

# Widespread inter-individual gene expression variability in *Arabidopsis thaliana*

Sandra Cortijo<sup>1</sup>, Zeynep Aydin<sup>1</sup>, Sebastian Ahnert<sup>1</sup>, James C.W. Locke<sup>1#</sup>

<sup>1</sup>The Sainsbury Laboratory, University of Cambridge, Cambridge, CB2 1LR, UK

Corresponding author: James Locke, The Sainsbury Laboratory, 47 Bateman St., Cambridge CB2 1LR, UK.

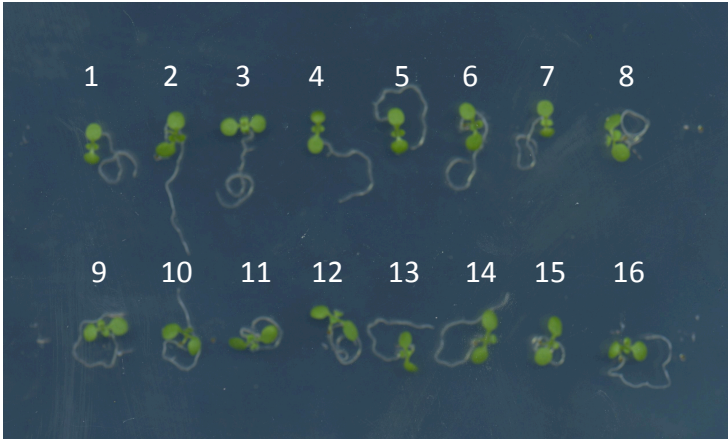
# Contact. Tel: +44(0)1223 761110. Email: [James.Locke@slcu.cam.ac.uk](mailto:James.Locke@slcu.cam.ac.uk)

