**Author for correspondence:**
Joy Bergelson
e-mail: jb7684@nyu.edu

## THE ROYAL SOCIETY
PUBLISHING

# Genome-wide association mapping within a local *Arabidopsis thaliana* population more fully reveals the genetic architecture for defensive metabolite diversity

Andrew D. Gloss[1,2], Amélie Vergnol[2], Timothy C. Morton[2], Peter J. Laurin[1,2], Fabrice Roux[3] and Joy Bergelson[1,2]

[1]Department of Biology and Center for Genomics and Systems Biology, New York University, New York, NY, USA
[2]Department of Ecology and Evolution, University of Chicago, Chicago, IL, USA
[3]LIPME, Université de Toulouse, INRAE, CNRS, Castanet-Tolosan, France

ADG, 0000-0003-3649-1836; TCM, 0000-0002-1310-3450; PJL, 0000-0003-0294-514X;
FR, 0000-0001-8059-5638; JB, 0000-0001-7893-7387

A paradoxical finding from genome-wide association studies (GWAS) in plants is that variation in metabolite profiles typically maps to a small number of loci, despite the complexity of underlying biosynthetic pathways. This discrepancy may partially arise from limitations presented by geographically diverse mapping panels. Properties of metabolic pathways that impede GWAS by diluting the additive effect of a causal variant, such as allelic and genetic heterogeneity and epistasis, would be expected to increase in severity with the geographical range of the mapping panel. We hypothesized that a population from a single locality would reveal an expanded set of associated loci. We tested this in a French *Arabidopsis thaliana* population (less than 1 km transect) by profiling and conducting GWAS for glucosinolates, a suite of defensive metabolites that have been studied in depth through functional and genetic mapping approaches. For two distinct classes of glucosinolates, we discovered more associations at biosynthetic loci than the previous GWAS with continental-scale mapping panels. Candidate genes underlying novel associations were supported by concordance between their observed effects in the TOU-A population and previous functional genetic and biochemical characterization. Local populations complement geographically diverse mapping panels to reveal a more complete genetic architecture for metabolic traits.

This article is part of the theme issue 'Genetic basis of adaptation and speciation: from loci to causative mutations'.

## 1. Introduction

Plants produce a vast array of secondary metabolites that collectively underpin many functions—from regulating growth and development to tolerating abiotic stresses, attracting pollinators and deterring enemies [1]. Illuminating the genetic architecture of secondary metabolism is not only integral to understanding plant physiology, adaptation and diversity across environments [2]; it also provides precise routes to create more durable and productive crops [3].

In recent years, genome-wide association studies (GWAS) have emerged as a tool of choice for elucidating the genotype-to-phenotype links that shape plant metabolic diversity [3–5]. GWAS involve tests for statistical associations between genetic variants and organismal phenotypes. Because they require only genotypic and phenotypic information across a panel of natural plant genotypes (accessions), GWAS offer a straightforward and efficient method for inferring links between millions of single-nucleotide polymorphisms (SNPs) across the

genome and thousands of metabolites, enabled by the parallel advances in genome sequencing and metabolomic profiling.

A paradoxical pattern emerging from GWAS of plant metabolic features is that only a few loci are associated with variation in the abundance of a given metabolite [5]. Indeed, an average of fewer than two significant loci per metabolite were discovered across four GWAS studies encompassing greater than 6500 metabolites in leaves and/ or seeds of *Arabidopsis*, rice and maize ($N = 305–529$ plant accessions per study) [6–9]. Such simple genetic architectures are surprising given that secondary metabolites are often the product of biosynthetic pathways that have many enzyme-catalyzed steps [10]. One potential explanation is that many genes in these pathways are subject to strong purifying selection [11,12], and thus lack polymorphisms to be mapped. However, this explanation does not fully explain the simple genetic architecture, as GWAS fails to replicate many functionally validated loci uncovered through other techniques for interrogating the genetic basis of metabolic variation [13].

Much attention has been paid to forces that reduce the efficacy of GWAS, and to both experimental designs and statistical approaches to mitigate them [14,15]. One relatively understudied factor is the composition of the mapping panel, especially the geographical distribution over which accessions are drawn [14,15]. This is an important consideration because GWAS mapping panels in plants have conventionally been assembled over broad geographical scales, such as the *Arabidopsis* Regional Mapping Population (RegMap) and 1001 Genomes Project (1001G), which are composed predominantly of accessions collected across the European continent [16,17]. This design ensures that a broad swath of the species' genetic diversity is included within the mapping panel, one of the main advantages of GWAS compared to quantitative trait locus (QTL) mapping. However, it also exposes analyses to a variety of geographically driven confounding forces.

The most popularized cause of confounding driven by geography concerns population structure [18,19]. False-positive associations arise at non-causal variants whose genotypes are correlated (i.e. in long-range linkage disequilibrium) with causal variants, and geographical population structure is a major source of these correlations [18]. The incorporation of a kinship matrix in the GWAS model can control these spurious associations [20,21], but at the cost of reducing power to detect causal variants whose geographical distribution tracks major axes of population structure [22,23].

However, even with effective control for the effects of long-range linkage disequilibrium, additional confounding factors are strengthened in geographically structured populations. Three processes in particular can dilute the strength of association at a causal variant. First, many alleles have geographically restricted distributions, causing the genetic basis of a trait to vary across regions (genetic heterogeneity) [14,24,25]. Because rare alleles in particular tend to be geographically restricted [26], mapping within local or regional panels would also have the benefit of elevating the frequencies of some rare alleles relative to their species-wide frequency [26], thus enhancing power to detect rare, informative SNPs. Second, a locus can have more than two functionally distinct haplotypes (allelic heterogeneity), especially in geographically broad mapping panels that have high genetic diversity [14,27]. Because GWAS typically interrogates biallelic SNPs, a variant's effect is diluted by averaging across the haplotypes tagged by each allele.

Third, population structure across multiple causal loci can produce different genotypic combinations in different geographical regions. GWAS is less powerful when a causal variant's effect is markedly weakened in some genetic backgrounds due to epistasis since standard GWAS models are formulated to detect average additive effects across genetic backgrounds [28,29]. All of these factors point to the benefit of mapping in local panels, provided that adequate phenotypic and genetic variation is present.

Glucosinolates (GSLs), the primary class of secondary defensive metabolites in *Arabidopsis* [30], are a well-studied example for which GWAS uncovers only a subset of genes in a complex biosynthetic pathway. As such, they offer a compelling opportunity to test the hypothesis that a local GWAS mapping population can better expose the genetic architecture of a complex trait than a geographically broad GWAS population. Glucosinolate biosynthesis requires a number of sequential enzyme-catalyzed reactions to produce a given aliphatic GSL (methionine-derived, 12–15 reactions) or indolic GSL (tryptophan-derived, 7–9 reactions) from their precursor amino acid [31]. Each step of the pathway has been functionally characterized through forward and reverse genetics approaches, leading to the identification of at least 45 genes involved [31]. Yet three GWAS of aliphatic GSL variation with large mapping populations ($N > 300$) spanning Europe have consistently described associations at only three biosynthetic loci [6,13,32], even though the causal polymorphisms underlying mapped QTL have been localized to additional biosynthetic genes [33].

Intriguingly, GSLs across the European distribution of *Arabidopsis* [13] exhibit all sources of confounding detailed above. Recurrent loss of function and gene conversion events have generated complex patterns of allelic heterogeneity, including rare variants, and the geographically restricted distributions of functionally defined haplotypes at a few major-effect loci implies strong genetic heterogeneity [13,32,34]. Higher-order epistatic interactions among major-effect loci determine which GSL molecules accumulate, resulting in GSL profiles that can be binned into qualitative 'chemotypes', defined by whether the gene(s) at each locus are functional [35]. Distributions of these epistatically defined chemotypes are also geographically biased, displaying regional or continental clines [13,32]. If similar patterns have arisen at other loci with more modest phenotypic effects, geographical confounding might hinder their detection through GWAS; at the very least, large effect epistasis has been documented for other GSL biosynthetic enzymes [33,36].

Here, we quantified variation in GSL profiles in a single local population of *Arabidopsis*, compared the genetic architecture revealed through GWAS in this local population and geographically broad mapping panels, and explored potential confounding factors underlying differences in the performance of the mapping populations. We focused on a population from Toulon-Sur-Arroux (TOU-A), France, which was collected along a fence line spanning only a few hundred metres [37]. Previous investigations found that the TOU-A population harbours less than 20% of the variants segregating at detectable frequencies in the 1001G, yet variants underlying heritable variation for a wide range of morphological, growth, defence, and fitness-related traits in TOU-A can be successfully mapped using GWAS [37,38]. We restricted our focus to genes with validated functions in GSL biosynthesis, broadly defined to include core structure

formation, side-chain elongation, and secondary modification [31]. Decades of research have compiled a near-exhaustive catalogue of the genes participating in these processes and their substrate specificities, providing functional data supporting novel associations that we uncovered at these loci. Overall, the expanded catalogue of natural polymorphisms shaping GSL variation in the TOU-A population suggests that GWAS in local mapping populations could complement and expand the genetic architecture for metabolic variation revealed from geographically broad mapping panels.

