Get Data** button, then the **Merge** button, leaving only **Accession** checked under the Taxa Properties section. Leave **Locality** and **Repetition** checked under the Environment Experiments Properties section. Data are merged as follows:

Environment Expe	riments	Studies	Localities	7	Genotypes	Phenotyp	es Co
Таха		Ta	xon Parents			Loci	
🕼 Merge	Merge Fu	nction:	Concat 🔻				
	2	7 Farms	Clayton_20.	. C	layton_20	Clayton_20	
38-11	65	5.0	92.0			75.0	
A272	53	3.0				84.3	
A441-5	68	3.0	92.0			73.0	
A554	52	2.0		62	2.0		
A6	78	3.0	104.0			87.25	
A619	54	k.0				67.0	
A632	55	5.0				67.0	
B103	53	3.0	73.0	74	1.0		
B104	67	'.O				71.0	
Ontology	1	axa (worki	ing list)		Environmer	nt Experiment	t
Cupules per rank	<b>A</b>	D02		-	38-11 Row	2000_Winte	<b></b>
Days to Pollen		D146		=	38-11 Row	2000 Sumr	=
Days to Silk		D940Y			38-11 Row	2000 Sumr	
Days to Silk		D940Y			38-11 Row	2000 Sumr	
Days to Tassel		D940Y			A272 Row	2000 Sumn	
Ear Diameter		D940Y			A272 Row	2000_Winte	
Ear Height		D940Y			A272 Row	2000_Sumn	
Ear Length		D940Y			A441-5 Rov	v_2000_Sun	
Ear Number	<b>_</b>	D940Y			A441-5 Rov	v_2000_Win	-
		D940Y		•	•		

#### 6.6.3 Importing GDPC data into TASSEL

Genotype and phenotype data must be loaded in separate steps. To load genotype data, first click on **GDPC** in **Data** mode. Then click on the **Genotypes** tab, followed by the **Load** button. The genotype data

is then loaded into TASSEL and labeled as "Genotype." To view the uploaded data, click on "Genotype" within the Data folder in TASSEL. Results will look as follows:

📓 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 2.1																	
File Tools Help <u>G</u> DPC																	
crowc crowc Data	Results Delete						C		0%	%				<b>?</b>			
🐏 🗸 GDPC 📃 File	🍸 Sites 🕨 Taxa 🗛:a Genot	ype	?+5 Tra	nsfor	m <b>p</b> +	+q Syn	onymi	zer	<b>ω</b> υ.	Join	٥	Join					
🛅 Data 🔨	Taxa	PZA	. PZA	. PZA	. PZA	. PZA	. PZA	PZA.	PZA	PZA	PZA	. PZA	. PZA	. PZA	PZA	. PZA	
Sequence		4	7	31	7	17	6	18	7	12	14	6	5	23	16	7	
😑 🚞 Polymorphisms	017	1018.	ID18.	ID18.	ID18.			ID18.	1018		ID18.		ID18.	ID18	ID18.	ID18.	
🖻 — 🫅 PZA03097.4.ID18566	DK2MA22	0.0	T.T	T.T	T.T	T.T	CIC.	TIT	T.T	CIC.	GIG	0.0	6.6	6.6	0.0	0.0	
I Genotype	DK4676A	TIT	TIC	TIG	T·T	TIC	6.6	TIC	0.0	4·G	4.6	GrG		6.6	TIC		-
Numerical	DK78002A	T:C	T:T	T:T	T:T	CC	G:G	T:T	T:T	A:A	G:G	CiC	G:G	G:G	C:C	A:A	-
Matrix	DK78004	TIT	T:C	T:G	T:T	T:C	-:-	CIC	T:C	-:-	-:-	CIC	-:-	G:G	TIC	A:A	-
Tree	DK78010	T:C	TIT	T:T	T:T	C:C	G:G	T:T	T:T	A:A	G:G	C:C	G:G	G:G	T:T	A:A	-
Fusions	DK78371A	T:T	C:C	T:G	-:-	T:T	G:G	T:T	C:C	G:G	A:A	G:G	G:G	G:G	C:C	G:G	-
Synonymizer	DKFAPW	T:T	T:T	T:T	T:T	C:C	G:G	T:T	C:C	A:A	G:G	G:G	G:G	G:G	C:C	G:G	-
Result	DKFBHJ	T:T	C:C	T:T	T:T	C:C	G:G	T:T	C:C	A:A	A:A	G:G	G:G	G:G	C:C	A:A	-
	DKHBA1	-:-	T:T	T:T	T:T	T:C	A:A	T:T	C:C	G:G	A:A	G:G	G:G	G:G	C:C	G:G	
	DKIBO14	C:C	C:C	T:T	T:T	C:C	A:A	T:T	T:T	G:G	A:A	C:C	G:G	G:G	C:C	G:G	
	DKIBO2	C:C	T:T	T:T	T:T	-:-	A:A	T:T	T:T	-:-	G:G	G:G	G:G	G:G	C:C	A:G	
	DKMBNA	T:T	C:C	-:-	T:T	T:T	-:-	-:-	-:-	G:G	A:A	-:-	G:G	-:-	C:C	A:A	
	DKMBST	T:T	C:C	T:T	T:T	-:-	-:-	T:T	C:C	-:-	G:G	-:-	G:G	-:-	C:C	G:G	
	DKMDF-13D	-:-	T:T	T:T	T:T	C:C	G:G	T:T	C:C	G:G	A:A	G:G	G:G	G:G	T:T	-:-	
	DKPB80	T:T	T:T	T:T	T:T	C:C	A:A	T:T	C:C	A:A	A:A	G:G	G:G	G:G	T:T	A:A	_
	<															<u>,</u>	