## Appendix

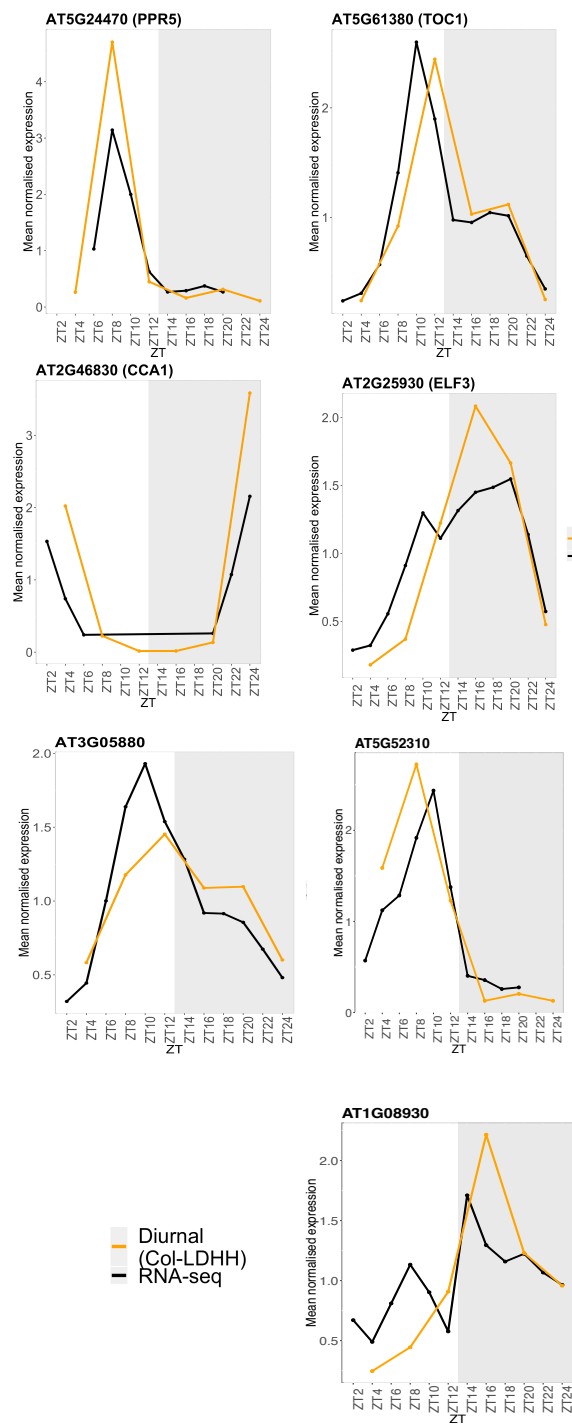
### Table of content

Appendix Figure S1 .....	p1
Appendix Figure S2 .....	p4
Appendix Figure S3 .....	p9
Appendix Figure S4 .....	p12
Appendix Figure S5 .....	p15
Appendix Figure S6 .....	p18
Appendix Figure S7 .....	p22