## 2. Material and methods

### (a) Plant growth

To minimize maternal effects, seeds were harvested from 294 TOU-A accessions grown at 22°C with a 16 : 8 h light : dark photoperiod, with 3 weeks vernalization at 4°C in 8 h : 16 h light : dark to synchronize flowering, in autumn 2017. For GSL profiling in mid-2019, seeds were sown on a 1 : 1 blend of nutrient retention (BM1) and seed germination (BM2) soil mixes (Berger, CA) in a complete randomized block design with four replicates per accession (i.e. $N = 1$ replicate per accession per complete block). After 4 days stratification at 4°C, growth trays were moved to a chamber with white LED light (180–200 µmol s$^{-1}$) at 20°C in 10 h : 14 h light : dark. Seedlings were thinned to one per cell 1 week after germination. Trays were rotated and bottom-watered every second day with fertilizer (15N-16P-17 K) solution at 100 ppm N until harvesting at 21 days.

### (b) Glucosinolates extraction and quantification

Glucosinolates were extracted and quantified for each harvested rosette individually, yielding $N = 4$ biological replicates per accession. All liquid preparation and storage steps throughout the following protocol were conducted in polypropylene 96-well plates sealed with silicone cap mats. Entire rosettes were first clipped from the root, weighed and directly submerged into 1.2 ml 80% methanol, which inhibits endogenous myrosinase activity [39]. After 2 days dark incubation at ambient temperature, samples were centrifuged for 1 min at 4000 × g, and the supernatant was transferred into a fresh plate and stored at −80°C. Immediately prior to GSL profiling, 240 µl was evaporated with a 96-pin air drier in a fresh plate and redissolved in 120 µl 25% methanol. This approach was chosen after favourable comparisons to alternative extraction methods with freezing and/or homogenization steps (see electronic supplementary material, Note).

GSL content was quantified with an Agilent 1200 Series HPLC machine coupled to an Agilent 6410 triple quadrupole mass spectrometer with parameters described in Humphrey et al. [40]. Samples were eluted with 0.1% formic acid in water (A) and 100% acetonitrile (B) using the following separation gradient: 3.5 min of 99% A followed by a gradient from 99% to 65% A (1 to 35% B) over 12.5 min, and a wash with 99% B for 4 min with 5 min post-run re-equilibration to 99% A. The mass spectrometer was run in precursor negative-ion electrospray mode, monitoring all parent ions from $m/z$ 350–520 with daughter ions of $m/z$ 97, which correspond to the sulfate moiety of the GSL analytes. External standards (sinigrin, every 12th sample; and a GSL extract from a mixture of TOU-A genotypes, every 24th sample) interspersed throughout each run were monitored to ensure consistency. Individual GSLs were identified based on their fragmentation pattern and retention time [32] (electronic supplementary material, table S1). Intensities for each molecule were integrated using MSnbase v. 2.8.3 [41] and xcms v. 3.4.4 [42] in R, using a customized approach that did not require delineating discrete peak boundaries and thus enabled increased sensitivity for low abundance molecules (see electronic supplementary material, Note).

### (c) Genotypes

Genotypes for the TOU-A population were obtained from Frachon et al. [37]. Genotype data for the RegMap [16] and 1001G [17] datasets were obtained from Arouisse et al. [43]. For the 1001G dataset, this consisted of SNPs that were directly genotyped through whole-genome resequencing (WGS). For the RegMap panel, this consisted of SNPs that were directly genotyped with a 250 K SNP chip and supported by WGS in resequenced accessions, and SNPs imputed by intersecting the RegMap chip genotypes and 1001G WGS genotypes. Of SNPs, 2.8 M with greater than 95% imputation accuracy were retained, which primarily excludes SNPs with low-frequency alleles.
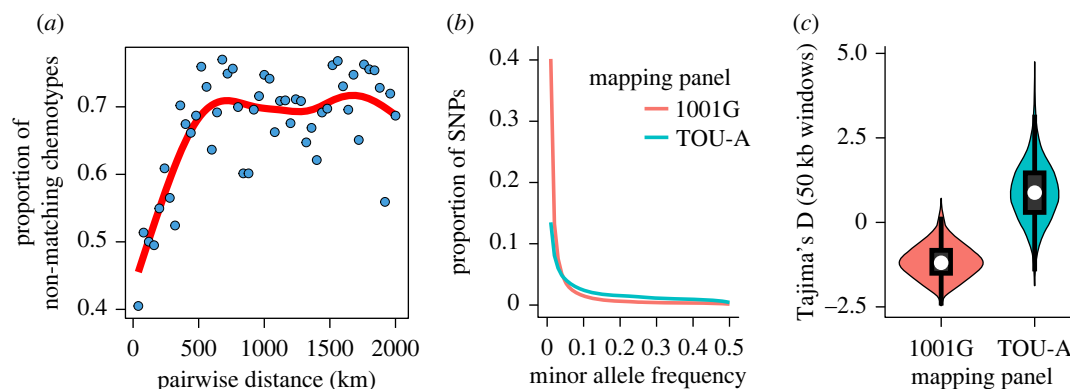
### (d) Broad-sense heritability of glucosinolates

We fitted linear mixed models for log-transformed ion counts per milligram of leaf tissue using lme4 [44], including random intercept effects for the plant accession identity and for the plate containing the sample during extraction and HPLC-MS/MS quantification. The models included all biological replicates per accession (mean $N = 3.93$ for TOU-A, see above; mean $N = 3.69$ for Brachi et al. [32]; mean $N = 2.00$ for Katz et al. [13]). Variance components were extracted from the model, and heritability was estimated as the proportion of total variance explained by accession identity, after excluding variance explained by sample plate identity. This approach leverages the independent biological replicates per accession to estimate variance components without the use of SNP genotypes or the matrix of genetic relatedness among accessions and thus is not biased by potential differences in patterns of population structure among mapping panels. Significance of accession identity was assessed by a likelihood ratio test with one degree of freedom. For published GSL measurements of RegMap [32] and 1001G [13] accessions, an identical model was implemented using GSL abundances scaled by sample weights as reported by the authors.

### (e) Genome-wide association mapping

To standardize comparisons across datasets, analyses were conducted identically for the TOU-A, 1001G and RegMap datasets. First, best unbiased linear predictors (BLUPs) were extracted from the linear mixed models above; for one dataset [6] that pooled biological replicates, abundances from the single technical replicate per accession were used directly. Values were converted to z-scores so that GWAS would produce effect size estimates in units of phenotypic standard deviations. Second, GWAS were implemented as linear mixed models in GEMMA v0.98.1 [45], including a centred genetic relatedness matrix (-gk 1) to account for population structure. Significance per SNP was assessed by Wald Tests (-lmm 1). Finally, to estimate the proportion of variance explained by a given SNP, we fitted a linear mixed model using lme4qtl v. 0.2.2 [46] with the centred genetic relatedness matrix as a random effect and the SNP genotype as the only fixed effect, and extracted the proportion of variance explained by fixed effects ($R_m^2$) using MuMIn v. 1.43.17 [47].

Traits that were modelled separately for GWAS included (i) abundances of each of the heritable GSL molecules, and (ii) log$_2$-transformed ratios of the abundances of pairs of molecules with precursor:product relationships (electronic supplementary material, figure S1). For indolic GSLs in TOU-A, we also implemented a multi-trait GWAS approach (multivariate linear mixed model, mvLMM [48]), which jointly models the relationships between the abundances of all detected molecules. Severe genomic inflation and/or algorithmic termination errors prevented the implementation of these models for other molecules and mapping panels.

**Figure 1.** Reduced genetic complexity within local *Arabidopsis* populations. (*a*) The proportion of non-matching GSL chemotypes, which reflect the joint genotype at three epistatically interacting loci (MAM, AOP, GS-OH), increases sharply and then plateaus as a function of geographical distance in pairwise comparisons among accessions. Points represent comparisons among European 1001G accessions in 40 km bins. (*b*) The allele frequency spectrum is skewed toward common alleles in TOU-A relative to European accessions in the 1001G. The plotted lines were produced by connecting points indicating the proportion of SNPs falling into 1% bins of minor allele frequency. (*c*) Tajima's D is also elevated in TOU-A, shown as a distribution of values across 50 kb genomic windows. The 1001G panel was downsampled to 192 individuals to match TOU-A, and both populations were downsampled to 100 individuals per SNP site, to avoid sample size and genotyping efficiency biases in panels (*b,c*). (Online version in colour.)