To load phenotype data from GDPC into TASSEL, first click on the **GDPC** button in **Data** mode. Then choose the **Phenotypes** tab, followed by the **Load** button. The phenotype data is then loaded into TASSEL and labeled as "4 traits/environ." To view the uploaded data, select "4 traits/environ" from the Phenotypes folder in TASSEL. Results will appear as follows:

📓 TASSEL (Trait Analysis by aSSociation, Evolution, and Linkage) 2.1					
File Tools Help GDPC					
Construction Const	Results Delete	pel ?+5 Transform P++ q Synonymizer @ ∪ Join ① ∩ Join			
Data Sequence Polymorphisms Genotype Wumerical Matrix Tree Fusions Synonymizer Result	Taxa           38-11           A272           A441-5           A554           A659           B103           B104           B104           B104           B103           B104           B105           B106           B107           C03           C1103           C1103           C1104           CM105           CM174           CM27           CML247           CM254           CML251           CM1261	DaysDaysDays           27 F Clay Clay           27 F Clay Clay           65.0         NaM           84.3         NaM           65.0         NaM           84.3         NaM           65.0         NaM           84.3         NaM           65.0         NaM           84.3         NaM           52.0         62.0           84.3         NaM           78.0         NaM           87.0         NaM           88.0         NaM           86.0         NaM           87.0         NaM           87.0         NaM           86.0         NaM           87.0         NaM           87.0         NaM           87.0         NaM           87.0         NaM           87.0         NaM      <			
Program Statur	× 207		×		

### 6.6.4 Saving GDPC Query Results

All query results, including both genotype and phenotype queries, can be saved as either Tab-delimited text files or XML files. Results are exported as tab-delimited text files by first choosing the **Query** Tab and then clicking on the **Export** button for the **Export**..., or by clicking the **Save As** button situations. Data in XML format to save results in XML format. Location and file name must be specified in both situations. Data in XML format can be imported back into GDPC by clicking on the **Open** button.

