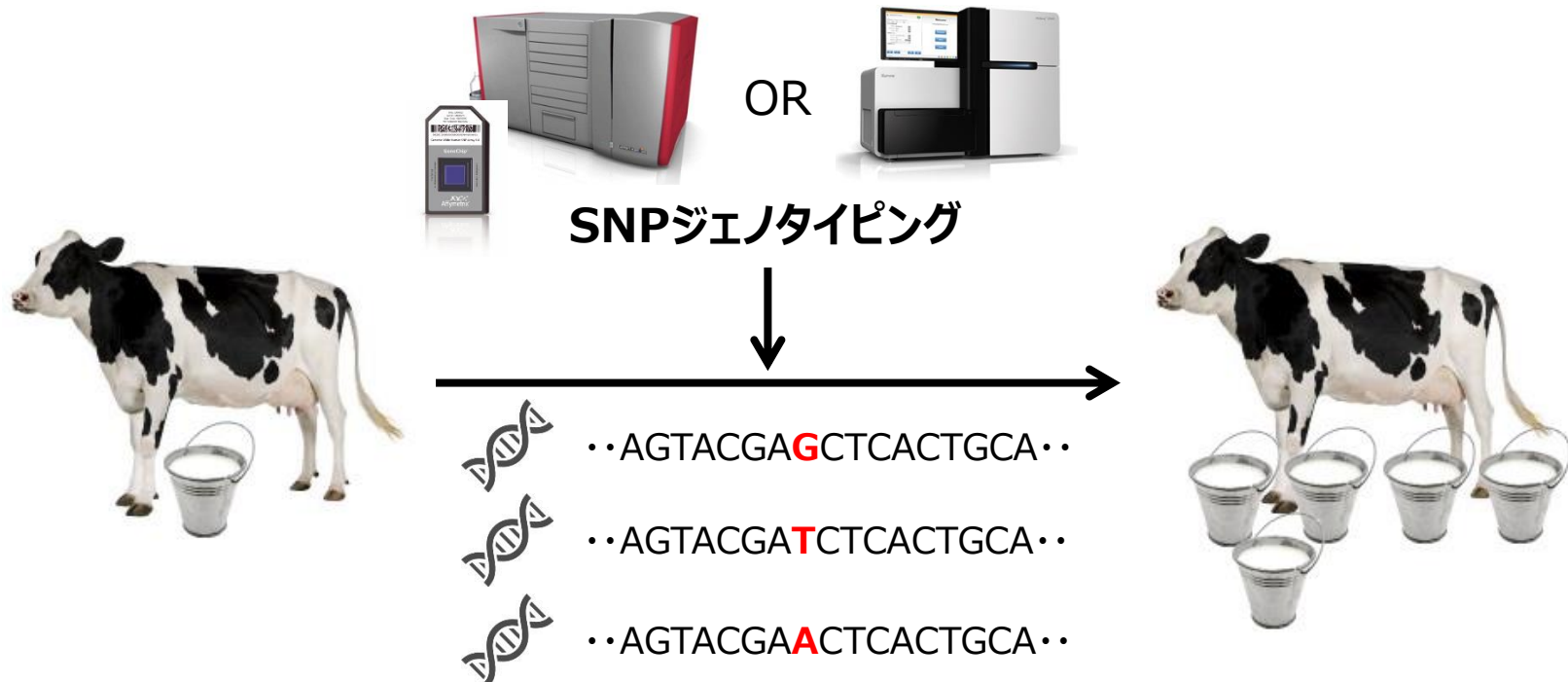


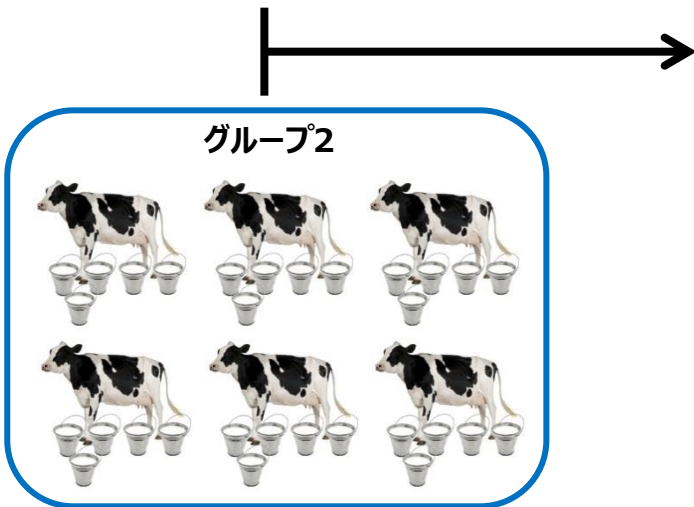
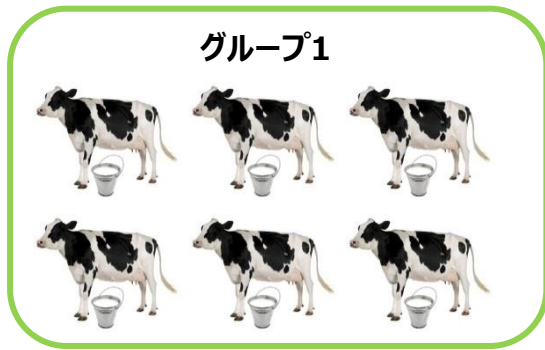
# GBS (Genotyping by Sequencing) によるアグリゲノミクス解析

フィルジエン株式会社 バイオサイエンス部  
(biosupport@filgen.jp)

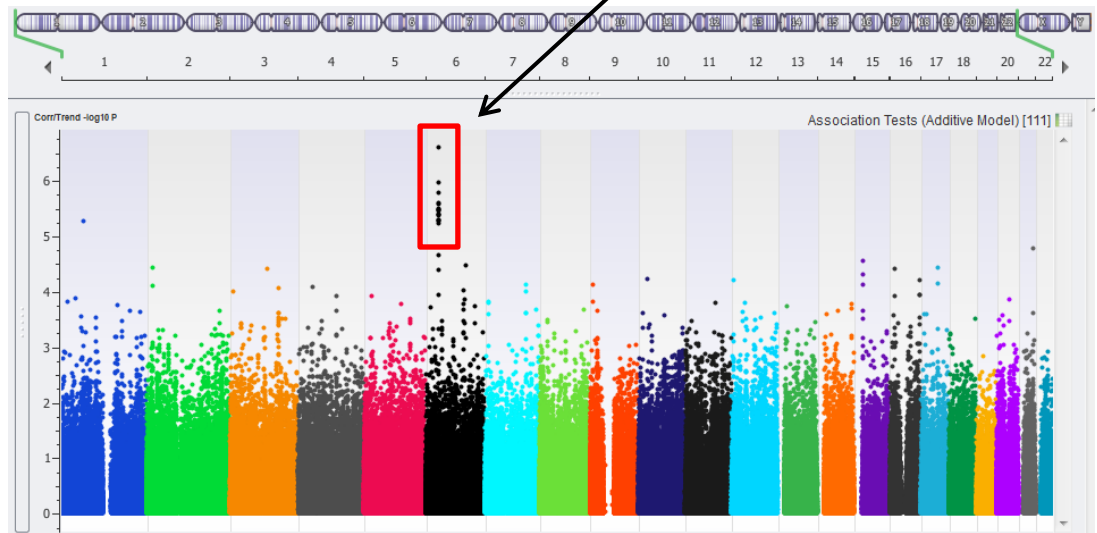
- 多数サンプルのSNPデータを統計的に比較し、表現型と関連するSNPを見つけ、農作物や家畜などの育種に役立てる。
- ハイスループットなSNPジェノタイピングには、従来はマイクロアレイが多く利用されていた。
- 現在では、コストの低下などにより、次世代シーケンサーを利用したDNAシーケンスによるジェノタイピング（Genotyping by Sequencing : GBS）も利用されるようになった。



- GBSデータから育種に有用なマーカーを見つけるために、**関連解析**という手法が使われる。
- 関連解析を行うことで、各サンプルのジェノタイプデータと表現型データを比較し、表現型のマーカーを見つけることができ、家畜や農作物の育種に応用できる。



表現型との関連が高いマーカー



## ゲノムワイド関連解析 (GWAS)

ゲノムレベルのデータ量を扱う、基本的な関連解析の手法。

- Chi-squared test
- Fisher's exact test
- Armitage test
- Correlation/Trend test
- Odd ratios
- Logistic or Linear regression ...etc

## 線形混合モデル解析

おもにサンプルの血縁関係による偏りを除外し、関連解析を行う場合に用いられる手法。育種分野で、近交系サンプルなどの解析に用いられる。

- Mixed Model GWAS using a single locus (EMMAX)
- Multi-locus mixed model GWAS (MLMM)
- Genomic Best Linear Unbiased Predictors (GBLUP)

- 次世代シーケンサーやマイクロアレイから作成された、大容量のジェノタイプデータや表現型データを扱うことができる。
- これらデータを処理するための、強力な統計解析機能が搭載されている。
- 一般的なモデル生物の他、多様な生物種のゲノムデータを扱うことができる。
- それら生物種の、各種アノテーションデータを利用できる。



## SNP & VARIATION SUITE

1. 各メーカー（Affymetrix, Illuminaなど）のSNP / CNV マイクロアレイデータ、および次世代シーケンサー解析のVCFファイルをインポート可能。

2. Golden Helix社のサーバーから、各種アノテーションデータを容易に取得可能。

3. 変異解析以外に、様々なアプリケーションを搭載。

- ゲノムワイド関連解析（GWAS）
- コピー数（CNV）解析
- 少数サンプルのNGS変異解析
- Genomic Prediction
- RNA-Seq解析
- 多数サンプルのNGS変異解析 ...など