Unless otherwise stated, all GWAS excluded SNPs with minor allele frequency (MAF) less than 0.05 or missing genotypes in greater than 5% of the accessions (relaxed to 10% for TOU-A, which had more uncalled sites). We excluded a small number of GWAS exhibiting systematic genomic inflation as determined from the $\chi^2$-test statistic corresponding to the median *p*-value ($\lambda_{GC} > 1.04$) [49] or an excess of associated SNPs (98th percentile of genome-wide *p*-values < 0.01).

To search for significant associations harbouring GSL biosynthetic loci, we used a recently compiled catalog of functionally validated genes in the aliphatic and indolic GSL biosynthetic pathways ([31]; categories: side chain elongation, core structure synthesis, side chain modification). Because peaks of association at known GSL biosynthetic loci in previous GWAS reside tens or even hundreds of kb from the causal genes [13,32,34]—which may arise from extended causal haplotypes [34], structural variants, or intergenic regulatory variants—we defined candidate SNPs as those within 30 kb of known biosynthetic genes. For the three loci with significant SNPs in our re-analysis of the 1001G and RegMap datasets, for which the causal genes are well-established, we further extended these windows in 10 kb increments until they captured 90% of the SNPs within 0.5 Mb of the known causal loci (AOP2/3, GS-OH, MAM1/3) that harboured significant associations with single GSL molecules or precursor:product ratios in those datasets.

## (f) Population genetic comparisons
Methods for all population genetic analyses are described in the electronic supplementary material, Methods.

## 3. Results

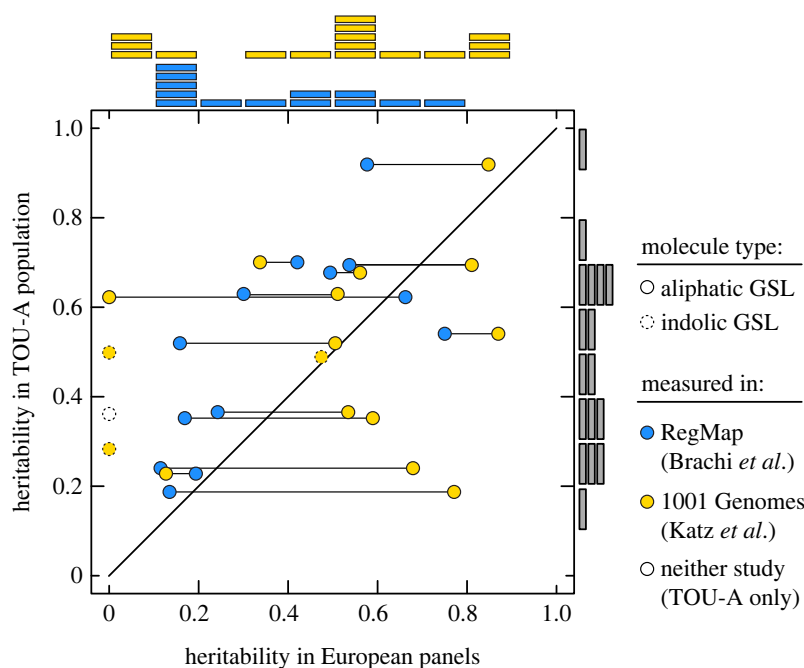### (a) A deficit of rare alleles in the local TOU-A population
A population genetic comparison between TOU-A and the European 1001G accessions revealed favourable conditions for GWAS relative to geographically broad mapping panels. First, for the particular example of glucosinolates, we found that epistatic variation increases rapidly with geographical distance (figure 1*a*). Second, despite reduced overall diversity (1.9 M SNPs in TOU-A versus 11.5 M SNPs in 1001G), the TOU-A population (1.3 M) and 1001G panel (2.2 M) had a relatively

comparable number of common variants (defined here as biallelic SNPs with MAF greater than 0.03). Indeed, a large fraction of common variants from the 1001G panel (2.2 M) were also common in TOU-A (0.83 M, 38%), indicating the reduced genetic diversity in TOU-A arises from a lessened contribution of rare variants. This was reflected in the allele frequency spectrum: after downsampling the 1001G to account for differences in sample size, the TOU-A population still displayed a less pronounced enrichment of rare relative to higher frequency variants (figure 1*b*), resulting in higher genome-wide values of Tajima's D (figure 1*c*). This strong reduction in both total and rare variants is expected to reduce confounding effects of genetic and allelic heterogeneity in TOU-A, while the presence of many common variants suggests this does not come at the expense of drastically culling the polymorphisms that can be interrogated through GWAS.

### (b) Heritable variation in glucosinolate profiles within the local TOU-A population
We quantified the relative concentrations of 13 major aliphatic and four indolic glucosinolates in 294 accessions from the TOU-A population under controlled growth chamber conditions. By contrast to broader geographical scales, where loss-of-function mutations within the glucosinolate biosynthetic pathway are pervasive, every TOU-A accession exhibited a fully functional GSL biosynthetic pathway. This was evidenced by abundant concentrations of the final products in the biosynthetic pathways for both short-chain aliphatic (hydroxyalkenyl) and indolic GSLs (electronic supplementary material, figure S2).

Genetic differences among individuals explained statistically significant portions of the between-accession variation in abundance for every GSL molecule: broad-sense heritabilities ranged from $0.19 < H^2 < 0.92$ (all $P_{Bonferroni} < 0.05$). In fact, analysis of GSL measurements from previous studies revealed systematically higher heritability estimates in TOU-A than the RegMap (Sign Test, median difference = 0.16 [95 %CI:0.04,0.31], *p* = 0.02) and no significant difference between TOU-A and the 1001G (median difference = 0.04 [−0.20,0.20], *p* = 0.46) (figure 2). Although the experimental design, tissue sampling,

**Figure 2.** Glucosinolate variation is highly heritable within the TOU-A local population. (*a*) Estimates of broad-sense heritability ($H^2$) for each GSL molecule in the TOU-A population are plotted against estimates in broader European mapping panels. Connected points indicate estimates of $H^2$ for the same molecule in different European panels. Points above the diagonal line exhibit higher $H^2$ in TOU-A. Histograms above and to the right of the plot indicate the distribution of $H^2$ values in each population. (Online version in colour.)

or data collection variables across studies could contribute to differences in heritability among the mapping populations, these data clearly indicate a high level of heritability for GSL traits within the TOU-A population, even in the absence of the loss-of-function alleles at biosynthetic loci that have dramatic effects on GSL profiles across broader geographical scales.

## (c) Genome-wide association studies within the local TOU-A population reveals known and novel variants shaping aliphatic glucosinolate profiles

For 192 phenotyped accessions with whole-genome sequences, we conducted GWAS using mixed models that controlled for confounding due to population structure by including a matrix of kinship among accessions as a random effect. We first focused on the abundances and relationships between 13 aliphatic GSLs.

### (i) Significant associations

The identity of associated loci in TOU-A depended on how GSL phenotypes were represented. First, we performed separate GWAS for the abundance of each molecule. This approach cumulatively uncovered significant associations at five biosynthetic loci (figure 3*a*). By contrast, only four cumulative associations (three per dataset) were recovered using the same approach in a re-analysis of three previous GWAS datasets, which consisted of mapping populations spanning the European continent ($N > 300$ accessions; electronic supplementary material, figure S3a).
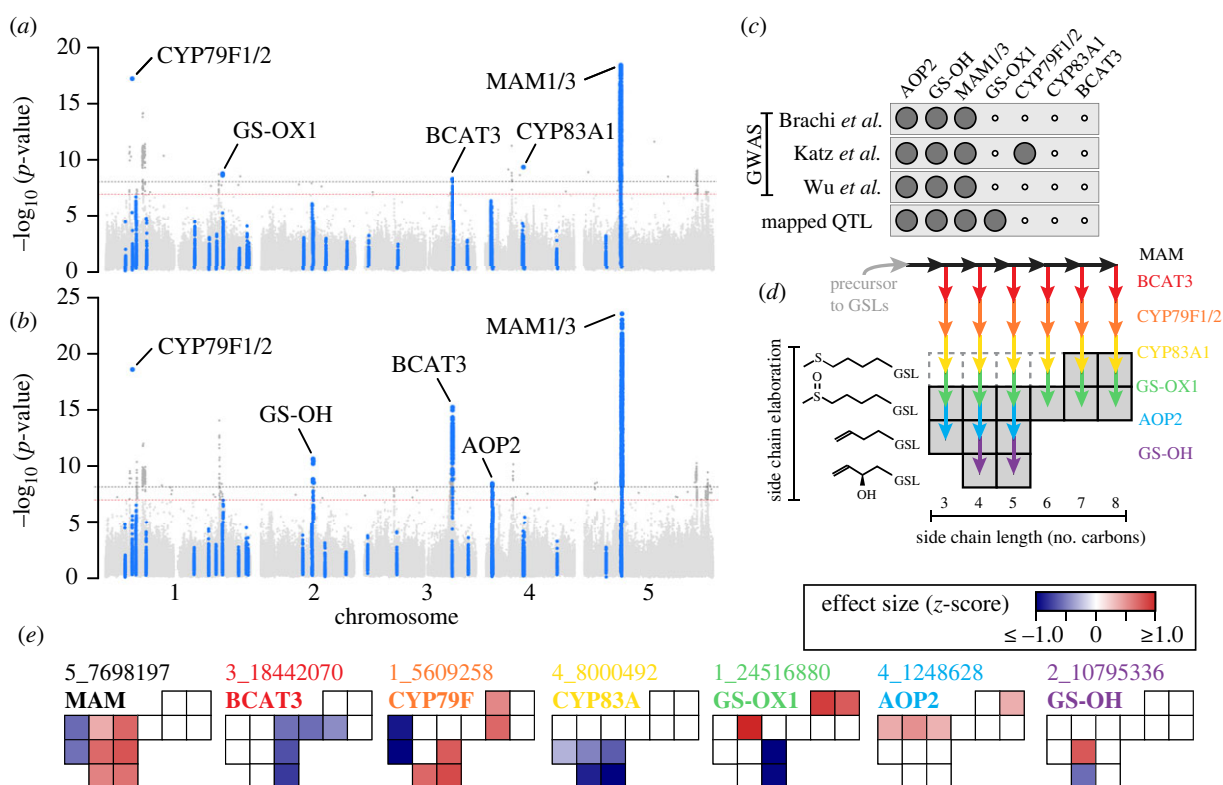
When traits are correlated, as expected for pleiotropic traits such as metabolites from the same biosynthetic pathway, GWAS approaches that use information from multiple traits in a single analysis may increase power [50]. Given the strong positive and negative genetic correlations among GSL molecules in the TOU-A population (electronic