# 7 Appendix

# 7.1 Nucleotide Codes (Derived from IUPAC)

Code	Meaning
Α	A:A
С	C:C
G	G:G
Т	T:T
R	A:G
Y	C:T
S	C:G
W	A:T
К	G:T
М	A:C
+	+:+ (insertion homozygous)
0	+:-
-	-:- (deletion homozygous)
Ν	Unknown

## 7.2 TASSEL Tutorial Data sets

The data set contains 9 files and can be downloaded at: http://www.maizegenetics.net/tassel/docs/TASSELTutorialData3.zip

File#	File name	Туре	Format
1	d8_sequence.phy	Genotype	Phylip Alignment
2	mdp_genotype.hmp.txt	Genotype	Hapmap Alignment
3 4	mdp_genotype.flpjk.geno mdp_genotype.flpjk.map	Genotype	Flapjack Alignment
5 6	mdp_genotype.plk.ped mdp_genotype.plk.map	Genotype	Plink Alignment
7	mdp_kinship.txt	Kinship	Numerical square matrix
8	mdp_population_structure.txt	Population structure	Numerical trait data
9	mdp_traits.txt	Phenotype	Numerical trait data

File #1 is the sequence of dwarf8 gene with 2466 sites on 91 maize inbred lines. The data was described by the paper on the association between Dwarf8 and flowering time<sup>26</sup>.

File #2-6 are 3093 SNPs on 281 maize association inbred lines. The data was presented in three formats (Hapmap, Plink and Flapjack). The data was created by the PANZEA project funded by NSF. Details of the data can be found at <u>http://www.panzea.org</u>.

File #3 and 4 are in pair for the format of Flapjack.

File #5 and 6 are in pair for the format of Plink.

File #7 is kinship created by Yu et al.<sup>9</sup>.

File #8 is population structure of 282 maize inbred line $^{27}$ .

File #9 is phenotype on three traits, including flowering time, on 282 maize inbred lines<sup>9</sup>.

# 7.3 Biography of TASSEL

2001	First public release
December, 2004	Score-able SNP Extractor Updated Main Panel
February, 2005	StepClade update
March, 2005	Fixed handling of ?s and non-standard characters Added Sliding Haplotype functionality Changed LD Fisher's Exact p-value to use two-sided p-value
April, 2005	Added Ability to visualize sequence quality scores "Synonymize"/match taxa names between data sets GLM analysis improvements
June, 2005	Code change preventing large data sets from being shown in JTable Update of GDPC which allows automatic restoration of last data source connection
October, 2005	Data transformation utilities added K-Nearest Neighbor Data Imputation added
January, 2006	Association analysis with Mixed Linear Model Taxa name " synonymizer" added Basic heterozygosity handling added Many ease-of-use improvements
March, 2006	<ul> <li>Fixed problem loading genotype data</li> <li>Mixed Linear Model changes: <ul> <li>Output NaN if non-converged</li> <li>Fixed problem loading genotype data</li> <li>Detection of duplicate ID in kinship</li> <li>Correction on progressive bar with MLM</li> <li>Starting values of NaN from previous marker are no longer used</li> </ul> </li> </ul>
September, 2006	MLM: Significant speed improvement (~10x faster) GLM: Added User-defined F-tests, Output taxa or marker means
October, 2006	Principle Components Analysis
September 2007	Architecture restructure and pipeline version for advanced users
April 2008	Genetic marker data numerical transformation
June 2008	MLM implemented P3D algorithm, increased speed in order of magnitude

	of at least ten times.
May 2009	EMMA implemented
November 2009	TASSEL Version 3 release (redesigned for large genomic data and large samples)
April 2010	Compression of MLM implemented

## 7.4 Frequently Asked Questions

## 1. What do I do if TASSEL misbehaves?

A: TASSEL is an open source software project hosted on SourceForge and has a bug tracking list at <a href="http://sf.net/projects/tassel">http://sf.net/projects/tassel</a> where you can notify the developer community of problems. In order for a bug to be fixed, we must be able to replicate the problem. Thus, it is important to document the steps that were taken that produced the error. If the data you are working with is not too sensitive, please include the files which were used in the faulty procedure. If you would rather not post your data file on SourceForge, you may email it to one of the software developers.

## 2. Where do I turn for more information?

A: If you are having difficulty with a certain aspect of TASSEL, you can either email one of the software developers listed at <u>www.maizegenetics.net</u> or you may check the TASSEL forum on SourceForge <u>http://sf.net/projects/tassel</u>), as another user may have already addressed a similar question. There is also a TASSEL discussion group at http://groups.google.com/group/tassel.

