# Network- and structure-informed prioritization of candidate genes associated with etofenprox phytotoxicity in soybean

**Anonymous Author(s)**
Affiliation
Address
email

## Abstract

Soybean cultivars can exhibit phytotoxic injury after exposure to the pyrethroid insecticide etofenprox, yet the molecular basis of cultivar-specific sensitivity remains unclear. We generated a controlled time-series RNA-seq dataset for two cultivars (15 libraries; control at 0/12/24 h; treatment at 12/24 h), yielding 469,847,012 raw reads, 96.1% retention after trimming, and a mean total mapping rate of 93.1% (88.7% uniquely mapped). Differential expression was strongest at 12 h (T_12h vs C_12h: 5,539 DEGs), and pathway-level inspection highlighted coordinated regulation in stress-associated processes. WGCNA identified a trait-associated co-expression module (salmon; MM–GS cor = 0.84), supporting hub-gene prioritization and integration with cultivar-specific variation for downstream validation.

## 1 Introduction

Soybean (Glycine max [L.] Merr.) is a globally important crop that supplies protein and oil for food, feed, and industrial uses. In modern production systems, chemical pest control is indispensable, and pyrethroid insecticides such as etofenprox are widely applied because of their efficacy and comparatively low mammalian toxicity. Nonetheless, field observations and controlled assays have reported phytotoxic symptoms in soybean after etofenprox exposure, including chlorosis, growth suppression, and delayed recovery, which can translate into yield penalties under intensive management [Kim et al., 2021]. Emerging evidence also suggests that metabolic activation of etofenprox may contribute to downstream injury processes in plants [Xu et al., 2025]. Understanding the molecular basis of etofenprox phytotoxicity is therefore essential for improving crop resilience and optimizing pesticide use.

Plant responses to xenobiotics are dynamic and involve coordinated regulation of detoxification, oxidative stress mitigation, hormone signaling, and metabolic reprogramming [Siminszky, 2006]. These processes evolve over hours to days, and single time-point transcriptomic snapshots often fail to capture the causal sequence of regulatory events. Time-series RNA sequencing provides a trajectory view of transcriptional changes, while weighted gene co-expression network analysis (WGCNA) can organize these changes into modules linked to physiological traits, enabling the discovery of regulatory hubs that govern tolerance or susceptibility [Langfelder and Horvath, 2008].

Genomic variation further shapes chemical sensitivity by altering enzyme activity, transport capacity, or signaling components. Whole-genome sequencing (WGS) and variant filtering can identify high-confidence polymorphisms between cultivars with contrasting responses, but functional interpretation remains a bottleneck. Recent advances in AI-based protein structure prediction provide a principled way to map candidate variants onto structural contexts, offering mechanistic hypotheses about how

sequence changes may influence protein stability or active sites and thereby modulate phytotoxic outcomes [Jumper et al., 2021].

Here, we investigate soybean responses to etofenprox using an integrated systems framework. We expose two cultivars with contrasting phytotoxicity phenotypes to a controlled etofenprox treatment, collect a time series of leaf tissues, and quantify physiological injury indices alongside transcriptomic profiles. We then construct WGCNA modules associated with injury traits, identify hub genes, and link these to cultivar-specific WGS variants. Finally, we apply AI-based protein structure analysis to prioritize variants with plausible functional impacts. This combined approach aims to clarify the regulatory networks and candidate genes that underlie soybean tolerance to etofenprox phytotoxicity.

## 2 Materials and Methods

### 2.1 Plant materials and experimental design

Two soybean cultivars with contrasting etofenprox sensitivity were grown in pots (9x7x7.5 cm) under controlled conditions: 26/20°C (day/night), 14/10 h (light/dark), and 60% humidity. At the V1 stage, plants were treated with etofenprox 20% EC (Sebero, KyungNong) diluted 1:1000 (v/v), applied at 10 mL per plant. Control plants received the carrier solution. Samples were collected to form 15 total RNA-seq libraries (3x5): Control at 0, 12, and 24 h, and Treatment at 12 and 24 h, with three biological replicates per condition.