supplementary material, figure S4), we reasoned that such an approach may reveal additional associations. Indeed, using ratios of the abundances of individual precursor versus product GSLs as the mapped traits cumulatively revealed significant associations at five biosynthetic loci in TOU-A, including two loci not recovered from GWAS using individual GSL abundances (figure 3*b*). The same approach in geographically broad European panels recovered only three cumulative associations (two to three per study; electronic supplementary material, figure S3a).

Overall, the significant associations in TOU-A included all three loci (GS-OH, AOP, MAM) that we also recovered in all broad European panels, and an additional locus (CYP79F1/F2) recovered in only one other panel (figure 3*c* and electronic supplementary material, figure S3a). Many of these same associations were reported in the authors' original analyses [6,13,32], although the CYP79F1/F2 polymorphism had not been previously reported. In addition, analyses in TOU-A uncovered three loci not found in GWAS of other mapping panels. The GS-OX locus, which was successfully mapped in biparental RILs, had not been mapped in the three GWAS with large mapping populations [33,51,52]. Further, we provide the first evidence for effects of natural polymorphisms in BCAT3 and CYP83A1. All of these associations had large phenotypic effects, with the leading SNP explaining an average of 24% of the phenotypic variation for its most strongly associated trait (range: 6–43%; electronic supplementary material, table S2).

### (ii) Mapping precision

Identifying candidate genes for functional validation is more efficient when GWAS associations narrowly peak at or near the causal mutations. Precise peaks of association arise when a causal variant recombines into different haplotypes. In populations with reduced genetic diversity, the presence of a given causal variant on fewer haplotypes could result in broader peaks of association, posing a drawback for

**Figure 3.** Seven biosynthetic loci are associated with aliphatic glucosinolate variation in the TOU-A local population. (*a,b*) The best *p*-value per SNP across individual GWAS, mapping either the abundance of individual GSL molecules (panel *a*, 13 traits) or the ratio of individual precursor versus product molecule abundances (panel *b*, 17 traits). SNPs assigned to known GSL biosynthetic loci (see Methods) are enlarged and coloured blue. Dotted lines indicate the Bonferroni genome-wide significance threshold for a single GWAS (red) or the full study (i.e. all individual GWAS across which *p*-values were merged; black). (*c*) For each locus associated with GSL variation in TOU-A, black circles indicate if the same locus was significant in GWAS in our re-analysis of GSL datasets from large ($N > 300$) European mapping populations [6,13,32] or was previously mapped as a QTL using biparental recombinant inbred lines (RILs) [33]. (*d*) A model for how these loci interact to generate variation in GSL profiles for the major aliphatic GSLs present in TOU-A plants (shaded boxes). Enzyme-catalyzed reactions from precursor to product are shown as coloured arrows. Dashed boxes indicate known intermediates that were not observed or quantifiable in TOU-A. (*e*) Effects on individual aliphatic GSLs for the minor allele of the leading SNP at each locus (identified as the SNP with the top association across any individual GWAS from panels (*a,b*), named as 'chromosome_position'). Boxes are oriented to represent the GSL molecules in panel (*d*). Effect sizes are shown for each single molecule GWAS with $p < 0.01$ for the focal SNP.

GWAS in local populations. Relative to European panels, GWAS peaks in TOU-A were indeed broader, although patterns varied among loci (electronic supplementary material, figure S5a). Nevertheless, the leading SNP for over half the associations fell within the transcribed or flanking intergenic regions of the biosynthetic gene (four of seven loci), and the cluster of significant SNPs overlapped with these regions in all but two cases (electronic supplementary material, figure S5a-b). Further, an inspection of the genome-wide *p*-value distribution revealed no systematic genomic inflation in any of the mapping populations (electronic supplementary material, figure S5c). Thus, GWAS in TOU-A retains the ability to narrowly pinpoint candidate genes.

### (iii) Effects on glucosinolate profiles

A model for how the putatively causal enzymes at the seven significant loci generate GSL profile variation in the TOU-A population emerges simply by overlaying the reaction catalyzed by each enzyme, from precursor to product molecules, onto a plot of the major aliphatic GSLs detected in TOU-A plants. This produces a visual map of the variable steps in the biosynthetic pathway (figure 3d). We sought to use these relationships, supplemented with GSL profiles from gene knockout mutants in previous studies, to validate

each locus by comparing them to the effects inferred in our GWAS. To do this, we identified the leading SNP (i.e. the SNP with the strongest experiment-wide *p*-value) at each locus, extracted its GWAS model-fitted effect on the abundance of each GSL molecule, and visualized the effects on the map of GSL molecular variation in TOU-A (figure 3e). In addition to offering further evidence supporting the hypothesized causal genes at each locus, this approach illuminates how these loci generate different aspects of GSL profile variation in the TOU-A population.

The effects of the BCAT3 locus in TOU-A suggest that this gene underlies a dimension of variation in GSL side-chain length previously undescribed in natural populations of *Arabidopsis*, distinct from effects of the well-characterized variation at the MAM locus. The BCAT3 locus affected the abundances of GSLs with intermediate-length side chains, mirroring effects previously observed in a BCAT3 knockout mutant (figure 3e and electronic supplementary material, figure S6). By contrast, functional genetic and biochemical assays have shown that the MAM1 and MAM2 enzymes primarily affect the abundance of GSLs with short side chains [53], similar to the inferred effect of the MAM locus in TOU-A, and MAM3 primarily affects the abundance of GSLs with long side chains (figure 3e and electronic supplementary material, figure S6).

Of two previously unreported associations at cytochrome P450 monooxygenases functioning downstream of MAM and BCAT3 in the biosynthetic pathway (figure 3d), the novel association at the paralogous CYP79F1 and CYP79F2 genes [54] is especially noteworthy. The leading SNP at this locus was associated with a larger magnitude of effect on some short-chain molecules in TOU-A than MAM or BCAT3 (figure 3e), with especially large effects on molecules with the shortest observed side-chain length. This is consistent with the finding that among all biosynthetic enzymes, CYP79F2 exerts the strongest effect on pathway flux, with an outsized effect on propyl GSLs (i.e. GSLs with 3C side-chain lengths) [12]. Functional polymorphism at a CYP79F gene also underlies a QTL affecting the propyl fraction of GSLs in *Brassica juncea* [55], and separately underlies adaptive variation in the proportion of GSLs derived from branched-chain amino acids relative to methionine in *Boechera stricta* [56]. The association at CYP79F paralogs was recovered in our re-analysis of one European *Arabidopsis* dataset (electronic supplementary material, figure S3), strengthening the evidence that CYP79F is a broadly important determinant of GSL profile variation across populations and species.

