

# Integrating co-expression networks with GWAS detects genes driving elemental accumulation in maize seeds

Robert J. Schaefer<sup>1</sup>, Jean-Michel Michno<sup>1,2</sup>, Joseph Jeffers<sup>3</sup>, Owen Hoekenga<sup>4</sup>, Brian Dilkes<sup>5</sup>,  
Ivan Baxter<sup>6,7\*</sup>, Chad L. Myers<sup>1,3\*</sup>

<sup>1.</sup> Biomedical Informatics and Computational Biology Graduate Program, University of Minnesota, Minneapolis, MN, USA

<sup>2.</sup> Department of Agronomy and Plant Genetics, University of Minnesota, St. Paul, MN, USA

<sup>3.</sup> Department of Computer Science, University of Minnesota, Minneapolis, MN, USA

<sup>4.</sup> Cayuga Genetics Consulting Group LLC, Ithaca, NY, USA

<sup>5.</sup> Department of Biochemistry, Purdue University, West Lafayette, IN, USA

<sup>6.</sup> Donald Danforth Plant Science Center, St. Louis, MO, USA

<sup>7.</sup> USDA-ARS Plant Genetics Research Unit, St. Louis, MO, USA

\* Corresponding Authors: Ivan Baxter, [ivan.baxter@ars.usda.gov](mailto:ivan.baxter@ars.usda.gov);

Chad L. Myers, [cmyers@cs.umn.edu](mailto:cmyers@cs.umn.edu)

## Abstract

GWAS have identified loci linked to hundreds of traits in many different species, yet, the causal genes often remain unknown. We developed a computational framework, Camoco, that integrates loci identified by GWAS with functional information derived from gene co-expression networks. We built networks from three distinct biological contexts and establish the precision of our method with simulated GWAS data. Using Camoco, we prioritized candidate genes from a large-scale GWAS examining the accumulation of 17 different elements in maize seeds, demonstrating the need to match GWAS datasets with biologically relevant co-expression networks. We show that taking the genes closest to GWAS loci will often lead to spurious results demonstrating the need for robust, integrated approaches. We validated a gene identified by our approach using mutants and annotate other candidates with ontological enrichment and curated literature support, resulting in a targeted set of candidate genes that drive elemental accumulation in maize grain.

## 37 Introduction

38 Genome-wide association studies (GWAS) are a powerful tool for understanding the genetic basis  
39 of trait variation. This approach has been successfully applied for hundreds of important traits in  
40 different species, including important yield-relevant traits in crops. Sufficiently powered GWAS  
41 often identify tens to hundreds of loci containing hundreds of single-nucleotide polymorphisms  
42 (SNPs) associated with a trait of interest(1). In *Zea mays* (maize) alone, GWAS have identified  
43 nearly 40 genetic loci for flowering time(2), 89 loci for plant height(3), 36 loci for leaf length(4),  
44 32 loci for resistance to southern leaf blight(5), and 26 loci for kernel protein(6). Despite an  
45 understanding of the overall genetic architecture and the ability to statistically associate many loci  
46 with a trait of interest, a major challenge has been the identification of causal genes and the  
47 biological interpretation of functional alleles associated with these loci.

48 Linkage disequilibrium (LD), which powers GWAS, acts as a major hurdle limiting the  
49 identification of causal genes. Genetic markers are identified by a GWAS, but often reside outside  
50 annotated gene boundaries(7) and can be relatively far from the actual causal mutation. Thus, a  
51 GWA “hit” can implicate many causal genes at each associated locus. In maize, LD varies between  
52 1 kb and over 1 Mb(8), and this range can be even broader in other crop species(9,10). Moreover,  
53 there is increasing evidence that gene regulatory regions play a significant role in functional  
54 variation, which means that causal variants will never fall within annotated gene boundaries(7,11).  
55 Several quantitative trait loci (QTLs) composed of non-coding sequences have been previously  
56 reported in maize(12–14). These challenging factors mean that even when a variant is strongly  
57 associated with a trait, many plausible candidate genes are equally implicated until a causal  
58 mutation is identified.

59 These issues are multiplied when studying complex traits involving the coordinated effects of  
60 many loci throughout the genome. Narrowing candidates to likely causal genes through prior  
61 knowledge is exacerbated in crop species, where gene annotation is largely incomplete. For  
62 example, in maize, only ~1% of genes have functional annotations based on mutant analyses(15).  
63 Thus, even when a list of potential candidate genes can be identified for a particular trait, there  
64 are very few sources of information that can help identify genes linked to a phenotype. The  
65 interpretation and narrowing of large lists of highly associated SNPs with complex traits are now  
66 the bottleneck in developing new mechanistic understanding of how genes influence traits.

Advanced mapping populations developed in crop species have enabled the rapid identification of hundreds of loci that characterize traits critical to important global issues such as worldwide food supply and crop nutritional quality, yet we lack the keys to understanding the wealth of information linking genotypic variation to phenotype, especially when the trait of interest involves many genes that have interactions that a GWAS cannot explicitly model.

One informative and easily measurable source of functional information is gene expression. Surveying gene expression profiles in different contexts, such as throughout tissue development or within different genetic backgrounds, helps establish how a gene's expression is linked to its biological function, including variation in phenotype. Comparing the similarity of two genes' expression profiles, or co-expression, quantifies the joint response of the genes to various biological contexts, and highly similar expression profiles can indicate shared regulation and function(16). Analysis of co-expression or co-expression networks has been used successfully for identifying functionally related genes, including in several crop species(17–23).

Because co-expression provides a global measure of functional relationships, it can serve as a powerful means for interpreting GWAS candidate loci. Specifically, we expect that variation in several different genes contributing to the same biological process would be associated with a given phenotype(24). Thus, if genetic variation driving the phenotype captured by GWAS is encoded by co-regulated genes, these datasets will non-randomly overlap. Systematic integration of candidate loci identified by GWAS with co-expression interactions provides an opportunity to prioritize candidate genes linked to GWAS SNPs based on putative functional information (captured by a gene co-expression network). Though not all functional relationships are captured using co-expression(25), these data still provide a highly informative, and sometimes the only, set of clues about genes that have otherwise not been studied. This principle has been used successfully with other types of networks, for example, protein-protein interactions(26), and co-expression has been used as a basis for understanding GWAS in mouse and human(27–29).

We developed a freely available, open-source computational framework called Camoco (**Co-analysis of molecular components**) designed specifically for integrating GWAS candidate lists with gene co-expression networks to prioritize individual candidate genes. Camoco evaluates candidate SNPs derived from a typical GWAS study, then identifies sets of high-confidence candidate genes with strong co-expression where multiple members of the set are associated with the phenotype of interest.

We applied this approach to maize, one of the most important agricultural crops in the world, yielding 15.1 billion bushels of grain in the United States alone in 2016(30). We specifically

focused on quantitative phenotypes measuring the accumulation of 17 different elements in the maize grain ionome (Al, As, B, Ca, Cd, Fe, K, Mg, Mn, Mo, Na, Ni, Rb, S, Se, Sr, and Zn). Plants must take up all elements except carbon and oxygen from the soil, making the plant ionome a critical component in understanding plant environmental response(31), grain nutritional quality(32), and plant physiology(33).

We evaluated the utility of three different types of co-expression networks for supporting the application of Camoco and demonstrate the efficacy of our approach by simulating GWAS to establish maize-specific SNP-to-gene mapping parameters as well as a robust null model for GWAS-network overlap. This approach does indeed confirm overlap between functional modules captured by co-expression networks and GWAS candidate SNPs for the maize grain ionome. We present high-confidence candidate genes identified for a variety of different ionomic traits, test single gene knockouts demonstrating the utility of this approach, and, more generally, highlight lessons about the connection between co-expression and GWAS loci from our study that are likely to generalize to other traits and other species.

## Results

### A framework for integrating GWAS results and co-expression networks

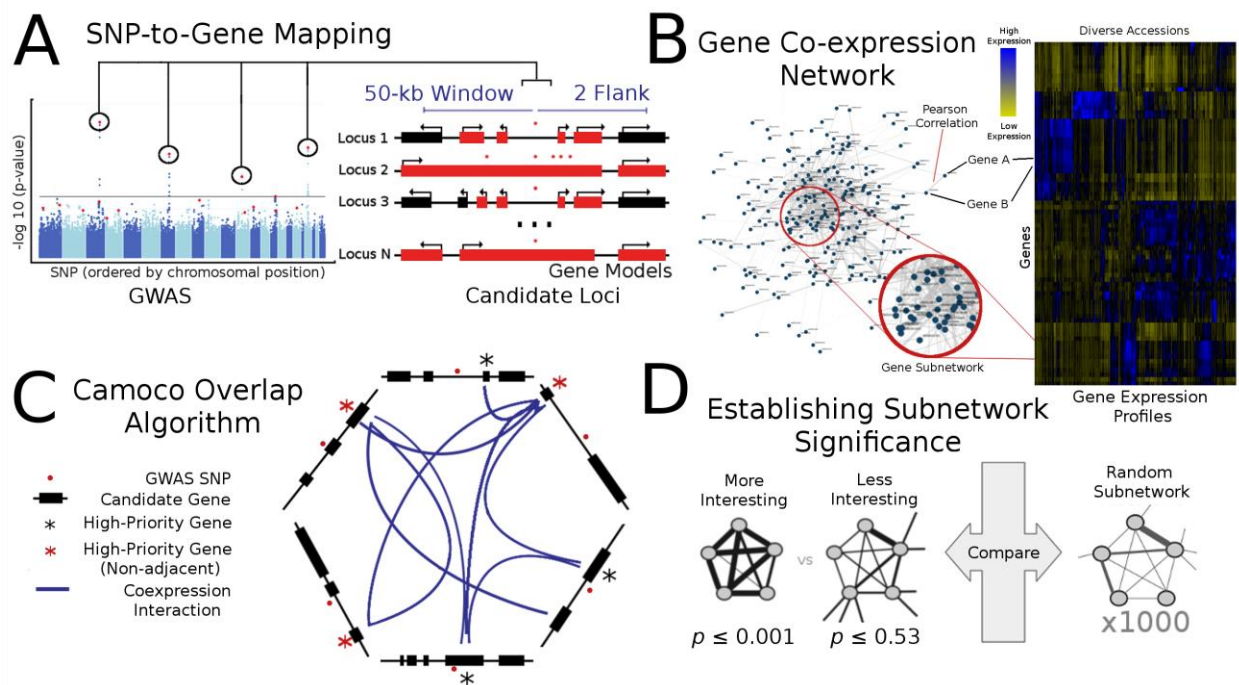
We developed a computational framework called Camoco that integrates the outputs of GWAS with co-expression networks to prioritize high-confidence causal genes associated with a phenotype of interest. The rationale for our approach is that genes that function together in a biological process that are identified by GWAS should also have non-random structure in co-expression networks that capture the same biological function. Our approach takes, as input, a list of SNPs associated with a trait of interest and a table of gene expression values and produces, as output, a list of high-priority candidate genes that are near GWAS peaks having evidence of strong co-expression.

There are three major components of the Camoco framework: a module for SNP-to-gene mapping (Figure 1A), tools for construction and analysis of co-expression networks (Figure 1B), and an "overlap" algorithm that integrates GWAS-derived candidate genes with the co-expression networks to identify high-priority candidate genes with strong co-expression support across multiple GWAS loci (Figure 1C) (see Materials and Methods for details on each component).

The overlap algorithm uses two network scoring metrics: subnetwork density and subnetwork locality (Eq. 1 and Eq. 2). Subnetwork density measures the average interaction strength between all pairwise combinations of genes near GWAS peaks. Subnetwork locality measures the

proportion of co-expression interactions among genes within a GWAS-derived subnetwork (local interactions) as compared to the number of global interactions with other genes in the genome. Density and locality were also calculated on a gene-specific basis (Eq. 3 and Eq. 4) (see Materials and Methods for details). For a given input GWAS trait and co-expression network, the statistical significance for both density and locality is determined by generating a null distribution based on randomly generated GWAS traits ( $n = 1,000$ ) with the same number of implicated loci and corresponding candidate genes. This null distribution is then used to derive a  $p$ -value for the observed subnetwork density and locality for all putative causal genes (Figure 1D). Thus, for a given input GWAS trait, Camoco produces a ranked list of candidate causal genes for both network metrics and a false discovery rate (FDR) that indicates the significance of the observed overlap between each candidate causal gene and the co-expression network. Using this integrated approach, the number of candidate genes prioritized for follow-up validation is reduced to those that have strong trait association and also are highly co-expressed with other GWAS-associated genes. Our method can be applied to any trait and species for which GWAS has been completed and sufficient gene expression data exist to construct a co-expression network.

Figure 1



## Schematic of the Camoco framework

The Camoco framework integrates genes identified by SNPs associated with complex traits with functional information inferred from co-expression

networks. **(A)** A typical GWAS result for a complex trait identifies several SNPs (circled) passing the threshold for genome-wide significance indicating a multigenic trait. SNP-to-gene mapping windows identify a varying number of candidate genes for each SNP. Candidate genes are identified based on user-specified window size and a maximum number of flanking genes surrounding a SNP (e.g., 50-kb and two flanking genes, designated in red). **(B)** Independently, gene co-expression networks identify interactions between genes uncovering an unbiased survey of putative biological co-function. Network interactions are identified by comparing gene expression profiles across a diverse set of accessions (e.g., experimental conditions, tissue, samples). Gene subnetworks indicate sets of genes with strongly correlated gene expression profiles. **(C)** Co-analysis of co-expression interactions among GWAS trait candidate genes identifies a small subset of genes with strong network connections. Blue lines designate genes that have similar co-expression patterns indicating co-regulation or shared function. Starred genes are potential candidate genes associated with GWAS traits based on SNP-to-gene mapping and co-expression evidence. Red stars indicate genes that are not the closest to the GWAS SNP (non-adjacent) that may have been missed without co-expression evidence. **(D)** Statistical significance of subnetwork interactions is assessed by comparing co-expression strength among genes identified from GWAS datasets to those from random networks containing the same number of genes. In the illustrated case, the more interesting subnetwork has both high density as well as locality.

## Generating co-expression networks from diverse transcriptional data

A co-expression network that is derived from the biological context generating the phenotypic variation subjected to GWAS is a key component of our approach. A well matched co-expression network will describe the most relevant functional relationships and identify coherent subsets of GWAS-implicated genes. We and others have previously shown that co-expression networks generated from expression data derived from different contexts capture different functional information(34,35). For example, experiments measuring changes in gene expression can explore environmental adaptation, developmental and organ-based variation, or variation in expression that arises from population and ecological dynamics (see (36) for review). For some species, published data contain enough experimental accessions to build networks from these different types of expression experiments (the term accession is used here to differentiate samples, tissues, conditions, etc.). We reasoned that these different sources of expression profiles likely have a strong impact on the utility of the co-expression network for interpreting genetic variation

captured by GWAS. Using this rationale, we constructed several different co-expression networks independently and assessed the ability of each to produce high-confidence discoveries using our Camoco framework.

Three co-expression networks representing three different biological contexts were built. The first dataset targeted expression variation that exists between diverse maize accessions built from whole-seedling transcriptomes on a panel of 503 diverse inbred lines from a previously published dataset characterizing the maize pan-genome(37) (called the ZmPAN network hereafter). Briefly, Hirsch et al. chose these lines to represent major heterotic groups within the United States, sweet corn, popcorn, and exotic maize lines and measured gene expression profiles for seedling tissue as a representative tissue for all lines. The second dataset examined gene expression variation from a previous study characterizing different tissues and developmental time points(38). Whole-genome RNA-Seq transcriptome profiles from 76 different tissues and developmental time points from the maize reference accession B73 were used to build a network representing a single-accession expression map (called the ZmSAM network hereafter). Finally, we created a third dataset as part of the ionomics GWAS research program. These data measure gene expression variation in the root, which serves as the primary uptake and delivery system for all the measured elements. Gene expression was measured from mature roots in a collection of 46 genotypically diverse maize inbreds (called the ZmRoot network hereafter). All datasets used here were generated from whole-genome RNA-Seq analysis, although Camoco could also be applied to microarray-derived expression data.

Table 1

	Number Significant ( $p \leq 0.01$ ) GO Terms (n = 1078)			
	Density	Locality	Both Scores	Either Score
ZmPAN	451 (41%)	539 (50%)	312 (29%)	678 (63%)
ZmSAM	365 (34%)	437 (40%)	234 (21%)	568 (53%)
ZmRoot	573 (53%)	331 (31%)	278 (26%)	626 (58%)

#### Significantly co-expressed GO terms

Co-expression was measured among genes within each GO term that had co-expression data in each network using both density (Eq. 1) and locality (Eq. 2). Significance of co-expression metrics was assessed by comparing values to 1,000 random gene sets of the same size.

Co-expression networks for each dataset were constructed from gene expression matrices using Camoco (see Materials and Methods for specific details on building these networks). Once built, several summary statistics were evaluated from interactions that arise from genes in the network



(Supp. Fig. 1–3). Co-expression was measured among genes within the same Gene Ontology (GO) term to establish how well density and locality captured terms with annotated biological functions.

Density and locality were measured for subnetworks consisting of the set of genes co-annotated to each GO term and compared to scores from 1,000 random sets of genes of the same size (see Table 1; Supp. Table 1 for full data). In total, 818 GO terms of the 1078 tested (76%) were composed of gene sets that were significantly co-expressed ( $p \leq 0.01$ ) in at least one network using density or locality relative to the randomized gene lists of the same size. Broken down by network as well by co-expression score, there was substantial co-expression among GO terms for both density and locality in each network. Density was significant for the most GO terms in the ZmRoot network, while locality performed best in ZmPAN (see Table 1). Considering terms captured by both scores or by either score, overlap between the two co-expression metrics was comparable. As previously reported(39), GO terms that exhibit strong co-expression between members often do so in only a subset of the networks (Supp. Table 1). Thus, both the biological context of the expression data and nature of the co-expression score influence the subset of GO terms with significantly co-expression. Overall, while density and locality recover different GO terms, there are substantially more co-expressed GO terms, for either score, than those found by size-matched randomly generated sets of genes (Supp. Table 1).

Table 2

	Network Clusters		
	Num Cluster: ( $10 \geq n > 100$ )	Num Clusters: ( $n \geq 100$ )	Num Clusters ( $n \geq 10$ ) Enriched for GO Terms ( $p \leq 0.01$ )
ZmPAN	76	18	71
ZmSAM	160	10	115
ZmRoot	150	10	106

#### Gene co-expression network cluster assignments

Gene clusters were calculated by running the Markov Cluster (MCL) algorithm on the co-expression matrix. Cluster values designate network specific gene clusters and are not compared across networks.

In addition to detecting strong co-expression among genes previously annotated by functional processes, unsupervised network clustering using the Markov Cluster algorithm(40) showed distinct modules within each network. A large number of clusters were significantly enriched for genes that are co-annotated for the same GO term (hypergeometric  $p$ -value  $\leq 0.01$ ; Supp. Table 3). Not all clusters identified previously annotated gene sets. Many strongly co-expressed clusters lacked any previously annotated function (Table 2; Supp. Table 3) potentially identifying novel

co-regulated biological processes. Additionally, all networks exhibited a truncated power law distribution in the number of significant interactions (degree) for genes in the network (Supp. Fig. 1–3), which is typical of biological networks(41).