4. 高度な統計学的計算アルゴリズムを多数搭載。

- カイ二乗検定、フィッシャー検定、トレンド検定
- ハーディ・ワインベルグ平衡の計算
- ハプロタイプブロックの検出
- 線形混合モデル(mixed linear model)
- CNV領域の検出および関連解析
- DE-SeqによるRNA-Seq発現解析
- 線形/ロジスティック回帰解析
- LD（連鎖不平衡）解析
- Runs of homozygosity (ROH)の検出
- Genomic BLUP (GBLUP)によるGenomic Prediction
- Collapsing Methodによるレアバリエント関連解析
- メタアナリシス ...など

5. 有償モジュールの追加により、家系情報に基づいた解析が可能。

6. 様々なグラフ表示機能。



# SNP & Variation Suite (SVS)

アノテーションデータ

Unsort		C 1	I 2	R 3	C 4	G 5	G 6	G 7	G 8	G 9	G 10
Map	sub	Phenotype 1	SBP	BMI	Gender	SNP_A-1909444	SNP_A-4303947	SNP_A-1886933	SNP_A-2116190	SNP_A-4291020	SNP_A-1902458
	Chromosome										
	Position					752566	779322	785989	1003629	1097335	1130727
	dbSNP RS ID					rs3094315	rs4040617	rs2980300	rs4075116	rs9442385	rs10907175
	Associated Gene					?	?	?	?	?	TTL10
	Cytoband					p36.33	p36.33	p36.33	p36.33	p36.33	p36.33
	Reference Alleles A/B					[C/T]	[A/G]	[A/G]	[A/G]	[G/T]	[A/C]
	Top Alleles					[G/A]	[A/G]	[T/C]	[T/C]	[G/T]	[A/C]
	Bottom Alleles					[C/T]	[T/C]	[A/G]	[A/G]	[C/A]	[T/G]
	Strand					-	+	-	-	+	+
	Strand Versus dbSNP					same	same	reverse	same	same	same
1	GSM233256_GSM233257	Case	128		M	T_T	A_A	G_G	A_A	G_G	A_A
2	GSM233258_GSM233259	Control	137		M	C_T	A_A	G_G	A_A	G_G	A_A
3	GSM233262_GSM233263	Control	117		M	C_T	A_G	A_G	A_G	G_T	A_A
4	GSM233264_GSM233265	Control	140		M	T_T	A_A	G_G	A_A	G_G	A_C
5	GSM233266_GSM233267	Case	113		M	T_T	A_A	G_G	G_G	G_G	A_A
6	GSM233270_GSM233271	Case	125		M	T_T	A_A	G_G	A_A	G_T	A_A
7	GSM233276_GSM233277	Control	88	20	M	T_T	A_A	G_G	A_A	G_T	A_A
8	GSM233278_GSM233279	Case	131	31.2	M	T_T	A_A	G_G	A_A	G_G	A_A
9	GSM233280_GSM233281	Case	124	26.2	M	T_T	A_A	A_G	A_A	G_T	A_A
10	GSM233284_GSM233285	Control	111	22.8	M	C_T	A_G	A_G	A_G	G_G	A_A
11	GSM233286_GSM233287	Control	114	26.1	F	T_T	A_A	G_G	A_A	G_G	A_A
12	GSM233288_GSM233289	Control	107	23.5	M	C_T	A_G	A_G	A_G	G_G	A_C
13	GSM233290_GSM233291	Case	146	37.4	M	C_C	?_?	A_A	A_G	G_G	A_C
14	GSM233292_GSM233293	Case	122	27.6	M	T_T	A_A	G_G	A_G	G_G	A_A
15	GSM233298_GSM233299	Case	130	24.4	F	T_T	A_A	G_G	A_A	G_T	A_A
16	GSM233300_GSM233301	Control	116	24.4	M	T_T	A_A	G_G	A_A	G_G	A_A
17	GSM233302_GSM233303	Control	125	30.5	M	T_T	A_A	G_G	A_A	?_?	A_A
18	GSM233304_GSM233305	Case	112	22.2	M	T_T	A_A	G_G	A_G	G_G	A_C
19	GSM233306_GSM233307	Case	130	29.7	M	T_T	A_A	G_G	A_G	G_G	A_A
20	GSM233308_GSM233309	Control	136	30.7	F	C_T	A_G	A_G	A_A	G_G	A_A
21	GSM233310_GSM233311	Case	104	28	M	C_C	?_?	A_A	A_G	G_G	A_A
22	GSM233314_GSM233315	Case	115	26.1	M	T_T	A_A	G_G	A_G	G_G	A_A
23	GSM233316_GSM233317	Case	143	33.6	M	T_T	A_A	G_G	A_G	G_G	A_A
24	GSM233320_GSM233321	Case	103	27.7	M	T_T	A_A	?_?	A_A	G_G	A_C
25	GSM233322_GSM233323	Control	144	33.1	F	T_T	A_A	G_G	A_A	G_T	A_A
26	GSM233324_GSM233325	Case	131	32.3	M	T_T	A_A	G_G	A_A	G_G	A_A

表現型データ

ジェノタイプデータ

- サンプルの表現型データとジェノタイプデータを統合表示し、各種データ解析を行う。

Genotype Association Tests

Case/Control dependent variable: Phenotype 1 - Binary (238 cases and 230 controls) Additive model: (dd) -> (Dd) -> (DD)

Classify alleles by allele frequency  Classify alleles by reference/alternate (Reference field in map: "Reference Alleles A/B")

Association Test Parameters **PCA Parameters** Overall Marker Statistics

Genetic Model or Tests Test Statistic or Method

Where D = minor allele, d = major allele

Basic allelic tests: D vs. d  
 Genotypic tests: (DD) vs. (dd) vs. (Dd)  
 Additive model: (dd) -> (Dd) -> (DD)  
 Dominant model: (DD, Dd) vs. (dd)  
 Recessive model: (DD) vs. (Dd, dd)

Correlation/Trend test  
 Cochran-Armitage test  
 Exact form of Cochran-Armitage test  
 Odds ratios: (Dd) vs. (dd) and (DD) vs. (Dd)  
 Logistic regression

Case/Control additive model:

Missing Values Multiple Testing Correction

Use missing values as predictors  
 Drop missing values

Bonferroni adjustment (on N SNPs)  
 False Discovery Rate (FDR)  
 Single value permutations  
 Full scan permutations

Additional Outputs

Output data for P-P/Q-Q plots  
 Output -log<sub>10</sub>(P)

Number of permutations:

Principal Components Analysis (PCA) Genomic Control of Output Data for Stratification

Correct for stratification with PCA

Show inflation factor (lambda), chi-squares, and corrected values  
 Correct using this inflation factor (lambda) instead:

Help Restore Options Save Options Run Cancel

Genotype Filtering by Marker

Case/Control dependent variable: Phenotype 1 - Binary (237 cases and 227 controls)

Classify alleles by allele frequency  Classify alleles by reference/alternate (Reference field in map: "Reference Alleles A/B")

Filter Genotype Columns

General Statistics Filtering

Drop if call rate   
 Drop if number of alleles   
 Drop if Minor Allele Frequency (MAF)   
 Drop if carrier count

Hardy Weinberg Equilibrium (HWE) Filtering

Perform HWE filtering based on:

Drop if Hardy Weinberg Equilibrium (HWE) P-Value   
 Drop if Fisher's exact test for HWE P-Value   
 Drop if signed HWE R (positive if more homozygous)

Actions

Inactivate genotype columns that meet above criteria for filtering  
 Output spreadsheet with marker statistics and 'Drop?' columns