Two distinct loci harbour paralogous GS-OX genes that catalyze the *S*-oxygenation of methylthioalkyl to methylsulfinylalkyl GSLs with broad substrate specificity. While natural variation in the locus containing GS-OX2, GS-OX3 and GS-OX4 had been detected through QTL mapping with biparental RILs [51,52], neither locus had been detected in the three large, European GWAS panels. In addition to harbouring a significant association when considering common variants (MAF > 0.05; figure 3a), GS-OX1 harboured the strongest genome-wide association for many molecules when slightly rarer variants were considered (MAF > 0.03; electronic supplementary material, figure S7). Although biases in our GWAS model can yield inflated or deflated signals of association for alleles below this threshold, the strength of the association for this variant is exceptional even among alleles of similar frequency (0.05 > MAF > 0.03). Intriguingly, the strongest associations at GS-OX1 did not involve methylthioalkyl GSL abundances individually or as a ratio compared to their derived methylsulfinylalkyl GSLs (electronic supplementary material, figure S7), suggesting that linkage disequilibrium with other loci (or an unexpected effect of GS-OX1) may contribute to this association. Nevertheless, the effect on its direct precursor and/or product molecules is sufficient to drive a significant association: we further performed GWAS for a principal component capturing opposing shifts in the abundance of long-chain methylthioalkyl versus methylsulfinylalkyl GSLs, and GS-OX1 harboured the strongest, statistically significant genome-wide association (electronic supplementary material, figure S7).

Finally, effects of the two remaining polymorphisms in TOU-A, at the AOP [57] and GS-OH [58] loci, differed from the effects of loss-of-function variants at these loci that segregate over broad geographical scales, which eliminate the production of their GSL products and generate qualitative presence/absence variation in GSL profiles [13]. In TOU-A, by contrast, both loci affected their precursor GSL abundances, with only GS-OH also oppositely affecting (but not abolishing) its product GSL abundances (figure 3e).

It is important to note that the predicted effects do not include epistatic interactions and that more subtle effects may not be discovered through GWAS. Accordingly, the effects described above should be interpreted only as the strongest, additive effects of each locus.

## (d) Genome-wide association studies within the local TOU-A population reveals known and novel variants shaping indolic glucosinolate profiles

### (i) Significant associations

We implemented the same association mapping approach for four indolic GSL molecules and were most successful when mapping traits that captured the relationships among abundances of different molecules. Our initial approach mapping the abundance of single molecules only recovered one association in both the TOU-A population and in the geographically broad European panel with high-quality indolic GSL data (electronic supplementary material, figure S3b). We recovered an additional association in TOU-A when mapping GSL precursor:product ratios (electronic supplementary material, figure S3b).

Multi-trait mixed models, which jointly model the relationships among two or more traits together, may further increase power by using the relationships among traits as additional information. Importantly, these models can recover genetic variants affecting both individual traits and the relationships among traits, which may have distinct genetic bases [48]. We employed a multi-trait GWAS jointly modelling the abundance of all four indolic GSLs detected in TOU-A. This recovered a third association, along with the two previously noted, in the TOU-A population (figure 4a). The model encountered algorithmic termination errors when applied to the geographically broad panel, preventing a comparison with TOU-A for this approach.

Overall, of these three loci recovered in TOU-A, two (both CYP81F loci) have been previously identified in GWAS [6] (figure 4b). One of these loci was also discovered through QTL mapping, and CYP81F2 was functionally validated as the causal gene [59,60]. The IGMT locus had not been linked to natural variation in GSL profiles previously.
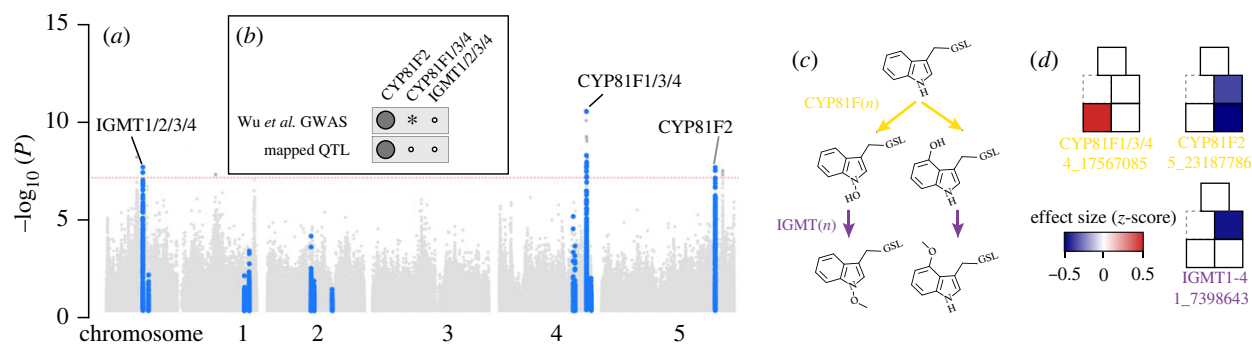
### (ii) Mapping precision

Similar to aliphatic GSLs, peaks of significantly associated SNPs varied from narrow to broad, depending on the locus (electronic supplementary material, figure S5). Notably, significance peaked directly within the tandem array of IGMT paralogs, again highlighting the ability of GWAS in TOU-A to pinpoint candidate genes.

### (iii) Effects on GSL profiles

Each putatively causal biosynthetic enzyme underlying the associations with indolic GSL variation in TOU-A has been functionally characterized through biochemical assays and gene knockout mutants. CYP81F paralogs collectively catalyze the first elaboration step at different sites of the indolic GSL ring structure [59,60], and IGMT paralogs collectively catalyze a subsequent elaboration step [61] (figure 4c). Using the effects of each locus extracted from our GWAS models (figure 4d), we looked for concordance with previous QTL mapping, functional genetic, and knockout mutant studies to inform how these loci shape GSL variation in TOU-A.

The CYP81F subfamily of cytochrome P450 monooxygenases are responsible for hydroxylation of indolyl-3-ylmethyl

**Figure 4.** Three biosynthetic loci are associated with indolic glucosinolate variation in the TOU-A local population. (a) p-values from a multi-trait GWAS (mvLMM) jointly modelling all indolic GSL abundances. The plot layout, colours, and significance thresholds are as described in figure 3a. (b) For each locus associated with GSL variation in TOU-A, black circles indicate if the same locus was significant in GWAS in our re-analysis of a GSL dataset from a large (N > 300) European mapping population [6] or was previously mapped as a QTL using biparental RILs [59]. An asterisk indicates a significant association in a published analysis that was not recovered in our standardized re-analysis. (c) The pathway for secondary modification of indole-3-ylmethyl GSL (top) through 1- or 4-hydroxylation (middle) and subsequent methoxylation (bottom). (d) Effects on individual indolic GSLs for the minor allele of the leading SNP at each locus, determined as in figure 3e. Boxes are oriented to represent the GSL molecules in panel (c).

(I3M) GSL [59,60], which can subsequently be methoxylated by other enzymes. The locus harbouring CYP81F2 affected two GSL molecules in TOU-A (4-hydroxy-I3M-GSL and its derivative, 4-methoxy-I3M-GSL), which also differentially accumulate due to the CYP81F2 locus in a previous QTL mapping experiment [59]. The locus harbouring CYP81F1, CYP81F3 and CYP81F4 paralogs affected the GSL that is methoxylated at a different site, 1-methoxy-I3M-GSL; the CYP81F-catalyzed product from which it derives, 1-hydroxy-I3M-GSL, is unstable and was not observable through our GSL profiling approach. These results further support evidence from previous mapping studies that paralogs at the two CYP81F loci affect different GSL molecules *in planta*, despite overlap in substrate specificities *in vitro* [59,60].

Four of the five indole glucosinolate O-methyltransferases (IGMT1-4) in *Arabidopsis* form a tandem array at the locus identified in our GWAS [61]. This locus had a strong effect on the abundance of its substrate, 4-hydroxy-I3M-GSL (figure 4d). Although IGMT1-4 enzymes cumulatively can methoxylate both 1- and 4-hydroxy-I3M-GSL in biochemical assays, our observation of effects restricted to 4-hydroxy-I3M-GSL methoxylation support a model previously inferred from the characterization of an IGMT5 knockout mutant, which retained functional copies of all four IGMT1-4 paralogs [61]. The mutant exhibited an absence of 1-methoxy-I3M-GSL but no reduction in 4-methoxy-I3M-GSL, suggesting the IGMT1-4 locus is responsible only for 4-methoxy-I3M-GSL's production *in planta*.

Taken together, our results more fully link the functional variation characterized in enzyme biochemical and gene knockout studies with the variation for indolic GSLs observed in natural populations, identifying loci acting at three of the four secondary modification steps that give rise to the major I3M-derived GSLs in the TOU-A population.