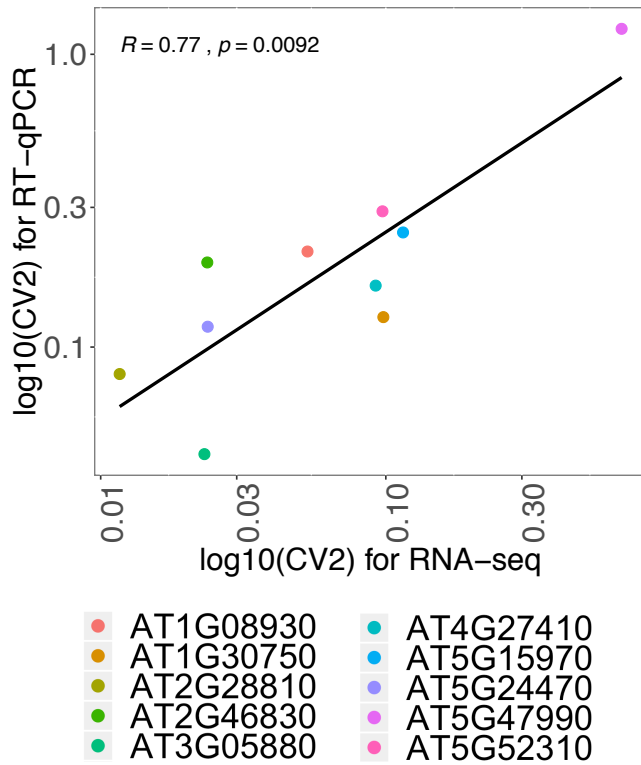
A



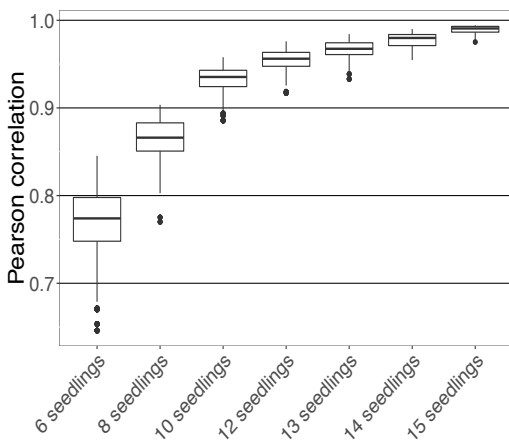
B



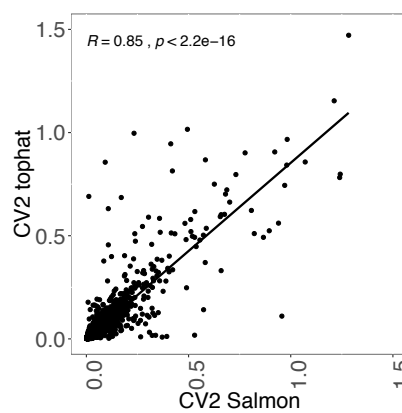
C



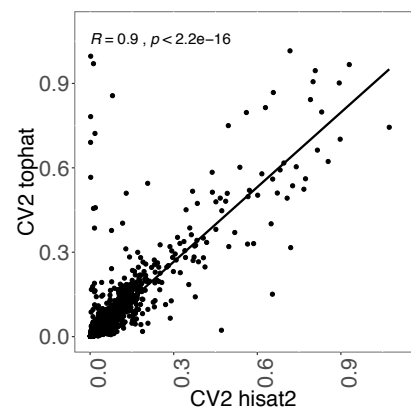