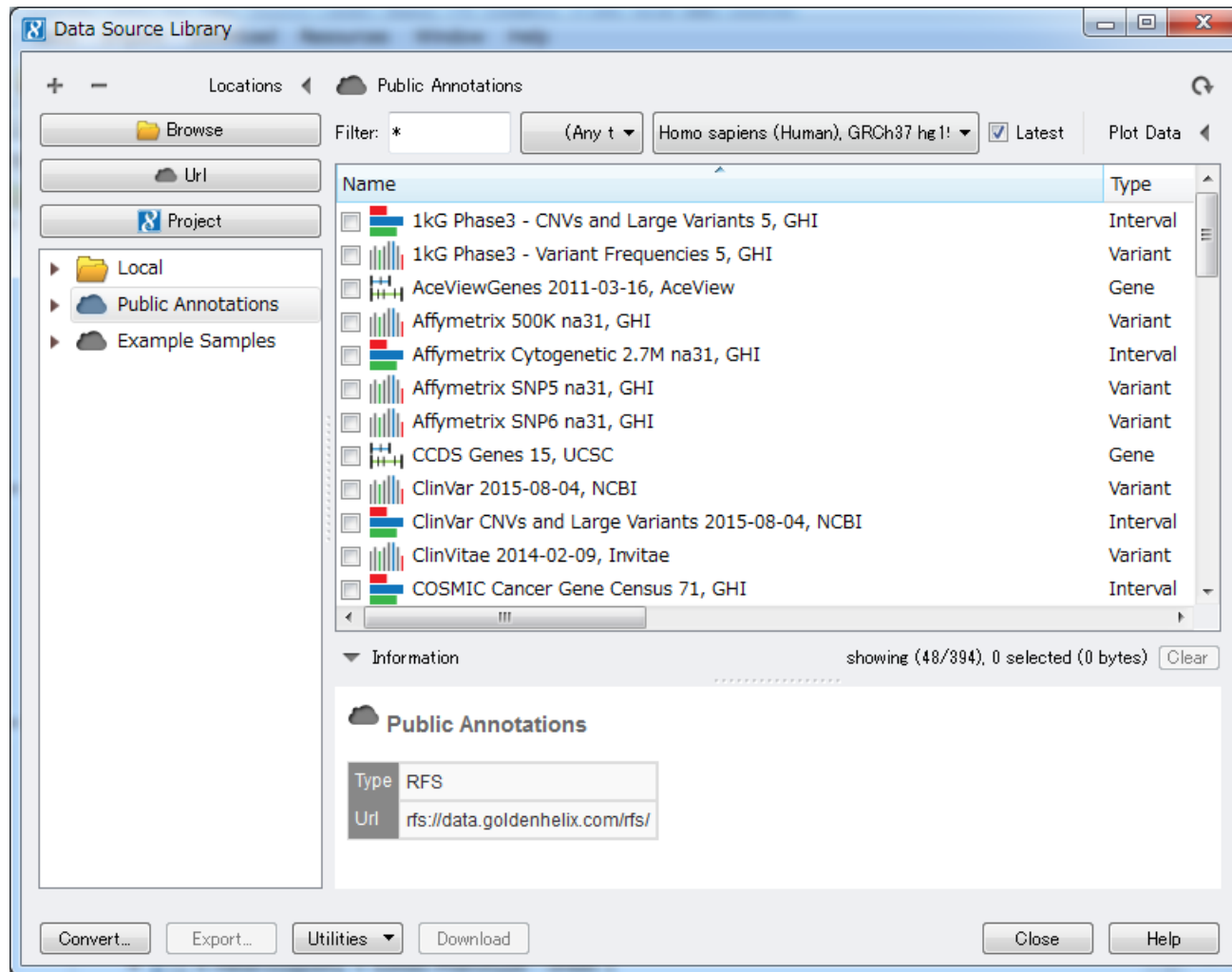
Additional Output

Output -log<sub>10</sub>(Value)

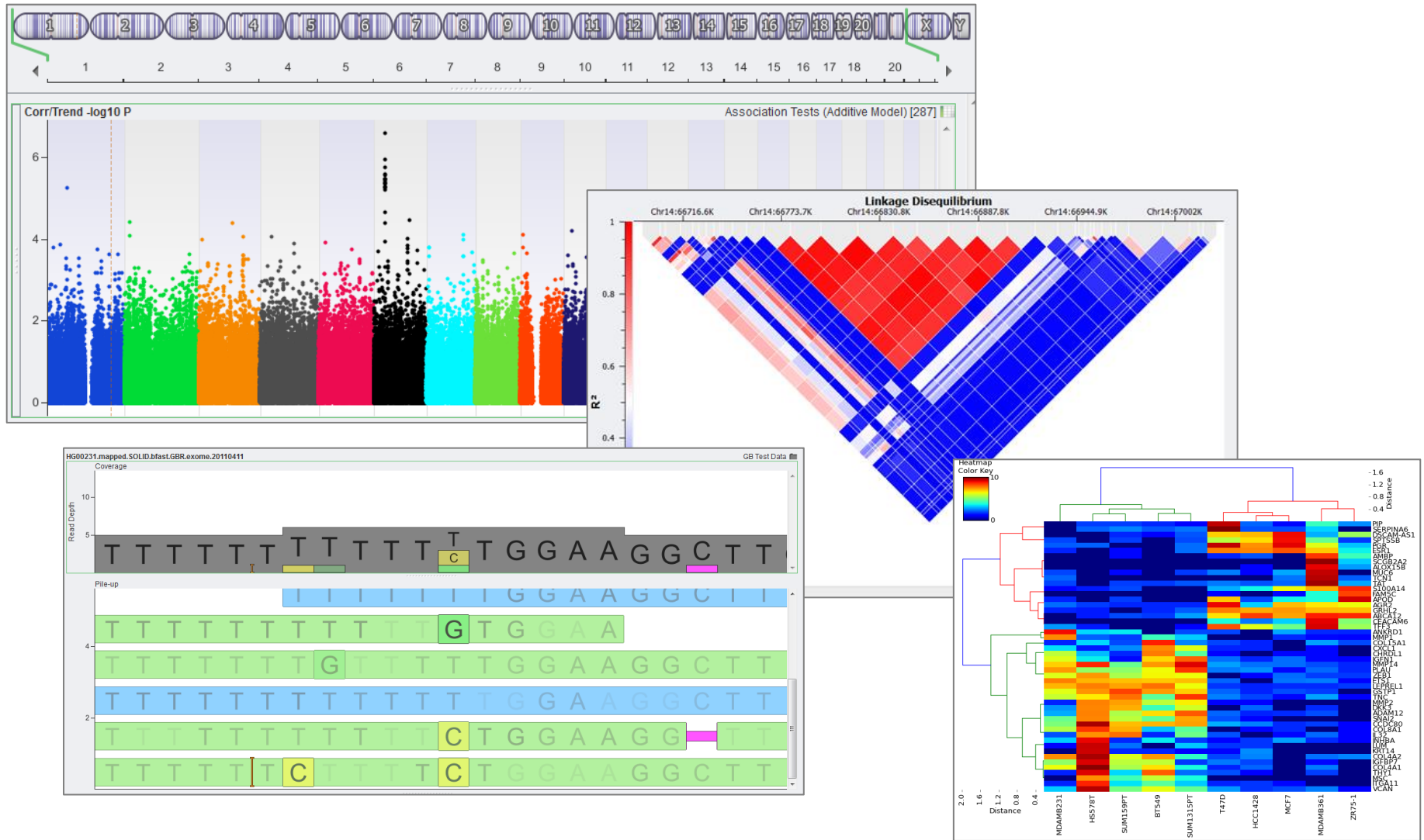
Help Restore Options Save Options Run Cancel

- 遺伝統計学の様々な解析アルゴリズムを搭載。





- 専用のデータ管理ツールを使用し、様々なデータソースのアノテーションデータを、容易にデータ解析に使用が可能。



- 解析データの様々なグラフ表示が可能。

## 対応生物種一覧

### 【哺乳動物】

- *Bos taurus* (ウシ)
- *Canis familiaris* (イヌ)
- *Capra hircus* (ヤギ)
- *Cricetulus griseus* (チャイニーズハムスター)
- *Equus caballus* (ウマ)
- *Felis catus* (ネコ)
- *Gallus gallus* (ニワトリ)
- *Heterocephalus glaber* (ハダカデバネズミ)
- *Homo sapiens* (ヒト)
- *Macaca mulatta* (アカゲザル)
- *Mus musculus* (マウス)
- *Nomascus leucogenys* (ホロジロテナガザル)
- *Ovis aries* (ヒツジ)
- *Rattus norvegicus* (ラット)
- *Sus scrofa* (ブタ)
- *Vicugna pacos* (アルパカ)

### 【魚類】

- *Danio retio* (ゼブラフィッシュ)
- *Medaka* (メダカ)
- *Oncorhynchus mykiss* (ニジマス)

### 【植物】

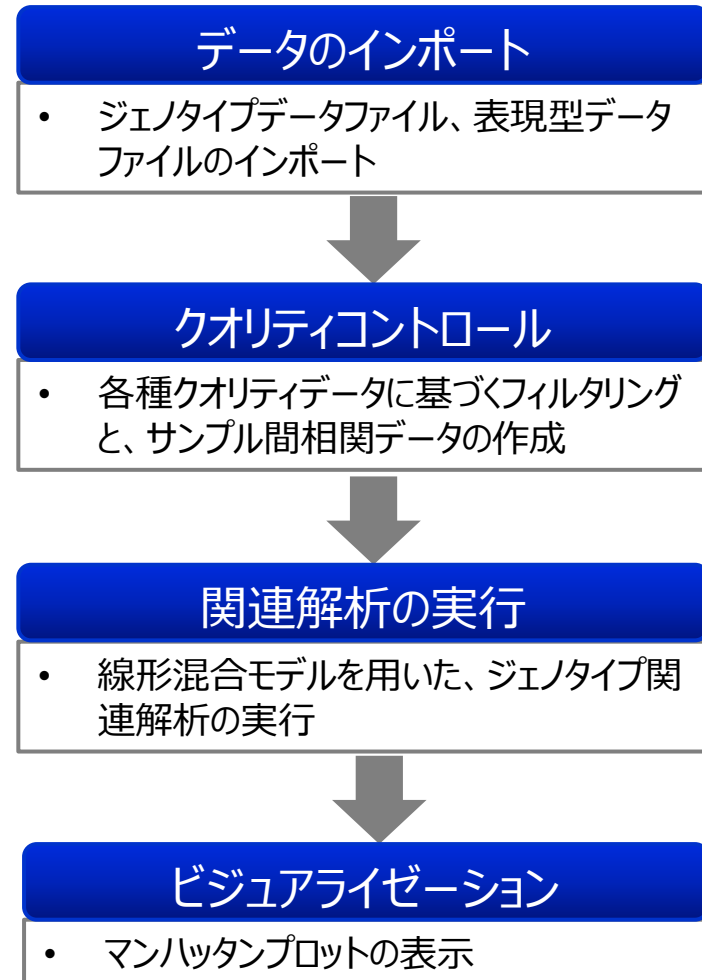
- *Arabidopsis thaliana* (シロイヌナズナ)
- *Brassica rapa* (アブラナ)
- *Capsicum annuum* (トウガラシ)
- *Carica papaya* (パパイヤ)
- *Citrullus lanatus* (スイカ)
- *Eucalyptus grandis* (ユーカリ)
- *Glycine max* (ダイズ)
- *Gossypium raimondii* (ワタ)
- *Oryza sativa* (イネ)
- *Prunus persica* (モモ)
- *Setaria italica* (アワ)
- *Solanum lycopersicum* (トマト)
- *Solanum tuberosum* (ジャガイモ)
- *Sorghum bicolor* (モロコシ)
- *Zea mays* (トウモロコシ)

