# GENETICS

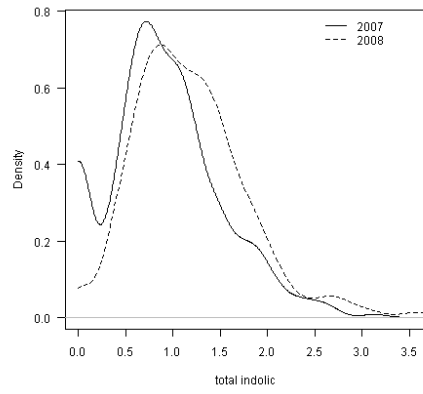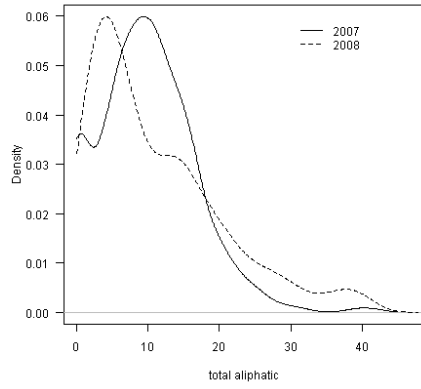## Understanding the Evolution of Defense Metabolites in *Arabidopsis thaliana* Using Genomewide Association Mapping

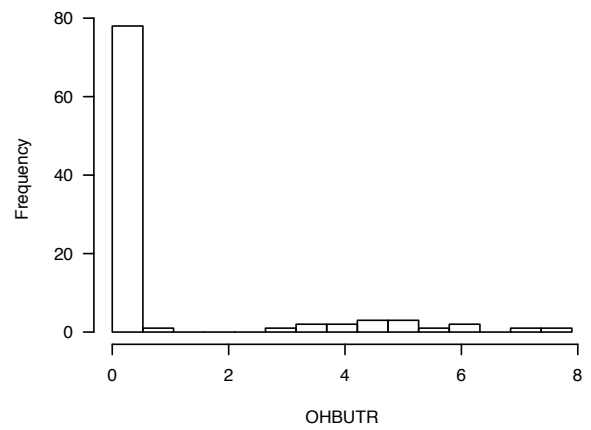Eva K. F. Chan, Heather C. Rowe and Daniel J. Kliebenstein
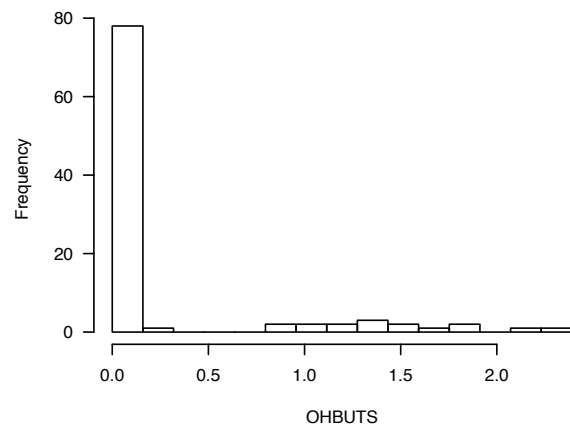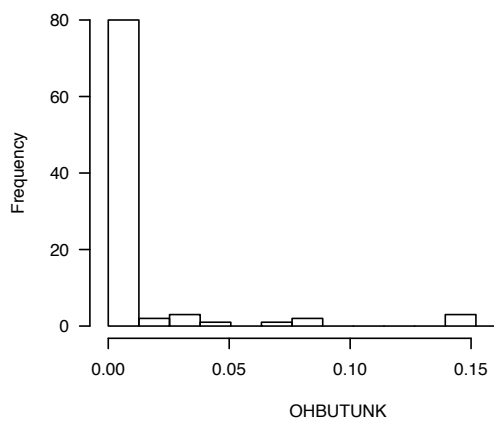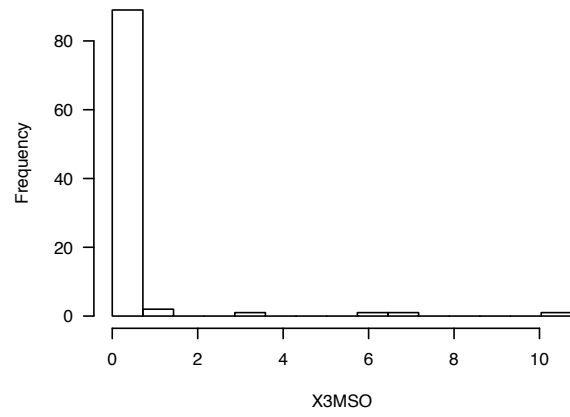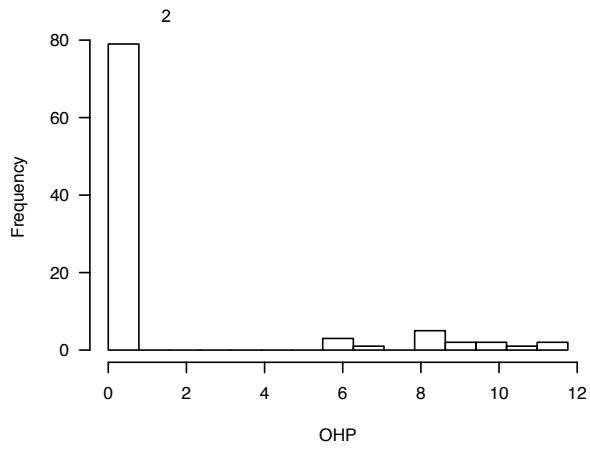
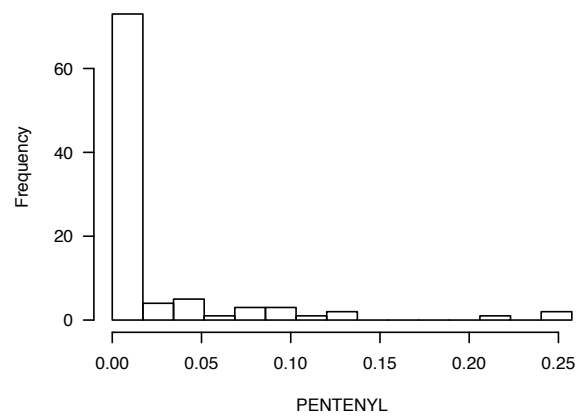E. K.F. Chan *et al.*