### 2.2 RNA sequencing, preprocessing, alignment, quantification, and differential expression

Total RNA was extracted from leaf tissue and quality-checked prior to sequencing. Raw reads were inspected using FastQC [Andrews, 2010], and adapter/low-quality sequences were removed using fastp [Chen et al., 2018]; reads shorter than 30 bp after trimming were discarded. Filtered reads were aligned to the soybean reference genome Wm82.a4.v1 (Phytozome) citepschmutz2010soybean using STAR [Dobin et al., 2013], and gene-level read counts were summarized with featureCounts [Liao et al., 2014] to generate a count matrix for downstream analyses. Differential expression analysis was conducted in DESeq2 [Love et al., 2014] using the gene-level count matrix, and genes were considered significant when |Log2FC| > 2 and p < 0.01. Where used, MultiQC was applied to summarize QC reports across samples [Ewels et al., 2016]. Software versions are available upon request.

### 2.3 WGCNA and hub gene prioritization

A variance-stabilized expression matrix was used to construct co-expression networks with the WGCNA R package [Langfelder and Horvath, 2008]. An appropriate soft-thresholding power was selected to approximate scale-free topology, followed by adjacency and topological overlap calculations. Modules were identified using dynamic tree cutting and merged based on eigengene similarity. Module-trait relationships were estimated by correlating module eigengenes with physiological injury metrics, including chlorosis scores, Fv/Fm, and MDA content. Hub genes were prioritized using high module membership and gene significance within trait-associated modules; network visualization was performed in Cytoscape when applicable.

### 2.4 Variant discovery and filtering

Genomic DNA from each cultivar was sequenced to high coverage and aligned to the reference genome using BWA-MEM [Li, 2013]. Variants were called with a haplotype-based caller following GATK best-practice recommendations [McKenna et al., 2010, Van der Auwera et al., 2013] and filtered using stringent hard filters on depth, quality by depth, strand bias, and mapping quality to obtain a high-confidence set of SNPs and indels. Functional annotation was performed to classify variants by genomic context and predicted effect; where applicable, common tools such as bcftools [Danecek et al., 2021] and annotation utilities (e.g., snpEff [Cingolani et al., 2012] or Ensembl VEP [McLaren et al., 2016]) were used.

## 2.5 Protein structure prediction

For candidate genes supported by both WGCNA and variant analyses, AI-based protein structure prediction tools were used where applicable to obtain 3D conformations and confidence scores (e.g., ESMFold or AlphaFold2) [Lin et al., 2023, Jumper et al., 2021]. Variants were mapped onto predicted structures to assess proximity to catalytic residues, ligand-binding pockets, or conserved motifs. Structural comparison and visualization were conducted using standard molecular graphics tools when needed to support qualitative interpretation.

# 3 Results

## 3.1 Read preprocessing and quality overview

Across the 15 RNA-seq libraries, the raw data contained a total of 469,847,012 reads, with a consistent mean read length of 151 bp per sample. After adapter/quality trimming, 451,394,808 reads were retained, corresponding to an overall retention of 96.1%. Trimming slightly reduced the average read length to 149.6 bp (mean across samples), indicating that only short low-quality or adapter-contaminated segments were removed while preserving the bulk of informative sequence. At the sample level, the proportion of retained reads ranged from 83.5% (C_12h-4) to 97.9% (C_24h-8), with most libraries clustering near the upper end of this range, supporting consistent preprocessing performance across the time-course dataset. Taken together, the high read retention and stable post-trimming read length distribution suggest that downstream alignment and quantification steps were performed on libraries of broadly comparable quality and complexity(Table 1).

Table 1: Trimming summary across 15 RNA-seq libraries

| Metric | Value |
|---|---|
| Total raw reads | 469,847,012 |
| Total retained reads | 451,394,808 |
| Overall retention | 96.1% |
| Mean read length (raw) | 151 bp |
| Mean read length (trimmed) | 149.6 bp |
| Retention range | 83.5% to 97.9% |

## 3.2 Alignment performance and mapping statistics

Trimmed reads were aligned to the reference genome, and alignment performance was summarized as uniquely mapped, multi-mapped, and unmapped fractions. Overall, the dataset showed robust alignment: the mean total mapping rate was 93.1%, with 88.7% of reads mapping uniquely and 4.4% mapping to multiple loci on average. The total mapping rate across samples ranged from 80.4% (12h_T-1_star) to 97.1% (0h-8_star), while the unmapped fraction ranged from 2.9% to 19.6%. Most libraries exhibited tightly grouped mapping profiles (typically >92% total mapped), indicating stable alignment behavior across experimental conditions. One library (12h_T-1_star) showed a comparatively lower mapping rate driven by an elevated unmapped fraction; however, the remaining libraries consistently achieved high unique mapping proportions, supporting reliable gene-level quantification for downstream comparative analyses. Collectively, these mapping statistics indicate that the majority of sequencing reads were successfully assigned to the reference, providing a strong basis for subsequent expression estimation and differential expression testing(Table 2).