### 【その他】

- *Anopheles gambiae* (ハマダラカ)
- *Caenorhabditis elegans* (線虫)
- *Drosophila melanogaster* (ショウジョウバエ)
- *E. coli* (大腸菌)
- *Leishmania infantum* JPCM5 (リーシュマニア寄生虫)
- *Mycobacterium tuberculosis* H37Rv (結核菌)
- *Plasmodium falciparum* 3D7 (マラリア)
- *Saccharomyces cerevisiae* (出芽酵母)
- *Schizosaccharomyces pombe* (分裂酵母)
- *Staphylococcus aureus* (黄色ブドウ球菌)

## 使用するジェノタイプデータ:

- 生物種 : トウモロコシ (Zea mays)
- NGSサンプルデータ数: 281例
- ファイルフォーマット : VCFファイル
- 変異数: 3,096個



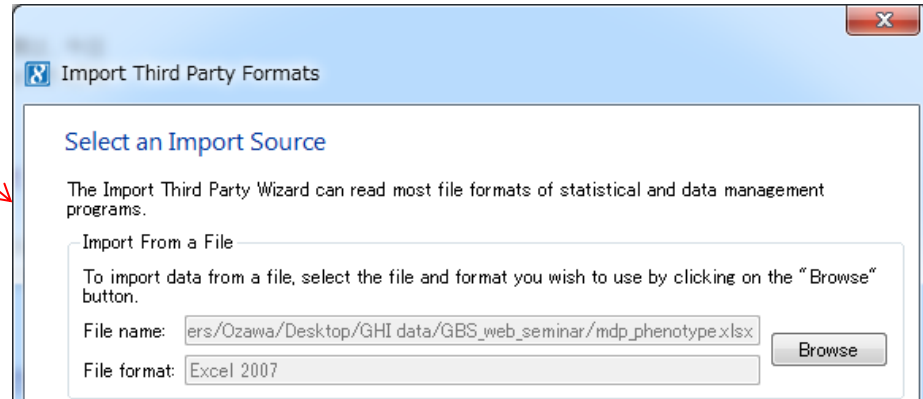
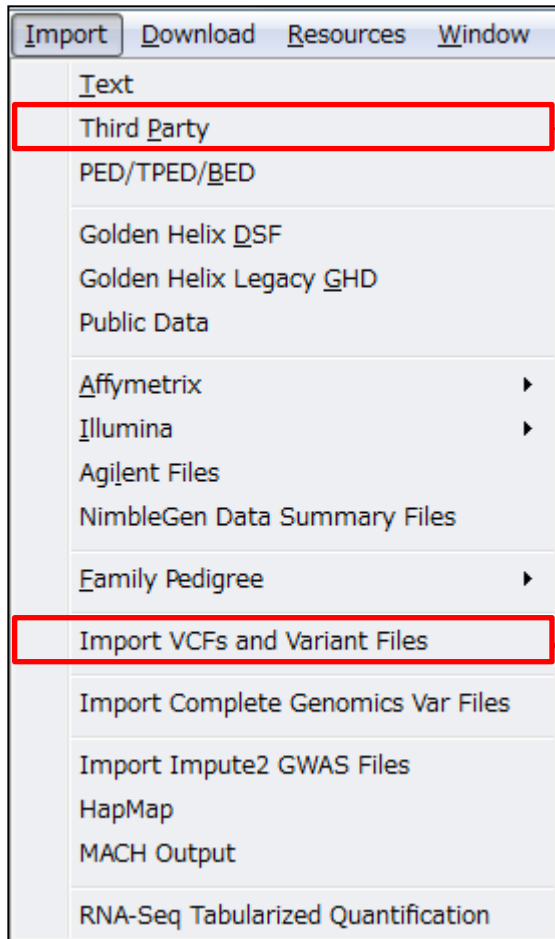
## 表現型データファイル (Excelファイル、Textファイル)

	A	B	C	D	E	F	G	H
1	Taxa	location	EarHT	dpool	EarDia	Q1	Q2	Q3
2	33-16	A	64.75	64.5		0.014	0.972	0.014
3	38-11	A	92.25	68.5	37.897	0.003	0.993	0.004
4	4226	A	65.5	59.5	32.21933	0.071	0.917	0.012
5	4722	A	81.13	71.5	32.421	0.035	0.854	0.111
6	A188	A	27.5	62	31.419	0.013	0.982	0.005
7	A214N	A	65	69	32.006	0.762	0.017	0.221
8	A239	A	47.88	61	36.064	0.035	0.963	0.002
9	A272	A	35.63	70		0.019	0.122	0.859
10	A441-5	A	53.5	67.5	35.008	0.005	0.531	0.464
11	A554	A	38.5	66	33.41775	0.019	0.979	0.002
12	A556	A	28	65	31.929	0.004	0.994	0.002
13	A6	A	109.5	80.5	31.5175	0.003	0.03	0.967
14	A619	A	36	61	40.63	0.009	0.99	0.001
15	A632	A	60	61	35.953	0.993	0.004	0.003
16	A634	A	54	59	35.601	0.897	0.1	0.003
17	A635	A	37	64	35.3005	0.825	0.171	0.004

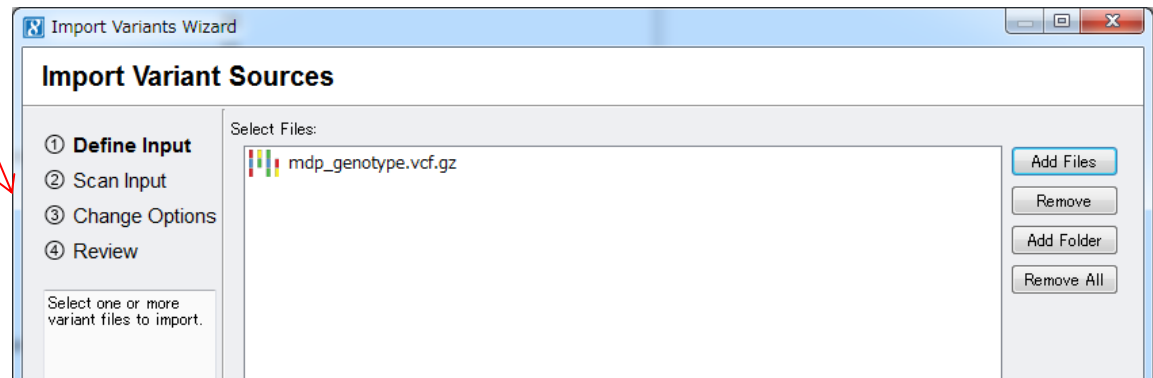
- 表現型データファイルには、疾患／正常などのグループ分類の他に、身長・体重などの連続値のデータも使用できる。
- ジェノタイプデータファイルは、バリアントコール用ツールなどで作成した、VCFファイルを使用する。

## ジェノタイプデータファイル (VCFファイル)

```
##fileformat=VCFv4.0
##Tasef=<ID=GenotypeTable,Version=5,Description="Reference allele is not known. The major allele was used as reference allele">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=AD,Number=,Type=Integer,Description="Allelic depths for the reference and alternate alleles in the order listed">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth (only filtered reads used for calling)">
##FORMAT=<ID=GQ,Number=1,Type=Float,Description="Genotype Quality">
##FORMAT=<ID=PL,Number=3,Type=Float,Description="Normalized, Phred-scaled likelihoods for AA,AB,BB genotypes where A=ref and B=alt; not applicable if site is not biallelic">
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=,Type=Float,Description="Allele Frequency">
#CHROM POS ID REF ALT QUAL FILTERINFO FORMAT 33-16 38-11 4226 4722 A188 A214N A239 A272 A441-5 A554 A556 A6
1 157104 PZB00859.1 C A . PASS . GT 0/0 0/0 0/0 0/0 1/1 0/0 1/1 1/1 0/0 0/0 0/0 1/1 0/0
1 1947984 PZA01271.1 G C . PASS . GT 1/1 0/0 1/1 0/0 1/1 1/1 1/1 1/1 1/1 0/0 1/1 1/1 0/0
1 2914066 PZA03613.2 T G . PASS . GT 1/1 1/1 1/1 1/1 1/1 0/0 0/0 0/0 1/1 0/0 1/1 0/0 1/1
1 2914171 PZA03613.1 T A . PASS . GT 0/0 0/0 0/0 0/0 0/0 1/1 0/0 0/0 0/0 0/0 0/0 0/0 0/0
1 2915078 PZA03614.2 G A . PASS . GT 0/0 0/0 0/0 0/0 0/0 0/0 1/1 1/1 0/0 1/1 0/0 1/1 0/0
1 2915242 PZA03614.1 T A . PASS . GT 0/0 0/0 0/0 0/0 0/0 1/1 1/1 1/1 0/0 0/0 0/0 1/1 0/0
1 2973508 PZA00258.3 C G . PASS . GT 1/1 0/0 0/0 1/0 0/0 0/0 0/0 1/1 0/0 0/0 ./ 0/0 0/0
1 3205252 PZA02962.13 T A . PASS . GT 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 ./ 0/0 1/1
1 3205262 PZA02962.14 C G . PASS . GT 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 ./ 0/0 1/1
1 3206090 PZA00599.25 T C . PASS . GT 1/1 0/0 1/1 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0
1 3706018 PZA02129.1 C T . PASS . GT 1/1 0/0 0/0 0/0 0/0 1/1 0/0 0/0 1/1 1/1 0/0 1/1 0/0
1 4175293 PZA00393.1 T C . PASS . GT 0/0 0/0 0/0 1/1 0/0 0/0 0/0 0/0 0/0 1/1 0/0 0/0 1/1
1 4429897 PZA02869.8 C T . PASS . GT 0/0 1/1 0/0 ./ 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0
1 4429927 PZA02869.4 C G . PASS . GT 0/0 0/0 0/0 ./ 1/1 0/0 0/0 0/0 1/1 0/0 0/0 0/0 0/0
1 4430055 PZA02869.2 T C . PASS . GT ./ 0/0 0/0 1/1 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0 0/0
1 4490461 PZA02032.1 A T . PASS . GT 0/0 1/1 0/0 0/0 0/0 1/1 0/0 0/0 0/0 0/0 1/1 0/0 0/0
```