1

E. K. F. Chan *et al.*

E. K. F. Chan *et al.*

E. K. F. Chan *et al.*

E. K. F. Chan *et al.*



FIGURE S1.—GSL distributions in 2007 and 2008. Shown on page 1 is a direct comparison for the distribution of the total aliphatic and indolic glucosinolates for the two datasets. The graphs on the remaining pages are individual histograms for all utilized glucosinolate phenotypes within each dataset: pages 2-8 are for the 2007 dataset and 9-15 are for the 2008 dataset.

FIGURE S2.—Genome-wide LD with SNPs in the *AOP* and *MAM* loci. Pairwise linkage disequilibrium between SNPs within coding regions of query genes (*AOP2, AOP3, MAM1, and MAM3*) and all ~240,000 SNPs is shown as a heatmap with significantly strong to weak LD indicated by green to white respectively. The top histogram represents the number of GSL phenotypes showing a significant association with the individual SNP in this region for both the 2007 and 2008 datasets. The right bar-graph shows the minor allele frequency of each SNP within the query genes. The cladogram on the left indicates the hierarchical clustering of SNPs within the query genes.

**MAM1**



**MAM1**



**MAM3**



**MAM3**

**AOP3**



**AOP3**



**AOP2**



**AOP2**

FIGURE S3.—LD decay around *MAM, AOP*, and recent selective sweeps. LD, as measured by $r^2$ and $-\log_{10}P(\chi^2)$, between SNPs within a query gene/region and SNPs within ±100kb around of the query gene is plotted against the physical distance (kilobases) between each SNP-pairs. The four query genes are *MAM1, MAM3, AOP2,* and *AOP3*, and the query regions are two previously proposed recent selective sweep (CLARK *et al.* 2007). The red line shows the fitted cubic smoothing spline (*spar*=1) between LD and distance. The top and bottom pink lines show the 99th and 95th percentiles of non-syntenic LD; i.e. LD between SNPs within the query gene and SNPs on chromosomes outside of the gene. The center of the graph is at the center of the query genes.

AOP2

E. K. F. Chan *et al.*

**AOP2**

AOP3

AOP3

FIGURE S4.—LD Structure at and around the *AOP,MAM* and sweep loci. Shown are cladograms of 96 accessions using SNPs within the *AOP2, AOP3, MAM1,* and *MAM* genes, as well as the two sweep regions. The main GSL structure within each accession is indicated by color of the accession name. Each gene is shown twice, the first figure displays the allelic values of the SNPs within the specific gene itself; blue represents the majority allele and yellow the minority allele. The second figure displays the allelic values of the SNPs that are encompassed by the region of extended LD surrounding the corresponding gene as defined using the 99th (or 98th or *MAM3*) percentile boundaries from Supplemental Figure 3.

