

July 18, 2023

Dear Dr. Sabot,

Thank you very much for your careful evaluation of our manuscript. We are grateful for the reviews received.

We have carefully revised our manuscript and have addressed the questions and suggestions by the reviewers. We provide a detailed explanation of the changes performed in our point-by-point response below.

We believe that our manuscript has improved considerably thanks to the reviewers' suggestions. Importantly, we have modified Supplementary Figure 2 to show no correlation between squared coefficient of variation and average gene expression; we have modified Supplementary Figure 3 (now 3 and 4) to address the reviewers' suggestions; and we have modified Supplementary Figure 4 (now 5) by adding a bootstrapping analysis to show differences between the distribution means. We hope that you find the new version of our manuscript suitable for recommendation.

We have posted the final version of our manuscript on biorxiv: <https://doi.org/10.1101/2022.04.26.489498> version 4. We also provide a tracked-changes pdf file (without figures, to fit within the file size limit) with the modified sections of the manuscript shown in red.

We are very pleased with our first experience submitting our work to PCI. Thank you very much for pushing this initiative forward.

Best regards,
Diego A. Hartasánchez
Arezki Boudaoud
Fraçoise Monéger

Point-by-point response to reviewers comments on “Expression of cell-wall related genes is highly variable and correlates with sepal morphology”

Reviewers comments in blue, our response in black, new or modified text in red

Reviewer #1: Anonymous

Minor comments:

- Why the clustering co-expression analysis was not performed with uniquely the HVGs? Did the authors make this analysis?

The clustering co-expression analysis was performed for the entire gene expression dataset after filtering, that is, for 14085 genes. By running the WGCNA package on this data, we obtained 16 modules of co-expressed genes. These modules were tested for correlation with morphological data. Gene ontology enrichment analysis was then performed on these modules and cell-wall related genes were found to be overrepresented in modules highly correlated with morphology. In parallel, cell-wall related genes happen to be highly variable. Given these results, at some point in our analysis, we were concerned that these two results could be confounded. Namely, we wondered if modules correlated with morphology had popped up because they were enriched in highly variable genes. As can be observed in Supplementary Figure 6, this is not the case. On the one hand, our top five correlated modules are not enriched in highly variable genes in general (it is true for the khaki and magenta modules, but not for the other three). Additionally, within modules, highly variable genes are not driving the correlation with morphology (their module membership and trait significance are not particularly high) nor are they driving the GO enrichment results. So, to answer the reviewer’s question, we did not perform the co-expression clustering analysis on our group of highly variable genes because, in our view, we would be filtering our dataset *a posteriori* based on criteria which resulted from the analysis itself. We believe that running WGCNA on HVGs would produce results that would be inevitably confounded. Additionally, WGCNA loses power to detect correlations for smaller datasets and for smaller modules so we do not expect more significant results *a priori*. We hope that this explanation convinces the reviewer. We consider that adding the analysis proposed by the reviewer to the current manuscript would not add relevant insight and would potentially make the reasoning within our manuscript circular, which we want to avoid.

- Figure 2C : add a color information for the variable contribution. codeLegend Supp figure1. Add a description of the color code D, E and F

We are not sure what the reviewer is referring to here. We believe that the reviewer is referring to adding the color code for plants D, E and F in the legend for figure 2B (and not 2C). We have clarified what the colors refer to by modifying the legends for Figure 2C, Figure 4A, and Supp Fig1.

- Supp Fig1: add color significance in the legend (Red, Blue, green)

Please refer to the answer to the previous comment. We have clarified what the colors refer to by modifying the legend of Supp Fig1.

- Supp Fig3: indicate HVG and LVG on the figure

Following this request and the suggestion of another reviewer, we have now modified Supp Fig3. We hope that the new figure is clearer.