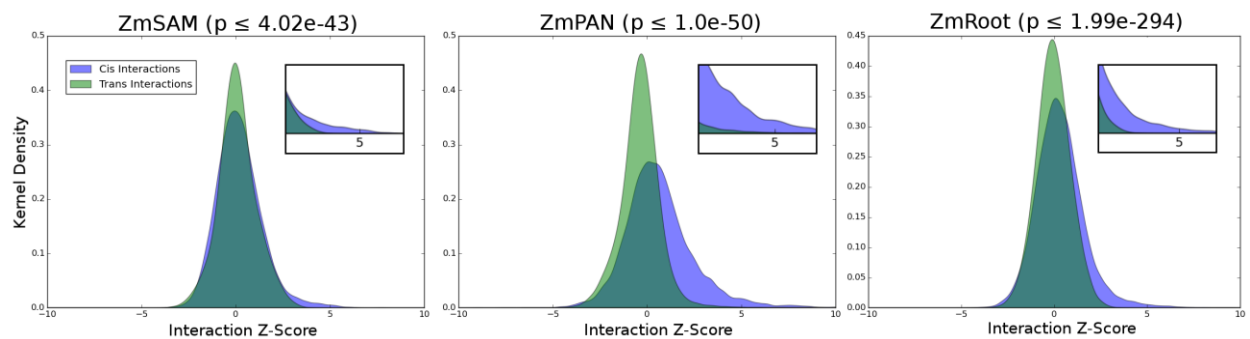
## Accounting for *cis* gene interactions

Camoco integrates GWAS candidates with co-expression interactions by directly assessing the density or locality of interactions among candidate genes near GWAS SNPs. However, the process of mapping SNPs to surrounding candidate genes has inherent complications that can strongly influence subnetwork co-expression calculations. While we assume that the majority of informative interactions among candidate genes are between GWAS loci, *cis*-regulatory elements and other factors can lead to co-expression between linked genes and produce skewed distributions in density and locality calculations, which can in turn bias co-expression statistics. Identifying significant overlap between GWAS loci and co-expression networks requires a distinction between co-expression among genes that are in close proximity to one another on a chromosome (*cis*) compared to those genes that are not (*trans*).

To assess the impact of *cis* co-expression, network interactions for genes located on different chromosome (*trans* interactions) were compared to *cis* interactions for pairs of genes less than 50 kb apart. The distributions of the two groups indicate that *cis* genes are more likely to have a strong co-expression interaction score than *trans* genes (Figure 2). This bias toward *cis* genes is especially pronounced for strong positive co-expression, where we observed substantially stronger enrichment for linked gene pairs compared to *trans* genes (e.g., z-score  $\geq 3$ ; see Figure 2 inset).

Figure

2



## *Cis* vs. *trans* co-expression network interactions

Comparing distributions of co-expression network interaction scores between *cis* and *trans* sets of genes. Distribution densities of *trans* gene pairs (green)

show interactions between genes on separate chromosomes. Distribution densities of *cis* gene pairs (blue) show interactions between genes with less than 50 kb intergenic distance. Inset figures show z-score values greater than 3. Non-parametric *p*-values were calculated between co-expression values taken from *cis* and *trans* distributions (Mann-Whitney U test).

The enrichment of significant co-expression among *cis* genes, likely due to shared *cis*-regulatory sequences, prompted us to remove *cis* interactions when examining co-expression relationships among candidate genes identified by GWAS SNPs in Camoco. To account for possible *cis* regulation within network metrics described here, only interactions that span different GWAS loci (*trans*) were included in density and locality calculations for GWAS-network overlap calculation (see Materials and Methods).

## Evaluation of the Camoco framework

To explore the limits of our approach, we examined factors that influence overlap detection between co-expression networks and genes linked to GWAS loci. In an idealized scenario, SNPs identified by GWAS map directly to true causal genes, all of which exhibit strong co-expression network interactions (Figure 3). But in practice, SNPs can affect regulatory sequences or be in linkage disequilibrium (LD) with the functionally important allele, leading to a large proportion of SNPs occurring outside of genic regions(7).

We evaluated two major challenges that influence SNP-to-gene mapping. The first is the total number of functionally related genes in a subnetwork, representing the fraction of genes involved in a biological process, that are simultaneously identified by GWAS. In cases where too few genes represent any one of the underlying causal processes, our proposed approach is not likely to perform well—for example, when GWAS identifies a single locus in a ten-gene biological process due to penetrance, limited allelic variation in the mapping population, or extensive gene-by-environment interactions. We refer to this source of noise as the *missing candidate gene rate* (*MCR*) or, in other words, the fraction of genes involved in the causal process not identified by the GWAS in question (Figure 3B; Eq. 6).

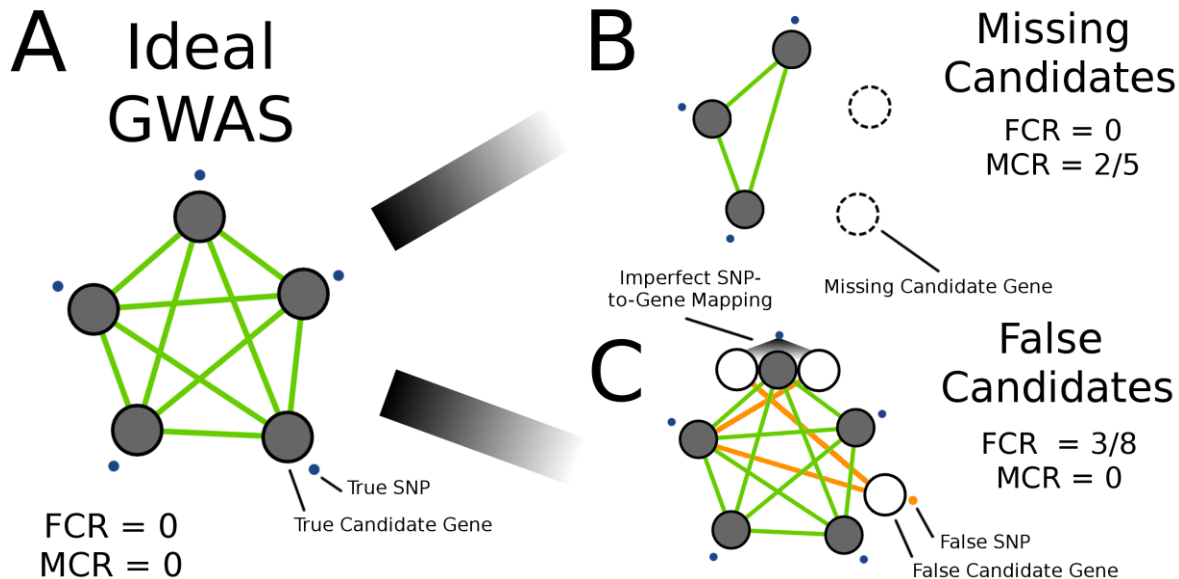
The second key challenge in identifying causal genes from GWAS loci is instances where associated SNPs each implicate a large number of candidate genes. Thus, in cases where the linked regions are large (i.e., imperfect SNP-to-gene mapping), the framework's ability to confidently identify subnetworks of highly co-expressed causal genes may be compromised. One would expect to find scenarios where the proposed approach does not work simply because there are too many non-causal genes implicated by linkage within each GWAS locus, such that the co-expression

signal among the true causal genes is diminished by the false candidates linked to those regions. We refer to this source of noise as the *false candidate gene rate (FCR)*, the fraction of all genes linked to GWAS loci that are not causal genes (Figure 3C; Eq. 7).

To explore the limits of our co-expression-based approach with respect to these factors, we simulated scenarios where we could precisely control both MCR and FCR. In practice, neither of these quantities can be controlled; MCR is a function of the genetic architecture of the phenotype as well as the degree of power within the study population of interest, and FCR is a function of recombination frequency in the GWAS population.

We evaluated the expected performance of the Camoco framework for a range of each of these parameters by simulating ideal GWAS scenarios using co-expressed GO terms ( $p \leq 0.05$ ; Table 1). These ideal cases were then subjected either to a subset of genes being replaced by random genes (i.e., to simulate MCR but conserve term size) or to functionally unrelated genes being added using SNP-to-gene mapping (i.e., to simulate FCR introduced by linkage). In both cases, simulated GWAS candidates (GO term set members) were subjected to varying levels of either FCR or MCR while tracking the number of GO terms that remained significantly co-expressed at each level. These simulations enabled us to explore a broad range of settings for these key parameters and establish whether our proposed approach had the potential to be applied in maize.

Figure 3



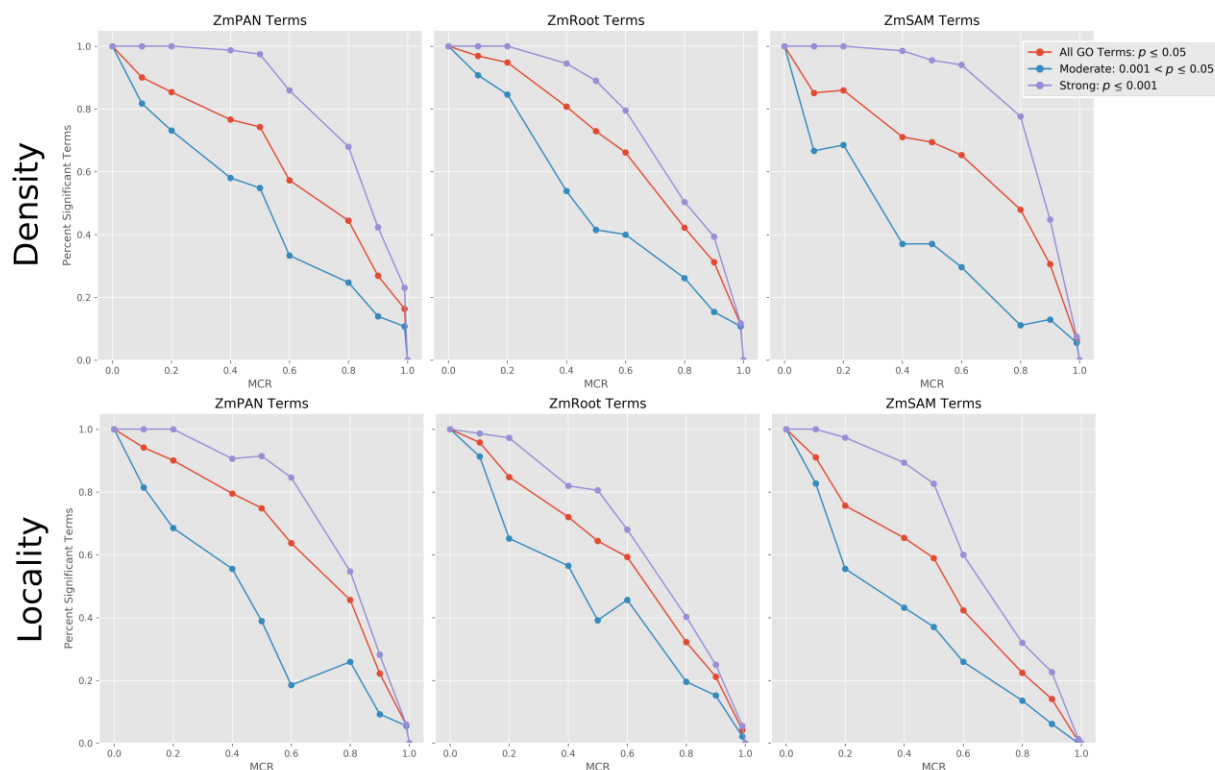
## Simulating GWAS-network overlap using GO terms

Several GWAS scenarios were simulated to assess the effect of noise on co-expression network overlap. Panel **(A)** shows an ideal GWAS, where SNPs (blue points) map directly to candidate genes within the same biological process (i.e., a GO term) and have strong co-expression (green lines). Signal is defined as the co-expression among the genes exclusive to the GO term. Noise in the overlap between GWAS and co-expression networks was introduced by varying two parameters: the missing candidate gene rate (MCR) and false candidate gene rate (FCR). Panel **(B)** demonstrates the effect of a large proportion of missing candidate genes (MCR = 2/5) on network signal. Likewise, panel **(C)** shows the effect of false candidate genes (FCR) on network overlap, either through false positive GWAS SNPs (orange points) or through imperfect SNP-to-gene mapping (FCR = 3/8). Orange lines designate the additional candidate genes that introduce co-expression noise that impedes the identification of network structure.

## Simulated GWAS datasets show robust co-expression signal to MCR and FCR

Subnetwork density and locality were measured for significantly co-expressed GO terms containing between 50 and 150 genes in each network at varying levels of MCR (see Supp. Table 4). At each MCR level, density and locality among the remaining genes were compared to 1,000 random sets of genes of the same size. The proportion of initial GO terms that remained significantly co-expressed was recorded for each network (see Figure 4, red curve; see Supp. Fig. 4A for absolute term numbers). GO terms were also split into two starting groups based on strength of initial co-expression: moderate ( $0.001 < p \leq 0.05$ ; blue curve) and strong ( $p \leq 0.001$ ; violet curve).

Figure 4



Strength of co-expression among GO terms at varying levels of MCR

Subnetwork density and locality were measured for all GO terms with strong initial co-expression ( $p \leq 0.05$ ) comparing co-expression in GO terms to 1,000 random networks of the same size. Co-expression density and locality were then compared again ( $n = 1,000$ ) with varying missing candidate rate (MCR), where a percentage of genes was removed from the term and replaced with random genes to conserve GO term size. Curves decline with increased MCR as the proportion of strongly co-expressed GO terms ( $p \leq 0.05$ ,  $n = 1,000$ ) decreases compared to the initial number of strongly co-expressed terms in each network (red curve). GO terms in each network were also split into two subsets based on initial co-expression strength: “strong,” (initial co-expression  $p \leq 0.001$ ; blue curve), and “moderate,” (initial co-expression  $0.001 < p \leq 0.05$ ; violet curve).

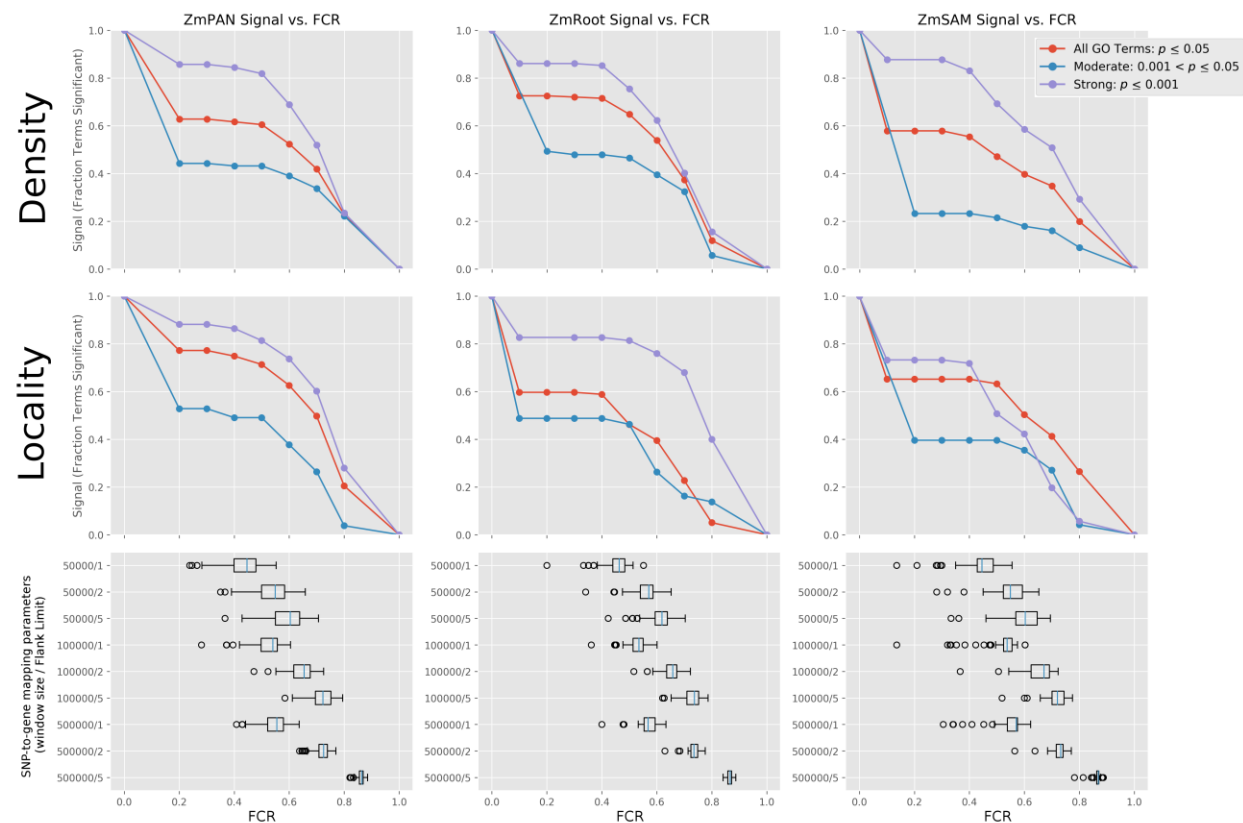
As expected, strength of co-expression among GO terms decreased as MCR increased. Figure 4 shows the decay in the proportion of GO terms that exhibit significant co-expression at increasing levels of MCR (red curve). In general, the decay of signal is similar between density and locality,

where signal initially decays slowly until approximately 60% MCR, when signal quickly diminishes.

In all three networks, GO terms with stronger initial co-expression were more robust to MCR. Signal among strongly co-expressed GO terms ( $p \leq 0.001$ ; violet curve) decayed at a substantially lower rate than moderately co-expressed GO terms, indicating that this approach is robust for GWAS datasets with moderate levels of missing genes when co-expression among true candidate genes is strong. Co-expression signal in relation to MCR was also compared between GO terms split by the number of genes within the term (see Supp. Fig. 4B–C), which did not influence the rate at which co-expression signal decayed.

Likewise, the effect of FCR was simulated. Significantly co-expressed GO terms of between 50 and 150 genes (MCR = 0) with significant co-expression ( $p \leq 0.05$ ; see Supp. Table 4) were selected. The nucleotide position of the starting base pair of each true GO term gene was used as input for our SNP-to-gene mapping protocol for identifying GWAS candidates (see Materials and Methods). Subnetwork density and locality were calculated for the simulated candidate genes corresponding to each SNP-to-gene mapping combination, in each network, to evaluate the decay of co-expression signal as FCR increases (Figure 5).

Figure 5



### Simulated GWAS: SNP-to-gene mapping density signal robustness

Strongly co-expressed GO terms (density or locality  $p$ -value  $\leq 0.05$ ) were used to simulate the effect of FCR on GWAS results. False candidates were added to GO terms by including flanking genes near true GO term genes according to SNP-to-gene mapping (window) parameters. Box plots show effective FCR of GO terms at each SNP-to-gene mapping parameter. Signal plots show the proportional number of GO terms that remain significant at  $\text{FCR} \geq x$  (red curve). GO terms in each network were also split into two subsets based on initial co-expression strength: “strong,” (initial co-expression  $p \leq 0.001$ ; blue curve), and “moderate,” (initial co-expression  $0.001 < p \leq 0.05$ ; violet curve).

Candidate genes were added by varying the window size for each SNP up to 50 kb, 100 kb, and 500 kb upstream and downstream and by varying the maximum number of flanking genes on each side to one, two, and five. Given the number of additional candidate genes introduced at each SNP-to-gene mapping combination, FCR was calculated for each GO term at each window size (see Figure 5 box plots).



Co-expression signal in relation to FCR was assessed by comparing subnetwork density and locality in each GO term at different SNP-to-gene mapping parameters for each of the three co-expression networks to random subnetworks with the same number of genes ( $n = 1,000$ ) (Figure 5, top). The proportion of significantly co-expressed GO terms decayed at higher levels of FCR (see Supp. Fig. 5A for absolute term numbers). The minimum FCR level for most GO terms was ~50% as the most stringent SNP-to-gene mapping (50 kb/one flank) approximately doubled the number of candidate genes. Two additional scenarios were considered in which signal was further split based on the initial co-expression strength: “moderate” ( $0.001 < p < 0.05$ ; blue curve) and “strong” ( $p \leq 0.001$ ; violet curve).