## 3. How do I join the fun: TASSEL on SourceForge?

A: TASSEL is an open source project distributed under the GNU general public license. This means that the source code is available and the user is free to modify the code to suit their particular needs. We welcome input from developers and those who wish to become involved in the improvement of this software. The project is hosted on SourceForge (<u>http://sf.net/projects/tassel</u>), thereby allowing anyone to access the most recent changes to the code. This setup makes it convenient for anyone to add special functionality to TASSEL if they so desire. It also serves as a good platform for anyone who wishes to become involved in a bioinformatics software development project.

#### 4. How do I change the amount of memory used? What do I do when the "Exception java.lang.OutOfMemoryError" appears?

A: If you are working with very large data sets or are running memory intensive procedures, there may be occasions when TASSEL runs out of memory. For most routine usage, however, TASSEL memory is sufficient. Memory issues usually result from attempting to execute a procedure like LD on a raw sequence alignment instead of selected SNPs. You may also experience a memory issue if you are not sufficiently specific when retrieving information through GDPC. By default, TASSEL is allocated up to 512 Mb of memory on your computer. If more is available on your computer, you can increase the amount allocated by downloading the "stand-alone" version of TASSEL and opening a command line window (in Windows use Start > Run and type in "cmd" or "command"). To run TASSEL from a command line, "cd" to go to the directory containing the stand-alone jar file then start TASSEL by typing the following:

java -Xms256M -Xmx768M -jar sTASSEL.jar

Where "-Xms###M" specifies the starting memory available and "-Xmx###M" specifies the maximum memory available to the Java Virtual Machine. You may set the values higher or lower as your hardware dictates. Alternatively, you can modify the start\_tassel.bat or start\_tassel.pl file that comes with the standalone distribution.

# 5. When I click on the most current version of TASSEL web start, a previous version appears. What should I do?

A: The previous version of TASSEL web start was cached in your machine. To replace it with the most current version, click the Start button in Windows, followed by Run. Type **javaws** and then click OK. In the window that opens, keep the most current version of TASSEL and delete the rest.

#### 6. What should I substitute for missing values in TASSEL?

A: For numerical data in version 3 format, use NA or NaN. For numerical data in version 2 format, use "-999" for missing values. For SNP data, use "N". For SSR data, use "?". Kinship does not allow missing values.

#### 7. Is it possible to change data names in the Data Tree?

A: Yes. Click on the desired data name in the Data Tree, wait for one second, and then click it again or immediately hit the F2 key. Rename the data set and then hit Enter to save the change.

## 8. How can I create a TASSEL icon on desktop?

A: Click "Start" on Microsoft Windows and select "Control Panel", then double click Java to show "java Control Panel". In "Temporary Internet Files" section, click "View" button show "Java Cache Viewer". Move mouse over TASSEL application and click right button and select "Install Shortcuts".

#### 9. Why do I get empty squares in MLM association analysis?

A: The empty square means null information. The major reasons include non-convergence in the estimation of variance componentsor that the statistic in question was not calculated. For example, marker F, p, and R<sup>2</sup> are not calculated when no marker is included in the model.

## 10. Why should I exclude one column of the population structure?

A: For some methods of calculating population structure, such as the software STRUCTURE, the population proportions sum to one. This produces linear dependence between the population covariates. While the algorithm used by GLM tolerates that dependency, MLM will fail because the design matrix will not be invertibleExcluding one column eliminates linear dependence between columns. Using PC axes to represent population structure does not result in linear dependency because all PC columns are guaranteed to be independent.

#### 11. Can kinship replace population structure?

A: Sometimes. For some traits and populations, the K-only model may be as good as or better than the Q+K model. For others, Q+K may be superior. The Q-only model is not as effective for controlling population structure as the alternatives. Unfortunately, no general guidelines exist for predicting which model will perform best. As a result, an investigator may wish to fit all three models and compare the results. If eliminating false positives is very important, then it may make sense to accept the most conservative model. However, if the objective is to identify candidates for further study and the cost of following up on a false lead is low, the most liberal model may be preferred.