表現型データファイル (Excelファイル) のインポート



ジェノタイプデータファイル (VCFファイル) のインポート

# データのインポート - データの統合

mdp\_phenotype Dataset - Sheet 1 [4]

Unsort	Taxa	C 1	R 2	R 3	R 4
1	33-16	A	64.75	64.5	?
2	38-11	A	92.25	68.5	37.897
3	4226	A	65.5	59.5	32.21933
4	4722	A	81.13	71.5	32.421
5	A188	A	27.5	62	31.419
6	A214N	A	65	69	32.006
7	A239	A	47.88	61	36.064
8	A272	A	35.63	70	?
9	A441-5	A	53.5	67.5	35.008
10	A554	A	38.5	66	33.41775

mdp\_genotype - Genotypes(G,T) - Sheet 1 [7]

Unsort	Map	G 4	G 5	G 6	G 7	G 8
1	33-16	C_C	C_C	G_G	T_T	G_G
2	38-11	C_C	G_G	G_G	T_T	G_G
3	4226	C_C	C_C	G_G	T_T	G_G
4	4722	C_C	G_G	G_G	T_T	G_G
5	A188	A_A	C_C	G_G	T_T	G_G
6	A214N	C_C	C_C	T_T	A_A	G_G
7	A239	A_A	C_C	T_T	T_T	A_A
8	A272	A_A	C_C	T_T	T_T	A_A
9	A441-5	C_C	C_C	G_G	T_T	G_G
10	A554	C_C	G_G	T_T	T_T	A_A

表現型データシート

ジェノタイプデータシート

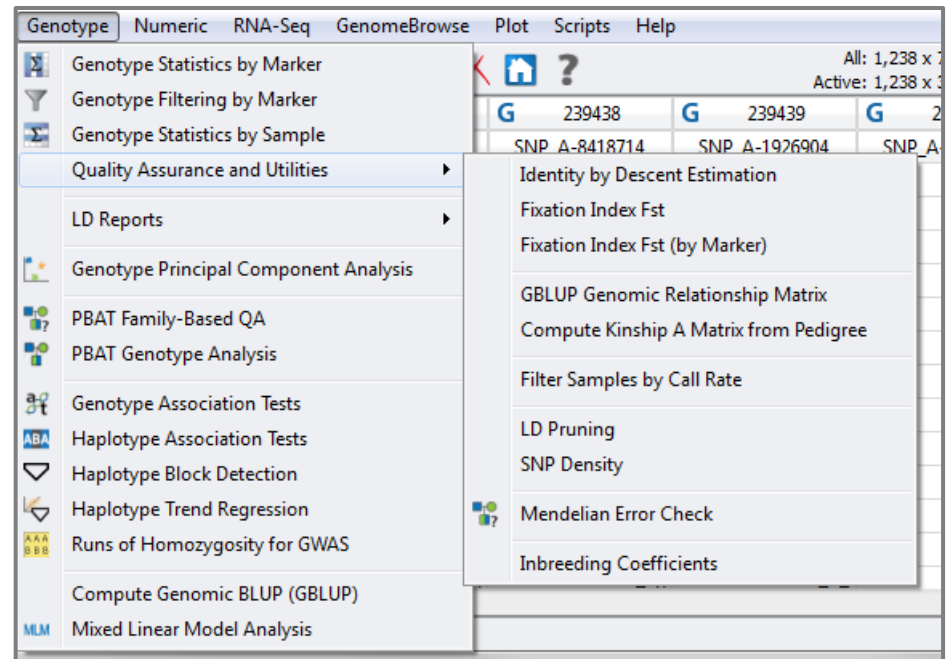
シートの結合

Editing: mdp\_phenotype Dataset + mdp\_genotype - Genotypes(G,T) - Sheet 1 [11]

	Taxa	C 1	R 2	R 3	R 4	G 5	G 6	G 7	G 8	G 9
1	33-16	A	64.75	64.5	?	C_C	C_C	G_G	T_T	G_G
2	38-11	A	92.25	68.5	37.897	C_C	G_G	G_G	T_T	G_G
3	4226	A	65.5	59.5	32.21933	C_C	C_C	G_G	T_T	G_G
4	4722	A	81.13	71.5	32.421	C_C	G_G	G_G	T_T	G_G
5	A188	A	27.5	62	31.419	A_A	C_C	G_G	T_T	G_G
6	A214N	A	65	69	32.006	C_C	C_C	T_T	A_A	G_G
7	A239	A	47.88	61	36.064	A_A	C_C	T_T	T_T	A_A
8	A272	A	35.63	70	?	A_A	C_C	T_T	T_T	A_A
9	A441-5	A	53.5	67.5	35.008	C_C	C_C	G_G	T_T	G_G
10	A554	A	38.5	66	33.41775	C_C	G_G	T_T	T_T	A_A

## SNP & Variation Suiteで使用可能なクオリティコントロール

- SNP Call Rateの検証
- Hardy Weinberg平衡 (HWE) の計算
- Minor Allele Frequency (MAF)に基づくフィルタリング
- 連鎖不平衡を示すSNPの除去
- 集団の階層化 (Population stractification)
- 性別誤認 (Gender misidentification) の検出
- メンデルエラーの検証
- 常染色体のヘテロ接合性
- Principal Component Analysis (PCA)
- Identity by Descent (IBD) の計算
- 多次元解析による異常値検出
- 染色体異常スクリーニング ...など





以下項目で、SNPのフィルタリングを実行

- Call Rate  
検出されたSNPの割合
- Number of allele  
検出されたアレル数
- Alternate allele frequency  
変異アレルの頻度
- Linkage disequilibrium (LD)  
SNP間の連鎖不平衡

Genotype Filtering by Marker

(No variable is set as dependent)

Classify alleles by allele frequency  Classify alleles by reference/alternate  
(Reference field in map: "Reference")

Filter Genotype Columns

General Statistics Filtering

Drop if call rate

Drop if number of alleles

Drop if alternate allele frequency

Drop if carrier count

Hardy Weinberg Equilibrium (HWE) Filtering

Perform HWE filtering based on:

Drop if Hardy Weinberg Equilibrium (HWE) P-Value

Drop if Fisher's exact test for HWE P-Value

Drop if signed HWE R (positive if more homozygous)