Despite high initial false candidate rates, co-expression signal among GO terms remained significant even at 60–70% FCR. Similar to the results with MCR, GO terms with stronger initial co-expression were more likely to remain significantly co-expressed at higher FCR levels. Co-expression signal in relation to FCR was also compared between GO terms split by the number of genes in the term (see Supp. Fig. 5B–C), which did not differentiate the rate at which co-expression signal decayed.

In cases where true candidate genes identified by GWAS were strongly co-expressed, as simulated here, a substantial number of false positive SNPs or an introduction of false candidate genes through uncertainty in SNP-to-gene mapping can be tolerated, and network metrics still detected the underlying co-expressed gene sets using our method. These results indicate that in GWAS scenarios where the majority of SNPs do not perfectly resolve to candidate genes, systematic integration with co-expression networks can efficiently filter out false candidates introduced by SNP-to-gene mapping if the underlying causative loci are strongly co-expressed. Moreover, in instances where several intervening genes exist between strongly associated SNPs in LD with each other and the true causative allele, true causal candidates can be detected using co-expression networks as a functional filter for candidate gene identification.

The potential for using this approach, however, is highly dependent on the LD of the organism in question, the genetic architecture of the trait being studied, and the degree of co-expression between causative loci. Simulations provide insight into the feasibility of using Camoco to evaluate overlap between co-expression networks and GWAS as well as a survey of the SNP-to-gene mapping parameters that should be used when using this approach (see Discussion for more details). In the context of maize, simulations performed here suggest that systematic integration of co-expression networks to interpret GWAS results will increase the precision with which causal genes associated with quantitative traits in true GWAS scenarios can be identified.

## Prioritizing causal genes driving elemental accumulation in maize grain

Identifying the biological processes underlying the elemental composition of plant tissues, also known as the ionome, can lead to a better understanding of plant adaptation as well as improved crops(42). High-throughput analytic approaches such as inductively coupled plasma mass spectrometry (ICP-MS) are capable of measuring elemental concentrations for multiple elements and are scalable to thousands of accessions per week. Using ICP-MS, we analyzed the accumulation of 17 elements in maize kernels described in depth by Ziegler et al.(43). Briefly, kernels from the nested association mapping (NAM) population were grown in four geographic locations(1). To reduce environmental-specific factors, the SNPs used in this study were from the GWAS performed on the all-location models. Approximately 30 million SNPs and small copy-number variants were projected onto the association panel and used to perform a GWAS for each of the 17 elements. SNPs were tested for significance of association for each trait using resampling model inclusion probability(44) ( $RMIP \leq 0.05$ ; see Materials and Methods). Significantly associated SNPs were used as input to Camoco to generate candidate genes from the maize filtered gene set (FGS;  $n = 39,656$ ) for each element using a range of SNP-to-gene mapping parameters: 50-kb, 100-kb, and 500-kb windows (up/downstream) limited each to one, two, or five flanking genes (up/downstream of SNP; see Figure 1A). In total, 4,243 statistically significant SNPs were associated with maize grain ionome traits. Summing the potential candidate genes across all 17 traits implicates between 5,272 and 22,927 unique genes depending on the SNP-to-gene mapping parameters used (between 13% and 57% of the maize FGS, respectively). On average, each trait's significantly associated SNPs identified 118 non-overlapping windows across the ten chromosomes of maize (i.e., effective loci; see Materials and Methods), and these implicate an average of 612 candidate genes per element (Materials and Methods).

Table 3

Name	GWAS SNPs	Effective Loci			Candidate Genes								
WindowSize	-	50KB	100KB	500KB	50KB			100KB			500KB		
FlankLimit	-	-	-	-	1	2	5	1	2	5	1	2	5
Ionome (Total)	4243	2279	1658	456	5272	7348	11612	7727	9664	13614	20024	20776	22927
Al27	176	149	140	98	239	336	417	350	523	699	804	1035	1684
As75	182	151	141	104	228	314	372	339	489	669	740	986	1657
B11	108	95	86	68	154	233	271	219	326	433	426	601	1007
Ca43	105	82	78	61	124	181	215	164	253	350	339	476	845
Cd111	630	471	418	251	869	1189	1395	1252	1786	2309	3159	3758	5283
Cu65	165	133	125	101	202	293	355	284	437	604	562	805	1431
Fe57	171	136	125	89	252	351	420	335	511	697	766	990	1546
K39	130	111	100	78	168	248	298	239	357	498	534	715	1176
Mg25	153	129	121	99	203	281	328	274	414	554	584	815	1398
Mn55	168	134	119	94	228	302	340	314	436	562	638	850	1364
Mo98	154	123	109	74	226	312	361	287	419	532	709	892	1354
Ni60	99	73	64	49	107	148	163	161	226	291	301	417	697
P31	123	101	91	70	159	223	260	210	312	424	485	643	1051
Rb85	135	105	93	78	168	223	251	245	335	414	409	590	1026
Se82	162	135	129	101	237	328	392	330	485	682	663	895	1563
Sr88	113	99	90	63	142	206	238	199	317	431	481	636	1009
Zn66	149	125	116	90	211	299	348	288	435	565	613	841	1419
Ionome (Average)	172	138	126	92	230	322	378	323	474	630	718	938	1501
		119			613								

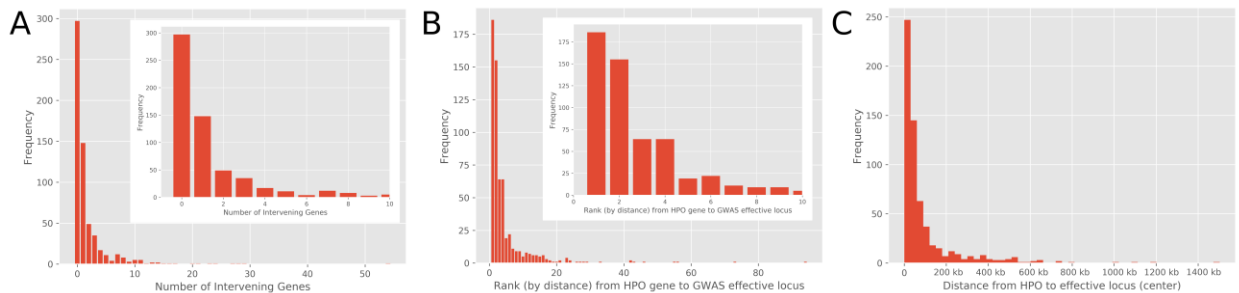
Maize grain ionome SNP-to-gene mapping results

Significant SNPs associated with the maize grain ionome were mapped to candidate genes by collapsing SNPs with overlapping windows down to effective SNPs, then taking genes upstream and downstream of the effective SNP up to the flank limit.

Camoco identifies high-priority candidate causal genes under ionomic GWAS loci

Given the large number of candidate genes associated with elemental accumulation, we used Camoco to integrate network co-expression with effective loci identified by GWAS for each of the 17 elemental traits separately. By combining candidate gene lists with the three gene expression datasets (ZmPAN, ZmRoot, and ZmSAM) and two co-expression network approaches (locality and density) high-priority candidate genes driving elemental accumulation in maize were identified (see Figure 1C). For each network-trait combination, Camoco identified a ranked list of prioritized candidate causal genes, each associated with an FDR that reflects the significance of co-expression connecting that candidate gene to genes near other loci associated with the same trait (Supp. Table 5). We defined a set of high-confidence discoveries by reporting candidates that were discovered at a  $FDR \leq 30\%$  in at least two SNP-to-gene mapping parameter settings (e.g., 50 kb/one flank and 100 kb/one flank), denoted as the high-priority overlap (HPO) set (see Supp. Table 6 and Materials and Methods).

Figure 6



Number of intervening genes between HPO gene and GWAS locus

The distribution of positional candidates and HPO genes. Panel **(A)** shows the distribution in the number of positional candidates between each of the 610 HPO genes and an effective locus (note: intervening gene could also be an HPO gene). Panel **(B)** shows candidate genes near GWAS SNPs, ranked by their absolute distance to effective loci. The distribution shows the rank of the absolute distance (either upstream or downstream) of HPO genes. In both panels, the inset plot shows the lower end of the distributions. Panel **(C)**

shows the distance between the center of HPO genes and the center of the effective locus identified by GWAS.

By these criteria, we found strong evidence of co-expression for 610 HPO genes that were positional candidates among the 17 ionic traits measured (1.5% maize FGS). The number of HPO genes discovered varied significantly across the traits we examined, with between 2 and 209 HPO genes for a given element considering either density or locality in any network (Table 4; Either:Any column). HPO genes discovered by Camoco were often non-adjacent to GWAS effective loci, either having genes intervening between the HPO candidate and the effective locus or having positional candidates that were closer either upstream or downstream of the GWAS locus (Figure 1C). Of the 610 HPO genes, 297 had zero intervening genes (Figure 6A). The remaining 313 HPO genes had between 1 and 54 intervening genes, though the majority (292 HPO genes) had 10 or fewer intervening genes (Figure 6; inset). Similar results were observed when considering candidate genes' absolute distance to the effective locus (Figure 6B), demonstrating that Camoco often identifies candidates with strong co-expression evidence that would not have been selected by choosing the closest positional candidate.

Table 4

Method Network	FDR 30%										
	Either	Density				Locality				Both	
	Any	ZmPAN	ZmRoot	ZmSAM	Any	ZmPAN	ZmRoot	ZmSAM	Any	Any	ZmRoot
Al	69	0	13	0	13	56	1	0	57	1	0
As	28	0	27	0	27	1	1	0	2	1	1
B	2	0	0	0	0	0	1	1	2	0	0
Ca	3	0	0	0	0	0	1	2	3	0	0
Cd	209	0	126	0	126	97	1	0	98	15	1
Cu	26	0	26	0	26	0	0	0	0	0	0
Fe	12	0	11	0	11	0	1	0	1	0	0
K	17	0	15	0	15	0	0	2	2	0	0
Mg	26	0	1	0	1	24	0	1	25	0	0
Mn	2	0	0	0	0	1	1	0	2	0	0
Mo	8	0	1	0	1	6	1	0	7	0	0
Ni	2	0	0	0	0	1	0	1	2	0	0
P	18	0	0	16	16	0	3	0	3	1	0
Rb	52	0	0	52	52	0	0	0	0	0	0
Se	105	0	76	0	76	34	0	1	35	6	0
Sr	60	0	58	0	58	4	0	0	4	2	0
Zn	49	0	8	0	8	43	0	0	4	2	0
Ionome	610	0	326	66	391	228	11	8	247	26	2

#### Maize grain ionome high-priority candidate genes

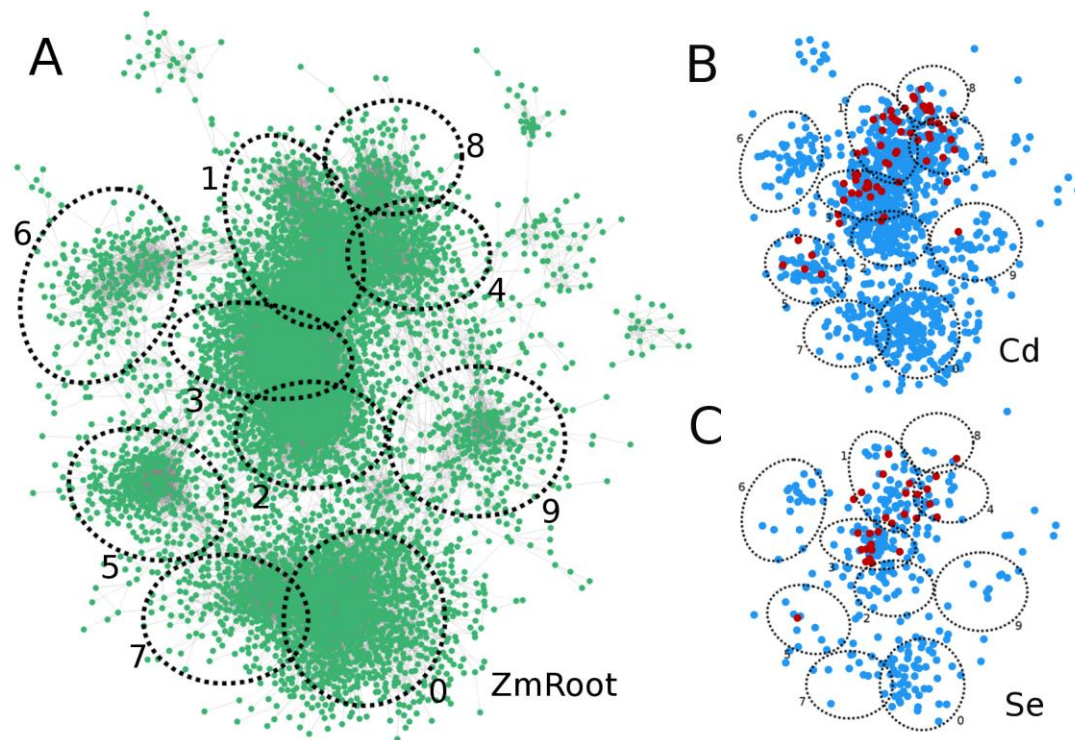
Gene-specific density and locality metrics were compared to ( $n = 1,000$ ) random sets of genes of the same size to establish a 30% FDR. Genes were considered candidates if they were observed at two or more SNP-to-gene mappings (i.e., HPO). Candidates in the "Either" column are HPO genes discovered by either density or locality in any network. The number of genes

discovered for each element is further broken down by co-expression method (density, locality, both) and by network (ZmPAN, ZmSAM, ZmRoot). Candidates in the "Both" column were discovered by density and locality in the same network or in different networks (Any). Note: zero elements had HPO genes using "Both" methods in the ZmPAN and ZmSAM networks.

Co-expression networks derived from variation across genotypically diverse accessions support stronger candidate gene discoveries

The variation in the number of genes discovered by Camoco depended on which co-expression network was used as the basis for discovery. The ZmRoot co-expression network proved to be the strongest input, discovering genes for 15 of the 17 elements (absent in Ni and Rb) for a total of 335 HPO genes, ranging from 1 to 126 per trait (Supp. Table 6). In contrast, the ZmSAM network, which was constructed based on a tissue and developmental expression atlas collected exclusively from the B73 accession, supported the discovery of candidate genes for only 8 elements (B, Ca, K, Mg, Ni, P, Rb, and Se) for a total of 74 HPO genes, ranging from 1 to 52 per trait (Supp. Table 6). The ZmPAN network, which was constructed from whole seedlings (pooled tissue) across 503 different accessions, provided intermediate results, supporting high-confidence candidate discoveries for 10 elements (Al, As, Cd, Mg, Mn, Mo, Ni, Se, Sr, and Zn) for a total of 228 HPO genes, ranging from 1 to 97 per trait (Supp. Table 6). The relative strength of the different networks for discovering candidate causal genes was consistent even at stricter FDR thresholds (e.g.,  $FDR \leq 0.10$ ; Supp. Table 6).

Figure 7



HPO genes for Cd and Se in the ZmRoot network

The strongest 100,000 interactions were used to visualize global clustering of genes ( $n = 7,844$ ) in the ZmRoot network. A force-directed algorithm positioned genes (**A**; green nodes) showing approximate boundaries (dotted black circles) of the top ten MCL clusters (Supp. Table 2). The ZmRoot network view was filtered to possible candidate genes (blue nodes) identified from SNP-to-gene mapping for Cd and Se (**B** and **C**, respectively). Network edges were removed from the visualization in panels (**B**) and (**C**), though MCL cluster boundaries were preserved. HPO genes for each element (highlighted in red) co-localize to specific clusters.

Figure 7 visualizes the discovery process for HPO genes in the ZmRoot network. Genes were organized in a global view containing the strongest 100,000 interactions using a force-directed layout algorithm to show high-level clustering (Figure 7A). For two elements, Cd and Se, a large number of possible candidate genes from SNP-to-gene mapping for each element (Figure 7B–C, blue nodes) spans many of the MCL clusters identified in the network (dotted ellipses). The HPO genes, in contrast, discovered by density and locality networks co-localize to a small number of MCL clusters (red nodes).

## Density and locality network metrics provide complementary information

Both density and locality were assessed on a gene-specific level to measure the strength of a given candidate causal gene's co-expression relationships with genes in other GWAS-identified loci (see Eq. 3 and Eq. 4). Gene-specific density measures the fraction of observed co-expression interactions to total possible co-expression relationships between the candidate gene and genes linked to other GWAS-identified loci, while gene-specific locality normalizes gene interactions to account for the proportion of interactions between the candidate gene and the rest of the genome (i.e., genes not near a GWAS locus). Overall, density identified more HPO genes than did locality. For example, across all traits and networks, 391 HPO candidate genes were discovered using density, while 247 HPO candidate genes were discovered using locality (see Table 4, Density:Any and Locality:Any). Interestingly, the high-confidence genes were largely complementary, in terms of both which traits and which network they produced results for. Among the two sets of genes (391 and 247 genes, respectively), 26 HPO genes in common were discovered (Table 4: Both:Any). While this overlap is statistically significant ( $p \leq 1.5e-13$ ; hypergeometric), the large number of uniquely discovered genes suggests that the two measures capture largely complementary biological information from co-expression subnetworks. Indeed, when we measured the direct correlation of gene-specific density and locality measures across several GWAS traits and GO terms, we observed very weak positive but significant correlations (Supp. Figure 6). Density was most effective at identifying HPO genes within the GWAS-linked loci when using the ZmRoot network (326 HPO genes using density vs. 11 HPO gene using locality). Locality provided stronger results on the ZmPAN network (228 HPO genes using locality and 0 HPO genes using density). We observed that the utility of the locality metric appeared to be linked to the number of accessions used to construct the network (Supp. Table 7), suggesting that the differences between networks in locality may simply reflect the number of accessions used to generate them (see Discussion).

## Most candidate causal genes are trait specific

One important question is the extent to which putative causal genes overlap across different ionic traits. It is plausible that some mechanisms affecting elemental accumulation are shared by multiple elements. We compiled the complete set of HPO genes discovered for each element and assessed overlap across the complete set of 17 elements (Table 5). Most of the discovered HPO genes are element specific, with relatively little overlap between elements (Table 5). However, a limited number of element pairs did exhibit statistically significant overlap, including Cd, which shared significant overlap with seven other elements (Al, Cu, K, Mg, Mo, Se, and Sr), and Se, which shared significant overlap with three other elements (As, Cd, and Mg), and Mo, which shared



significant overlap with two other elements (Al and Cd). These candidate genes represent important potential modulators of elemental composition and are particularly worthy of further study (Supp. Table 8).

Table 5

	Al	As	B	Ca	Cd	Cu	Fe	K	Mg	Mn	Mo	Ni	P	Rb	Se	Sr	Zn
Al	69	0	0	0	14	0	1	0	1	0	2	0	0	0	3	0	1
As	1	28	0	0	2	0	0	0	0	0	0	0	0	0	4	0	0
B	1	1	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Ca	1	1	1	3	0	0	0	0	0	0	0	0	0	0	0	0	0
Cd	0	0.056	1	1	209	6	2	3	4	0	4	0	0	1	12	9	3
Cu	1	1	1	1	1E-06	26	0	1	0	0	0	0	0	0	0	2	0
Fe	0.053	1	1	1	0.011	1	12	0	0	0	0	0	0	0	0	1	0
K	1	1	1	1	0.002	0.029	1	17	0	0	0	0	0	0	0	1	0
Mg	0.112	1	1	1	4E-04	1	1	1	26	0	0	0	0	0	3	2	0
Mn	1	1	1	1	1	1	1	1	1	2	0	0	0	0	1	0	0
Mo	6E-04	1	1	1	2E-06	1	1	1	1	1	8	0	0	0	1	0	0
Ni	1	1	1	1	1	1	1	1	1	1	1	2	0	0	0	0	0
P	1	1	1	1	1	1	1	1	1	1	1	1	18	2	0	0	0
Rb	1	1	1	1	0.514	1	1	1	1	1	1	1	0.002	52	0	0	0
Se	0.012	4E-05	1	1	0	1	1	1	7E-04	0.014	0.054	1	1	1	105	2	3
Sr	1	1	1	1	0	0.005	0.046	0.065	0.005	1	1	1	1	1	0.065	60	0
Zn	0.2	1	1	1	0.03	1	1	1	1	1	1	1	1	1	0.005	1	49

#### Element HPO candidate gene overlap

Overlap between the 610 HPO genes discovered between different elements by either density or locality and in any network. The diagonal shows the number of HPO genes discovered for each element. Values in the upper triangular region (green) show the number of genes that overlap between elements. The values in the lower triangle designate the  $p$ -values (hypergeometric) for overlap between the two sets of HPO genes. Red cells indicate significance with Bonferroni correction.