Table 2: Mapping summary across 15 RNA-seq libraries

| Metric | Value |
|---|---|
| Mean total mapping | 93.1% |
| Mean uniquely mapped | 88.7% |
| Mean multi-mapped | 4.4% |
| Mean unmapped | 6.9% |
| Total mapping range | 80.4% to 97.1% |
| Unmapped range | 2.9% to 19.6% |

## 3.3 Differential expression overview across contrasts

Differential expression analysis was conducted across eight pairwise contrasts representing time-course changes within control or treatment conditions and treatment-control comparisons at matched time points. Using the predefined significance criteria ($|\log2FC| > 2$ and adjusted $p < 0.01$), the number of differentially expressed genes (DEGs) varied substantially by contrast. Among control time comparisons, C_24h vs C_12h yielded the largest DEG set (5,310 up-regulated; 3,909 down-regulated; total 9,219), followed by C_24h vs C_0h (2,360 up; 3,395 down; total 5,755), whereas C_12h vs C_0h showed a smaller but still notable shift (1,497 up; 2,090 down; total 3,587). For treatment-related comparisons, T_12h vs C_12h showed a pronounced transcriptional difference (3,242 up; 2,297 down; total 5,539), while T_12h vs C_0h and T_24h vs C_0h exhibited comparable DEG magnitudes (1,944 up / 2,892 down; total 4,836 and 1,772 up / 2,732 down; total 4,504, respectively). The within-treatment contrast T_24h vs T_12h produced a moderate DEG set (513 up; 924 down; total 1,437). In contrast, T_24h vs C_24h produced a markedly smaller DEG set (178 up; 6 down; total 184), indicating minimal differential expression at this matched time point under the applied thresholds. Across most contrasts, down-regulated DEGs were numerically dominant, while a subset of contrasts showed a stronger up-regulated component. Collectively, these results confirm that DEG magnitude and directionality are contrast-dependent across the dataset and provide a quantitative basis for subsequent enrichment and network analyses(Fig 1).
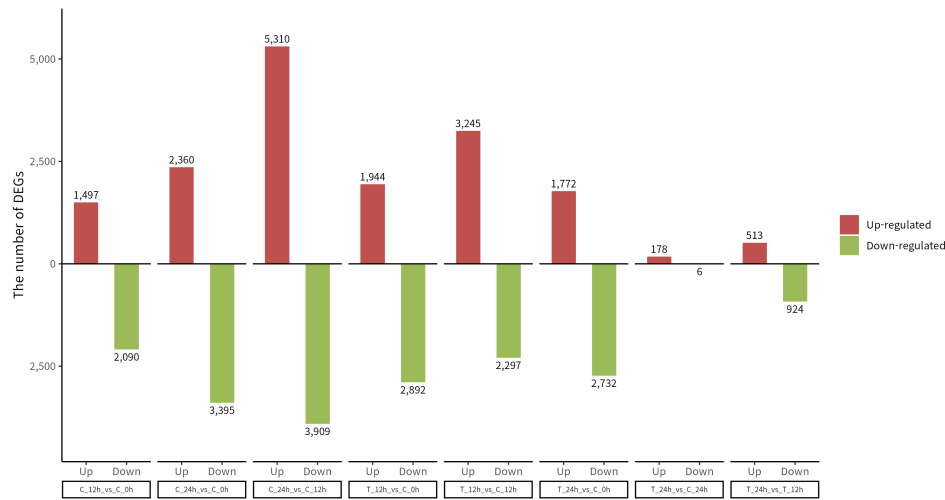


Figure 1: Numbers of up- and down-regulated DEGs across eight contrasts under the applied thresholds ($|\log2FC| > 2$, adjusted $p < 0.01$).