## (e) Reduced population structure is unlikely to underlie improved performance of genome-wide association studies for glucosinolate profiles in the local TOU-A population

GSL profiles, and some of the large effect loci that underlie them, show strong geographical clines within and across Europe [13,32]. This raises the possibil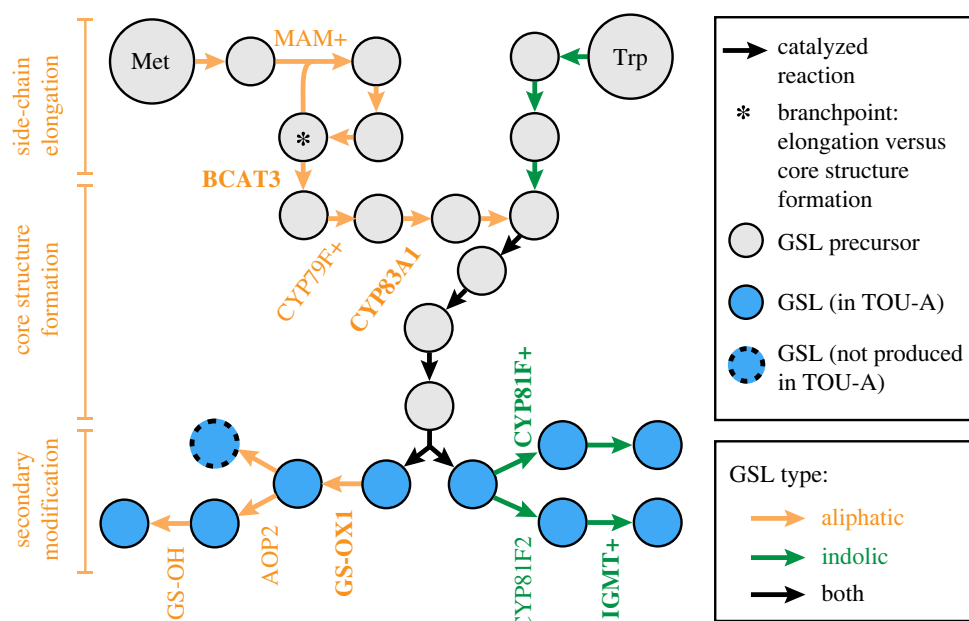ity that methods to control for population structure in GWAS could weaken signals of association with GSLs at loci whose genotypes are strongly correlated with population structure. To investigate this, we used ADMIXTURE to infer subgroups (k = 5) contributing to population structure separately within the TOU-A and the 1001G accessions. Focusing on the 10 glucosinolate biosynthetic loci recovered by GWAS in TOU-A, we found that among-group variation in allele frequency was not elevated in the 1001G relative to TOU-A (electronic supplementary material, figure S8). This suggests that the efficacy of GWAS for GSLs in TOU-A is unlikely to be the product of weaker population structure at causal loci, and may instead arise from differences in other confounding factors that are exaggerated in geographically broad mapping panels.

## 4. Discussion

As one of the best-studied secondary metabolite pathways in plants—with a wealth of functional genetic knowledge from GWAS and QTL mapping of natural variation, characterization of genetic mutant lines, and enzyme biochemical assays [30]—GSLs offered a compelling opportunity to investigate the performance of GWAS using a local mapping population. The expanded genetic architecture revealed for GSLs in the TOU-A population highlights the benefits of this approach. A modest mapping panel (N = 192 accessions) led not only to the discovery of variants that were absent in geographically broad mapping panels with 1.5–4.0× more accessions but also to novel loci whose contribution to natural variation was unknown despite numerous QTL mapping studies (albeit typically with relatively small sample sizes) previously conducted for GSLs. These associations spanned each major portion of the pathway (figure 5): the MAM-catalyzed reaction loop for side-chain elongation in GSL precursor molecules, sequential steps for synthesis of the GSL core structure, and every level of secondary modification subsequent to the formation of a functional GSL molecule [31]. Thus, GWAS within a local population can offer a deep catalogue of functional polymorphism within a biosynthetic pathway.

The simplest explanation for the effectiveness of GWAS in TOU-A may be the observed reduction in genetic diversity

**Figure 5.** An overview of glucosinolate biosynthetic loci associated with GSL variation in the TOU-A population. The diagram shows each enzyme-catalyzed step, beginning with the amino acid precursor (Met or Trp). Genes harbouring significant GWAS associations in TOU-A are listed at the biosynthetic step they catalyze. Bolded genes are novel associations, defined as those significantly associated in TOU-A but not in our re-analysis of three datasets with geographically broad European mapping panels. A '+' indicates that multiple paralogous genes at a locus could contribute to the association (e.g. CYP79F1 and CYP79F2 are represented as CYP79F+). The pathway and enzyme positions are based on Harun *et al.* [31]. Note that additional steps producing GSLs that accumulate only at very low levels in leaves are omitted. (Online version in colour.)

relative to the broader European population. Theory predicts that allelic heterogeneity, which poses a major obstacle for GWAS, will be more pervasive in more genetically diverse populations. Further, the fact that diversity was reduced in TOU-A primarily through a relative deficit of rare variants, as expected if rare variants are geographically restricted and therefore locally more common [26], likely provides an additional benefit. Rare variants are not only poorly detected through GWAS, but their presence can obscure true associations at causal loci [62]. Consistent with this, GWAS has uncovered more associations and a broader (albeit largely unvalidated) functional repertoire of underlying candidate genes—including biosynthetic enzymes, transcription factors, and transporters—across cultivars of *Brassica napus* than in European panels of *Arabidopsis* [63–65]. *Brassica napus* cultivars are less genetically diverse and have an excess of common variants (reflected in elevated Tajima's D) relative to *Arabidopsis* [17,65,66], which may have been further exaggerated at glucosinolate-related genes by the diversity-reducing effects of directional selection during the breeding process [66].

While the general benefits of reduced geography-driven confounding in local populations should extend to GWAS for a variety of traits, our findings also illustrate properties of local populations likely to be especially beneficial when studying metabolite diversity specifically. In particular, the confounding effects of loss-of-function polymorphisms were absent from the major loci (MAM, AOP, GS-OH) that segregate such mutations over broad geographical scales. Loss-of-function mutations produce a particularly severe form of allelic heterogeneity. Many different mutations can produce analogous loss-of-function alleles at a gene, resulting in a high gene-wide mutation rate, such that many loss-of-function polymorphisms involve multiple haplotypes with parallel loss-of-function mutations [27]. Furthermore, loss-

of-function mutations underlie dramatic epistatic effects, which may dilute additive effects modelled by GWAS. An extreme example involves the GS-OH locus that catalyzes the final secondary modification in the biosynthetic pathway (figure 5): loss of function alleles at upstream enzymes fully mask the effect of GS-OH on GSL variation in the majority of genetic backgrounds in *Arabidopsis*, and GS-OH itself segregates numerous loss-of-function alleles [13]. Of the three major large-effect loci mapped in other GWAS of aliphatic GSLs, only GS-OH has failed to consistently yield associations across previous analyses [6,13,32,34].

Although statistical approaches exist to mitigate geographically driven confounding factors, they cannot entirely control for them. For example, GWAS models can be extended to include epistatic interactions alongside, or instead of, additive effects [67]. However, the immense number of possible pairwise interactions across the genome creates computational challenges and a severe multiple testing burden [68]. Other confounding factors can be lessened by altering genotype information rather than the GWAS models themselves. One simple yet powerful approach involves collapsing all predicted loss-of-function variants at a gene into a single allele, reducing their contribution to allelic heterogeneity [69]. Nevertheless, this approach requires genotyping to be conducted through whole-genome sequencing, and even then, many cases of abolished or altered gene function are difficult to annotate from DNA sequence data alone. Furthermore, while this approach can improve power to discover associations at loci with heterogeneous loss-of-function variants, it does not address their confounding epistatic effects on other loci. Even in cases where various genotyping and statistical approaches do largely succeed in mitigating specific confounding factors, integrating them to address many factors simultaneously is challenging. For many research questions, the use of local mapping populations in which these

confounding factors are lessened offers an attractive alternative to these more tailored GWAS approaches.

Despite their benefits, GWAS in local populations are certainly not ideal for every research question. GWAS of GSLs in different mapping populations illustrate this clearly: integrating population genomic analyses with GWAS using *Arabidopsis* accessions sampled throughout Europe revealed how GSL profiles have been shaped by adaptation and demography across the species range [13,32,34], which would be impossible to infer from a single local population. Meanwhile, GWAS using the TOU-A population implicated more loci in natural phenotypic variation than could be detected in broader mapping panels. Complementary GWAS in local and geographically broad mapping panels thus provide an exciting avenue toward a fuller understanding of the genetic variation and evolutionary processes that shape phenotypic diversity in nature.

Authors' contributions. A.G.: conceptualization, data curation, formal analysis, investigation, methodology, project administration, supervision, validation, visualization, writing—original draft, writing—review and editing; A.V.: data curation, formal analysis, investigation, methodology, visualization, writing—review and editing; T.C.M.: investigation, methodology, project administration, resources, supervision, validation, writing—review and editing; P.L.: data curation, formal analysis, investigation, methodology, visualization, writing—review and editing; F.R.: conceptualization, resources, supervision; J.B.: conceptualization, funding acquisition, investigation, project administration, resources, supervision, validation, writing—original draft, writing—review and editing.

All authors gave final approval for publication and agreed to be held accountable for the work performed therein.