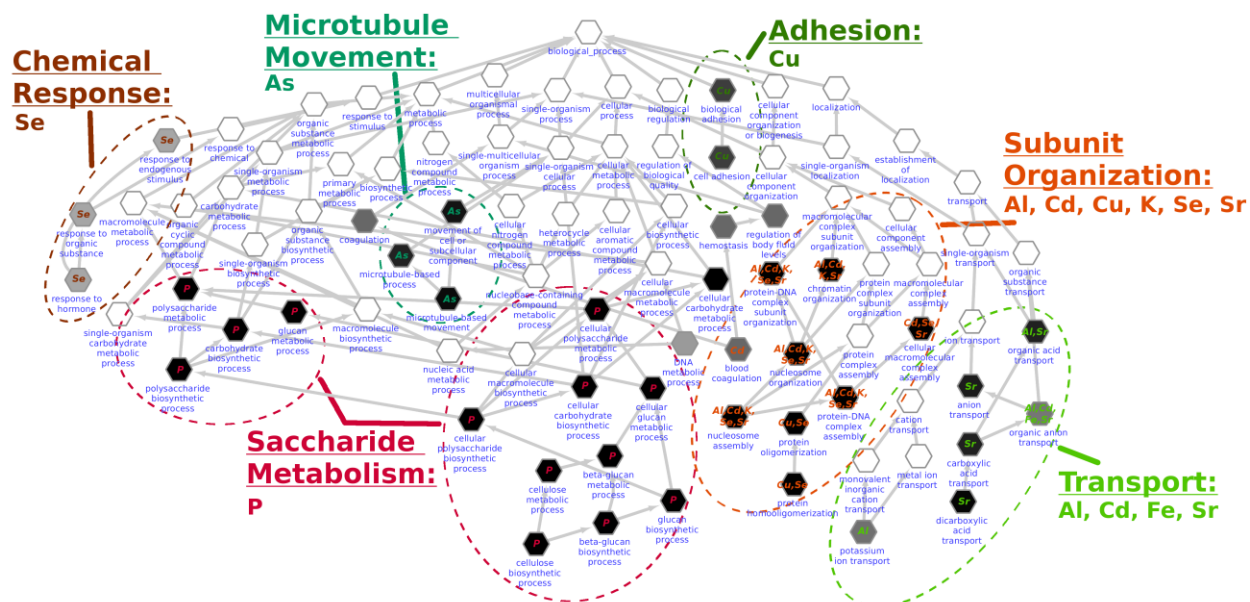
#### Enrichment analysis of putative causal genes

To explore the broader biological processes represented among HPO genes, we performed GO enrichment analysis on the candidate lists, revealing enrichments for five elements (Supp. Table 9). For example, Sr was enriched for anion transport (GO:0006820;  $p \leq 0.008$ ) and metal ion transmembrane transporter activity (GO:0046873;  $p \leq 0.015$ ). Possibly due to insufficient functional annotation of the maize genome, these enrichment results were limited, and zero elements passed a strict multiple-test correction (Bonferroni). To compensate for the sparsity of annotations, we used the HPO gene set discovered for each trait to identify the set of highly connected co-expression network neighbors, designated the HPO+ sets. Inclusion in HPO+ was determined by a gene's aggregate connectedness to the HPO set (see Materials and Methods). The HPO+ sets for several of the ionic traits showed strong GO enrichments, many of which had terms that passed strict multiple-test correction, including Al, As, Cd, Cu, Fe, K, P, Se, Sr, and Zn



(Supp. Table 10). Several of the enriched GO terms were common across HPO+ sets for different elements (Figure 8). For example, we found enrichment for a collection of GO terms related to ion transport (GO:0006811), including anion transport (GO:0006820), potassium ion transport (GO:0006813), and others (GO:0015849, GO:0015711, GO:0046942, GO:0006835), which were supported by enrichments from multiple elements (Al, Cd, Fe, Sr) (see Figure 8; “Transport” cluster). We also observed a set of six elements whose HPO+ sets (Al, Cd, Cu, K, Se, Sr) were enriched for GO terms related to chromatin organization (e.g., GO:0006325, GO:0071824, GO:0034728, GO:0006334; see Figure 8, “Subunit Organization” cluster). This may result from changes in cell cycle or endoreduplication control in roots, which is expected to alter the accumulation of multiple elements(45).

Figure 8



GO biological process enrichment for the ionome

The HPO+ gene sets were analyzed for GO enrichment in the “biological process” namespace. Each node represents a GO term organized hierarchically in a tree with directed edges designating parent terms. Shaded terms were enriched for HPO+ genes ( $p \leq 0.05$ ; hypergeometric). Dotted ovals represent curated functional terms describing the enriched nodes in different clades of the tree. Each clade is annotated with the ionomic terms that were represented in the GO enrichment.

Several of the observed GO enrichments were trait specific, including collections of GO terms reflecting “chemical response” (Se), “microtubule movement” (As), “adhesion” (Cu), and

“saccharide metabolism” (P). For example, the “saccharide metabolism” collection of GO term enrichments was driven by five HPO+ genes for P, one of which was *tgd1* (GRMZM2G044027; see Supp. Table 10). Mutations in the *Arabidopsis thaliana* ortholog for *tgd1* caused the accumulation of triacylglycerols and oligogalactolipids and showed a decreased ability to incorporate phosphatidic acid into galactolipids(46), which may alter P accumulation directly or via phosphatidic acid signaling(47). TGD1 is an ATP-binding cassette (ABC) transporter known to transport other substrates, including inorganic and organic cations and anions(48). The *tgd1* gene was present in the HPO set, and the four other genes were identified as strongly connected neighbors (HPO+) in the co-expression network. Two genes, GRMZM2G018241 and GRMZM2G030673, are of unknown function, and the other two, GRMZM2G122277 and GRMZM2G177631, are involved in cellulose synthesis. We should note that these enriched GO terms demonstrated idiosyncrasies in automated annotation approaches. Terms related to “blood coagulation” and “regulation of body fluid levels” were recovered, which were likely due to annotations translated to maize genes on the basis of protein sequence homology in humans. While, at face value, these term descriptions are not applicable to plant species, the fact that these terms contained HPO genes as well as strong network co-expression suggests that annotations assigned through orthology might be capturing underlying biological signals for which the accepted name is inappropriate (see Discussion).

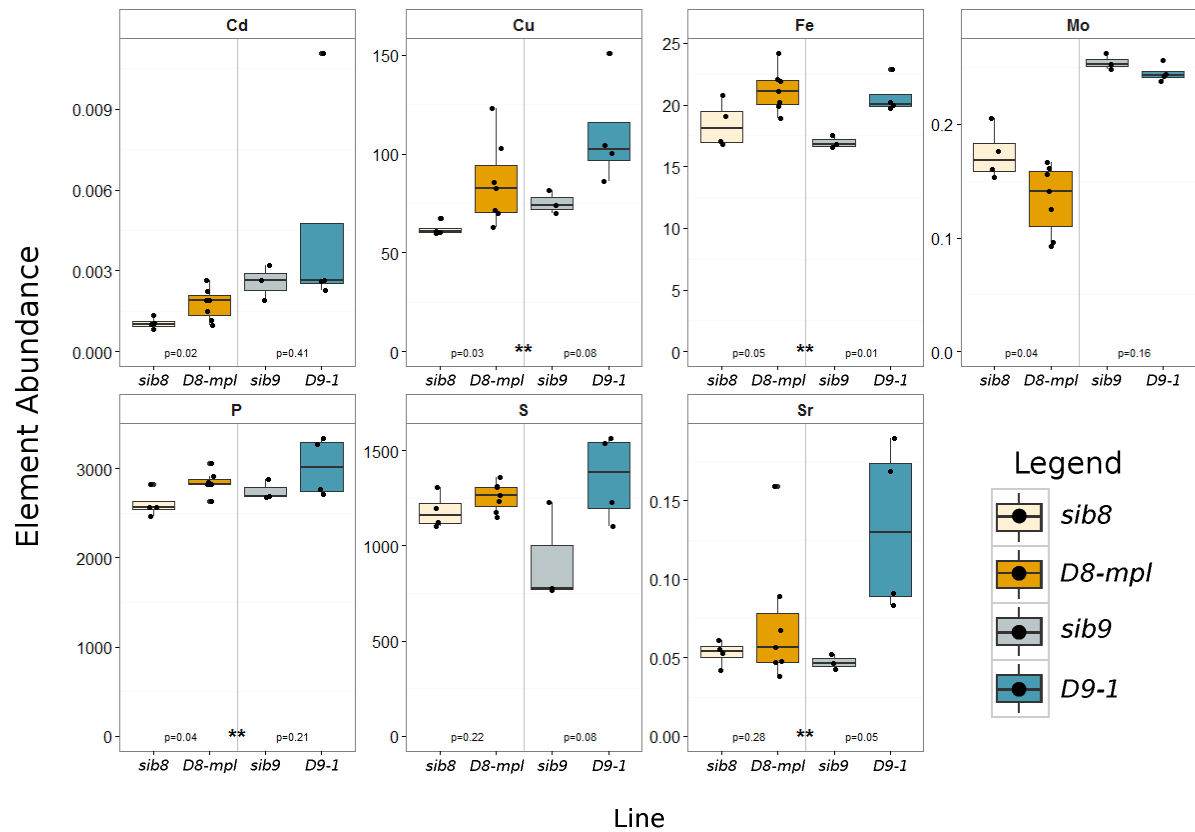
In general, using co-expression networks to expand the neighborhood of the high-confidence candidate causal genes and then assessing the entire set for functional coherence through GO enrichment is a productive strategy for gaining insight into what processes are represented. Yet this approach is particularly challenging in the annotation-sparse maize genome, where only ~1% of genes have mutant phenotypes(49). GO terms were too broad or insufficiently described to distinguish causal genes. However, the terms discovered here contain genes that act in previously described pathways known to impact elemental traits. With greater confidence that subnetworks containing HPO genes contained coherent biological information, we refined our analysis by curating HPO genes for their involvement in specific biological processes, namely, those that are known or suspected to affect the transport, storage, and utilization of elements.

## GA-signaling DELLA domain transcription factors influence the ionome of maize

One of the high-confidence candidate genes, which appeared in the HPO sets comparing Cd and the ZmRoot network, is the gibberellin (GA)-signaling component and DELLA and GRAS domain transcription factor *dwarf9* (GRMZM2G024973; *d9*(50)). *d9* is one of two DELLA paralogs in the maize genome, the other being *dwarf8* (GRMZM2G144744; *d8*); both can be mutated to

dominant-negative forms that display dwarf phenotypes and dramatic suppression of GA responses(51). Camoco ranked *d9* among the high-confidence candidates for Cd but not *d8*, though both are present in the root-based co-expression network (ZmRoot). There was only moderate, but positive, co-expression between D8 and D9 (ZmRoot:  $z = 1.03$ ; ZmPAN:  $z = 1.04$ ). Given the indistinguishable phenotypes of the known dominant mutants of *d8* and *d9*, the most likely explanation for this result is that there was allelic variation for *d9* but not *d8* in the GWAS panel. Moreover, the GA biosynthetic enzyme ent-kaurene synthase (GRMZM2G093603) encoding the *dwarf5* locus(52) affected the concentration of seed Cd and appeared among the HPO genes for Sr in the ZmRoot network. This gene is required for the biosynthesis of bioactive GA via ent-kaurene, strongly suggesting that GA signaling in the roots shapes the ionome and alters the accumulation of Cd in seeds, with potential impacts on human health.

Figure 9



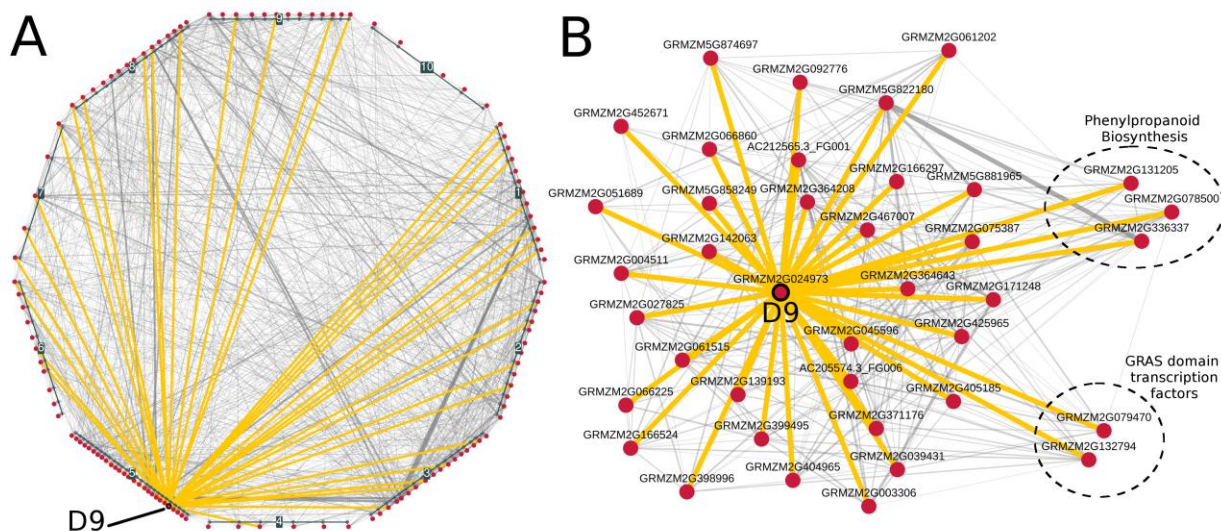
Ionomic profiles of *D8-mpl* and *D9-1* mutants

Box plots displaying ICP-MS values for *D8-mpl* and *D9-1* along with null segregating siblings (*sib8* and *sib9*). Embedded *p*-values indicate statistical differences between mutants and wild-type siblings, while asterisks (\*\*) indicate significant differences.

indicate significant differences in a joint analysis between dwarf and wild-type.

To test for an impact of GA signaling on the ionome and provide single-locus tests, we grew the dominant GA-insensitive mutants *D9-1* and *D8-mpl* and their null segregating siblings (*sib9* and *sib8*). The dominant *D8-mpl* and *D9-1* alleles have nearly equivalent effects on above-ground plant growth and similar GA insensitivity phenotypes in the shoots(50). Both mutants were obtained from the maize genetics co-op and crossed three times to inbred B73 to generate BC2F1 families segregating 1:1 for the dwarf phenotype. Ears from phenotypically dwarf and phenotypically wild-type siblings were collected and processed for single-seed ionomic profiling using ICP-MS (Figure 9). Both dwarf lines had significantly different elemental compositions compared to their wild-type siblings. A joint analysis by *t*-tests between least-squared means comparing dwarfs and wild-types revealed that Cu, Fe, P, and Sr were higher in the dwarf than wild-type seeds (designated with two asterisks in Figure 9). Dominant mutants of *d8* are expressed at lower levels than *d9* in the root but at many fold higher levels in the shoot (qteller.com(53)). *D8-mpl* was also significantly different from its sibling in Cd and Mo accumulation. It is possible that *D8-mpl* has a shoot-driven effect on Mo accumulation in the seed, but we note that previous work(54) identified a large-effect QTL affecting Mo and containing the *Mot1* gene a mere 22 Mb away from *d8*. As the allele at *Mot1* is unknown in the original *D8-mpl* genetic background, linkage drag carrying a *Mot1* allele cannot be ruled out. This dominant-negative allele of *D9-1* did not recapitulate the Cd accumulation effect of the linked GWAS QTL that was the basis for its discovery as a high-confidence candidate gene by Camoco. However, the *D8-mpl* allele did recapitulate the accumulation effect, and our data demonstrate that both *D8* and *D9* have broad effects on other ionomic phenotypes.

Figure 10



Co-expression network for D9 and cadmium HPO genes

Co-expression interactions among high-priority candidate (HPO) genes were identified in the ZmRoot network for Cd and visualized at several levels. Panel (A) shows local interactions among the 126 cadmium HPO genes (red nodes). Genes are grouped and positioned based on chromosomal location. Interactions among HPO genes and D9 (GRMZM2G024973) are highlighted in yellow. Panel (B) shows a force-directed layout of D9 with HPO neighbors. Circled genes show sets of genes with previously known roles in elemental accumulation.

Genes co-expressed with D9 were investigated to determine which were associated with ionic traits, in particular, seed Cd levels. In the ZmRoot network, D9 had strong co-expression interactions with 38 other HPO genes (Figure 10A). Among these were the maize Shortroot paralog (GRMZM2G132794) and a second GRAS domain transcription factor (GRMZM2G079470). Both of these, as well as the presence of many cell-cycle genes among the co-expressed genes and ionomics traits affecting genes, raised the possibility that, like in *Arabidopsis*(55), DELLA-dependent processes, which are responsive to GA, shape the architecture of the root and the maize ionome. In *Arabidopsis*, DELLA expression disrupts Fe uptake, and loss of DELLA prevents some Fe-deficiency-mediated root growth suppression. Our finding that constitutive DELLA activity in the roots results in excess Fe, as determined by the *D9-1* and *D8-mpl* mutants, points to a conserved role for the DELLA domain transcription factors and GA signaling for Fe homeostasis in maize, a plant with an entirely different Fe uptake system than *Arabidopsis*. However, the direction of the effect was opposite to that observed in

*Arabidopsis*. Future research into the targets of the DELLA proteins in maize will be required to further address these differences.

Remarkably, the HPO co-expression network associated with D9 in the roots contained three genes with expected roles in the biosynthesis and polymerization of phenylpropanoids(56). The genes encoding enzymes that participate in phenylpropanoid biosynthesis, *ccr1* (GRMZM2G131205), the maize *ligB* paralog (GRMZM2G078500), and a laccase paralog, were co-expressed with D9 (GRMZM2G336337). The gene, *ligB*, which in angiosperms such as *Arabidopsis* is only known to be required for the formation of a pioneer specialized metabolite of no known function, was linked to QTL for multiple ions including Cd, Mn, Zn, and Ni. The gene, *ccr1*, however, was only in the HPO set for Cd. The *laccase-12* gene (GRMZM2G336337) was also a multi-ionomic hit with linked SNPs affecting Cd, Fe, and P. Genes co-expressed with D9 also were identified in the ZmPAN network. Consistent with the hypothesis that maize DELLAs regulated the type II iron uptake mechanism used by grasses, the *nicotianamine synthase3* gene (GRMZM2G439195, ZmPAN-Cd), which is required for making the type II iron chelators, was both a Cd GWAS hit and substantially co-expressed with D9 in the ZmPAN network, such that it contributed to the identification of d9 as an HPO gene for Cd.