Reviewer #2: Sandra Cortijo

Below are some minor changes I suggest:

- The authors indicate they do not see a correlation between CV2 and average expression level since they removed the lowly expressed genes. However, this is very hard to see in the Supplementary Figure 2 because of the different colors. Moreover, the fact that lowly variable genes (LVG) have a higher expression than highly variable genes (HVG) as shown in Supplementary Figure 2 (top) contradicts this claim. To see if there is indeed a correlation or not, I recommend the authors to do one or several of the following suggestions: (i) add a regression line to the plot; (ii) aid visualization by showing the point density; (iii) If the authors do find a correlation, I suggest to use the method described in Brennecke et al, 2013 (<https://doi.org/10.1038/nmeth.2645>) to identify HVG.

We thank Dr. Cortijo for her suggestion. Indeed, following Dr. Cortijo's work in Cortijo et al. 2019, we had already performed a correction to account for the correlation between CV2 and average gene expression during an initial phase of the project when our dataset was slightly larger and included some lowly expressed genes. We then decided to impose stronger filtering criteria which made said correction irrelevant in practical terms (for example, it only modified one of the 718 genes in the top 5%). We, hence, removed the correction, and clarified so in the main text: "We then classified genes according to the CV² values of their expression. Given the higher variability (both technical and biological) found for lowly expressed genes, one can correct CV² to account for this association (Cortijo et al., 2019). However, the gene expression threshold that we have set has eliminated most lowly expressed genes from our dataset and essentially, there is no correlation between CV² values and average gene expression (Supplementary Figure 2)."

We agree that this is not clear from the data, in particular Supplementary Figure 2. We have now included a 2d density visualization as suggested by Dr. Cortijo and added the results of a Pearson correlation test ($r=-0.00189$, $p\text{-value}=0.823$). We have added the following to the figure legend: "Point density in the 2 dimensional space (CV² against μ) is shown in the bottom plot with the Pearson correlation coefficient (r) and the corresponding p-value shown highlighting no correlation between the two variables."

- The Supplementary Figure 3 is very difficult to understand. I would recommend the authors to explain a bit more what they did in the text and the legend. Also, an alternative to Venn diagrams, like the one proposed by the UpSetR package, might help understand the figure.

We agree with Dr. Cortijo that this figure was difficult to understand. We explored the possibility of using the UpSetR package but found that it was not ideal for our purposes since the comparisons that we are more interested in are difficult to highlight with the UpSetR approach.

We have, instead, first separated Supp Figure 3 in Supp Figure 3 for highly variable genes and Supp Figure 4 for lowly variable genes. Second, we have separated each of the Venn diagrams of the previous version into two complementary Venn diagrams. Third, we have written a more clarifying figure legend and added a more precise explanation on what we did in the main text.

The text now reads: “We verified that the identification of LVGs and HVGs was not exclusively due to differences in gene expression between different plants (Supplementary Figures 3 & 4) by extracting highly and lowly variable genes independently for each plant and comparing the lists obtained between plants and their overlap with the LVG and HVG lists obtained considering all 27 sepal samples.”

The figure legend for Supplementary Figure 3 now reads (equivalent for Supplementary Figure 4): “Supplementary Figure 3: Venn diagrams for the bottom 5% lowly variable genes (LVGs) and for the bottom 15% LVGs according to their gene expression CV² values. LVGs were extracted in four different ways: considering only sepals from plant D, only sepals from plant E, only sepals from plant F, and all 27 sepals (DEF). The left Venn diagrams show the overlap between genes found for each plant independently. The number of genes in each red, green and blue circle adds up to 704 genes (top) and up to 2,112 genes (bottom) corresponding to 5% and 15% of our complete 14,085 gene set, respectively. Numbers shown in parentheses are the expected number of common genes found in each intersection if genes had been chosen randomly for each plant. The right Venn diagrams show how the LVG set obtained by considering all 27 sepal samples compares to the LVG sets obtained for each plant. Genes within the gray ellipse also add up to 704 (top) and 2,112 (bottom). All genes in the intersection between the three plants are also present within the gray ellipse (89 out of 89, top; 564 out of 564, bottom), whereas only a few genes (82, top; 61, bottom) are found in the DEF set and not for any plant independently. Circle and ellipse sizes and intersections are drawn to aid visualization but their sizes are not proportional to the number of genes found within them.”