D



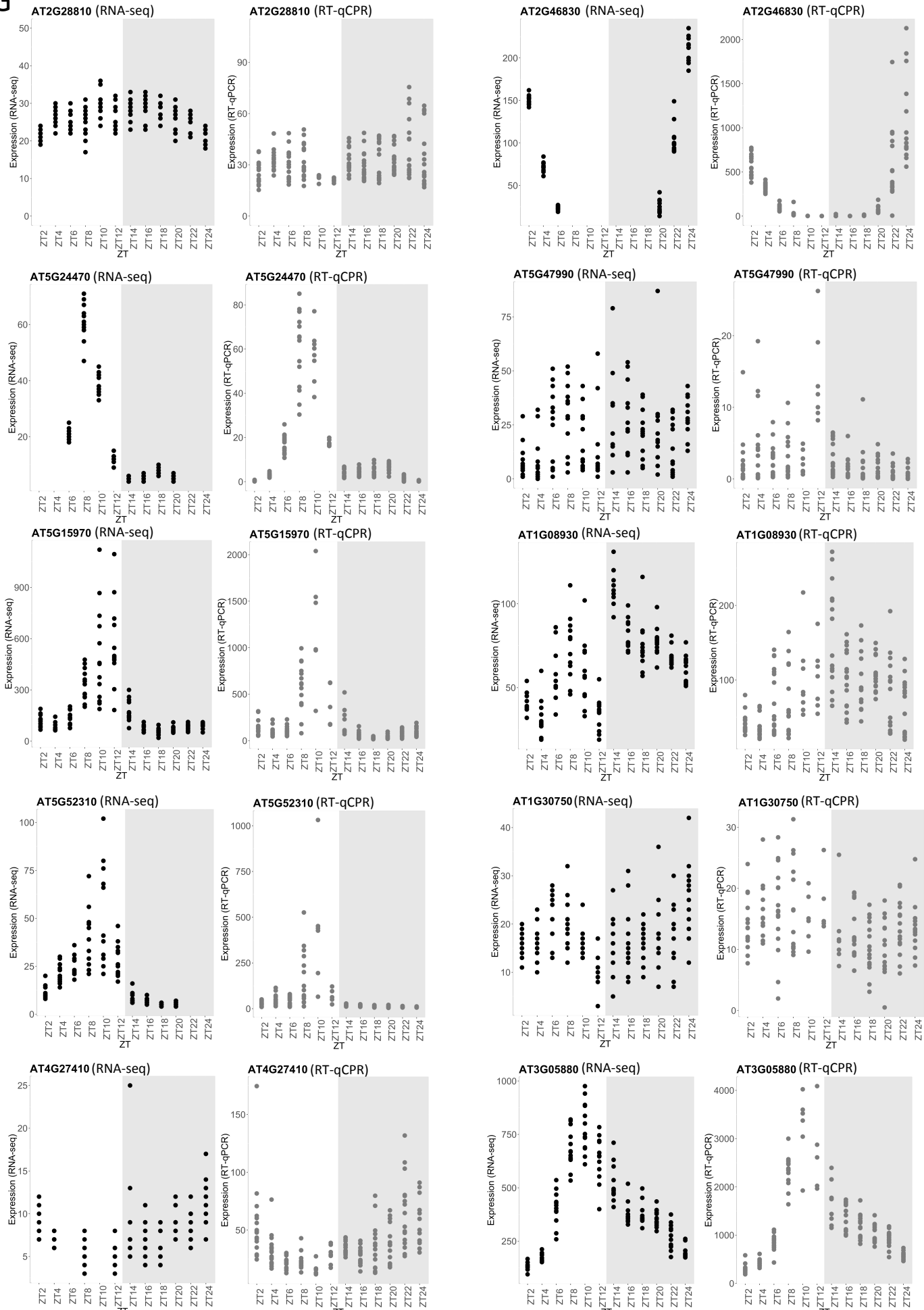
E



F



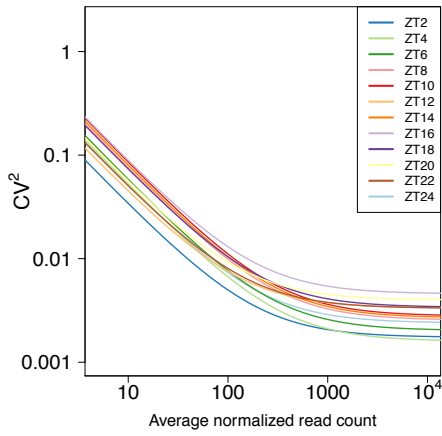
G



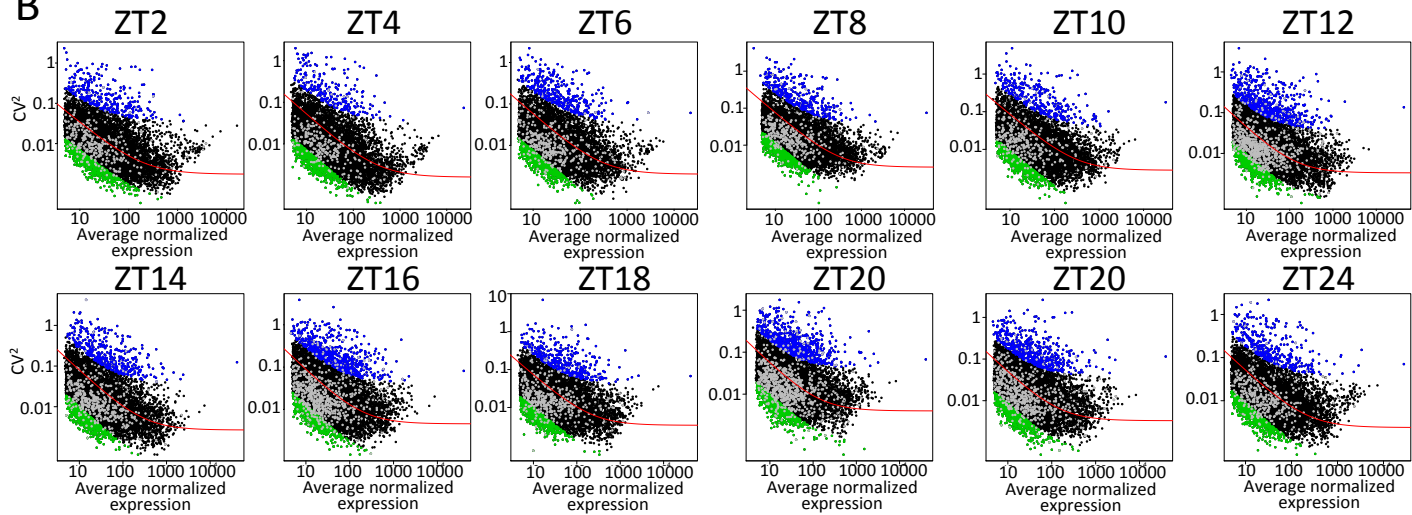
## Appendix Figure S1:

- A.** Photo of seedlings grown in exactly the same conditions used in the project, at the time they would be harvested. Seedling number is indicated in white.
- B.** Comparison of the average expression profiles (of the 14 seedlings) over the 24 hour time course (with 12 time-points) from our RNA-seq dataset and an already available diurnal expression dataset (<http://diurnal.mocklerlab.org>). Data from the condition Col-LDHH from the diurnal dataset was used, as its growing conditions are the closest to ours. Expression data for the first 24 hour were used; with a total of 6 time-points (every 4 hours). Each dot is the mean normalised expression level for the available diurnal expression data (orange) or for the average over all 14 seedlings from the RNAseq dataset (black dot).
- C.** Comparison of the  $CV^2$  measured for 10 genes over the 24 hour time course (with 12 time-points) from the RNA-seq dataset and a replicate experiment carried out by RT-qPCR. Each dot is the average  $CV^2$  over the time course for one of the 10 genes. The Spearman correlation between the  $CV^2$  measured in the RNA-seq and in the RT-qPCR experiments is 0.77, with a p-value of 0.0092.
- D.** Comparison of the corrected  $CV^2$  calculated from a subset of the data, from 6 to 15 seedlings, with the corrected  $CV^2$  measured from 16 seedlings at ZT6.
- E.** Comparison of the  $CV^2$  of all genes at ZT24 when mapping is performed with tophat or with salmon (Spearman correlation of 0.85, p-value < 2.2e-16).
- F.** Comparison of the  $CV^2$  of all genes at ZT24 when mapping is performed with tophat or with hisat2 (Spearman correlation of 0.9, p-value < 2.2e-16).
- G.** Comparison of the expression profiles in the 14 seedlings over a 24 hour time course (with 12 time-points) in the RNA-seq and in a replicate experiment done by RT-qPCR for 10 genes with different expression and variability profiles. Each dot is the mean normalised expression level for a single seedling from the RNAseq (left plot) or from the RT-qPCR (right plot).