#### Camoco identifies GWAS candidates for ion accumulation in maize seeds

In addition to the mutant analysis of HPO genes identified by our approach, we manually examined literature support for the association of candidate genes with ionomic traits. Complementing genes with known roles in elemental homeostasis, HPO gene sets for some ionomic traits included multiple genes encoding known members of the same pathway or protein complex. This suggests that biological signal was enriched by our novel combination of expression level polymorphisms and GWAS and provided evidence of novel associations between multiple pathways and elemental homeostasis.

For example, one gene with highly pleiotropic effects on the maize kernel ionome is *sugary1* (*su1*; GRMZM2G138060)(57). Genetic polymorphisms that affect seed compartment proportions or the major storage constituents are expected to contribute disproportionately to variation in seed ionomic content. Within the NAM population, functional variation for *su1* can be found in the B73 x IL14H subpopulation. For this reason, six IL14H recombinant inbred lines (RILs) that were still segregating for the recessive *su1* allele were previously tested for ionomic effects(57). This demonstrated that segregation for a loss of function allele at *su1*, on the cob, affected the levels of P, S, K, Ca, Mn, Fe, As, Se, and Rb in the seed(57). The *su1* gene was present among the HPO genes for Se accumulation (Supp. Table 6) based on the root co-expression network (ZmRoot-Se).

The *su1* locus was only identified in the HPO set for the element Se, but was linked to significant NAM GWAS SNPs for the elements P, K, and As. Thus, of the eight elements that were identified as co-segregating with the *su1* allele in the IL14H RIL population and measured in the NAM panel, four were associated with *su1* variation in the association panel. It is possible that *su1*, which is expressed in multiple plant compartments including the roots (qteller.com(53)), might also have effects throughout the seed ionome beyond a dramatic loss of seed starch. This may result from coordinate regulation of the encoded isoamylase and other root-expressed determinants of S and Se metabolism, or from unexpected coordination between root and seed expression networks. The finding that HPO network neighbors for P were enriched for carbohydrate biosynthetic enzymes favors the former of these two hypotheses (see Figure 8).

Our combined analysis of loci-linked GWAS SNPs and gene co-expression networks identified a large number of HPO genes associated with Se accumulation. Several genes with known effects on the ionome, or known to be impacted by the ionome, were identified within this HPO set. For example, GRMZM2G327406, encodes an adenylyl-sulfate kinase (*adenosine-5'-phosphosulfate [APS] kinase 3*), which is a key component of the sulfur and selenium assimilation pathway and plays a role in the formation of the substrate for protein and metabolite sulfation (ZmRoot-Se). At another locus, Camoco identified a cysteine desulfurase (GRMZM2G581155), critical for the metabolism of sulfur amino acids and the biosynthesis of the 21st amino acid selenocysteine, as an HPO gene (ZmRoot-Se).

Based on the work of Chao et al. in *Arabidopsis*, alterations in cell size and cell division in the root are expected to have effects on K accumulation in leaves(45). Two of the four subunits of the polycomb repressive complex 2 (PRC2), known to act on the cell cycle via the retinoblastoma-related proteins (RBRs), were identified as HPO genes for the K analog Rb. Both *msi1* (GRMZM2G090217; ZmSAM-Rb) and *fie2* (GRMZM2G148924; ZmSAM-Rb), members of the PRC2, are co-expressed in the ZmSAM network. The RBR-binding E2F-like transcription factor encoded by GRMZM2G361659 (ZmSAM-Rb) was also found, a further indication that cell-cycle regulation via these proteins' interactions could provide a common mechanism for these associations. Histone deacetylases from the RPD3 family are known to interact with RBR proteins as well. The RPD3-like *histone deacetylase 2* from maize was identified in the same HPO set (GRMZM2G136067; ZmSAM-Rb). The *Arabidopsis* homologs of both MSI and this histone deacetylase have known roles as histone chaperones, and the latter directly binds histone H2B. Remarkably, *histone H2B* (GRMZM2G401147; ZmSAM-Rb) was also an HPO hit. Lastly, an actin utilizing SNF2-like *chromatin regulator18* (GRMZM2G126774; ZmSAM-Rb) was identified as yet another SAM-Rb hit. This mirrors the similar finding of GO enrichment for chromatin regulatory

categories in the HPO+ enrichment analysis presented above. Taken together, these demonstrate a strong enrichment for known protein-protein interactors important for chromatin regulation and cell-cycle control among the HPO set for the K analog Rb.

A number of transporters with known roles in ionome homeostasis were also identified among the HPO genes. Among these were a P-type ATPase transporter of the ACA P2B subfamily 4 (GRMZM2G140328; ZmRoot-Sr) encoding a homolog of known plasma membrane localized Ca transporters in multiple species(58), an ABC transporter homolog of the family involved in organic acid secretion in the roots from the As HPO set (GRMZM2G415529; ZmRoot-As)(59), and a pyrophosphate energized pump (GRMZM2G090718; ZmPAN-Cd). Several annotated transporters were identified in the HPO sets for multiple elements: a sulfate transporter (GRMZM2G444801; ZmRoot-K), a cationic amino acid transporter (AC207755.3\_FG005; ZmPAN-Cd, ZmPAN-Mo), and an inositol transporter (GRMZM2G142063; ZmRoot-Fe, ZmRoot-Cd, ZmRoot-Sr).

Cadmium is well measured by ICP-MS and affected by substantial genetic variance(43). We detected the largest number of HPO candidate genes for Cd (209 genes; see Table 4). Among these were the maize *glossy2* gene (GRMZM2G098239; ZmPAN-Cd), which is responsible for a step in the biosynthesis of hydrophobic barriers(60). This implicates the biosynthesis and deposition of hydrophobic molecules in accumulation of ions and may point to root processes, rather than epicuticular waxes deposition, as the primary mode by which these genes may affect water dynamics. An ARR1-like gene, GRMZM2G067702, was also an HPO gene associated with Cd (ZmRoot). Previous work has shown that ARR genes from *Arabidopsis* are expressed in the stele, where they regulate the activity of HKT1(61). This gene was expressed at the highest level in the stele at 3 days after sowing.

Integrating GWAS data with co-expression networks resulted a set of 610 HPO genes that are primed for functional validation (1.5% of the maize FGS). The further curated subset of genes described above all have previously demonstrated roles in elemental accumulation, yet represent only a small proportion of the HPO genes discovered by Camoco. Functional validation is expensive and time consuming. Combining data-driven approaches such as network integration with expert biological curation is an extremely efficient means for the prioritization of genes driving complex traits like elemental accumulation.



## Discussion

### The effects of linkage disequilibrium

Our approach addresses a challenging bottleneck in the process of translating large sets of statistically associated loci into shorter lists based on a more mechanistic understanding of these traits. Marker SNPs identified by a GWAS provide an initial lead on a region of interest, but due to linkage disequilibrium, the candidate region can be quite broad and implicate many potentially causal genes. In addition to LD, many SNPs identified by GWAS studies lie in regulatory regions quite far from their target genes(12–14). Previous studies in maize found that while LD decays rapidly in maize (~1 kb), the variance can be large due to the functional allele segregating in a small number of lines(7). Additionally, Wallace et al. showed that the causal polymorphism is likely to reside in regulatory regions, that is, outside of exonic regions.

These factors can result in a very large (upward of 57% of all genes here) and ambiguous set of candidate genes. Until we precisely understand the regulatory landscape in the species being studied, even the most powerful GWAS will identify polymorphisms that implicate genes many base pairs away. Here, we find that the large majority of HPO genes were often not the closest genes to the identified SNPs. These genes would likely not have been identified using the common approach of prioritizing the genes closest to each marker SNP (Figure 6).

A common approach to interpreting such a locus is through manual inspection of the genome region of interest with a goal of identifying candidate genes whose function is consistent with the phenotype of interest. This can introduce bias into the discovery process and necessarily ignores uncharacterized genes. For non-human and non-model species, like maize, this manual approach is especially ineffective because the large majority of the genome remains functionally uncharacterized. Camoco leverages the orthogonal use of gene expression data, which can now be readily collected for most species of interest, to add an unbiased layer of relevant biological context to the interpretation of GWAS data and the prioritization of potentially causal variants for further experimental validation. We evaluated our framework under simulated conditions as well as applied to a large scale GWAS in order to define different co-expression metrics and networks, biases such as *cis* co-expression, and network parameters needed to be considered in order to identify co-expression signal.

Camoco successfully identified subsets of genes linked to candidate SNPs that also exhibit strong co-expression with genes near other candidate SNPs. The resulting prioritized gene sets (HPO genes) reflect groups of co-regulated genes that can potentially be used to infer a broader

biological process in which genetic variation affects the phenotype of interest. Indeed, using Camoco, we found strong evidence for HPO gene sets in 13 of the 17 elemental accumulation phenotypes we examined (with 5 or more HPO genes). These high-priority sets of genes represent a small, targeted subset of the candidates implicated by the GWAS for each phenotype (see Supp. Table 5 and Table 4).

### Establishing performance expectations of Camoco

It is important to note caveats to our approach. For example, phenotypes caused by genetic variation in a single or small number of genes or, alternatively, caused by a diverse set of otherwise functionally unrelated genes are not good candidates for our approach. The core assumption underpinning Camoco is that there are multiple variants in different genes, each driving phenotypic variation by virtue of their involvement in a common biological. We expect that this assumption holds for many phenotypes (supported by the fact that we have discovered strong candidates for the most traits examined), but we expect there are exceptional traits and causal genes that will violate this assumption. For these traits and genes, Camoco cannot be applied. Additionally, expression data used to build networks do not fully overlap with genomic data included in GWAS. For example, of the 39,656 genes in the maize filtered gene set, 11,718 genes did not pass quality control qualifications and were absent from the three co-expression networks analyzed here; they thus could not be analyzed despite the possibility there were potentially significant GWAS SNPs nearby.

### Camoco-discovered gene sets are as coherent as GO terms

In evaluating the expected performance of our approach, we simulated the effect of imperfect SNP-to-gene mapping by assuming that GO terms were identified by a simulated GWAS trait. Neighboring genes (encoded nearby on the genome) were added to simulate the scenario where we could not resolve the causal gene from linked neighboring genes. This analysis was useful, as it established the boundaries of possibility for our approach, that is, how much noise in terms of false candidate genes can be tolerated before our approach fails. As described in Figure 5, this analysis suggests a sensitivity of ~40% using a  $\pm 500$ -kb window to map SNPs to genes (two flanking genes maximum), or a tolerance of nearly 75% false candidates due to SNP-to-gene mapping. Therefore, if linkage regions implicated by GWAS extend so far as to include more than 75% false candidates, we would not be likely to discover processes as coherent as GO terms.

At the same window/flank parameter setting noted above, we were able to make significant discoveries (genes with  $FDR \leq 0.30$ ) for 7 of 17 elements (41%) using the density metric in the ZmRoot network. This success rate is remarkably consistent with what was predicted by our GO

simulations at the same window/flanking gene parameter setting. Intriguingly, HPO gene sets alone were not significantly enriched for GO term genes, indicating that while the HPO gene sets and GO terms exhibited strikingly similar patterns of gene expression, the gene sets they described do not significantly overlap. It was not until the HPO gene sets were supplemented with co-expression neighbors (HPO+) that gene sets exhibited GO term enrichment, though the resulting terms were not very specific. We speculate that this is due to discovery bias in the GO annotations that were used for our evaluation, which were largely curated from model species and assigned to maize through orthology. There are likely a large number of maize-specific processes and phenotypes that are not yet annotated in ontologies such as GO, yet have strong co-expression evidence and can be given functional annotations through GWAS.

Our analysis shows that loci implicated by ionomic GWAS loci exhibit patterns of co-expression as strong as many of the maize genes co-annotated to GO terms. Additionally, gene sets identified by Camoco have strong literature support for being involved in elemental accumulation despite not exhibiting GO enrichment. Indeed, one of the key motivations of our approach was that crop genomes like maize have limited species-specific gene ontologies, and this result emphasizes the extent of this limitation. Where current functional annotations, such as GO, rely highly on orthology, future curation schemes could rely on species-specific data obtained from GWAS and co-expression.

Beyond highlighting the challenges of a genome lacking precise functional annotation, these results also suggest an interesting direction for future work. Despite maize genes' limited ontological annotations, many GWAS have been enabled by powerful mapping populations (e.g., NAM(1)). Our results suggest that these sets of loci, combined with a proper mapping to the genes they represent using co-expression, could serve as a powerful resource for gene function characterization. Systematic efforts to curate the results from such GWAS using Camoco and similar tools, then providing public access in convenient forms, would be worthwhile. Maize is exceptional in this regard due to its excellent genomic tools and powerful mapping populations. There are several other crop species with rich population genetic resources but limited genome functional annotation that could also benefit from this approach.

#### Co-expression context matters

Using our approach, we evaluated 17 ionomic traits for overlap with three different co-expression networks. Two of the co-expression networks were generated from gene expression profiles collected across a diverse set of individuals (ZmRoot, ZmPAN) and performed substantially better than the ZmSAM network, which was based on a large collection of expression profiles across

different tissues and developmental stages derived from a single reference line (B73). We emphasize that this result is not a reflection of the data quality or even the general utility of the co-expression network used to derive the tissue/developmental atlas. Evaluations of this network showed a similar level of enrichment for co-expression relationships among genes involved in the same biological processes (Table 1) and had very similar network structure (Table 2). Instead, our results indicate that the underlying processes driving genotypic variation associated with traits captured by GWAS are better captured by transcriptional variation observed across genetically diverse individuals. Indeed, despite networks having similar levels of GO term enrichment (Table 1), the actual GO terms that drove that enrichment are quite different (Supp. Table 1), which is consistent with our previous analysis demonstrating that the experimental context of co-expression networks strongly influences which biological processes it captures(34).

Between the two co-expression networks based on expression variation across genotypically diverse individuals, we also observed differences depending on which tissues were profiled. Our co-expression network derived from sampling of root tissue across a diverse set of individuals (ZmRoot) provided the best performance at the FDR we analyzed (Table 4), producing a total of 335 (326 from density and 11 from locality, 2 in both) HPO candidate genes as compared to 228 (all from locality) HPO candidate genes produced by the ZmPAN network, which was derived from expression profiles of whole seedlings. This result affirms our original motivation for collecting tissue-specific gene expression profiles: we expected that processes occurring in the roots would be central to elemental accumulation phenotypes, which were measured in kernels. However, the difference between the performance of these two networks was modest and much less significant than the difference between the developmental/tissue atlas-derived network and the diverse genotype-derived network. Furthermore, we expect neither the ZmRoot nor the ZmPAN network to fully describe elemental accumulation processes. While ions are initially acquired from the soil via the root system, we do not directly observe their accumulation in the seed. The datasets presented here could be further complemented by additional tissue-specific data, such as genotypically diverse seed or leaf networks.

The performance of the ZmRoot versus the ZmPAN network was also quite different depending on which network metric we used. Specifically, HPO gene discovery in the ZmRoot network was driven by the density metric, while performance of the ZmPAN network relied on the locality metric (Table 4). However, locality and density were positively correlated in both networks (Supp. Figure 6), implying that these two metrics are likely complementary. Indeed, this relationship was also seen for density and locality of GO terms. Table 1 shows that both metrics had similar overall performance, each capturing ~40% of GO terms in each network; however, only ~25% was

captured by both metrics, indicating that there are certain biological processes where one metric is more appropriate than the other. In addition to the tissue source differing between the ZmRoot and ZmPAN networks, the number of experimental accessions drastically differed as well (503 accessions in ZmPAN and 48 in ZmRoot), and this influenced the performance of network metrics. We showed that locality was sensitive to the number of accessions used to calculate co-expression (Supp. Table 7) and thus could explain the bias between network metrics and the number of input accessions. This result also suggests that the 46 accessions in ZmRoot did not saturate this approach for co-expression signal and that expanding the ZmRoot dataset to include 503 accessions would result in greater power to detect overlap and the identification of more true positives using locality.

In general, our results strongly suggest that co-expression networks derived from expression experiments profiling genetically diverse individuals, as opposed to deep expression atlases derived from a single reference genotype, will be more powerful for interpreting candidate genetic loci identified in a GWAS. Furthermore, our findings suggest that where it is possible to identify relevant tissues for a phenotype of interest, tissue-specific expression profiling across genetically diverse individuals is an effective strategy. Identifying the best co-expression context for a given GWAS has important implications for data generation efforts in future studies.

## Conclusion

We developed a tool, Camoco, which integrates co-expression networks with GWAS data in order to better identify functionally relevant causal variants. We used Camoco to examine loci associated with elemental accumulation in maize grain. To do this, we built three different co-expression networks and simulated their ability to detect co-expression using GO terms. We then used these networks to identify patterns of co-expression in a set of GWAS traits measuring seed accumulation for 17 different elements, resulting in the discovery of 610 high-confidence candidate causal genes. These candidate gene sets were enriched for bioprocesses related to the ionome. Although the large majority of the high-confidence candidate genes are uncharacterized and worth further study, we did find linkage between ionomic traits and alleles at genes that have previously been demonstrated to affect the plant ionome. We validated our approach using genes and pathways not previously demonstrated to affect the ionome in maize and demonstrated that GA signaling through the DELLA domain transcription factors broadly impacted the plants' elemental profiles. Our approach successfully prioritizes causal genes underlying GWAS-identified loci based solely on gene expression data and establishes a basis for functional interpretation of otherwise uncharacterized genes associated with complex traits.

## Materials and Methods

### Software implementation of Camoco

Camoco (Co-analysis of molecular components) is a python library that includes a suite of command line tools to inter-relate and co-analyze different layers of genomic data. Specifically, it integrates genes present near GWAS loci with functional information derived from gene co-expression networks. Camoco was developed to build and analyze co-expression networks from gene transcript expression data (i.e., RNA-Seq), but it can also be utilized on other expression data such as metabolite, protein abundance, or microarray data.

This software implements three main routines: (1) construction and validation of co-expression networks from a counts or abundance matrix, (2) mapping SNPs (or other loci) to genes, and (3) an algorithm that assesses the *overlap* of co-expression among candidate genes near significant GWAS peaks.

Camoco is open source and freely available under the terms of the MIT license. Full source code, software examples, as well as instructions on how to install and run Camoco are available at <http://github.com/schae234/Camoco>. Camoco version 0.5.0 (DOI:10.5281/zenodo.1049133) was used for this article.

### Construction quality control of co-expression networks

#### Camoco Parameters

All networks were built (using the CLI) with the following Camoco QC parameters:

- min\_expr\_level: 0.001 (expression [FPKM] below this is set to NaN)
- max\_gene\_missing\_data: 0.3 (genes missing expression data more than this percent were removed from analysis)
- max\_accession\_missing\_data: 0.08 (Accessions missing expression data in more than this percent were removed from analysis)
- min\_single\_sample\_expr: 1.0 (genes must have at least this amount of expression [FPKM] in one accession)

### ZmPAN: A genotypically diverse, PAN genome co-expression network

Camoco was used to process the fragments per kilobase per million reads (FPKM) table reported by Hirsh et al. and to build a co-expression network. The raw gene expression data were passed through the quality control pipeline in Camoco. After QC, 24,756 genes were used to build the network. For each pairwise combination of genes, a Pearson correlation coefficient (PCC) was

calculated across FPKM profiles to produce ~306 million network edge scores (Supp. Fig. 1A), which were then mean centered and standard normalized (z-score hereafter) to allow cross network comparison (Supp. Fig. 1B). A global significance threshold of  $z \geq 3$  was set on co-expression interactions in order to calculate gene degree and other conventional network measures.