# References

1. Weng J-K, Philippe RN, Noel JP. 2012 The rise of chemodiversity in plants. *Science* **336**, 1667–1670. (doi:10.1126/science.1217411)

2. Fernie AR, Tohge T. 2017 The genetics of plant metabolism. *Annu. Rev. Genet.* **51**, 287–310. (doi:10.1146/annurev-genet-120116-024640)

3. Pott DM, Durán-Soria S, Osorio S, Vallarino JG. 2021 Combining metabolomic and transcriptomic approaches to assess and improve crop quality traits. *CABI Agric. Biosci.* **2**, 1. (doi:10.1186/s43170-020-00021-8)

4. Luo J. 2015 Metabolite-based genome-wide association studies in plants. *Curr. Opin. Plant Biol.* **24**, 31–38. (doi:10.1016/j.pbi.2015.01.006)

5. Fang C, Luo J. 2019 Metabolic GWAS-based dissection of genetic bases underlying the diversity of plant metabolism. *Plant J.* **97**, 91–100. (doi:10.1111/tpj.14097)

6. Wu S et al. 2018 Mapping the *Arabidopsis* metabolic landscape by untargeted metabolomics at different environmental conditions. *Mol. Plant* **11**, 118–134. (doi:10.1016/j.molp.2017.08.012)

7. Chen W et al. 2014 Genome-wide association analyses provide genetic and biochemical insights into natural variation in rice metabolism. *Nat. Genet.* **46**, 714–721. (doi:10.1038/ng.3007)

8. Chen W et al. 2016 Comparative and parallel genome-wide association studies for metabolic and agronomic traits in cereals. *Nat. Commun.* **7**, 11767. (doi:10.1038/ncomms12767)

9. Wen W et al. 2014 Metabolome-based genome-wide association study of maize kernel leads to novel biochemical insights. *Nat. Commun.* **5**, 3438. (doi:10.1038/ncomms4438)

10. Kliebenstein DJ. 2014 Synthetic biology of metabolism: using natural variation to reverse engineer systems. *Curr. Opin. Plant Biol.* **19**, 20–26. (doi:10.1016/j.pbi.2014.03.008)

11. Wright KM, Rausher MD. 2010 The evolution of control and distribution of adaptive mutations in a metabolic pathway. *Genetics* **184**, 483–502. (doi:10.1534/genetics.109.110411)

12. Olson-Manning CF, Lee C-R, Rausher MD, Mitchell-Olds T. 2013 Evolution of flux control in the glucosinolate pathway in *Arabidopsis thaliana*. *Mol. Biol. Evol.* **30**, 14–23. (doi:10.1093/molbev/mss204)

13. Katz E et al. 2021 Genetic variation, environment and demography intersect to shape *Arabidopsis* defense metabolite variation across Europe. *Elife* **10**, e67784. (doi:10.7554/eLife.67784)

14. Korte A, Farlow A. 2013 The advantages and limitations of trait analysis with GWAS: a review. *Plant Methods* **9**, 1–9. (doi:10.1186/1746-4811-9-29)

15. Brachi B, Morris GP, Borevitz JO. 2011 Genome-wide association studies in plants: the missing heritability is in the field. *Genome Biol.* **12**, 232. (doi:10.1186/gb-2011-12-10-232)

16. Horton MW et al. 2012 Genome-wide patterns of genetic variation in worldwide *Arabidopsis thaliana* accessions from the RegMap panel. *Nat. Genet.* **44**, 212–216. (doi:10.1038/ng.1042)

17. 1001 Genomes Consortium. 2016 1,135 Genomes reveal the global pattern of polymorphism in *Arabidopsis thaliana*. *Cell* **166**, 481–491. (doi:10.1016/j.cell.2016.05.063)

18. Vilhjálmsson BJ, Nordborg M. 2013 The nature of confounding in genome-wide association studies. *Nat. Rev. Genet.* **14**, 1–2. (doi:10.1038/nrg3382)

19. Sul JH, Martin LS, Eskin E. 2018 Population structure in genetic studies: confounding factors and mixed models. *PLoS Genet.* **14**, e1007309. (doi:10.1371/journal.pgen.1007309)

20. Kang HM, Zaitlen NA, Wade CM, Kirby A, Heckerman D, Daly MJ, Eskin E. 2008 Efficient control of population structure in model organism association mapping. *Genetics* **178**, 1709–1723. (doi:10.1534/genetics.107.080101)

21. Segura V, Vilhjálmsson BJ, Platt A, Korte A, Seren Ü, Long Q, Nordborg M. 2012 An efficient multi-locus mixed-model approach for genome-wide association studies in structured populations. *Nat. Genet.* **44**, 825–830. (doi:10.1038/ng.2314)

22. Atwell S et al. 2010 Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana* inbred lines. *Nature* **465**, 627–631. (doi:10.1038/nature08800)

23. Liu X, Huang M, Fan B, Buckler ES, Zhang Z. 2016 Iterative usage of fixed and random effect models for powerful and efficient genome-wide association studies. *PLoS Genet.* **12**, e1005767. (doi:10.1371/journal.pgen.1005767)

24. Lopez-Arboleda WA, Reinert S, Nordborg M, Korte A. 2021 Global genetic heterogeneity in adaptive traits. *Mol. Biol. Evol.* **38**, 4822–4831. (doi:10.1093/molbev/msab208)

25. Lander ES, Schork NJ. 1994 Genetic dissection of complex traits. *Science* **265**, 2037–2048. (doi:10.1126/science.8091226)

26. Biddanda A, Rice DP, Novembre J. 2020 A variant-centric perspective on geographic patterns of human allele frequency variation. *Elife* **9**, e60107. (doi:10.7554/eLife.60107)

27. Monroe JG, McKay JK, Weigel D, Flood PJ. 2021 The population genomics of adaptive loss of function. *Heredity* **126**, 383–395. (doi:10.1038/s41437-021-00403-2)

28. Eaves LJ. 1994 Effect of genetic architecture on the power of human linkage studies to resolve the contribution of quantitative trait loci. *Heredity* **72**, 175–192. (doi:10.1038/hdy.1994.25)

29. Platt A, Vilhjálmsson BJ, Nordborg M. 2010 Conditions under which genome-wide association studies will be positively misleading. *Genetics* **186**, 1045–1052. (doi:10.1534/genetics.110.121665)

30. Jensen LM, Halkier BA, Burow M. 2014 How to discover a metabolic pathway? An update on gene identification in aliphatic glucosinolate biosynthesis, regulation and transport. *Biol. Chem.* **395**, 529–543. (doi:10.1515/hsz-2013-0286)

31. Harun S, Abdullah-Zawawi M-R, Goh H-H, Mohamed-Hussein Z-A. 2020 A comprehensive gene inventory for glucosinolate biosynthetic pathway in *Arabidopsis thaliana*. *J. Agric. Food Chem.* **68**, 7281–7297. (doi:10.1021/acs.jafc.0c01916)

32. Brachi B, Meyer CG, Villoutreix R, Platt A, Morton TC, Roux F, Bergelson J. 2015 Coselected genes determine adaptive variation in herbivore resistance throughout the native range of *Arabidopsis thaliana*. *Proc. Natl Acad. Sci. USA* **112**, 4032–4037. (doi:10.1073/pnas.1421416112)

33. Kliebenstein DJ. 2009 A quantitative genetics and ecological model system: understanding the aliphatic glucosinolate biosynthetic network via QTLs. *Phytochem. Rev.* **8**, 243–254. (doi:10.1007/s11101-008-9102-8)

34. Chan EKF, Rowe HC, Kliebenstein DJ. 2010 Understanding the evolution of defense metabolites in *Arabidopsis thaliana* using genome-wide association mapping. *Genetics* **185**, 991–1007. (doi:10.1534/genetics.109.108522)

35. Kerwin R et al. 2015 Natural genetic variation in *Arabidopsis thaliana* defense metabolism genes modulates field fitness. *Elife* **4**, e05604. (doi:10.7554/eLife.05604)

36. Rowe HC, Hansen BG, Halkier BA, Kliebenstein DJ. 2008 Biochemical networks and epistasis shape the *Arabidopsis thaliana* metabolome. *Plant Cell* **20**, 1199–1216. (doi:10.1105/tpc.108.058131)

37. Frachon L et al. 2017 Intermediate degrees of synergistic pleiotropy drive adaptive evolution in ecological time. *Nat. Ecol. Evol.* **1**, 1551–1561. (doi:10.1038/s41559-017-0297-1)

38. Aoun N, Desaint H, Boyrie L, Bonhomme M, Deslandes L, Berthomé R, Roux F. 2020 A complex network of additive and epistatic quantitative trait loci underlies natural variation of *Arabidopsis thaliana* quantitative disease resistance to *Ralstonia solanacearum* under heat stress. *Mol. Plant Pathol.* **21**, 1405–1420. (doi:10.1111/mpp.12964)

39. Doheny-Adams T, Redeker K, Kittipol V, Bancroft I, Hartley SE. 2017 Development of an efficient glucosinolate extraction method. *Plant Methods* **13**, 17. (doi:10.1186/s13007-017-0164-8)