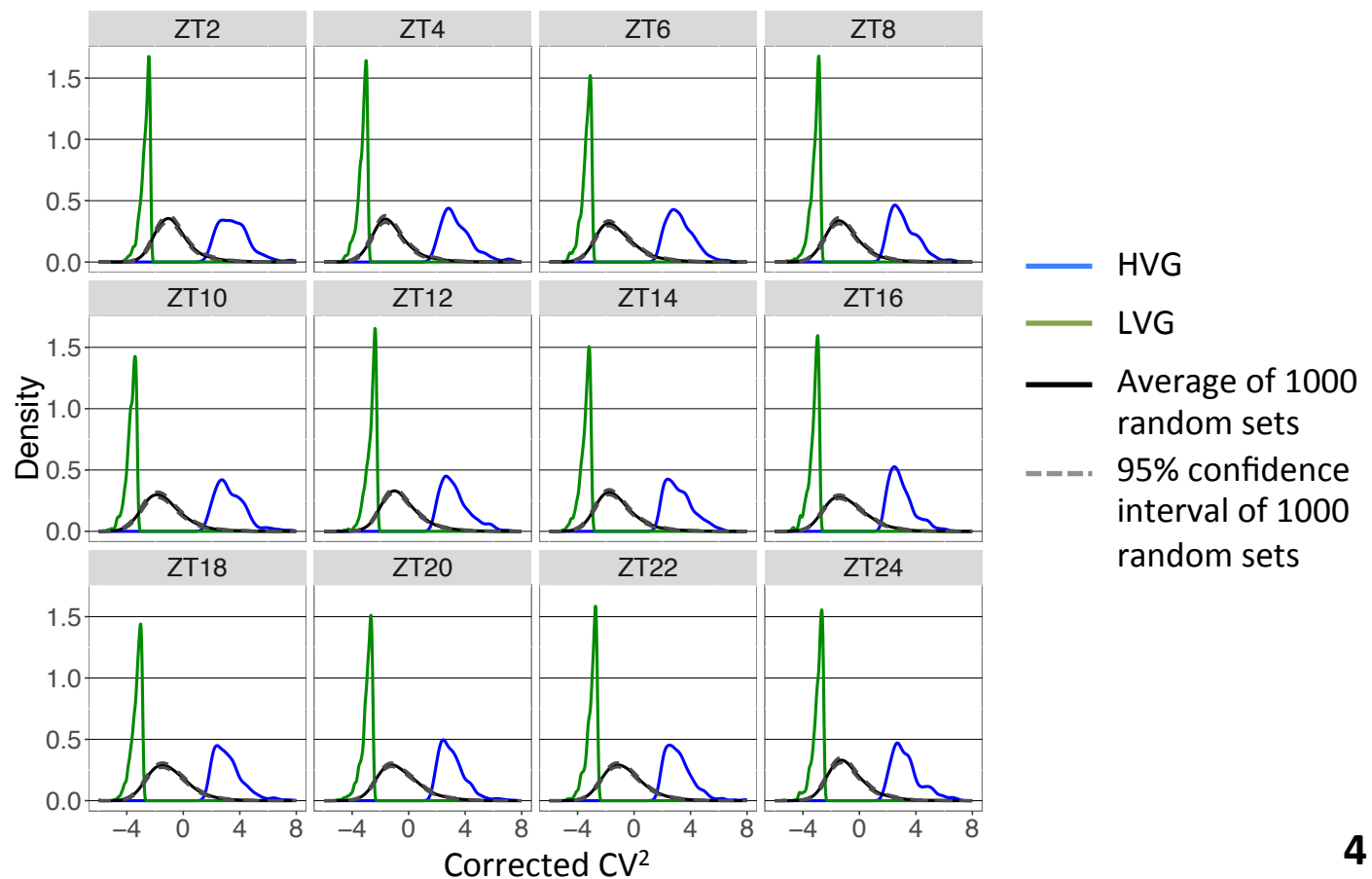
A

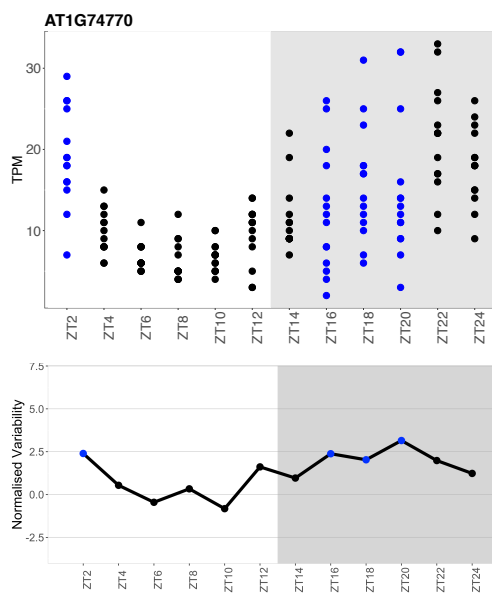
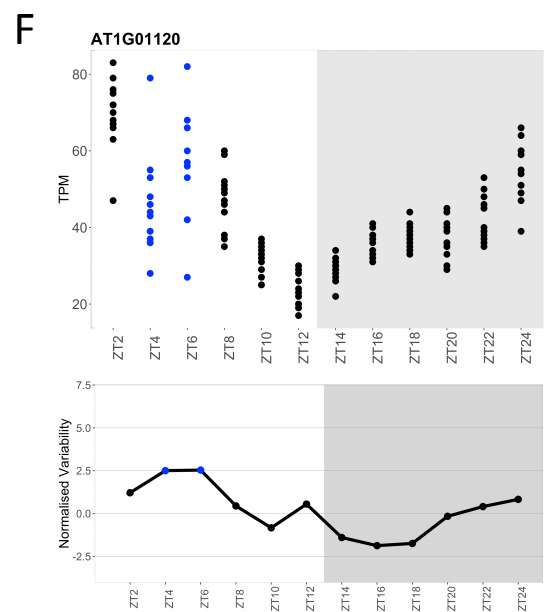