To assess overall network health, several approaches were taken. First, the z-scores of edges between genes co-annotated in the maize gene ontology (GO) terms were compared to edges in 1,000 random terms containing the same number genes. Supp. Fig. 1C shows the distribution of  $p$ -values compared to empirical z-score of edges within a GO term. With a nominal  $p$ -value cutoff of 0.05, the PAN co-expression network had 11.9-fold more GO terms than expected with  $p \leq 0.05$ , suggesting that edges within this co-expression network capture meaningful biological variation. Degree distribution is also as expected within the network. Supp. Fig. 1D shows empirical degree distributions compared to the power law, exponential, and truncated power law distributions. Typically, the degree distributions of biological networks are best fit by a truncated power law distribution, which is consistent with the ZmPAN genome co-expression network(41).

#### ZmSAM: A maize single accession map co-expression network

Publicly available gene expression data were generated from Stelpflug et al(38). In total, 22,691 genes passed quality control metrics. Similar to the ZmPAN network described above, gene interactions were calculated between each pairwise combination of genes to produce ~257 million network edges. A global significance threshold of  $z \geq 3$  was set on co-expression interactions in order to differentiate significantly co-expressed gene pairs.

Supp. Fig. 2A shows the distribution of edge scores before they were mean centered and standard normalized (Supp. Fig. 2B). The ZmSAM network shows a 10.8-fold enrichment for strong edge scores ( $p \leq 0.05$ ) between genes annotated to the same GO terms (Supp. Fig. 2C). A final network health check shows that the empirical degree distribution of the ZmSAM network is consistent with previously characterized biological networks (Supp. Fig. 2D).

#### ZmRoot: A genotypically diverse maize root co-expression network

Root RNA was extracted and sequenced from 48 diverse maize lines using TruSeq Stranded RNA Library Prep and Illumina HiSeq 100-bp paired-end RNA sequencing (RNA-Seq) reads. Raw reads were deposited into the short read archive (SRA) under project number PRJNA304663. Raw reads were passed through quality control using the program AdapterRemoval(62), which collapses overlapping reads into high-quality single reads while also trimming residual PCR

adapters. Reads were then mapped to the maize 5b reference genome using BWA(63,64), PCR duplicates were detected and removed, and then realignment was performed across detected insertions and deletions, resulting in between 14 and 30 million high-quality, unique nuclear reads per accession. Two accessions were dropped due to low coverage, bringing the total number to 46.

Quantification of gene expression levels into FPKM was done using a modified version of HTSeq that quantifies both paired- and unpaired-end reads(65), available at <http://github.com/schae234/MixedHTSeq>. Raw FPKM tables were imported into Camoco and passed through the quality control pipeline. After QC steps, 25,260 genes were included in co-expression network construction containing ~319 million interactions. Supp. Fig. 3A shows raw PCC scores, while Supp. Fig. 3B shows z-scores after standard normal transformation. Similar to ZmPAN and ZmSAM, co-expression among GO terms was compared to random gene sets of the same size as GO terms (1,000 instances) showing a 13.5-fold enrichment for significantly co-expressed GO terms (Supp. Fig. 3C). The degree distribution of the ZmRoot network closely follows a truncated power law similar to the other networks built here (Supp. Fig. 3D).

#### SNP-to-gene mapping and effective loci

Two parameters are used during SNP-to-gene mapping: candidate window size and maximum number of flanking genes. Windows were calculated both upstream and downstream of input SNPs. SNPs having overlapping windows were collapsed down into *effective loci* containing the contiguous genomic intervals of all overlapping SNPs, including windows both upstream and downstream of the effective locus' flanking SNPs (e.g., locus 2 in Figure 1A). Effective loci were cross referenced with the maize 5b functional gene set (FGS) genome feature format (GFF) file ([http://ftp.maizesequence.org/release-5b/filtered-set/ZmB73\\_5b\\_FGS.gff.gz](http://ftp.maizesequence.org/release-5b/filtered-set/ZmB73_5b_FGS.gff.gz)) to convert effective loci to candidate gene sets containing all candidate genes within the interval of the effective SNP and also including up to a certain number of flanking genes both upstream and downstream from the effective SNP. For each candidate gene identified by an effective locus, the number of intervening genes was calculated from the middle of the candidate gene to the middle of the effective locus. Candidate genes were ranked by the absolute value of their distance to the center of their parental effective locus. Algorithms implementing SNP-to-gene mapping used here are accessible through the Camoco command line interface.

#### Calculating subnetwork density and locality

Co-expression was measured among candidate genes using two metrics: density and locality. Subnetwork *density* is based off a z-score statistic and is formulated as the average interaction



1087 strength between *all* (un-thresholded) pairwise combinations of input genes, normalized for the  
1088 total number of input gene pairs:

1089 Eq. 1

$$1090 \quad \text{Subnetwork Density} = \frac{\bar{X} - E(X)}{\sigma(X)/\sqrt{N}}$$

1091 where  $\bar{X}$  is the calculated, mean subnetwork interaction score and  $N$  is  
1092 the number of interactions in the subnetwork. As the interaction data were  
1093 standard normalized, the expected network interaction score,  $E(X)$ , is 0, and  
1094 the standard deviation of network interactions,  $\sigma(X)$ , is 1.

1095 Network *locality* assesses the proportion of significant co-expression interactions ( $z \geq 3$ ) that are  
1096 locally connected to other subnetwork genes compared to the number of global network  
1097 interactions. To quantify network locality, both local and global degree are calculated for each  
1098 gene within a subnetwork. To account for degree bias, where genes with a high global degree are  
1099 more likely to have more local interactions, a linear regression is calculated on local degree using  
1100 global degree (designated:  $\text{local} \sim \text{global}$ ), and regression residuals for each gene are analyzed:

1101 Eq. 2

$$1102 \quad \text{Subnetwork Locality} = \text{mean}(\text{residual}(\text{local\_degree} \sim \text{global\_degree}))$$

1103 Gene-specific density is calculated by considering subnetwork interactions on a per-gene basis:

1104 Eq. 3

$$1105 \quad \text{Gene-Specific Density} = \frac{\sum \text{subnetwork\_interaction\_score}(\text{gene})}{\text{number\_of\_genes} - 1}$$

1106

1107 Gene locality residuals can be interpreted independently to identify gene-specific locality:

1108 Eq. 4

$$1110 \quad \text{Gene-Specific Locality} = \text{residual}(\text{local\_degree} \sim \text{global\_degree})$$

1109

1111 Interactions among genes that originate from the same effective GWAS locus (i.e., *cis*  
1112 interactions) were removed from density and locality calculations due to biases in *cis* co-  
1113 expression. During SNP-to-gene mapping, candidate genes retained information containing a

reference back to the parental GWAS SNP. A software flag within Camoco allows for interactions derived from the same parental SNP to be discarded from co-expression score calculations.

Statistical significance of subnetwork density and locality was assessed by comparing subnetwork scores to 1,000 random sets of candidate genes, conserving the number of input genes.

#### Simulating GWAS using Gene Ontology (GO) terms

GO(66) annotations were downloaded for maize genes from [http://ftp.maizesequence.org/release-4a.53/functional\\_annotations/](http://ftp.maizesequence.org/release-4a.53/functional_annotations/). Co-annotated genes within a GO term were treated as true causal genes identified by a hypothetical GWAS. Terms between 50 and 100 genes were included to simulate the genetic architecture of a multi-genic trait. In each co-expression network, terms having genes with significant co-expression ( $p$ -value  $\leq 0.05$ ; density or locality) were retained for further analysis. Noise introduced by imperfect GWAS was simulated using two different methods to decompose how noise affects significantly co-expressed networks.

#### Missing Candidate Rate

Eq. 6

$$MCR = 1 - \frac{\# \text{ True\_Candidate\_Genes }}{\# \text{ Candidate\_Genes }}$$

#### False Candidate Rate

Eq. 7

$$FCR = \frac{\# \text{ Candidate\_Genes } - \# \text{ True\_Candidate\_Genes }}{\# \text{ Candidate\_Genes }}$$

#### Simulating missing candidate gene rate (MCR)

The effects of MCR were evaluated by subjecting GO terms with significant co-expression ( $p \leq 0.05$ ; described above) to varying levels of missing candidate rates. True GO term genes were replaced with random genes at varying rates (MCR: 0%, 10%, 20%, 50%, 80%, 90%, 100%). The effect of MCR was evaluated by assessing the number of GO terms that retained significant co-expression (compared to 1,000 randomizations) at each level of MCR.

#### Adding false candidate genes by expanding SNP-to-gene mapping parameters

To determine how false candidates due to imperfect SNP-to-gene mapping affected the ability to detect co-expressed candidate genes linked to a GWAS trait, significantly co-expressed GO terms

were reassessed after incorporating false candidate genes. Each gene in a GO term was treated as a SNP and remapped to a set of candidate genes using the different SNP-to-gene mapping parameters (all combinations of 50 kb 100 kb, 500 kb and one, two, or five flanking genes). Effective FCR at each SNP-to-gene mapping parameter setting was calculated by dividing the number of true GO genes with candidates identified after SNP-to-gene mapping. Since varying SNP-to-gene mapping parameters changes the number of candidate genes considered within a term, each term was considered independently for each parameter combination.

## Maize ionome GWAS

Elemental concentrations were measured for 17 different elements in the maize kernel using inductively coupled plasma mass spectrometry (ICP-MS) as described in Ziegler et al.(43) Outliers were removed from single-seed measurements using median absolute deviation(67). Basic linear unbiased predictors (BLUPs) for each elemental concentration were calculated across different environments and used to estimate variance components(68). Joint-linkage analysis was run using TASSEL version 3.0(69) with over 7,000 SNPs obtained by genotype by sequencing (GBS)(70). An empirical  $p$ -value cutoff was determined by performing 1,000 permutations in which the BLUP phenotype data were shuffled within each NAM family before joint-linkage analysis was performed. The  $p$ -value corresponding to a 5% false discovery rate was used for inclusion of a QTL in the joint-linkage model.

Genome-wide association was performed using stepwise forward regression implemented in TASSEL version 4.0 similar to other studies(4,6,7). Briefly, genome-wide association was performed on a chromosomal-by-chromosome basis. To account for variance explained by QTLs on other chromosomes, the phenotypes used were the residuals from each chromosome calculated from the joint-linkage model fit with all significant joint-linkage QTLs except those on the given chromosome. Association analysis for each trait was performed 100 times by randomly sampling, without replacement, 80% of the lines from each population.

The final input SNP dataset contained 28.9 million SNPs obtained from the maize HapMap1(8), the maize HapMap2(71), as well as an additional ~800,000 putative copy-number variants from analysis of read depth counts in HapMap2(7,71). These ~30 million markers were projected onto all 5,000 lines in the NAM population using low-density markers obtained through GBS. A cutoff  $p$ -value value ( $p \leq 1e-6$ ) was used from inclusion in the final model. SNPs associated with elemental concentrations were considered significant if they were selected in more than 5 of the 100 models (resample model inclusion probability [RMIP])(44).

## Identifying ionome high-priority overlap (HPO) genes and HPO+ genes

Gene-specific density and locality were calculated for candidate genes identified from the 17 ionome GWAS traits as well as for 1,000 random sets of genes of the same size. Gene-specific metrics were converted to the standard normal scale (z-score) by subtracting the average gene-specific score from the randomized set and dividing by the average randomized standard deviation. A false discovery rate was established by incrementally evaluating the number of GWAS candidates discovered at a z-score threshold compared to the average number discovered in the random sets. For example, if ten GWAS genes had a gene-specific z-score of 3 and an average of 2.5 randomized genes (in the 1,000 random sets) had a score of 3 or above, the FDR would be 25%.

High-priority overlap (HPO) candidate genes for each element were identified by requiring candidate genes to have a co-expression  $FDR \leq 30\%$  in two or more SNP-to-gene mapping scenarios in the same co-expression network using the same co-expression metric (i.e., density or locality).

HPO+ candidate gene sets were identified by taking the number of HPO genes discovered in each element ( $n$  genes) and querying each co-expression network for the set of ( $n$ ) genes that had the strongest aggregate co-expression. For example, of the 18 HPO genes for P, an additional 18 genes (36 total) were added to the HPO+ set based on co-expression in each of the networks. Genes were added based on the sum of their co-expression to the original HPO set.

## Reduced-accession ZmPAN networks

Both the ZmPAN and ZmRoot networks were rebuilt using only the 20 accessions in common between the 503 ZmPAN and 46 ZmRoot experimental datasets. The ZmPAN network was also built using the common set of 20 accessions as well as 26 accessions selected from the broader set of 503 to simulate the number of accessions used in the ZmRoot network. Density and locality were assessed in these reduced-accession networks using the same approach as the full datasets.

## Acknowledgements

We would like to thank Ben VanderSluis, Henry Ward, and Joanna Dinsmore for their helpful comments and feedback in writing this article. We would also like to thank Abby Cabunoc-Mayes and other members of the Mozilla Science Lab for their mentorship and help in making Camoco a free and open scientific resource.

## Competing Interests

The authors declare no competing interests.

## References

1. McMullen MD, Kresovich S, Villeda HS, Bradbury P, Li H, Sun Q, et al. Genetic properties of the maize nested association mapping population. *Science* [Internet]. 2009 Aug 7 [cited 2012 Oct 29];325(5941):737–40. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/19661427>
2. Buckler ES, Holland JB, Bradbury PJ, Acharya CB, Brown PJ, Browne C, et al. The genetic architecture of maize flowering time. *Science* [Internet]. 2009 Aug 7 [cited 2012 Oct 29];325(5941):714–8. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/19661422>
3. Peiffer JA, Romay MC, Gore MA, Flint-Garcia SA, Zhang Z, Millard MJ, et al. The Genetic Architecture of Maize Height. *Genetics* [Internet]. 2014 Feb 10 [cited 2014 Mar 19]; Available from: <http://www.ncbi.nlm.nih.gov/pubmed/24514905>
4. Tian F, Bradbury PJ, Brown PJ, Hung H, Sun Q, Flint-Garcia S, et al. Genome-wide association study of leaf architecture in the maize nested association mapping population. *Nat Genet* [Internet]. 2011 Feb [cited 2012 Oct 29];43(2):159–62. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/21217756>
5. Kump KL, Bradbury PJ, Wisser RJ, Buckler ES, Belcher AR, Oropeza-Rosas M a, et al. Genome-wide association study of quantitative resistance to southern leaf blight in the maize nested association mapping population. *Nat Genet* [Internet]. 2011;43(2):163–8. Available from: <http://www.nature.com/doifinder/10.1038/ng.747>
6. Cook JP, McMullen MD, Holland JB, Tian F, Bradbury P, Ross-Ibarra J, et al. Genetic architecture of maize kernel composition in the nested association mapping and inbred association panels. *Plant Physiol* [Internet]. 2012 Feb [cited 2012 Oct 5];158(2):824–34. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/22135431>
7. Wallace JG, Bradbury PJ, Zhang N, Gibon Y, Stitt M, Buckler ES. Association mapping across numerous traits reveals patterns of functional variation in maize. *PLoS Genet* [Internet]. 2014 Dec 4 [cited 2015 Sep 24];10(12):e1004845. Available from: <http://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1004845>
8. Gore M a, Chia J-M, Elshire RJ, Sun Q, Ersoz ES, Hurwitz BL, et al. A first-generation haplotype map of maize. *Science* [Internet]. 2009 Nov 20 [cited 2011 Jun 10];326(5956):1115–7. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/19965431>
9. Morrell PL, Toleno DM, Lundy KE, Clegg MT. Low levels of linkage disequilibrium in wild barley (*Hordeum vulgare* ssp. *spontaneum*) despite high rates of self-fertilization. *Proc Natl Acad Sci U S A*. 2005;102(7):2442–7.
10. Caldwell KS, Russell J, Langridge P, Powell W. Extreme population-dependent linkage disequilibrium detected in an inbreeding plant species, *Hordeum vulgare*. *Genetics*. 2006;172(1):557–67.
11. Wray GA. The evolutionary significance of cis-regulatory mutations. *Nat Rev Genet* [Internet]. 2007

- 1243 Mar [cited 2014 Jul 11];8(3):206–16. Available from: <http://dx.doi.org/10.1038/nrg2063>
- 1244 12. Clark RM, Wagler TN, Quijada P, Doebley J. A distant upstream enhancer at the maize  
1245 domestication gene *tb1* has pleiotropic effects on plant and inflorescent architecture. *Nat Genet*  
1246 [Internet]. 2006 May [cited 2015 Aug 6];38(5):594–7. Available from:  
1247 <http://www.ncbi.nlm.nih.gov/pubmed/16642024>
- 1248 13. Castelletti S, Tuberosa R, Pindo M, Salvi S. A MITE transposon insertion is associated with  
1249 differential methylation at the maize flowering time QTL *Vgt1*. G3 (Bethesda) [Internet]. 2014 May  
1250 [cited 2015 Sep 15];4(5):805–12. Available from:  
1251 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4025479&tool=pmcentrez&rendertype=abstract>  
1252
- 1253 14. Louwers M, Bader R, Haring M, van Driel R, de Laat W, Stam M. Tissue- and expression level-specific  
1254 chromatin looping at maize *b1* epialleles. *Plant Cell* [Internet]. 2009 Mar [cited 2015 Sep  
1255 13];21(3):832–42. Available from:  
1256 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2671708&tool=pmcentrez&rendertype=abstract>  
1257
- 1258 15. Andorf CM, Cannon EK, Portwood JL, Gardiner JM, Harper LC, Schaeffer ML, et al. MaizeGDB  
1259 update: new tools, data and interface for the maize model organism database. *Nucleic Acids Res*  
1260 [Internet]. 2015;gkv1007. Available from:  
1261 <http://nar.oxfordjournals.org/content/early/2015/10/01/nar.gkv1007.full>
- 1262 16. Eisen MB, Spellman PT, Brown PO, Botstein D. Cluster analysis and display of genome-wide  
1263 expression patterns. *Proc Natl Acad Sci* [Internet]. 1998 Dec 8;95(25):14863–8. Available from:  
1264 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=24541&tool=pmcentrez&rendertype=abstract>  
1265 <http://www.pnas.org/cgi/doi/10.1073/pnas.95.25.14863>
- 1266 17. Schaefer RJ, Briskine R, Springer NM, Myers CL. Discovering functional modules across diverse  
1267 maize transcriptomes using COB, the co-expression browser. *PLoS One*. 2014;9(6).
- 1268 18. Mochida K, Uehara-Yamaguchi Y, Yoshida T, Sakurai T, Shinozaki K. Global landscape of a co-  
1269 expressed gene network in barley and its application to gene discovery in Triticeae crops. *Plant Cell*  
1270 *Physiol* [Internet]. 2011 May [cited 2011 Aug 15];52(5):785–803. Available from:  
1271 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3093127&tool=pmcentrez&rendertype=abstract>  
1272
- 1273 19. Obayashi T, Okamura Y, Ito S, Tadaka S, Aoki Y, Shirota M, et al. ATTED-II in 2014: Evaluation of  
1274 Gene Coexpression in Agriculturally Important Plants. *Plant Cell Physiol* [Internet]. 2014;55(1):e6–  
1275 e6. Available from: <http://pcp.oxfordjournals.org/cgi/doi/10.1093/pcp/pct178>
- 1276 20. Sarkar NK, Kim Y-K, Grover A. Coexpression network analysis associated with call of rice seedlings  
1277 for encountering heat stress. *Plant Mol Biol*. 2014 Jan;84(1–2):125–43.
- 1278 21. Zheng Z-L, Zhao Y. Transcriptome comparison and gene coexpression network analysis provide a  
1279 systems view of citrus response to “*Candidatus Liberibacter asiaticus*” infection. *BMC Genomics*.  
1280 2013;14(1):27.
- 1281 22. Ozaki S, Ogata Y, Suda K, Kurabayashi A, Suzuki T, Yamamoto N, et al. Coexpression analysis of  
1282 tomato genes and experimental verification of coordinated expression of genes found in a  
1283 functionally enriched coexpression module. *DNA Res* [Internet]. 2010 Apr [cited 2016 Apr

1284 16];17(2):105–16. Available from:  
1285 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2853382&tool=pmcentrez&rendert](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2853382&tool=pmcentrez&rendertype=abstract)  
1286 [ype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2853382&tool=pmcentrez&rendertype=abstract)

1287 23. Swanson-Wagner R, Briskine R, Schaefer R, Hufford MB, Ross-Ibarra J, Myers CL, et al. Reshaping  
1288 of the maize transcriptome by domestication. PNAS [Internet]. 2012 Jul 17 [cited 2016 Jun  
1289 7];109(29):11878–83. Available from: <http://www.pnas.org/cgi/doi/10.1073/pnas.1201961109>

1290 24. Wolfe CJ, Kohane IS, Butte AJ. Systematic survey reveals general applicability of “guilt-by-  
1291 association” within gene coexpression networks. BMC Bioinformatics [Internet]. 2005 Jan [cited  
1292 2016 Apr 8];6:227. Available from:  
1293 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1239911&tool=pmcentrez&rendert](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1239911&tool=pmcentrez&rendertype=abstract)  
1294 [ype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1239911&tool=pmcentrez&rendertype=abstract)

1295 25. Ritchie MD, Holzinger ER, Li R, Pendergrass SA, Kim D. Methods of integrating data to uncover  
1296 genotype–phenotype interactions. Nat Rev Genet [Internet]. 2015 Jan 13;16(2):85–97. Available  
1297 from: <http://dx.doi.org/10.1038/nrg3868>

1298 26. Li M, Chen J, Wang J, Hu B, Chen G. Modifying the DPCLUS algorithm for identifying protein  
1299 complexes based on new topological structures. BMC Bioinformatics [Internet]. 2008 Jan 25 [cited  
1300 2016 Apr 28];9(1):398. Available from:  
1301 <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-9-398>

1302 27. Calabrese GM, Mesner LD, Stains JP, Tommasini SM, Horowitz MC, Rosen CJ, et al. Integrating  
1303 GWAS and Co-expression Network Data Identifies Bone Mineral Density Genes SPTBN1 and MARK3  
1304 and an Osteoblast Functional Module. Cell Syst [Internet]. 2017;4(1):46–59.e4. Available from:  
1305 <http://dx.doi.org/10.1016/j.cels.2016.10.014>

1306 28. Bunyavanich S, Schadt EE, Himes BE, Lasky-Su J, Qiu W, Lazarus R, et al. Integrated genome-wide  
1307 association, coexpression network, and expression single nucleotide polymorphism analysis  
1308 identifies novel pathway in allergic rhinitis. BMC Med Genomics [Internet]. 2014;7(1):48. Available  
1309 from: <http://bmcmmedgenomics.biomedcentral.com/articles/10.1186/1755-8794-7-48>

1310 29. Taşan M, Musso G, Hao T, Vidal M, Macrae C a, Roth FP. Selecting causal genes from genome-wide  
1311 association studies via functionally coherent subnetworks. 2014;12(2).