Actions

Inactivate genotype columns that meet above criteria for filtering

Output spreadsheet with marker statistics and 'Drop?' columns

Additional Output

Output  $-\log_{10}(\text{Value})$

Buttons: Help, Restore Options, Save Options, Run, Cancel

LD Pruning

Window Size

Window Increment

LD Statistic   $r^2$    $D'$

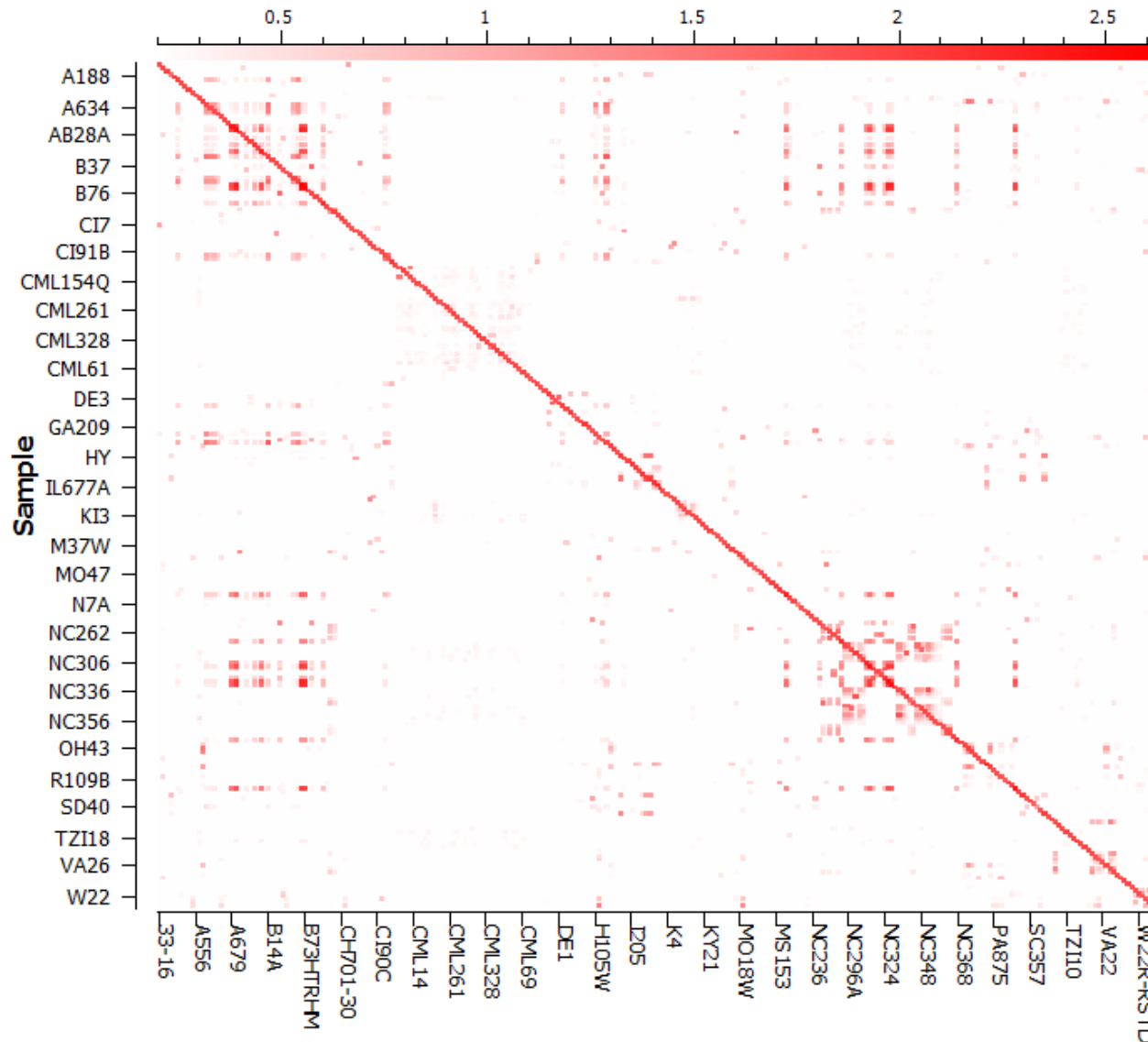
LD Threshold

LD Computation Method  CHM  EM

Buttons: OK, Cancel, Help

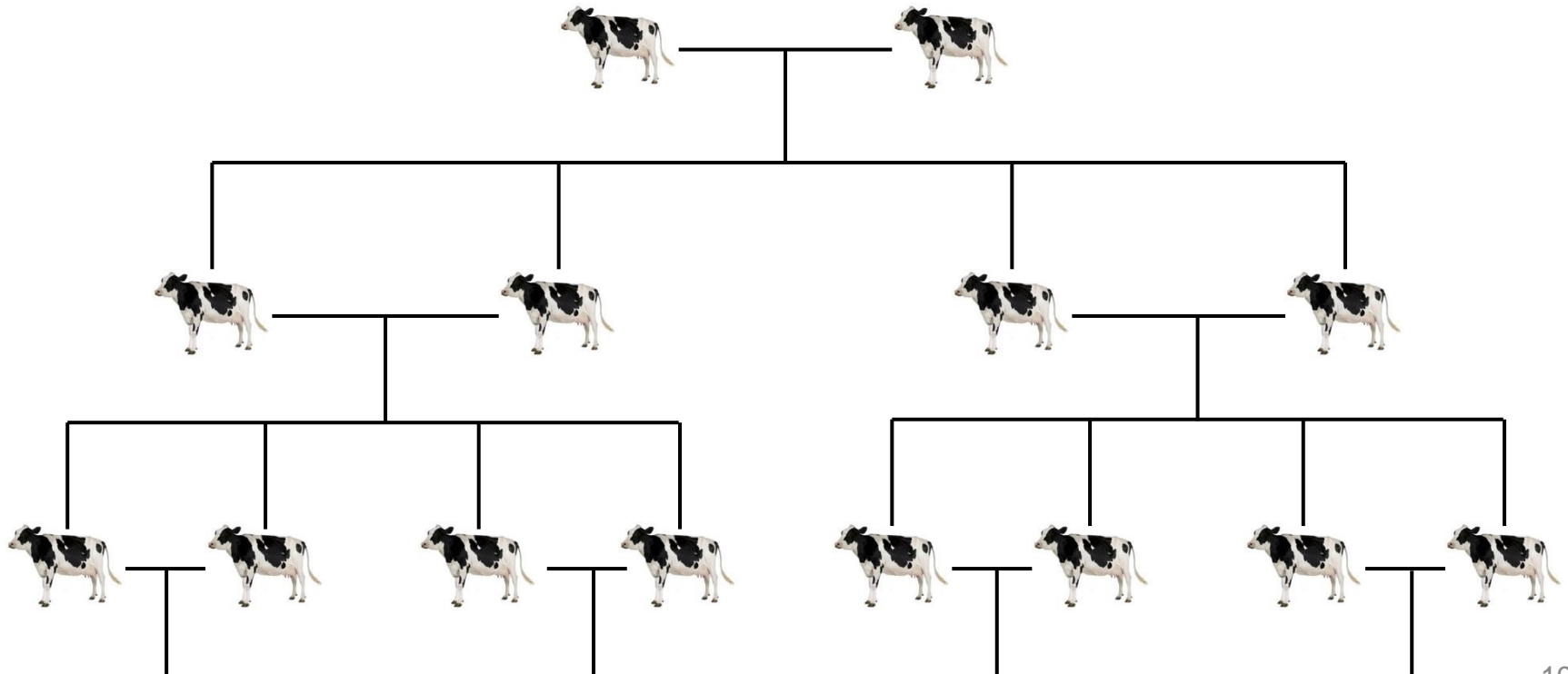
mdp\_phenotype Dataset + mdp\_genotype - Genotypes(G\_T) - Sheet 2 [18]

Unsort		G 238	G 239	G 240	G 241
Map	Taxa	1:150749370-SNV	1:161072169-SNV	1:161472433-SNV	1:161472477-REF
1	33-16	G_G	?_?	C_C	G_G
2	38-11	G_G	C_C	C_C	G_G
3	4226	G_G	T_T	C_C	G_G
4	4722	?_?	T_T	C_T	G_G
5	A188	G_G	C_C	C_C	G_G
6	A214N	G_G	T_T	C_C	G_G
7	A239	G_G	T_T	C_C	G_G
8	A272	G_G	C_C	?_?	G_G
9	A441-5	A_A	T_T	C_C	G_G
10	A554	A_A	T_T	C_C	G_G



- フィルタリングを行ったSNPデータを使用し、GBLUPモデルでサンプル間の相関を計算。

- Inbred lines (近交系) サンプルを解析に使用する場合は、関連解析実行時に、血縁関係によるバイアスを取り除く必要がある。
- 今回使用する線形混合モデルでは、サンプル間の相関データを使って、血縁関係にあるデータを補正することができる。
- 血縁関係の他、民族の違いによるバイアスも補正が可能。



## Mixed Model GWAS using a single locus (EMMAX)

- ジェノタイプデータによるサンプル間の相関データを用いて、血縁関係の偏りを補正する。
- 1か所のSNPごとに表現型との関連を計算する。

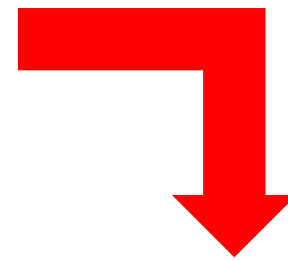
## Multi-locus mixed model GWAS (MLMM)

- ジェノタイプデータによるサンプル間の相関データを用いて、血縁関係の偏りを補正する。
- 複数か所のSNPをまとめて、表現型との関連を計算する。

## Genomic Best Linear Unbiased Predictors (GBLUP)

- ジェノタイプデータによるサンプル間の相関データを用いて、血縁関係の偏りを補正する。
- サンプルごとのランダム効果と、SNPごとのアレル代替効果を計算する。
- 農学分野における、育種価の計算に用いられる。

Map	Taxa	location	1	R	2	R	3	R	4	G	5	G	6	G	7
					EarHT		dpoll		EarDia	1:157104-SNV		1:1947984-SNV		1:2914066-SNV	
1	33-16	A			64.75		64.5		?	C_C		C_C		G_G	
2	38-11	A			92.25		68.5		37.897	C_C		G_G		G_G	
3	4226	A			65.5		59.5		32.21933	C_C		C_C		G_G	
4	4722	A			81.13		71.5		32.421	C_C		G_G		G_G	
5	A188	A			27.5		62		31.419	A_A		C_C		G_G	
6	A214N	A			65		69		32.006	C_C		C_C		T_T	
7	A239	A			47.89		61		36.064	A_A		C_C		T_T	



- 最初に、検定に使用するサンプルの表現型データの種類を指定する。
- 線形混合モデルのパラメータで、クオリティコントロールで計算しておいた、サンプル間相関データを選択する。

Mixed Linear Model Analysis

MLM Parameters      Additional Outputs

Regression Model(s) To Use

- Linear regression (fixed effects only)
- Mixed Model GWAS
- Single-locus mixed model GWAS (EMMAX)
- Multi-locus mixed model GWAS (MLMM)

Number of steps to use: 10

- Use Pre-Computed Kinship Matrix (Cov. Matrix of Random Effects)
  - GBLUP Genomic Relationship Matrix

Correct for Additional Covariates

Genetic Model and Imputation

Genetic model to use:

- Additive     Dominant     Recessive

Impute missing data as:

- Homozygous major allele     Numerically as average value

Correct For Hemizygous Males

Choose Sex Column:      

Chromosome that is hemizygous for males: X

P-Values from Linear Regression [46]