**TABLE S1**

**ANOVA of accession and population structure for GSL 2007**

Table S1 is available for download as an Excel file at http://www.genetics.org/cgi/content/full/genetics.109.108522/DC1.

For each GSL trait the following model was examined: $y_{sar} = \beta_0 + S_s + A(S)_{sa} + R_r + \varepsilon_{sar}$ where s = 1,…, 8 (NORDBORG *et al.* 2002; NORDBORG *et al.* 2005); r = 1,…4; and a = 1,…,95. The main effects are denoted as S, A and R to represent structure, accession, and replicate block respectively, and error, $\varepsilon_{sar} \sim N(0, \sigma_\varepsilon^2)$. F: F-statistic of the model, $df_{num}$: numerator degrees of freedom, $df_{denom}$: denominator d.f., P(F): nominal P-value of the F-statistic, $R^2$: fraction of total variance explained by the model; $\bar{x}$: observed mean, $\beta_0$: model mean, $\eta^2$: partial $R^2$ of the corresponding predictor variable, and P: P-value of the corresponding predictor variable. Description shows the chemical structure of the glucosinolate or the algebraic approximation of the trait.

**TABLE S2**

**ANOVA of accession and population structure for GSL 2008**

Table S2 is available for download as an Excel file at http://www.genetics.org/cgi/content/full/genetics.109.108522/DC1.

For each GSL trait the following model was examined: $y_{sar} = \beta_0 + S_s + A(S)_{sa} + R_r + \varepsilon_{sar}$ where s = 1,…, 8 (NORDBORG *et al.* 2002; NORDBORG *et al.* 2005); r = 1,…4; and a = 1,…,95. The main effects are denoted as S, A and R to represent structure, accession, and replicate block respectively, and error, $\varepsilon_{sar} \sim N(0, \sigma_\varepsilon^2)$. F: F-statistic of the model, $df_{num}$: numerator degrees of freedom, $df_{denom}$: denominator d.f., P(F): nominal P-value of the F-statistic, $R^2$: fraction of total variance explained by the model; $\bar{x}$: observed mean, $\beta_0$: model mean, $\eta^2$: partial $R^2$ of the corresponding predictor variable, and P: P-value of the corresponding predictor variable. Description shows the chemical structure of the glucosinolate or the algebraic approximation of the trait.

**TABLE S3**

**Average glucosinolate values per accession for both experiments**

Table S3 is available for download as an Excel file at http://www.genetics.org/cgi/content/full/genetics.109.108522/DC1.

All metabolite values are in nmol per mg fresh weight tissue. This is the average of four independent plants per accession.

**TABLE S4**

**Correlation of glucosinolate phenotypes**

Table S4 is available for download as an Excel file at http://www.genetics.org/cgi/content/full/genetics.109.108522/DC1.

Shown are the significant spearman rank correlations for all glucosinolate phenotypes used in this study.

**TABLE S5**

**Gene to Trait Associations as identified from GWA**

Table S5 is available for download as an Excel file at http://www.genetics.org/cgi/content/full/genetics.109.108522/DC1.

Logical table indicating whether each of 31,505 genes is significantly associated to each of the 43 traits. For each trait, a gene is significantly associated if at least two SNPs within ±1kb flanking the coding region has a P-value in the bottom 0.1 percentile of the P-value distribution. T: is significant in either 2007 or 2008 set; F: not significant in both sets; -: not significant in either 2007 or 2008 set and untested in the reverse set.

**TABLE S6**

**Known and putative genes involved in the GSL biosynthetic pathway**

Table S6 is available for download as an Excel file at http://www.genetics.org/cgi/content/full/genetics.109.108522/DC1.

List of genes either known or predicted to play a role in GSL metabolism and regulation. AGI:  shows the AGI (Arabidopsis Genome Initiative) code for each gene; Pathway: specific part of the GSL metabolic system the gene is thought to function, this includes the core biosynthetic pathways as well as pathways required for the production of immediate precursors such as methionine, tryptophan, etc.; Pseudogene: whether or not the gene is predicted to be a pseudogene; Evidence: experimental evidence (Genetic or Biochemical) or sequence evidence base on homology to validated GSL gene (Homology).