1312 30. USDA. Crop Production 2015 Summary. 2016.

1313 31. Baxter I. Ionomics: The functional genomics of elements. Brief Funct Genomics [Internet]. 2010  
1314 Mar [cited 2012 Oct 29];9(2):149–56. Available from:  
1315 <http://www.ncbi.nlm.nih.gov/pubmed/20081216>

1316 32. Guerinot M Lou, Salt DE. Fortified Foods and Phytoremediation . Two Sides of the Same Coin 1.  
1317 2017;3755.

1318 33. Baxter IR, Vitek O, Lahner B, Muthukumar B, Borghi M, Morrissey J, et al. The leaf ionome as a  
1319 multivariable system to detect a plant’s physiological status. Proc Natl Acad Sci U S A [Internet].  
1320 2008 Aug 19 [cited 2015 Oct 2];105(33):12081–6. Available from:  
1321 <http://www.pnas.org/content/105/33/12081.abstract>

1322 34. Schaefer RJ, Briskine R, Springer NM, Myers CL. Discovering functional modules across diverse  
1323 maize transcriptomes using COB, the co-expression browser. PLoS One. 2014;9(6):99193.

- 1324 35. Swanson-Wagner R, Briskine R, Schaefer R, Hufford MB, Ross-Ibarra J, Myers CL, et al. Reshaping  
1325 of the maize transcriptome by domestication. *Proc Natl Acad Sci U S A*. 2012;109(29).
- 1326 36. Schaefer RJ, Michno J-M, Myers CL. Unraveling gene function in agricultural species using gene co-  
1327 expression networks. *Biochim Biophys Acta - Gene Regul Mech*. 2016;
- 1328 37. Hirsch CN, Foerster JM, Johnson JM, Sekhon RS, Muttoni G, Vaillancourt B, et al. Insights into the  
1329 maize pan-genome and pan-transcriptome. *Plant Cell [Internet]*. 2014 Jan 31 [cited 2014 Jul  
1330 14];26(1):121–35. Available from:  
1331 <http://www.plantcell.org/content/early/2014/01/31/tpc.113.119982.abstract>
- 1332 38. Stelpflug SC, Rajandeeep S, Vaillancourt B, Hirsch CN, Buell CR, Leon N De, et al. An expanded maize  
1333 gene expression atlas based on RNA-sequencing and its use to explore root development. *Plant*  
1334 *Genome*. 2015;(608):314–62.
- 1335 39. Schaefer RJ, Briskine R, Springer NM, Myers CCL, Wei H, Persson S, et al. Discovering functional  
1336 modules across diverse maize transcriptomes using COB, the co-expression browser. Börnke F,  
1337 editor. *PLoS One [Internet]*. 2014 Jun 12 [cited 2016 Jun 7];9(6):99193. Available from:  
1338 <http://dx.plos.org/10.1371/journal.pone.0099193>
- 1339 40. Dongen S van. MCL: A Cluster Algorithm for Graphs. Center for Information Workshop; 2000.
- 1340 41. Ghazalpour A, Doss S, Zhang B, Wang S, Plaisier C, Castellanos R, et al. Integrating genetic and  
1341 network analysis to characterize genes related to mouse weight. Gibson G, editor. *PLoS Genet*  
1342 *[Internet]*. 2006 Aug 18 [cited 2014 Apr 29];2(8):e130. Available from:  
1343 <http://dx.plos.org/10.1371/journal.pgen.0020130>
- 1344 42. Baxter I, Dilkes BP. Elemental profiles reflect plant adaptations to the environment. *Science*  
1345 *[Internet]*. 2012 Jun 29 [cited 2015 Oct 4];336(6089):1661–3. Available from:  
1346 <http://www.sciencemag.org/content/336/6089/1661.abstract>
- 1347 43. Ziegler G, Kear PJ, Wu D, Ziyomo C, Lipka AE, Gore M, et al. Elemental Accumulation in Kernels of  
1348 the Maize Nested Association Mapping Panel Reveals Signals of Gene by Environment Interactions.  
1349 *bioRxiv*. 2017;(May).
- 1350 44. Valdar W, Holmes CC, Mott R, Flint J. Mapping in structured populations by resample model  
1351 averaging. *Genetics [Internet]*. 2009 Aug 1 [cited 2015 Aug 6];182(4):1263–77. Available from:  
1352 <http://www.genetics.org/content/182/4/1263.long>
- 1353 45. Chao D-Y, Gable K, Chen M, Baxter I, Dietrich CR, Cahoon EB, et al. Sphingolipids in the Root Play  
1354 an Important Role in Regulating the Leaf Ionome in *Arabidopsis thaliana*. *Plant Cell [Internet]*.  
1355 2011;23(3):1061–81. Available from: <http://www.plantcell.org/cgi/doi/10.1105/tpc.110.079095>
- 1356 46. Fan J, Zhai Z, Yan C, Xu C. *Arabidopsis* TRIGALACTOSYLDIACYLGLYCEROL5 Interacts with TGD1,  
1357 TGD2, and TGD4 to Facilitate Lipid Transfer from the Endoplasmic Reticulum to Plastids. *Plant Cell*  
1358 *[Internet]*. 2015;27(October):tpc.15.00394. Available from:  
1359 <http://www.plantcell.org/lookup/doi/10.1105/tpc.15.00394>
- 1360 47. Katagiri T, Ishiyama K, Kato T, Tabata S, Kobayashi M, Shinozaki K. An important role of  
1361 phosphatidic acid in ABA signaling during germination in *Arabidopsis thaliana*. *Plant J*.  
1362 2005;43(1):107–17.
- 1363 48. Roston RL, Gao J, Murcha MW, Whelan J, Benning C. TGD1, -2, and -3 proteins involved in lipid



- 1364 trafficking form ATP-binding cassette (ABC) transporter with multiple substrate-binding proteins. J  
1365 Biol Chem. 2012;287(25):21406–15.
- 1366 49. Lawrence CJ, Dong Q, Polacco ML, Seigfried TE, Brendel V. MaizeGDB, the community database for  
1367 maize genetics and genomics. Nucleic Acids Res [Internet]. 2004 Jan 1 [cited 2012 Oct  
1368 30];32(Database issue):D393-7. Available from:  
1369 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=308746&tool=pmcentrez&renderty](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=308746&tool=pmcentrez&rendertype=abstract)  
1370 [pe=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=308746&tool=pmcentrez&rendertype=abstract)
- 1371 50. Winkler RG, Freeling M. Physiological genetics of the dominant gibberellin-nonresponsive maize  
1372 dwarfs, Dwarf8 and Dwarf9. Planta. 1994;193:341–8.
- 1373 51. Lawit SJ, Wych HM, Xu D, Kundu S, Tomes DT. Maize dwarf proteins dwarf plant8 and dwarf plant9  
1374 as modulators of plant development. Plant Cell Physiol. 2010;51(11):1854–68.
- 1375 52. Fu J, Ren F, Lu X, Mao H, Xu M, Degenhardt J, et al. A Tandem Array of *ent*-Kaurene Synthases in  
1376 Maize with Roles in Gibberellin and More Specialized Metabolism. Plant Physiol [Internet].  
1377 2016;170(2):742–51. Available from:  
1378 <http://www.plantphysiol.org/lookup/doi/10.1104/pp.15.01727>
- 1379 53. Wang X, Elling AA, Li X, Li N, Peng Z, He G, et al. Genome-Wide and Organ-Specific Landscapes of  
1380 Epigenetic Modifications and Their Relationships to mRNA and Small RNA Transcriptomes in Maize.  
1381 Plant Cell Online [Internet]. 2009;21(4):1053–69. Available from:  
1382 <http://www.plantcell.org/cgi/doi/10.1105/tpc.109.065714>
- 1383 54. Asaro A, Ziegler G, Ziyomo C, Hoekenga O, Dilkes B, Baxter I. The Interaction of Genotype and  
1384 Environment Determines Variation in the Maize Kernel Ionome. G3&#58;  
1385 Genes|Genomes|Genetics [Internet]. 2016;6(December):4175–83. Available from:  
1386 <http://g3journal.org/cgi/doi/10.1534/g3.116.034827>
- 1387 55. Wild M, Davi??re JM, Regnault T, Sakvarelidze-Achard L, Carrera E, Lopez Diaz I, et al. Tissue-  
1388 Specific Regulation of Gibberellin Signaling Fine-Tunes Arabidopsis Iron-Deficiency Responses. Dev  
1389 Cell. 2016;37(2):190–200.
- 1390 56. Monaco MK, Sen TZ, Dharmawardhana PD, Ren L, Schaeffer M, Naithani S, et al. Maize Metabolic  
1391 Network Construction and Transcriptome Analysis. Plant Genome [Internet]. 2013;6(1):0.  
1392 Available from:  
1393 <https://www.crops.org/publications/tpg/abstracts/6/1/plantgenome2012.09.0025>
- 1394 57. Baxter IR, Ziegler G, Lahner B, Mickelbart M V., Foley R, Danku J, et al. Single-kernel ionomic profiles  
1395 are highly heritable indicators of genetic and environmental influences on elemental accumulation  
1396 in maize grain (*Zea mays*). PLoS One. 2014;9(1).
- 1397 58. Baxter I, Tchieu J, Sussman MR, Boutry M, Palmgren MG, Gribskov M, et al. Genomic Comparison  
1398 of P-Type ATPase Ion Pumps in Arabidopsis and Rice 1. 2003;132(June):618–28.
- 1399 59. Badri D V., Loyola-Vargas VM, Broeckling CD, De-la-Pena C, Jasinski M, Santelia D, et al. Altered  
1400 Profile of Secondary Metabolites in the Root Exudates of Arabidopsis ATP-Binding Cassette  
1401 Transporter Mutants. Plant Physiol [Internet]. 2007;146(2):762–71. Available from:  
1402 <http://www.plantphysiol.org/cgi/doi/10.1104/pp.107.109587>
- 1403 60. Tacke E, Korfhage C, Michel D, Maddaloni M, Motto M, Lanzini S, et al. Transposon tagging of the  
1404 maize Glossy2 locus with the transposable element En/Spm. Vol. 8, The Plant Journal. 1995. p.

- 1405 907–17.
- 1406 61. Mason MG, Jha D, Salt DE, Tester M, Hill K, Kieber JJ, et al. Type-B response regulators ARR1 and  
1407 ARR12 regulate expression of AtHKT1;1 and accumulation of sodium in Arabidopsis shoots. *Plant*  
1408 *J.* 2010;64(5):753–63.
- 1409 62. Lindgreen S. AdapterRemoval: easy cleaning of next-generation sequencing reads. *BMC Res Notes*  
1410 [Internet]. 2012 Jan [cited 2015 Sep 4];5(1):337. Available from:  
1411 <http://www.biomedcentral.com/1756-0500/5/337>
- 1412 63. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
1413 *Bioinformatics* [Internet]. 2009 Jul 15 [cited 2014 Jul 9];25(14):1754–60. Available from:  
1414 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2705234&tool=pmcentrez&rendertype=abstract>  
1415
- 1416 64. Schubert M, Ermini L, Der Sarkissian C, Jónsson H, Ginolhac A, Schaefer R, et al. Characterization of  
1417 ancient and modern genomes by SNP detection and phylogenomic and metagenomic analysis  
1418 using PALEOMIX. *Nat Protoc* [Internet]. 2014;9(5):1056–82. Available from:  
1419 <http://www.ncbi.nlm.nih.gov/pubmed/24722405>
- 1420 65. Anders S, Pyl PT, Huber W. HTSeq - A Python framework to work with high-throughput sequencing  
1421 data. *Bioinformatics* [Internet]. 2014 Sep 25 [cited 2014 Sep 29];31(2):166–9. Available from:  
1422 <http://bioinformatics.oxfordjournals.org/content/31/2/166>
- 1423 66. Harris M a, Clark J, Ireland a, Lomax J, Ashburner M, Foulger R, et al. The Gene Ontology (GO)  
1424 database and informatics resource. *Nucleic Acids Res* [Internet]. 2004 Jan 1 [cited 2012 Oct  
1425 10];32(Database issue):D258–61. Available from:  
1426 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=308770&tool=pmcentrez&rendertype=abstract>  
1427
- 1428 67. Davies L, Gather U. The Identification of Multiple Outliers. *J Am Stat Assoc* [Internet]. 2012 Feb 27  
1429 [cited 2015 Nov 20]; Available from:  
1430 <http://www.tandfonline.com/doi/abs/10.1080/01621459.1993.10476339>
- 1431 68. Hung H-Y, Browne C, Guill K, Coles N, Eller M, Garcia A, et al. The relationship between parental  
1432 genetic or phenotypic divergence and progeny variation in the maize nested association mapping  
1433 population. *Heredity (Edinb)* [Internet]. 2012 May [cited 2015 Nov 20];108(5):490–9. Available  
1434 from:  
1435 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3330692&tool=pmcentrez&rendertype=abstract>  
1436
- 1437 69. Bradbury PJ, Zhang Z, Kroon DE, Casstevens TM, Ramdoss Y, Buckler ES. TASSEL: software for  
1438 association mapping of complex traits in diverse samples. *Bioinformatics* [Internet]. 2007 Oct 1  
1439 [cited 2014 Jul 12];23(19):2633–5. Available from:  
1440 <http://www.ncbi.nlm.nih.gov/pubmed/17586829>
- 1441 70. Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, et al. A robust, simple  
1442 genotyping-by-sequencing (GBS) approach for high diversity species. *PLoS One* [Internet]. 2011 Jan  
1443 [cited 2014 Jul 9];6(5):e19379. Available from:  
1444 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3087801&tool=pmcentrez&rendertype=abstract>  
1445

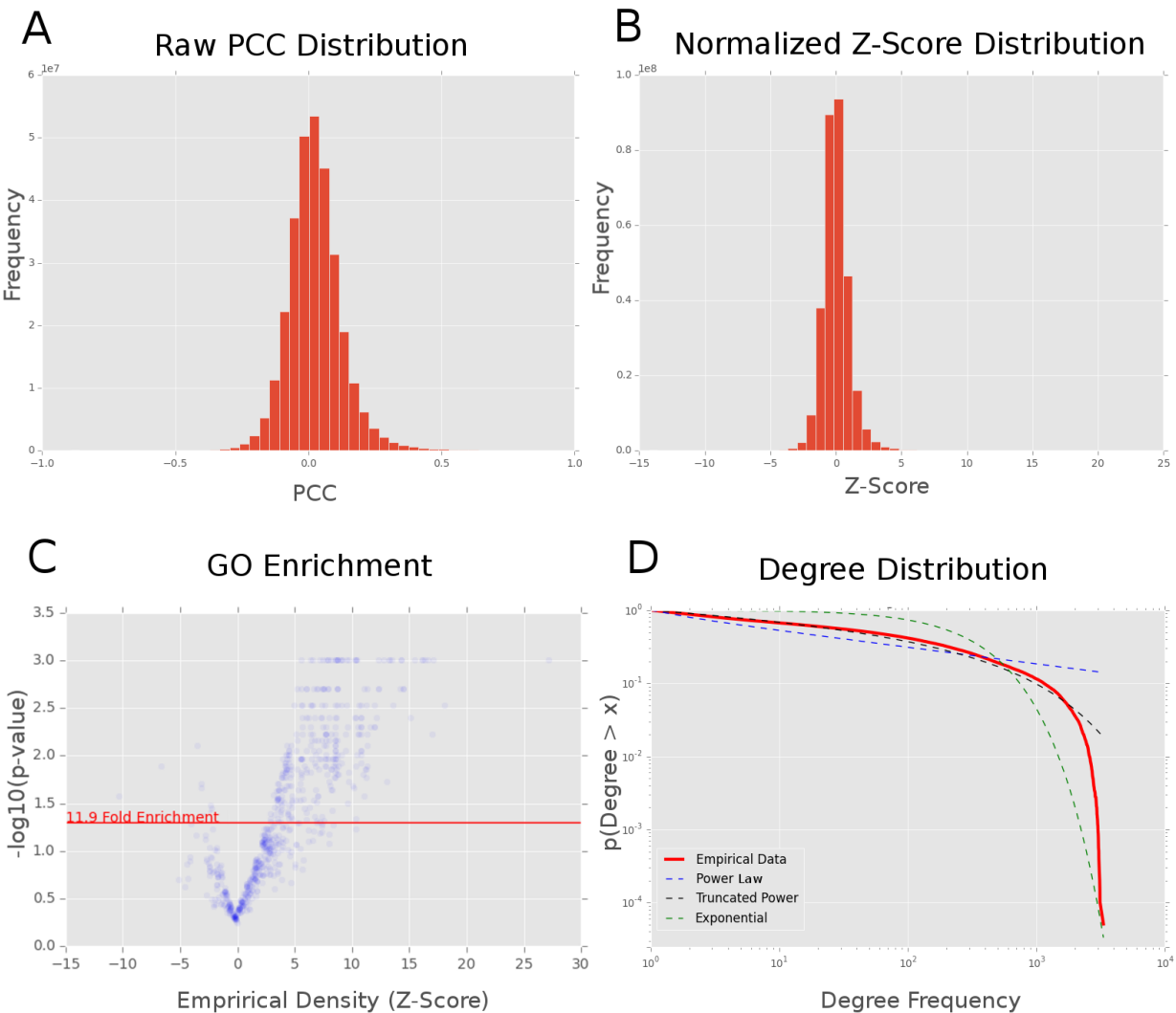
1446 71. Chia J-M, Song C, Bradbury PJ, Costich D, de Leon N, Doebley J, et al. Maize HapMap2 identifies  
 1447 extant variation from a genome in flux. Nat Genet [Internet]. 2012 Jul [cited 2012 Oct 9];44(7):803–  
 1448 7. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/22660545>