## 3.4 GO and KEGG enrichment summary for T_12h vs C_12h DEGs

To summarize functional signals associated with the treatment effect at 12 h, Gene Ontology (GO) and KEGG pathway enrichment analyses were performed separately for up- and down-regulated DEGs in the T_12h vs C_12h contrast. In the GO analysis, enriched biological process terms among up-regulated DEGs included vesicle- and secretion-related categories (e.g., exocytic process and vesicle docking involved in exocytosis) together with broader metabolic and repair-associated terms. Enriched molecular function terms for up-regulated DEGs were dominated by nucleotide-binding categories (e.g., ATP binding and related purine/ribonucleotide binding terms). For down-regulated DEGs, enriched GO biological process terms included biosynthetic process and macromolecule biosynthetic process, along with hormone-related response categories and electron transport chain. Enriched cellular component terms prominently included ribosome and ribosomal subunit, and molecular function terms included ribosome-associated functions and multiple transporter/oxidoreductase-related annotations. Consistent with the GO patterns, KEGG enrichment for up-regulated DEGs highlighted Endocytosis (93 genes), Spliceosome (81 genes), mRNA surveillance pathway (60 genes), and Circadian rhythm - plant (73 genes), alongside additional genetic information processing pathways. For down-regulated DEGs, the largest KEGG category was Ribosome (174 genes),

4

accompanied by Protein processing in endoplasmic reticulum (49 genes), Oxidative phosphorylation (35 genes), and additional metabolism-related pathways. Overall, these enrichment outputs provide a structured summary of GO and KEGG categories observed for up- and down-regulated DEGs in the T_12h vs C_12h comparison(Fig 2, Table 3).
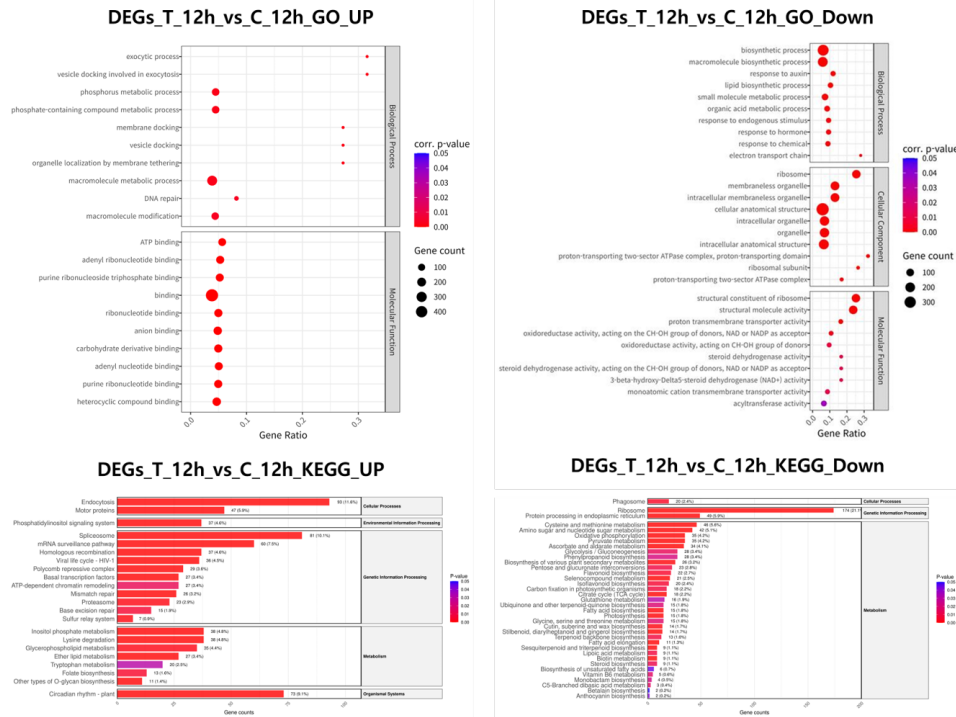


Figure 2: GO and KEGG enrichment results for up- and down-regulated DEGs in the T_12h vs C_12h contrast. Dot plots indicate GO enrichment; bar charts indicate KEGG pathway enrichment with gene counts.