#### 12. Why do TASSEL and SPAGeDi give different kinship estimates?

A: First, many algorithms exist to calculate kinship and their estimates will differ from one another. Secondly, the algorithm in TASSEL treats each genotype as a haplotype. It is not recommended that TASSEL be used to generate a kinship matrix from heterozygous genotype. In the near future, the TASSEL kinship algorithm will be modified to handle heterozygous diploids.

### 13. Can I get Marker R square using SAS Proc Mixed or TASSEL MLM?

A: SAS Proc Mixed does not produce an  $R^2$  statistic. MLM in TASSEL does. The user manual describes how it is calculated.

#### 14. Does MLM find more associations than GLM?

A: Sometimes. MLM has higher statistical power than GLM and may detect more true associations.. When the tested genetic markers are confounded with kinship structure, GLM does not correct for that as effectively as MLM and may produce more false positives

#### 15. Do I need multiple test correction for the p value from Tassel?

A: Yes.

#### 16. Can TASSEL handle diploid genotype data?

A: While TASSEL accepts most common sequence alignment formats which handle polyploid genotype data including haploid and diploid, some analyses are not appropriate for heterozygous data. GLM or MLM fit SNPs one at a time, treating each distinct genotype as a separate class. This has the effect of fitting an additive plus dominance model. Separating the two effects is under consideration. Because handling heterozygotes as a third marker class is not appropriate for kinship or LD those analyses should not be used for that type of data at the present time. Work to improve handling heterzygotes is ongoing.

#### 17. How to cite TASSEL?

A: The paper that describes  $TASSEL^1$  as a software package and the papers that introduce specific methods implemented in TASSEL should be cited as appropriate, such as the unified ("Q+K") approach, EMMA, compression of mixed linear model and P3D. For example,:

- A. Linkage disequilibrium (D',  $R^2$  and P value) were calculated by TASSEL<sup>1</sup>.
- B. Association analyses were performed with the mixed linear model approach<sup>9</sup> implemented by TASSEL<sup>1</sup>.
- C. GWAS was performed with the compressed mixed linear model approach<sup>4,9</sup> carried by TASSEL<sup>1</sup> which also implemented the EMMA<sup>3</sup> and P3D<sup>4</sup> algorithms to reduce computing time.

## REFERENCES

- 1. Bradbury, P.J. et al. TASSEL: software for association mapping of complex traits in diverse samples. *Bioinformatics* **23**, 2633-2635 (2007).
- Zhang, Z., Buckler, E.S., Casstevens, T.M. & Bradbury, P.J. Software engineering the mixed model for genome-wide association studies on large samples. *Brief Bioinform* **10**, 664-75 (2009).
- 3. Kang, H.M. et al. Efficient Control of Population Structure in Model Organism Association Mapping. *Genetics* **178**, 1709-1723 (2008).
- 4. Zhang, Z. et al. Mixed linear model approach adapted for genomewide association studies. *Nat Genet* **42**, 355-60 (2010).
- Kang, H.M. et al. Variance component model to account for sample structure in genome-wide association studies. *Nat Genet* 42, 348-54 (2010).
- 6. Thornsberry, J.M. et al. Dwarf8 polymorphisms associate with variation in flowering time. *Nature Genetics* **28**, 286-289 (2001).
- Pritchard, J.K., Stephens, M., Rosenberg, N.A. & Donnelly, P. Association mapping in structured populations. *American Journal* of Human Genetics 67, 170-181 (2000).
- 8. Zhao, K. et al. An Arabidopsis example of association mapping in structured samples. *PLoS Genet* **3**, e4 (2007).
- 9. Yu, J.M. et al. A unified mixed-model method for association mapping that accounts for multiple levels of relatedness. *Nature Genetics* **38**, 203-208 (2006).
- 10. Casstevens, T.M. & Buckler, E.S. GDPC: connecting researchers with multiple integrated data sources. *Bioinformatics* **20**, 2839-2840 (2004).
- 11. Ware, D. et al. Gramene: a resource for comparative grass genomics. *Nucleic Acids Research* **30**, 103-105 (2002).
- 12. Ware, D.H. et al. Gramene, a tool for grass Genomics. *Plant Physiology* **130**, 1606-1613 (2002).
- 13. Jaiswal, P. et al. Gramene: development and integration of trait and gene ontologies for rice. *Comparative and Functional Genomics* **3**, 132-136 (2002).