File Edit Select DNA-Seq Genotype Numeric RNA-Seq GenomeBrowse Plot Scripts Help

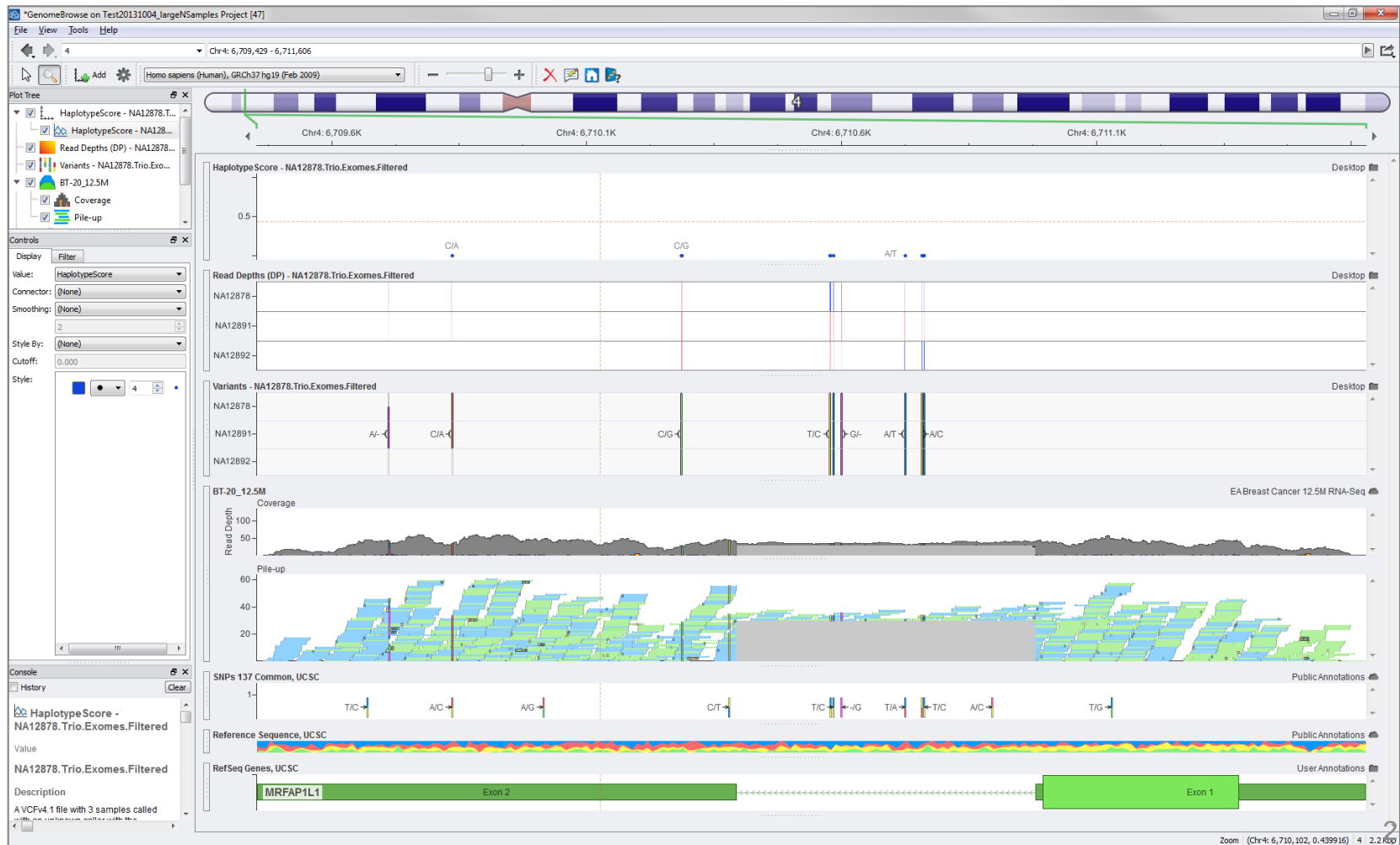
All: 2,107 x 13  
Active: 2,107 x 13

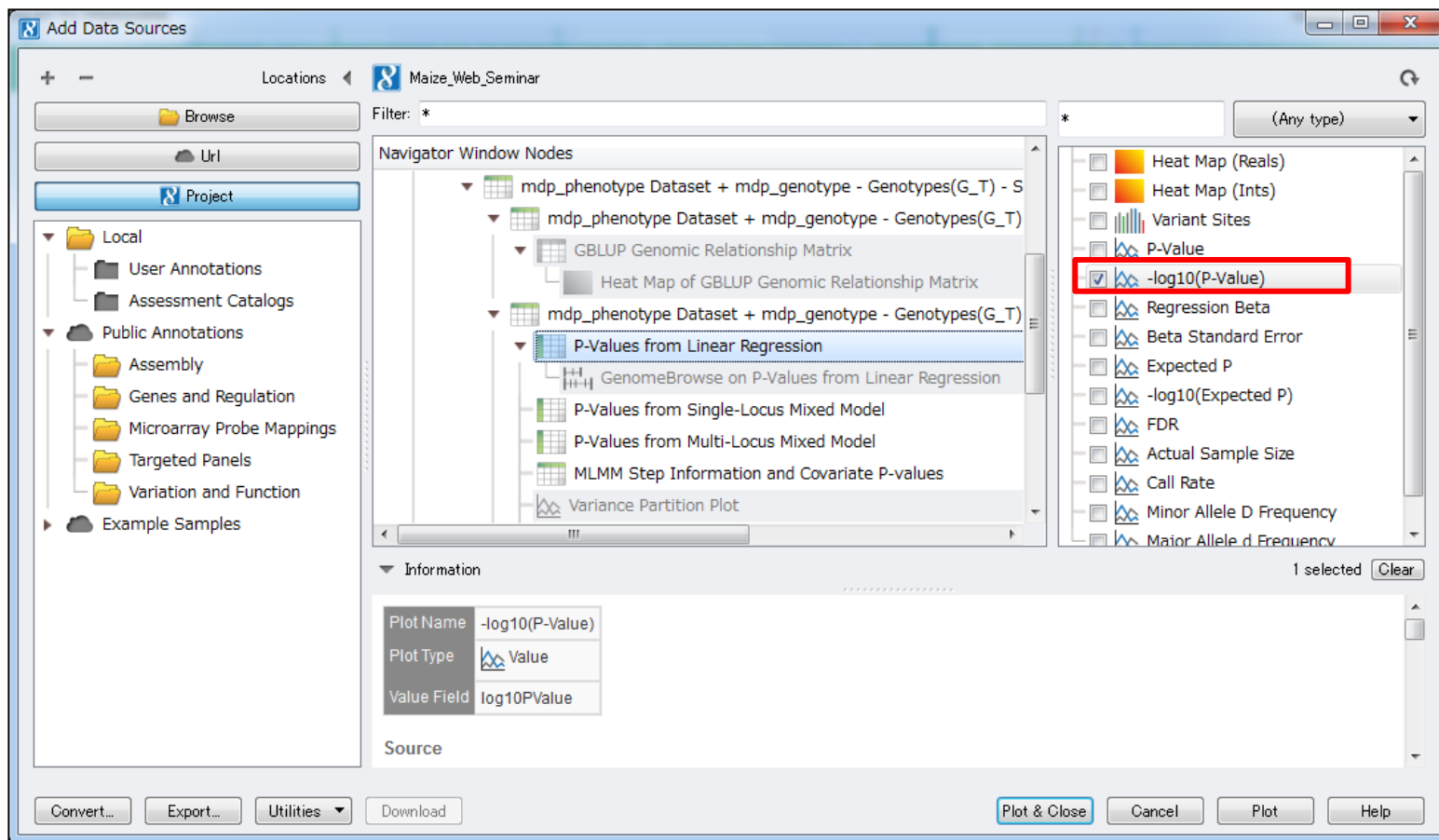
Unsort							R	1	R	2	R	3	R	4
Map	Marker	Chromosome	Position	Identifier	Reference	Alternates		P-Value		-log10(P-Value)		Regression Beta		Beta Standard Error
1	1:157104-SNV	1	157104	PZB00859.1	C	A		0.00682301335410328		2.16602377871164		-4.21360034760133		1.54514328023671
2	1:1947984-SNV	1	1947984	PZA01271.1	G	C		0.154230814494175		0.811828847894587		-1.79482856153314		1.25611161532419
3	1:2914066-SNV	1	2914066	PZA03613.2	T	G		0.456226211335439		0.340819767039785		-1.03499020156775		1.38704714709467
4	1:2914171-SNV	1	2914171	PZA03613.1	T	A		0.672174830549305		0.172517753631196		-0.627906392147472		1.48216434747834
5	1:2915078-SNV	1	2915078	PZA03614.2	G	A		0.558670458773075		0.252844292540025		0.738363000635729		1.26093650356777
6	1:2915242-SNV	1	2915242	PZA03614.1	T	A		0.406213581319227		0.391245560514189		1.05342311059235		1.26626038838676
7	1:2973508-SNV	1	2973508	PZA00258.3	C	G		0.0139014017823995		1.85694140437427		3.53690209154289		1.42821114162166
8	1:3205262-SNV	1	3205262	PZA02962.14	C	G		0.829032399857615		0.0814284962236546		0.586983935742978		2.71552558472251
9	1:3206090-SNV	1	3206090	PZA00599.25	T	C		0.912551826761436		0.0397424611867301		0.22738573259291		2.06852538723326
10	1:3706018-SNV	1	3706018	PZA02129.1	C	T		0.379069033633547		0.421281691889495		1.10703287545261		1.25641551969684
11	1:4175293-SNV	1	4175293	PZA00393.1	T	C		0.00337594548441879		2.47160457508171		-4.6341742612538		1.5664918817349
12	1:4430055-SNV	1	4430055	PZA02869.2	T	C		0.265438465415753		0.576036142107076		-1.99081710820301		1.78386324436588
13	1:4490461-SNV	1	4490461	PZA02032.1	A	T		0.435796435297428		0.360716326492215		1.29212305544475		1.65548955352549
14	1:5353319-SNV	1	5353319	PZB00919.1	C	A		0.312470154856717		0.50519145735998		-1.54329900839054		1.52499407297895
15	1:5562502-SNV	1	5562502	PHM2244.142	G	C		0.113334743117698		0.945636935403036		-2.75697682119204		1.73536951468481
16	1:8075572-SNV	1	8075572	PZA03093.10	G	C		0.722612474672404		0.141094545349473		-0.450242259926535		1.26703487457908
17	1:8367944-SNV	1	8367944	PZA00528.1	C	A		0.142904076869023		0.844955381170794		-2.29687963873737		1.56306274210936
18	1:8510027-SNV	1	8510027	PZA00175.2	T	A		0.0141144381762963		1.85033640421847		3.11896507842557		1.2622830965203
19	1:9024005-SNV	1	9024005	PZA00447.8	T	C		0.000515142803057097		3.28807236322522		-4.35707470090457		1.23902884211047
20	1:9029842-SNV	1	9029842	PZB01915.1	G	T		0.307995658742144		0.511455404920396		1.54786658855038		1.51540017941652
21	1:9084948-SNV	1	9084948	PZA03128.1	T	C		0.0313172932904999		1.50421578045181		-2.99197342180964		1.38219364118203
22	1:9084979-SNV	1	9084979	PZA03128.3	C	G		0.049168390152931		1.30831401093181		3.34002367673365		1.6900295422596
23	1:9273299-SNV	1	9273299	PZA02284.1	T	C		0.266074719903578		0.574996386357766		-1.41827328568897		1.27253440911744