Table 3: Selected KEGG pathways enriched in T_12h vs C_12h up- and down-regulated DEGs (gene counts summarized alongside the KEGG bar chart).

| Category / Pathway | Gene count |
| --- | --- |
| 'Up': 'Endocytosis' | 93 |
| 'Up': 'Spliceosome' | 81 |
| 'Up': 'mRNA surveillance pathway' | 60 |
| 'Up': 'Circadian rhythm - plant' | 73 |
| 'Up': 'Motor proteins' | 47 |
| 'Down': 'Ribosome' | 174 |
| 'Down': 'Protein processing in endoplasmic reticulum' | 49 |
| 'Down': 'Oxidative phosphorylation' | 35 |
| 'Down': 'Glycolysis/Gluconeogenesis' | 28 |
| 'Down': 'Phenylpropanoid biosynthesis' | 28 |

## 3.5 Co-expression network analysis identifies key modules associated with response traits

We next used weighted gene co-expression network analysis (WGCNA) to summarize time-course transcriptional dynamics into co-expressed gene modules and to connect these modules to phenotypic response traits. Prior to network construction, sample-level clustering was inspected to verify that global expression profiles were coherent with the experimental design and that no outlier libraries dominated downstream module detection (see the sample clustering figure below).

Module-trait correlation analysis highlighted a small number of modules with strong associations to the response phenotype. In particular, the salmon module showed the highest positive relationship with the sensitivity-related trait and a concurrent association with the time variable, indicating that this module captures a coordinated transcriptional program that tracks injury severity over the course of exposure (see the module-trait relationship heatmap below). To evaluate whether trait association within this module reflected coherent intramodular organization, we compared module membership (MM) to gene significance (GS). Genes in the salmon module exhibited a strong positive MM-GS relationship (cor = 0.84; p = 7.5e-185), consistent with a module in which highly connected genes also show the strongest trait relevance (see the MM-GS scatter plot below). These results support prioritizing salmon-module hub genes as candidates for downstream pathway inspection and integrative variant/structure-informed interpretation.
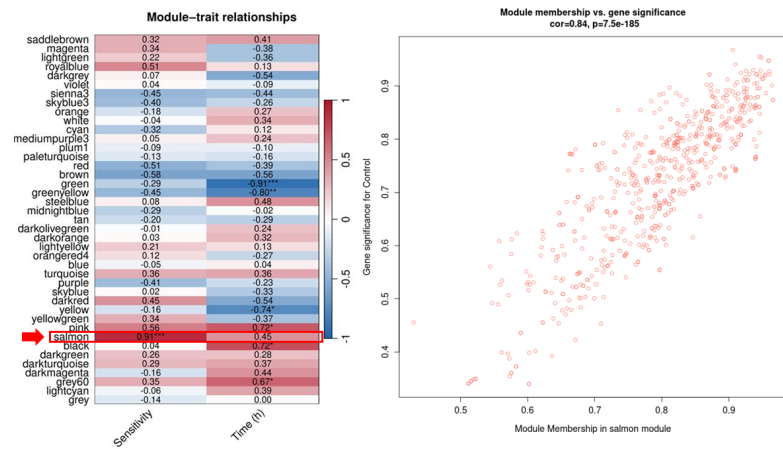


Figure 3: Module-trait relationships and intramodular evidence for the salmon module. Left: correlations between module eigengenes and traits. Right: relationship between module membership and gene significance within the salmon module (cor = 0.84; p = 7.5e-185).

## 3.6  Pathway-level inspection supports coordinated regulation of representative pathways

To complement term-level enrichment and module association results, we inspected KEGG pathway maps for representative signaling programs that can contextualize coordinated transcriptional regulation at the network level. We focused on two pathways commonly implicated in rapid stress signaling in plants, MAPK signaling (gmx04016) and Plant-pathogen interaction (gmx04626). Both pathways provide a compact view of upstream signal perception, kinase cascades, and downstream transcriptional outputs that can be compared directly to the directionality observed in the RNA-seq contrast.