- 14. Yamazaki, Y. & Jaiswal, P. Biological ontologies in rice databases. An introduction to the activities in gramene and oryzabase. *Plant and Cell Physiology* **46**, 63-68 (2005).
- 15. Zhao, W. et al. Panzea: a database and resource for molecular and functional diversity in the maize genome. *Nucleic Acids Research* **34**, D752-D757 (2006).
- 16. Canaran, P., Stein, L. & Ware, D. Look-Align: an interactive webbased multiple sequence alignment viewer with polymorphism analysis support. *Bioinformatics* **22**, 885-886 (2006).
- 17. Du, C.G., Buckler, E. & Muse, S. Development of a maize molecular evolutionary genomic database. *Comparative and Functional Genomics* **4**, 246-249 (2003).
- 18. SAS, I.I. SAS. Statistical Analysis Software for Windows, 9.0 ed. *Cary, NC. USA.* (2002.).
- 19. Hardy, O.J. & Vekemans, X. SPAGEDi: a versatile computer program to analyse spatial genetic structure at the individual or population levels. *Molecular Ecology Notes* **2**, 618-620 (2002).
- 20. Cover, T. & Hart, P. Nearest neighbor pattern classification. *Proc IEEE Trans Inform Theory* **13**(1967).
- 21. Weir. Genetic Data Analysis II. Sunderland, MA. (1996).
- 22. Farnir, F. et al. Extensive genome-wide linkage disequilibrium in cattle. *Genome Res* **10**, 220-7 (2000).
- 23. Henderson, C.R. Best Linear Unbiased Estimation and Prediction under a Selection Model. *Biometrics* **31**, 423-447 (1975).
- 24. Kang, H.M. et al. Efficient control of population structure in model organism association mapping. *Genetics* **178**, 1709-23 (2008).
- 25. Laird, N.M. & Ware, J.H. Random-Effects Models for Longitudinal Data. *Biometrics* **38**, 963-974 (1982).
- 26. Thornsberry, J.M. et al. Dwarf8 polymorphisms associate with variation in flowering time. *Nat Genet* **28**, 286-9 (2001).
- 27. Flint-Garcia, S.A. et al. Maize association population: a highresolution platform for quantitative trait locus dissection. *Plant J* **44**, 1054-64 (2005).
- 28. Anderson, M.J. & Ter Braak, C.J.F. Permutations tests for multifactorial analysis of variance. Journal of Statistical Computation and Simulation 73, 85-113 (2003)
## INDEX

2D Plot, 34 Analysis Mode, 25 Annotated alignment, 14 BLOBs, 12 Chart, 36 Cladogram, 27 Collapse Non Major Alleles, 19 compressed MLM, 31 Compression, 31 compression level, 31 Data Mode, 10 data tree, 38 Diversity, 25 EM algorithm, 30 EMMA, 30 expectation and maximization algorithm, 30 File Menu, 38 Flapjack, 13 GDPC, 10, 54 General Linear Model, 28 Genome-Wide Association Study, 53 Genotype Numericalization, 18 GLM, 47 Hapmap, 12 Henderson. See MLM Heritability, 30 Impute Phenotype, 20 Impute SNPs, 18 Intersection Join, 24 Join, 23 Kinship, 15, 28, 30, 46

LD Plot, 35 Linkage Disequilibrium, 26 menus, 38 Mixed Linear Model, 30 MLM, 51 Numerical data, 14 Open source code, 8 P3D, 31 Panels, 8 PCA, 21 Plink, 12 Population parameters previously determined, 31 Principal component analysis, 21 Principal Component Analysis, 42 REML, 30 Restricted Maximum Likelihood, 30 Result Mode, 33 Sites, 16 SNP Extract, 27 Specified number of rows, columns, and labels. See Kinship Square Numerical Matrix, 15 Stand-alone, 8 Synonymize Taxa Names, 21 Table, 33 Taxa, 17 Traits, 17 Transform, 18 Tree Plot, 33 Union Join, 23 Web start, 7