40. Humphrey PT, Gloss AD, Frazier J, Nelson-Dittrich AC, Faries S, Whiteman NK. 2018 Heritable plant phenotypes track light and herbivory levels at fine spatial scales. *Oecologia* **187**, 427–445. (doi:10.1007/s00442-018-4116-4)

41. Gatto L, Gibb S, Rainer J. 2021 MSnbase, efficient and elegant R-based processing and visualization of raw mass spectrometry data. *J. Proteome Res.* **20**, 1063–1069. (doi:10.1021/acs.jproteome.0c00313)

42. Smith CA, Want EJ, O'Maille G, Abagyan R, Siuzdak G. 2006 XCMS: processing mass spectrometry data for metabolite profiling using nonlinear peak alignment, matching, and identification. *Anal. Chem.* **78**, 779–787. (doi:10.1021/ac051437y)

43. Arouisse B, Korte A, van Eeuwijk F, Kruijer W. 2020 Imputation of 3 million SNPs in the *Arabidopsis* regional mapping population. *Plant J.* **102**, 872–882. (doi:10.1111/tpj.14659)

44. Bates D, Mächler M, Bolker B, Walker S. 2014 Fitting linear mixed-effects models using lme4. *arXiv* [stat.CO]. (doi:10.18637/jss.v067.i01)

45. Zhou X, Stephens M. 2012 Genome-wide efficient mixed-model analysis for association studies. *Nat. Genet.* **44**, 821–824. (doi:10.1038/ng.2310)

46. Ziyatdinov A, Vázquez-Santiago M, Brunel H, Martinez-Perez A, Aschard H, Soria JM. 2018 lme4qtl: linear mixed models with flexible covariance structure for genetic studies of related individuals. *BMC Bioinf.* **19**, 1–5. (doi:10.1186/s12859-018-2057-x)

47. Barton K. 2020 *MuMIn: multi-model inference*. R Package 1.43.

48. Zhou X, Stephens M. 2014 Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nat. Methods* **11**, 407–409. (doi:10.1038/nmeth.2848)

49. Bacanu SA, Devlin B, Roeder K. 2000 The power of genomic control. *Am. J. Hum. Genet.* **66**, 1933–1944. (doi:10.1086/302929)

50. Stephens M. 2013 A unified framework for association analysis with multiple related phenotypes. *PLoS ONE* **8**, e65245. (doi:10.1371/journal.pone.0065245)

51. Li J, Hansen BG, Ober JA, Kliebenstein DJ, Halkier BA. 2008 Subclade of flavin-monooxygenases involved in aliphatic glucosinolate biosynthesis. *Plant Physiol.* **148**, 1721–1733. (doi:10.1104/pp.108.125757)

52. Hansen BG, Kliebenstein DJ, Halkier BA. 2007 Identification of a flavin-monooxygenase as the S-oxygenating enzyme in aliphatic glucosinolate biosynthesis in *Arabidopsis*. *Plant J.* **50**, 902–910. (doi:10.1111/j.1365-313X.2007.03101.x)

53. Textor S, de Kraker J-W, Hause B, Gershenzon J, Tokuhisa JG. 2007 MAM3 catalyzes the formation of all aliphatic glucosinolate chain lengths in *Arabidopsis*. *Plant Physiol.* **144**, 60–71. (doi:10.1104/pp.106.091579)

54. Chen S et al. 2003 CYP79F1 and CYP79F2 have distinct functions in the biosynthesis of aliphatic glucosinolates in *Arabidopsis*. *Plant J.* **33**, 923–937. (doi:10.1046/j.1365-313X.2003.01679.x)

55. Sharma M, Mukhopadhyay A, Gupta V, Pental D, Pradhan AK. 2016 BjuB.CYP79F1 regulates synthesis of propyl fraction of aliphatic glucosinolates in oilseed mustard *Brassica juncea*: functional validation through genetic and transgenic approaches. *PLoS ONE* **11**, e0150060. (doi:10.1371/journal.pone.0150060)

56. Prasad KVSK et al. 2012 A gain-of-function polymorphism controlling complex traits and fitness in nature. *Science* **337**, 1081–1084. (doi:10.1126/science.1221636)

57. Kliebenstein DJ, Lambrix VM, Reichelt M, Gershenzon J, Mitchell-Olds T. 2001 Gene duplication in the diversification of secondary metabolism: tandem 2-oxoglutarate-dependent dioxygenases control glucosinolate biosynthesis in *Arabidopsis*. *Plant Cell* **13**, 681–693. (doi:10.1105/tpc.13.3.681)

58. Hansen BG, Kerwin RE, Ober JA, Lambrix VM, Mitchell-Olds T, Gershenzon J, Halkier BA, Kliebenstein DJ. 2008 A novel 2-oxoacid-dependent dioxygenase involved in the formation of the goiterogenic 2-hydroxybut-3-enyl glucosinolate and generalist insect resistance in *Arabidopsis*. *Plant Physiol.* **148**, 2096–2108. (doi:10.1104/pp.108.129981)

59. Pfalz M, Vogel H, Kroymann J. 2009 The gene controlling the indole glucosinolate modifier1 quantitative trait locus alters indole glucosinolate structures and aphid resistance in *Arabidopsis*. *Plant Cell* **21**, 985–999. (doi:10.1105/tpc.108.063115)

60. Pfalz M, Mikkelsen MD, Bednarek P, Olsen CE, Halkier BA, Kroymann J. 2011 Metabolic engineering in *Nicotiana benthamiana* reveals key enzyme functions in *Arabidopsis* indole glucosinolate modification. *Plant Cell* **23**, 716–729. (doi:10.1105/tpc.110.081711)

61. Pfalz M, Mukhaimar M, Perreau F, Kirk J, Hansen CIC, Olsen CE, Agerbirk N, Kroymann J. 2016 Methyl transfer in glucosinolate biosynthesis mediated by indole glucosinolate o-methyltransferase 5. *Plant Physiol.* **172**, 2190–2203. (doi:10.1104/pp.16.01402)

62. Mathieson I, McVean G. 2012 Differential confounding of rare and common variants in spatially structured populations. *Nat. Genet.* **44**, 243–246. (doi:10.1038/ng.1074)

63. Kittipol V, He Z, Wang L, Doheny-Adams T, Langer S, Bancroft I. 2019 Genetic architecture of glucosinolate variation in *Brassica napus*. *J. Plant Physiol.* **240**, 152988. (doi:10.1016/j.jplph.2019.06.001)

64. Liu S, Huang H, Yi X, Zhang Y, Yang Q, Zhang C, Fan C, Zhou Y. 2020 Dissection of genetic architecture for glucosinolate accumulations in leaves and seeds of *Brassica napus* by genome-wide association study. *Plant Biotechnol. J.* **18**, 1472–1484. (doi:10.1111/pbi.13314)

65. Wei D, Cui Y, Mei J, Qian L, Lu K, Wang Z-M, Li J, Tang Q, Qian W. 2019 Genome-wide identification of loci affecting seed glucosinolate contents in *Brassica napus* L. *J. Integr. Plant Biol.* **61**, 611–623. (doi:10.1111/jipb.12717)

66. Lu K et al. 2019 Whole-genome resequencing reveals *Brassica napus* origin and

genetic loci involved in its improvement. *Nat. Commun.* **10**, 1154. (doi:10.1038/s41467-019-09134-9)

67. Ritchie MD, Van Steen K. 2018 The search for gene–gene interactions in genome-wide association studies: challenges in abundance of methods, practical considerations, and biological interpretation. *Ann. Transl. Med.* **6**, 157. (doi:10.21037/atm.2018.04.05)

68. Crawford L, Zeng P, Mukherjee S, Zhou X. 2017 Detecting epistasis with the marginal epistasis test in genetic mapping studies of quantitative traits. *PLoS Genet.* **13**, e1006869. (doi:10.1371/journal.pgen.1006869)

69. Barboza L, Effgen S, Alonso-Blanco C, Kooke R, Keurentjes JJB, Koornneef M, Alcázar R. 2013 *Arabidopsis* semidwarfs evolved from independent mutations in GA20ox1, ortholog to green revolution dwarf alleles in rice and barley. *Proc. Natl Acad. Sci. USA* **110**, 15 818–15 823. (doi:10.1073/pnas.1314979110)

70. Gloss AD, Vergnol A, Morton TC, Laurin PJ, Roux F, Bergelson J. 2022 Data from: Genome-wide association mapping within a local *Arabidopsis thaliana* population more fully reveals the genetic architecture for defensive metabolite diversity. Dryad Digital Repository. (doi:10.5061/dryad.4mw6m90b6)

71. Gloss AD, Vergnol A, Morton TC, Laurin PJ, Roux F, Bergelson J. 2022 Genome-wide association mapping within a local *Arabidopsis thaliana* population more fully reveals the genetic architecture for defensive metabolite diversity. FigShare. (https://doi.org/10.6084/m9.figshare.c.5958817)