In the MAPK signaling pathway map (gmx04016), WRKY33 was highlighted as an induced node, positioned downstream of MAPK cascade branches that connect to canonical defense outputs, including camalexin biosynthesis (via PAD3) and late defense gene induction (via PR1). In the Plant-pathogen interaction map (gmx04626), WRKY22 and PBS1 were highlighted, spanning signaling axes that link pattern-triggered immunity components to effector-triggered responses and hypersensitive response (HR)-associated outputs. Together, these maps provide a pathway-level, qualitative view that multiple upstream signaling routes converge on transcriptional regulation and defense-associated response programs under the T_12h vs C_12h comparison. This pathway inspection is consistent with the enrichment profile reported above, where up-regulated genes were enriched for pathways such as Endocytosis and Spliceosome, while down-regulated genes were dominated by categories such as Ribosome and Oxidative phosphorylation(Fig 5).
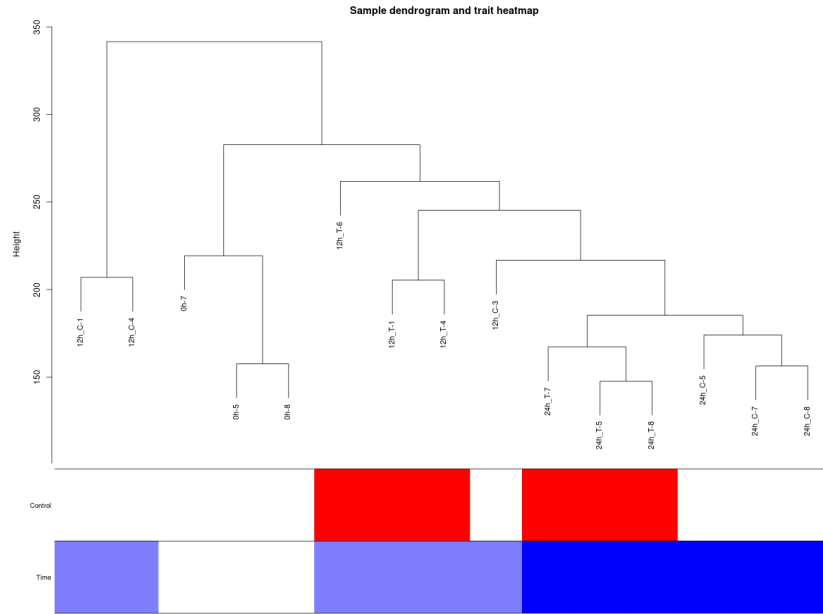
6

Figure 4: Sample clustering and trait heatmap used for WGCNA quality inspection. Hierarchical clustering of libraries and corresponding trait annotations indicate consistency with the experimental design.
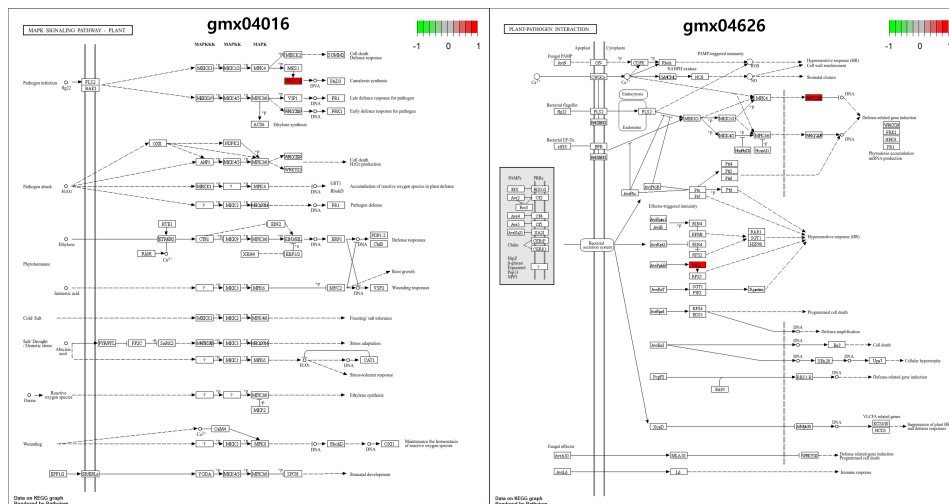


Figure 5: Pathway-level inspection of representative KEGG maps. Left: MAPK signaling pathway (gmx04016). Right: Plant-pathogen interaction pathway (gmx04626). Colored nodes indicate transcript-level directionality for the T_12h vs C_12h comparison (red, higher; green, lower; scale shown).

## 3.7 Candidate gene prioritization integrating network evidence and sequence/structure context

Finally, we integrated network-level evidence with sequence-level variant filtering and structure-based context to prioritize candidates for follow-up. Candidate selection was performed in a layered manner: (i) trait-associated co-expression modules (with an emphasis on the salmon module) defined a network-relevant search space; (ii) cultivar- contrasting variants were filtered to retain high-confidence polymorphisms within expressed transcripts; and (iii) a small subset of protein-altering variants was mapped into predicted 3D structures to provide qualitative context on whether sequence changes might localize to regions plausibly relevant to stability or ligand interaction. Importantly, this step is used to contextualize candidates rather than to assert causality.