- In the Supplementary Figure 4, the authors do a Wilcoxon test to compare the CV² of two groups of genes of very different size (718 genes vs 14085 genes). In order to actually be able to compare the CV² of these two groups the number of genes should be similar. For this, I recommend to use a bootstrap strategy with 1000 random sets of 718 genes taken in the 14085 genes set and to compare the average and the 95% confidence interval obtained for these 1000 random sets with the CWRG list.

We thank Dr. Cortijo for the suggestion. We have modified Supplementary Figure 4 (now Supplementary Figure 5) to include the results from the bootstrap strategy. Gene expression CV² mean values for the cell-wall related gene list is clearly above the 95% confidence interval of the distribution of mean CV² values for 1000 random subsets of the whole 14085 gene set.

We have added this paragraph in the main text: “Gene expression CV² values of these 718 genes is higher compared to those of the entire gene dataset (Supplementary Figure 5). We confirmed this through a bootstrapping approach by calculating gene expression CV² means for 1000 random samples of 718 genes from the entire dataset to have gene sets of equal size. The gene expression CV² mean of our CWRG list falls clearly above the 95% confidence interval of the distribution corresponding to the entire dataset (Supplementary Figure 5).”

We have also modified the Supplementary Figure 5 legend accordingly: “Supplementary Figure 5: Left: violin plots of gene expression CV² values for our entire 14,085-gene set (purple)

compared to that of the 718 cell-wall related gene (CWRG) list found in the dataset (blue). Grey box plots show median, 1st and 3rd quartiles of the distributions. Right: histogram of gene expression CV² mean values from 1000 random subsets of 718 genes from the 14,085-gene set. Pink vertical dashed lines delimit the 95% confidence interval of the distribution. Red asterisk shows the gene expression CV² mean value of the CWRG list, clearly above the 95% confidence interval of the random subset gene expression CV² distribution.”

- In the Figures 5 and 6, the authors use a trait significance. However, they never explain in the text how it was measured and what it means.

We thank Dr. Cortijo for pointing this out. We had thought of using “trait significance” instead of “gene significance for a trait” but have now realized that it is not clear and might serve as a confusion. We have reverted back to “gene significance” across the text and figures. We have additionally defined both module membership and gene significance in the Materials & Methods WGCNA section with the following text: “Each gene within a module has a module membership value, calculated as the correlation between that gene’s expression and the expression of the module eigengene across all samples. It is representative of the gene’s intramodular connectivity. Each gene also has a gene significance for any trait, calculated as the correlation between the gene’s expression values and the morphological parameter values across all samples.”

- About the discussion on cell-wall related genes, I would like to add that this category of GO is also enriched among HVG detected between plants in Cortijo et al. 2019. The authors’ observations might thus be more general and not specific to sepals, making it a phenomenon of wide interest.

We thank Dr. Cortijo for pointing this out. We have modified the last paragraph of the discussion and added Dr. Cortijo’s observation: “We extend this interpretation to cell-wall related genes, which also show significant variability of expression in sepals. Considering that reproducibility in morphology should be achieved in part thanks to underlying regulatory and compensatory mechanisms, our results indicate that cell-wall related genes could be fundamental for these mechanisms. The possible role of cell-wall related genes in determining sepal organ morphology could be applicable to other organs since cell-wall related genes exhibit high variability across plants also in Arabidopsis seedlings (Cortijo et al., 2019). In general, our work sheds light on the links between expression variability, gene regulatory networks, and developmental robustness and thereby, opens the way to informed functional analysis.”