**TABLE S7**

**GSL genes as a group not preferentially identified in 2007 or 2008**

Table S7 is available for download as an Excel file at http://www.genetics.org/cgi/content/full/genetics.109.108522/DC1.

The Kolmogorov-Smirnov D-statistics and corresponding P-values are shown for each trait and for each of 2007 & 2008 datasets from the GWA results. Each KS-test assessed the null hypothesis that: the GWA P-value distribution corresponding to SNPs within ±1kb flanking region of known GSL genes is identical to that of non-GSL SNPs.

**TABLE S8**

**Using known GSL genes to estimate call thresholds**

|  | Dataset | ≥1 SNP per gene | ≥2 SNP per gene | ≥20% SNP per gene |
|---|---|---|---|---|
| # of Sig Genes | 2007 | 5333 | 1056 | 1135 |
|  | 2008 | 4703 | 893 | 960 |
| GSL Genes in Sig Genes | 2007 | 0.7% | 1.2% | 0.8% |
|  | 2008 | 0.5% | 1.5% | 0.5% |
| Sig GSL Genes | 2007 | 25.3% | 8.7% | 6.0% |
|  | 2008 | 16.7% | 8.7% | 3.3% |

The number of genes significantly associated with at least one trait in the two datasets (2007 and 2008) using three different call thresholds are presented. The percentage of significant genes overlapping with the list of known and putative GSL genes is provided as well as the percentage of all known GSL genes exceeding each threshold.

**FILE S1**

**GWA Results**

File S1 is available for download as a compressed (.zip) file at http://www.genetics.org/cgi/content/full/genetics.109.108522/DC1.

This dataset contains three files. One is the SNP information (snp.info.txt) which includes the chromosome and genomic position of each SNPs, as well as the minor allele frequencies in this dataset. The two other are tables of P-values of association tests between each GSL trait and SNPs; one file for each experiment (2007 & 2008).

**FILE S2**

**Supporting Method: Modifications to R/EMMA**

This file details the changes made to three functions of the R/EMMA package, version 1.0.7 (KANG *et al.* 2008).

**emma.eigen.L.w.Z()**

changed:
```
      eig <- eigen(K%*%crossprod(Z,Z),symmetric=FALSE,EISPACK=TRUE)
```
to:
```
      eig <- eigen(K%*%crossprod(Z,Z),symmetric=F,EISPACK=F)
```

**emma.eigen.R.w.Z()**

changed:
```
      eig <- eigen(K%*%crossprod(Z,SZ),symmetric=FALSE,EISPACK=TRUE)
```
to:
```
      eig <- eigen(K%*%crossprod(Z,SZ),symmetric=F,EISPACK=F)
```

**emma.ML.LRT()**

changed:
```
    for(i in 1:g) {
      vrows <- !is.na(ys[i,])
      if ( is.null(Z) ) {
        ML0[i] <-
emma.MLE(ys[i,vrows],X0[vrows,,drop=FALSE],K[vrows,vrows],NULL,ngrids,l
lim,ulim,esp)$ML
      }
      else {
        vids <- colSums(Z[vrows,]>0)
        ML0[i] <-
emma.MLE(ys[i,vrows],X0[vrows,,drop=FALSE],K[vids,vids],Z[vrows,vids],n
grids,llim,ulim,esp)$ML
      }
    }
```
to:
```
    for(i in 1:g) {      # for each trait
       vrows <- !is.na(ys[i,])
       if( var(ys[i,vrows],na.rm=T)>0 ) {
          if ( is.null(Z) ) {
          ML0[i] <-
emma.MLE(ys[i,vrows],X0[vrows,,drop=FALSE],K[vrows,vrows],NULL,ngrids,l
lim,ulim,esp)$ML
          }
          else {
            #vids <- colSums(Z[vrows,]>0)
            vids <- as.logical(colSums(Z[vrows,]>0))
            ML0[i] <-
emma.MLE(ys[i,vrows],X0[vrows,,drop=FALSE],K[vids,vids],Z[vrows,vids],n
grids,llim,ulim,esp)$ML
          }
       }
    }
```

E. K.F. Chan *et al.*

## FILE S3

### Genotypes of 248,584 SNP for the 96 *A. thaliana* accessions

File S3 is available for download as a compressed (.zip) file at http://www.genetics.org/cgi/content/full/genetics.109.108522/DC1.

Genotypes of 248,584 SNP for the 96 A. thaliana accessions were obtained using the SNP chip described by Kim et al. (2007) and were generated by the groups of J. Bergelson, J. O. Borevitz, and M. Nordborg.