1449

1450 **Supplementary Figures**

1451 Supp. Fig. 1

# ZmPAN Network Stats



1452

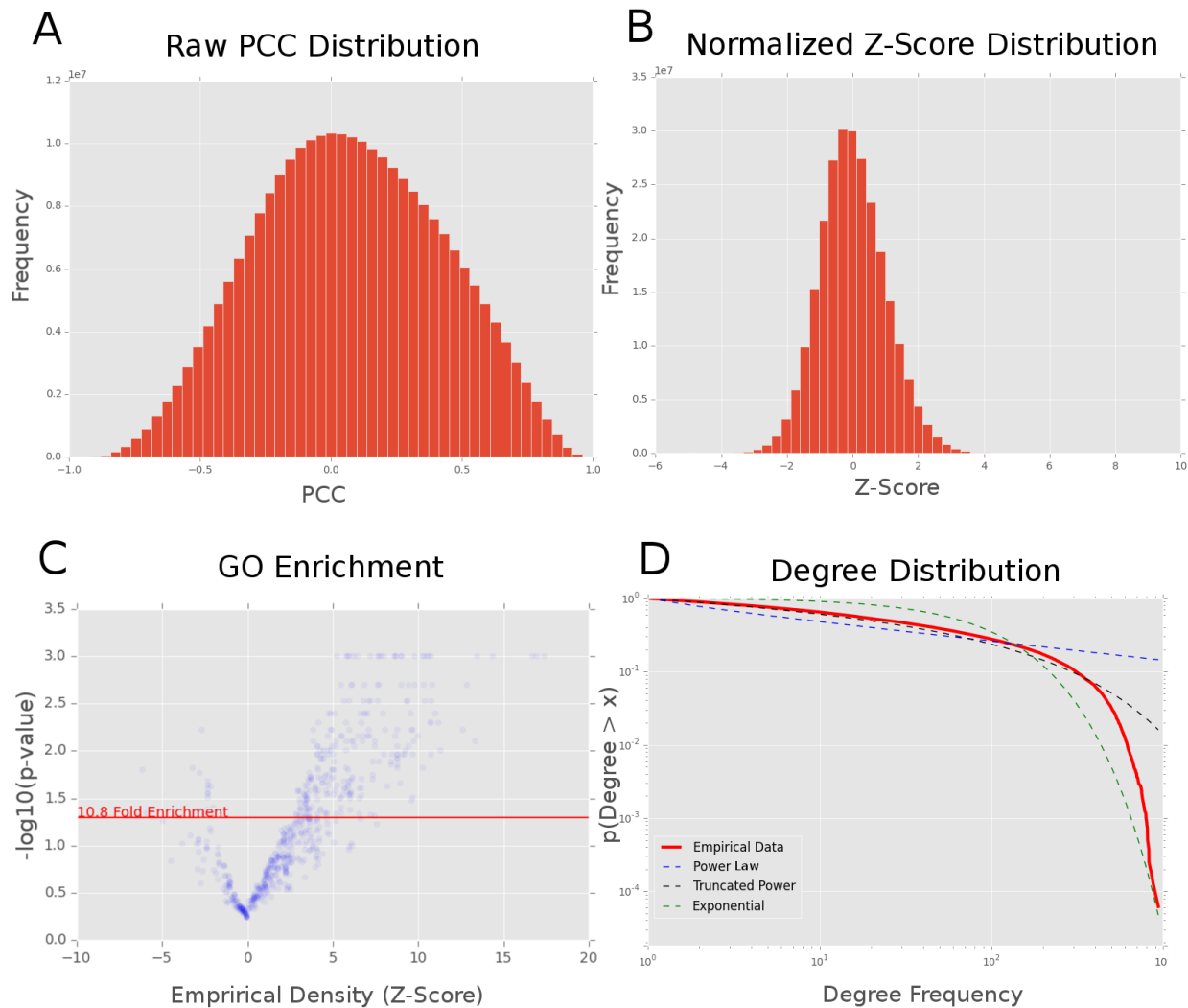
1453 **ZmPAN network health**

1454 Global network health of the maize PAN genome (ZmPAN) co-expression  
 1455 network. **(A)** Raw Pearson correlation coefficient distribution of all co-

expression interactions. **(B)** Fisher-transformed, variance-stabilized, and mean centered network interactions. **(C)** A volcano plot showing empirical density for genes in each GO term compared to the corresponding  $p$ -value derived from measuring density in 1,000 random gene sets of the same size. **(D)** Degree distribution of ZmPAN genome co-expression network compared to power law, exponential, and truncated power law distributions.

Supp. Fig. 2

## ZmSAM Network Stats



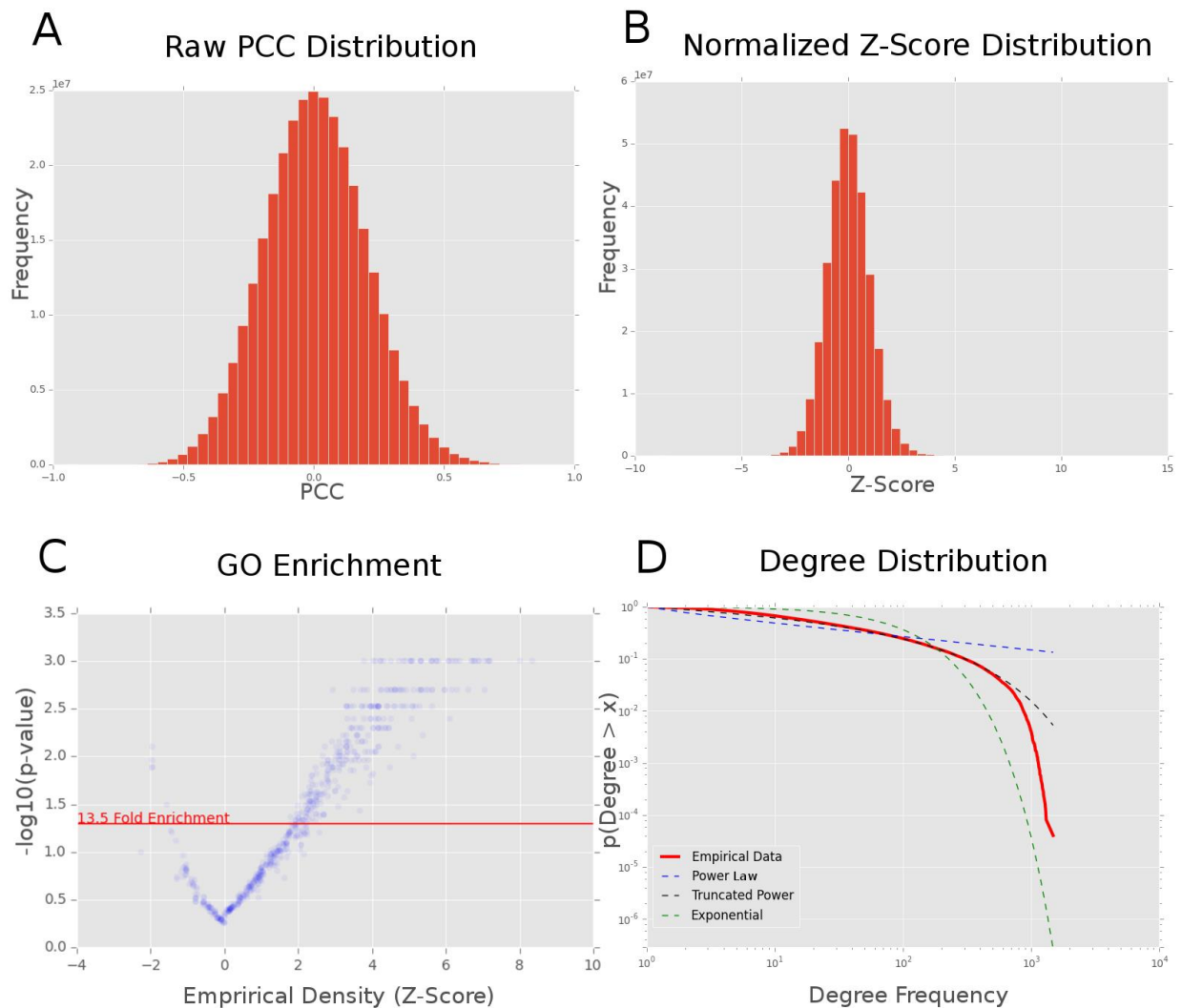
ZmSAM network health

Global network health of the maize ZmSAM co-expression network. **(A)** Raw Pearson correlation coefficient distribution of all co-expression interactions.

**(B)** Variance-stabilized and mean centered network interactions. **(C)** A volcano plot showing empirical density for genes in each GO term compared to the corresponding  $p$ -value derived from measuring density in 1,000 random gene sets of the same size. **(D)** Degree distribution of tissue/developmental co-expression network compared to power law, exponential, and truncated power law distributions.

Supp. Fig. 3

## ZmRoot Network Stats



1476 ZmRoot network health

1477 Global network health of the maize ZmRoot co-expression network. **(A)** Raw  
1478 Pearson correlation coefficient distribution of all co-expression interactions.  
1479 **(B)** Variance-stabilized and mean centered network interactions. **(C)** A  
1480 volcano plot showing empirical density for genes in each GO term compared  
1481 to the corresponding  $p$ -value derived from measuring density in 1,000 random  
1482 gene sets of the same size. **(D)** Degree distribution of ZmRoot co-expression  
1483 network compared to power law, exponential, and truncated power law  
1484 distributions.

1485



## MCR supplemental figure

Panel **(A)** shows the absolute number of GO terms that remain significantly co-expressed at varying levels of MCR in each network. Red curves show all GO terms with an initial co-expression  $p$ -value  $\leq 0.05$ . Blue and violet curves show GO terms with either moderate or strong initial co-expression (at MCR = 0). Panels **(B-C)** show the percent and absolute number of GO terms that remain significantly co-expressed at varying levels of MCR. The red curves show small GO terms ( $50 \leq n < 65$ ), the blue curves show medium sized GO terms ( $65 \leq n < 80$ ), and the violet curves show large terms ( $80 \leq n < 100$ ).



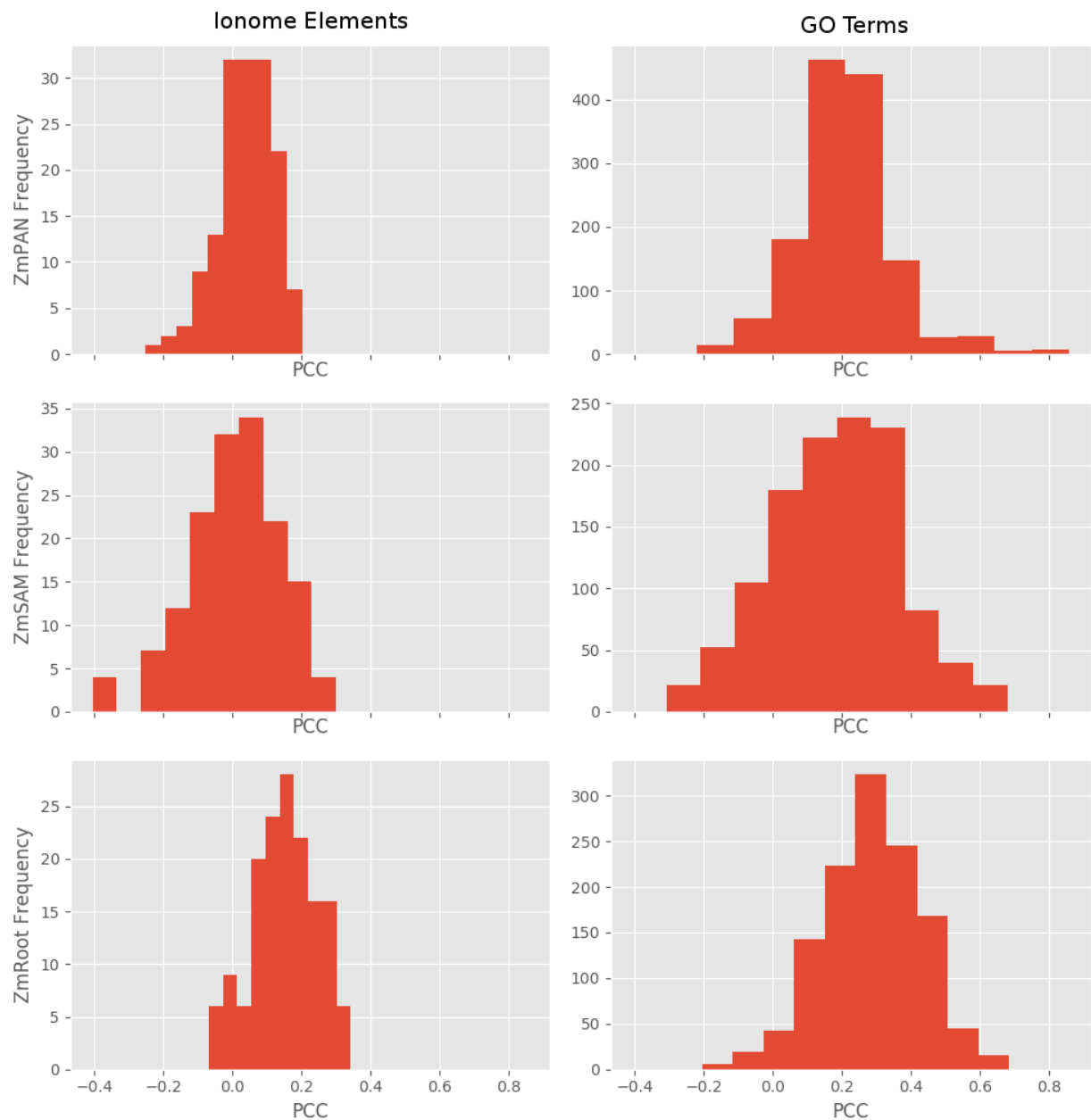
1497    Supp. Fig. 5

1498

## FCR supplemental figure

Panel **(A)** shows the absolute number of GO terms that remain significantly co-expressed at varying levels of FCR in each network. Red curves show all GO terms with an initial co-expression  $p$ -value  $\leq 0.05$ . Blue and violet curves show GO terms with either moderate or strong initial co-expression. Panels **(B-C)** show the percent and absolute number of GO terms that remain significantly co-expressed at varying levels of FCR. The red curves show small GO terms ( $50 \leq n < 65$ ), the blue curves show medium sized GO terms ( $65 \leq n < 80$ ) and the violet curves show large terms ( $80 \leq n < 100$ ).

Correlation distributions of gene-specific density vs. locality



1509

1510 Distribution of Pearson correlation coefficients between gene-specific density and locality

1511        Pearson correlation was measured between gene-specific density and locality

1512        in each network for both ionome elements and GO terms. PCCs between

1513        metrics were calculated by grouping sets of genes in either ionome elements

1514        (e.g., Al, Fe) or GO terms at the same SNP-to-gene mapping parameters (50-

1515        , 100-, and 500-kb window size and one, two, and five gene flank limits). The

distribution shows the PCCs between the metrics aggregated across all SNP-to-gene mapping parameters.

## Supplementary Tables

### Supp. Table 1

#### Full gene ontology term density and locality *p*-values

Density and locality scores were measured between genes within each GO term. Subnetwork *p*-values were generated for both density and locality by comparing each term's metric to 1,000 randomized gene sets of the same size.

### Supp. Table 2

#### Network MCL cluster gene assignments

Clusters in all three networks were identified using the MCL algorithm. Genes in each network were assigned to cluster IDs. Lower cluster IDs have a larger number of genes.

### Supp. Table 3

#### Network MCL cluster GO enrichment

Enrichment of genes co-annotated for GO terms in each MCL cluster. Significance of enrichment was calculated using the hypergeometric test with a Bonferroni corrected *p*-value of  $\leq 0.05$ .

### Supp. Table 4

#### Network signal of GO terms with various levels of MCR/FCR

Co-expression among genes co-annotated to GO terms was compared to random gene sets of the same size to generate *p*-values. Noise was introduced by varying the missing candidate rate (MCR) or false candidate rate (FCR). Missing candidates were removed in proportion to the values in the table, while false candidates were introduced using SNP-to-gene mapping values (see WindowSize and FlankLimit columns). FCR values are reported as averages across 10% quantiles (see Figure 5).

1544 Supp. Table 5

1545 Maize grain ionome GWAS network overlap candidate genes

1546 Candidate genes were identified in each co-expression network (ZmSAM,  
1547 ZmPAN, or ZmRoot) using SNP-to-gene mapping for each element (using  
1548 WindowSize and FlankLimit). Co-expression (density or locality) among all  
1549 genes within a subnetwork was compared to randomized gene sets of the  
1550 same size to establish *p*-values. Gene-specific z-scores were computed by  
1551 comparing the empirical gene-specific density (Eq. 3) or locality (Eq. 4) to the  
1552 average density or locality observed in randomized gene sets, then correcting  
1553 for standard deviation. False discovery rates (FDRs) were calculated for  
1554 candidate genes with positive gene-specific co-expression values by  
1555 comparing the number of genes discovered at a z-score cutoff to the average  
1556 number of genes discovered in randomized sets.

1557 Supp. Table 6

1558 Maize grain ionome GWAS high-priority overlap (HPO) candidate genes

1559 High-priority overlap (HPO) genes were identified by calculating gene-specific  
1560 density or locality (Method column) for each element at different SNP-to-gene  
1561 mapping parameters (see WindowSize and FlankLimit columns). At an FDR  
1562 cutoff of 30%, genes were defined as HPO if they were observed at two or  
1563 more SNP-to-gene mapping parameters.

1564 Supp. Table 7

1565 HPO genes discovered with networks built from accessions subsets

1566 The number of HPO genes discovered in full ZmPAN (503 accessions) and  
1567 ZmRoot (46 accessions) networks was compared to networks built with a  
1568 subset of accessions. Both ZmPAN and ZmRoot networks were re-built using  
1569 a common set of 20 accessions. The ZmPAN network was re-built using 46  
1570 accessions consisting of the 20 common accessions and either 26 random or  
1571 26 CML biased accessions to simulate the number used in the full 46  
1572 accession ZmRoot network. Each network was analyzed for HPO genes in the  
1573 17 GWAS elements.

1574 Supp. Table 8

1575 Multiple element HPO gene list

1576 The number of commonly discovered HPO genes, hypergeometric *p*-values of  
1577 set overlap, and GRMZM IDs across multiple elements.

1578    Supp. Table 9

1579    Element gene ontology enrichment

1580            HPO genes for each element were tested for enrichment among genes co-

1581            annotated for gene ontology (GO) terms (hypergeometric test). Bonferroni

1582            correction is included as a column, treating each GO term as an independent

1583            test.

1584    Supp. Table 10

1585    HPO plus neighbors gene ontology enrichment

1586            Elemental HPO gene sets were supplemented with an additional set of highly

1587            connected neighbors equal to the number of genes in the HPO set. These

1588            HPO+ gene sets were tested for enrichment among genes annotated for GO

1589            terms (hypergeometric test).

1590