P-Values from Linear Regression

- 計算が終了すると、使用した計算手法ごとに別のデータとして、各SNPを評価した P-Valueなどのデータが出力される。

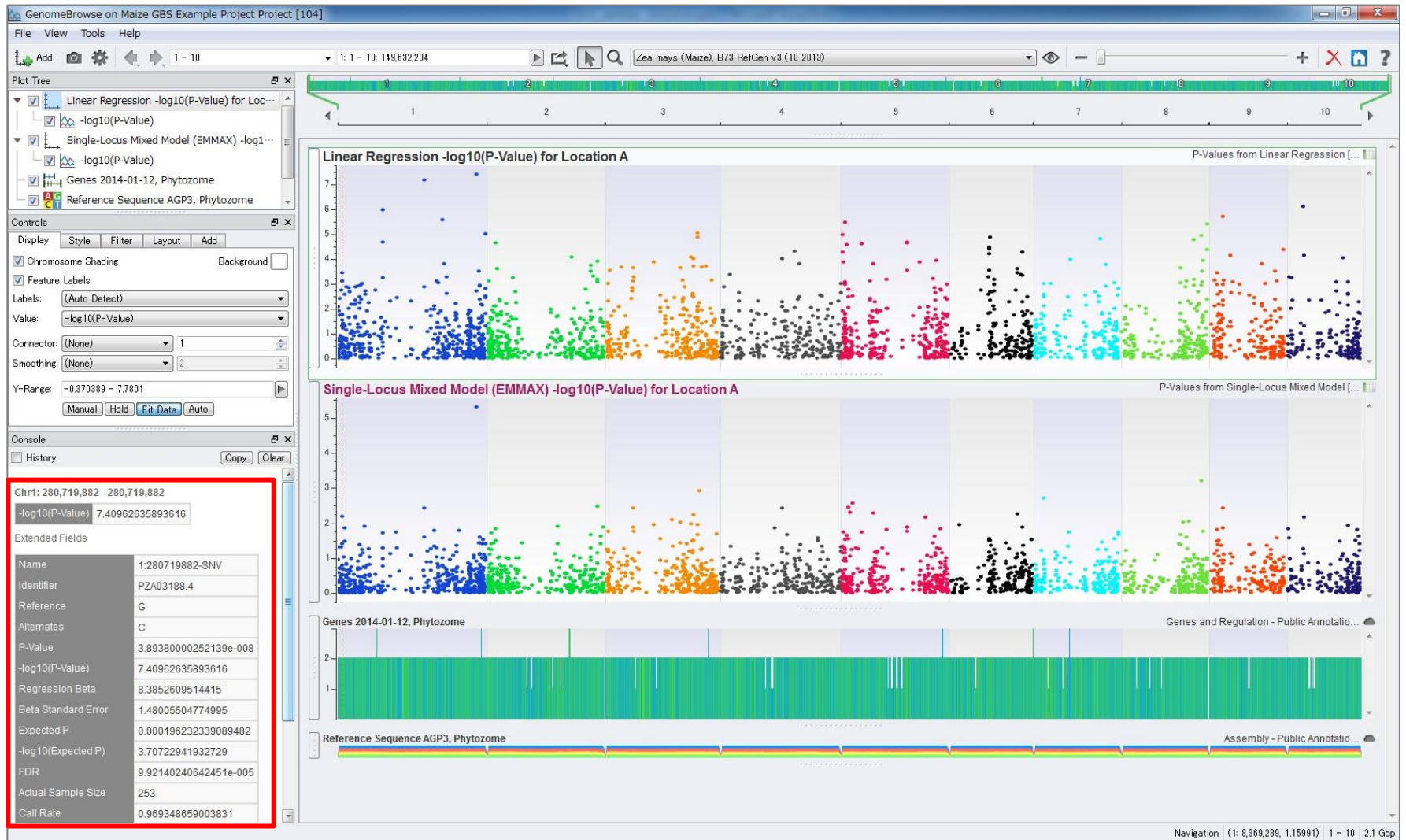
- Golden Helix社より無償で提供されているゲノムブラウザ「GenomeBrowse®」が、SVSに組み込まれており、BAMファイルデータ、VCFファイルデータ、各種数値データやアノテーションデータなどを統合表示が可能。





- 関連解析で計算したシートの、各SNPごとの「 $-\log_{10}(\text{P-Value})$ 」を選択してプロットする。
- 同時に、各種データベースのアノテーションデータや、ユーザー作成データのプロットも可能。





- プロット上の各ポイントをクリックすると、画面左下にSNPの詳細情報が表示される。

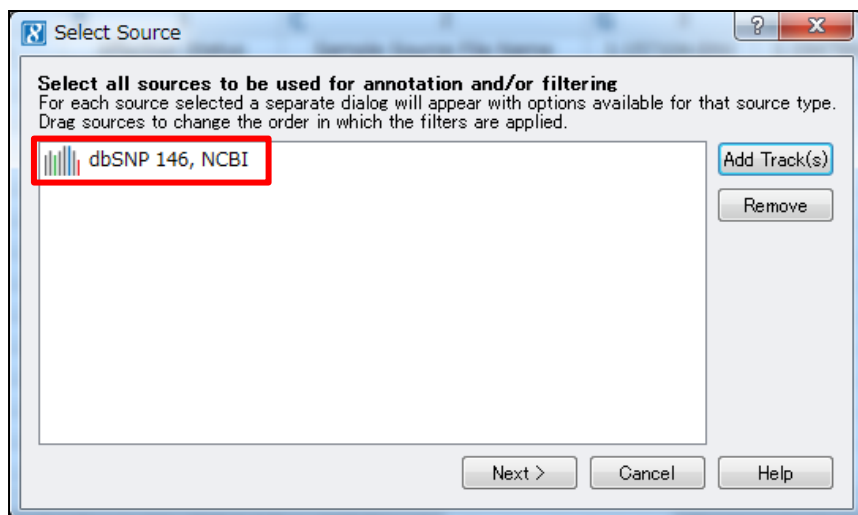
Coding Variant Classification - Loss of Function and Missense Mutations [74]

File Edit Select DNA-Seq Genotype Numeric RNA-Seq GenomeBrowse Plot Scripts Help

All: 49 x 7  
Active: 49 x 7

Unsort	1	2	3	4	5	6	7	
Map	Variant	Classification	Priority	Gene 1	Transcript 1	Exon 1	HGVS Coding 1	HGVS Protein 1
1	1:12208569-SNV	Nonsyn SNV	3	GRMZM2G096252	GRMZM2G096252_T01	1	c.500C>T	p.Ala167Val
2	1:21457474-SNV	Nonsyn SNV	3	GRMZM2G004259	GRMZM2G004259_T01	3	c.170G>T	p.Ser57Ile
3	1:208143724-SNV	Nonsyn SNV	3	GRMZM2G479318	GRMZM2G479318_T01	2	c.49C>G	p.Arg17Gly
4	1:212244300-SNV	Nonsyn SNV	3	GRMZM2G422720	GRMZM2G422720_T01	2	c.239C>T	p.Thr80Ile
5	1:213323599-SNV	Nonsyn SNV	3	GRMZM2G311550	GRMZM2G311550_T01	3	c.602G>A	p.Gly201Asp
6	1:255148812-SNV	Nonsyn SNV	3	GRMZM2G081953	GRMZM2G081953_T01	2	c.92C>T	p.Pro31Leu
7	1:255148947-SNV	Nonsyn SNV	3	GRMZM2G081953	GRMZM2G081953_T01	2	c.227C>T	p.Ser76Phe
8	1:278690966-SNV	Splicing	9	GRMZM2G147596	GRMZM2G147596_T01	2	c.635-2A>T	?
9	1:293632755-SNV	Nonsyn SNV	3	GRMZM2G124805	GRMZM2G124805_T01	3	c.657C>G	p.Asp219Glu
10	1:293632805-SNV	Nonsyn SNV	3	GRMZM2G124805	GRMZM2G124805_T01	3	c.607C>A	p.Gln203Lys
11	2:10429605-SNV	Nonsyn SNV	3	GRMZM2G429762	GRMZM2G429762_T01	2	c.133T>G	p.Phe45Val
12	2:11678938-SNV	Nonsyn SNV	3	GRMZM2G072614	GRMZM2G072614_T01	3	c.320C>T	p.Pro107Leu
13	2:170899280-SNV	Nonsyn SNV	3	GRMZM2G074687	GRMZM2G074687_T01	13	c.1207G>A	p.Asp403Asn
14	2:220397345-SNV	Nonsyn SNV	3	AC183504.4_FG003	AC183504.4_FGT003	1	c.326A>G	p.Gln109Arg

Coding Variant Classification - Loss of Function and Missense Mutations



- 多数サンプルデータによる統計処理以外にも、各種データベースのアノテーションデータを利用した解析も可能。
- タンパク質アミノ酸配列の変化による、非同義変異の抽出や、コモンSNPの除去などが可能。

ソフトウェアの詳細は、以下の弊社Webサイトをご覧ください。

SNP & Variation Suite (SVS):

<http://www.filgen.jp/Product/BioScience21-software/goldenhelix/index.html>

お問い合わせ先: フィルジェン株式会社

TEL 052-624-4388 (9:00~17:00)

FAX 052-624-4389

E-mail: [biosupport@filgen.jp](mailto:biosupport@filgen.jp)