As an example of the structure-informed layer, we generated 3D models for Glyma.01G117900.27 and compared predicted conformations between CMJ 047 and CMJ 213. The overall fold was visualized to verify that both sequences yielded interpretable globular conformations and to identify regions showing notable geometric differences. We further visualized a representative docking pose (ligand shown as sticks) to document the presence and location of a plausible pocket-like region in the predicted structures. These visual checks provide a structural sanity layer for prioritization by indicating whether candidate sequence differences coincide with or lie near pocket-adjacent regions, while reserving mechanistic claims for downstream validation.
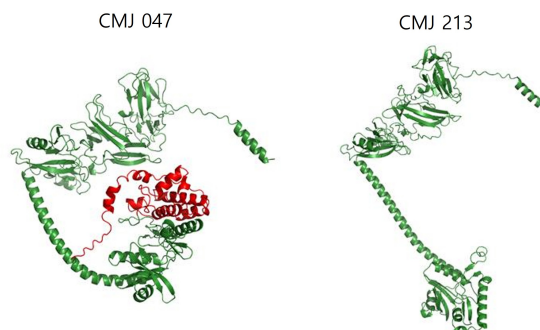


Figure 6: Predicted protein structure comparison for Glyma.01G117900.27 between CMJ 047 and CMJ 213. Ribbon diagrams are shown to document overall fold and major domain organization.
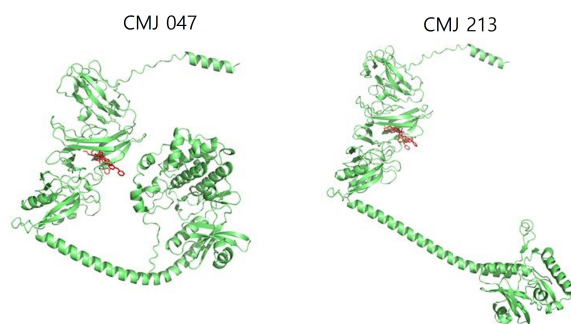


Figure 7: Representative docking visualization on the predicted structures of Glyma.01G117900.27. The ligand pose (sticks) is shown to document pocket location in each cultivar-specific model.

## 4 Discussion

Plant responses to xenobiotics typically involve a staged detoxification process that includes oxidative transformation and conjugation (often mediated by cytochrome P450s and GSTs), followed by compartmentalization or transport of modified metabolites [Siminszky, 2006]. A soybean cultivar-specific

8

phytotoxic response to etofenprox has been reported in Korea, with symptoms such as leaf deformation and necrosis and genetic evidence consistent with a single major locus controlling sensitivity [Kim et al., 2021]. Evidence for metabolic activation of etofenprox suggests additional biochemical layers that may modulate injury outcomes [Xu et al., 2025]. This prior observation supports the biological plausibility of our study design, where we combine time-resolved transcriptome dynamics with genotype contrasts to narrow down candidate determinants of etofenprox sensitivity.

A central contribution of our analysis is the use of network-level evidence to reduce the candidate search space beyond what is achievable with differential expression alone. WGCNA provides a principled approach to summarize coordinated transcriptional programs into modules and to prioritize intramodular hub genes that are more likely to reflect core regulatory control points. In the context of chemical stress responses, this is particularly relevant because visible injury phenotypes can emerge after a delay (e.g., 12–24 h), when downstream defense, oxidative stress, and signaling cascades have propagated through the transcriptome. Consistent with this, our pathway-level inspection highlighted coherent regulation of representative stress-related pathways, including signaling modules commonly associated with defense activation. Rather than interpreting individual pathway nodes as definitive causal drivers, we treat these pathway patterns as a contextual layer that helps explain why certain co-expression modules and hub genes are plausible follow-up targets.

We further incorporated sequence/structure context as a conservative "sanity layer" for candidate prioritization. Recent advances in AI-based protein structure prediction enable rapid generation of plausible 3D conformations directly from sequence, providing a qualitative check on whether cultivar-specific variants are compatible with stable folding or may localize near pocket-adjacent regions relevant to ligand interaction [Jumper et al., 2021, Lin et al., 2023]. In this study, predicted models and representative docking visualizations were used strictly to document fold-level plausibility and pocket localization across cultivar-specific sequences, not to claim a finalized mechanism. This distinction is important for research integrity: structural inspection can down-select candidates and guide targeted biochemical validation, but does not replace experimental confirmation of enzymatic activity, metabolite profiles, or binding kinetics.

Several limitations should be acknowledged. First, the genomic contrast (one sensitive vs. eight insensitive individuals) is well suited to detect private/enriched variants but may miss polygenic effects or background-dependent modifiers; expanded sampling and independent validation populations will improve resolution. Second, transcriptome-based prioritization depends on tissue, developmental stage, and treatment conditions, so replication across environments and time-points is needed. Finally, structure/docking analyses are assumption-dependent and should be viewed as hypothesis-supporting annotations rather than proof. Despite these constraints, our multi-layer prioritization integrating network evidence, pathway coherence, and conservative structure-context checks provides a practical framework to generate a compact, testable candidate list for downstream functional assays and for studies on the mechanism of etofenprox-induced injury.

# References

Simon Andrews. Fastqc: a quality control tool for high throughput sequence data, 2010. URL https://www.bioinformatics.babraham.ac.uk/projects/fastqc/.

Shifu Chen, Yanqing Zhou, Yaru Chen, and Jia Gu. fastp: an ultra-fast all-in-one fastq preprocessor. *Bioinformatics*, 34(17):i884–i890, 2018. doi: 10.1093/bioinformatics/bty560.

Pablo Cingolani, Adrian Platts, Lili Wang, et al. A program for annotating and predicting the effects of single nucleotide polymorphisms, snpeff. *Fly*, 6(2):80–92, 2012. doi: 10.4161/fly.19695.

Petr Danecek, James K. Bonfield, Jennifer Liddle, et al. Twelve years of samtools and bcftools. *GigaScience*, 10(2):giab008, 2021. doi: 10.1093/gigascience/giab008.

Alexander Dobin, Carrie A. Davis, Felix Schlesinger, et al. Star: ultrafast universal rna-seq aligner. *Bioinformatics*, 29(1):15–21, 2013. doi: 10.1093/bioinformatics/bts635.

Philip Ewels, Måns Magnusson, Sverker Lundin, and Max Käller. Multiqc: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics*, 32(19):3047–3048, 2016. doi: 10.1093/bioinformatics/btw354.

John Jumper, Richard Evans, Alexander Pritzel, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596:583–589, 2021. doi: 10.1038/s41586-021-03819-2.

Kyung-Hye Kim, Jungmin Ha, Taeklim Lee, Jinho Heo, Jiyeong Jung, Juseok Lee, and Sungteag Kang. Identification of a novel trait associated with phytotoxicity of an insecticide etofenprox in soybean. *Journal of Pesticide Science*, 2021. doi: 10.1584/jpestics.D20-073.

Peter Langfelder and Steve Horvath. Wgcna: an r package for weighted correlation network analysis. *BMC Bioinformatics*, 9:559, 2008. doi: 10.1186/1471-2105-9-559.

Heng Li. Aligning sequence reads, clone sequences and assembly contigs with bwa-mem. *arXiv*, page arXiv:1303.3997, 2013.

Yang Liao, Gordon K. Smyth, and Wei Shi. featurecounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics*, 30(7):923–930, 2014. doi: 10.1093/bioinformatics/btt656.

Zeming Lin, Halil Akin, Roshan Rao, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023. doi: 10.1126/science.ade2574.

Michael I. Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome Biology*, 15:550, 2014. doi: 10.1186/s13059-014-0550-8.

Aaron McKenna, Matthew Hanna, Eric Banks, et al. The genome analysis toolkit: a mapreduce framework for analyzing next-generation dna sequencing data. *Genome Research*, 20(9):1297–1303, 2010. doi: 10.1101/gr.107524.110.

William McLaren, Laurent Gil, Sarah E. Hunt, et al. The ensembl variant effect predictor. *Genome Biology*, 17:122, 2016. doi: 10.1186/s13059-016-0974-4.

Balazs Siminszky. Plant cytochrome p450-mediated herbicide metabolism. *Phytochemistry Reviews*, 5:445–458, 2006. doi: 10.1007/s11101-006-9017-7.

Geraldine A. Van der Auwera, Mauricio O. Carneiro, Christopher Hartl, et al. From fastq data to high-confidence variant calls: the genome analysis toolkit best practices pipeline. *Current Protocols in Bioinformatics*, 43:11.10.1–11.10.33, 2013. doi: 10.1002/0471250953.bi1110s43.

Jun Xu et al. Metabolic activation of etofenprox and downstream oxidative responses in plants. *Environmental Science & Technology*, 2025.