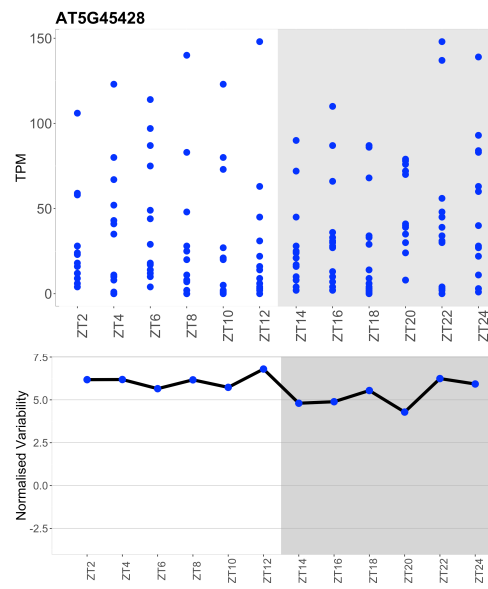
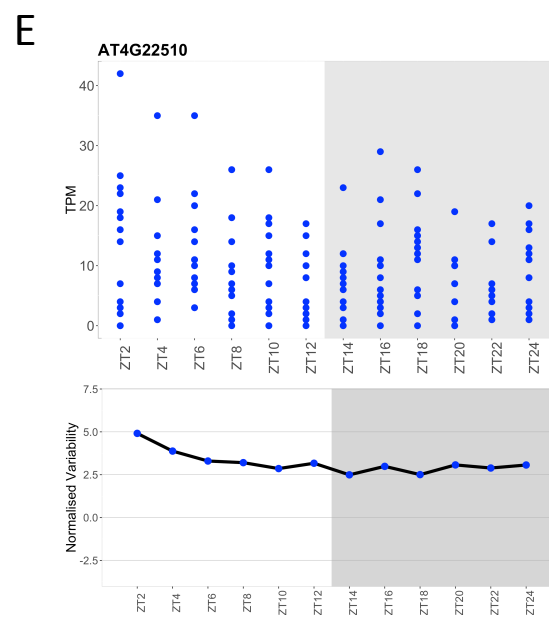
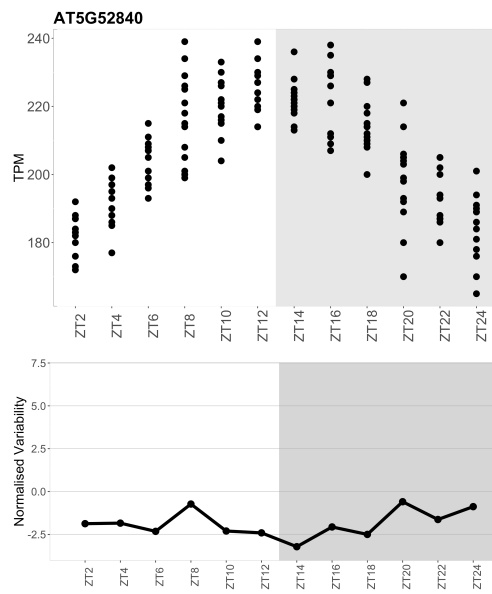
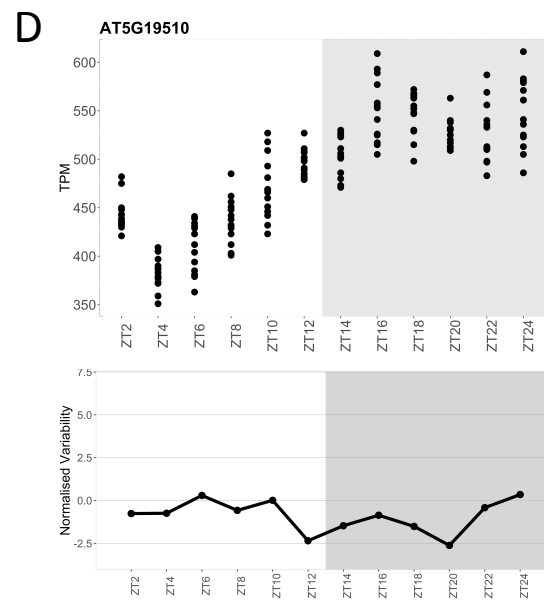


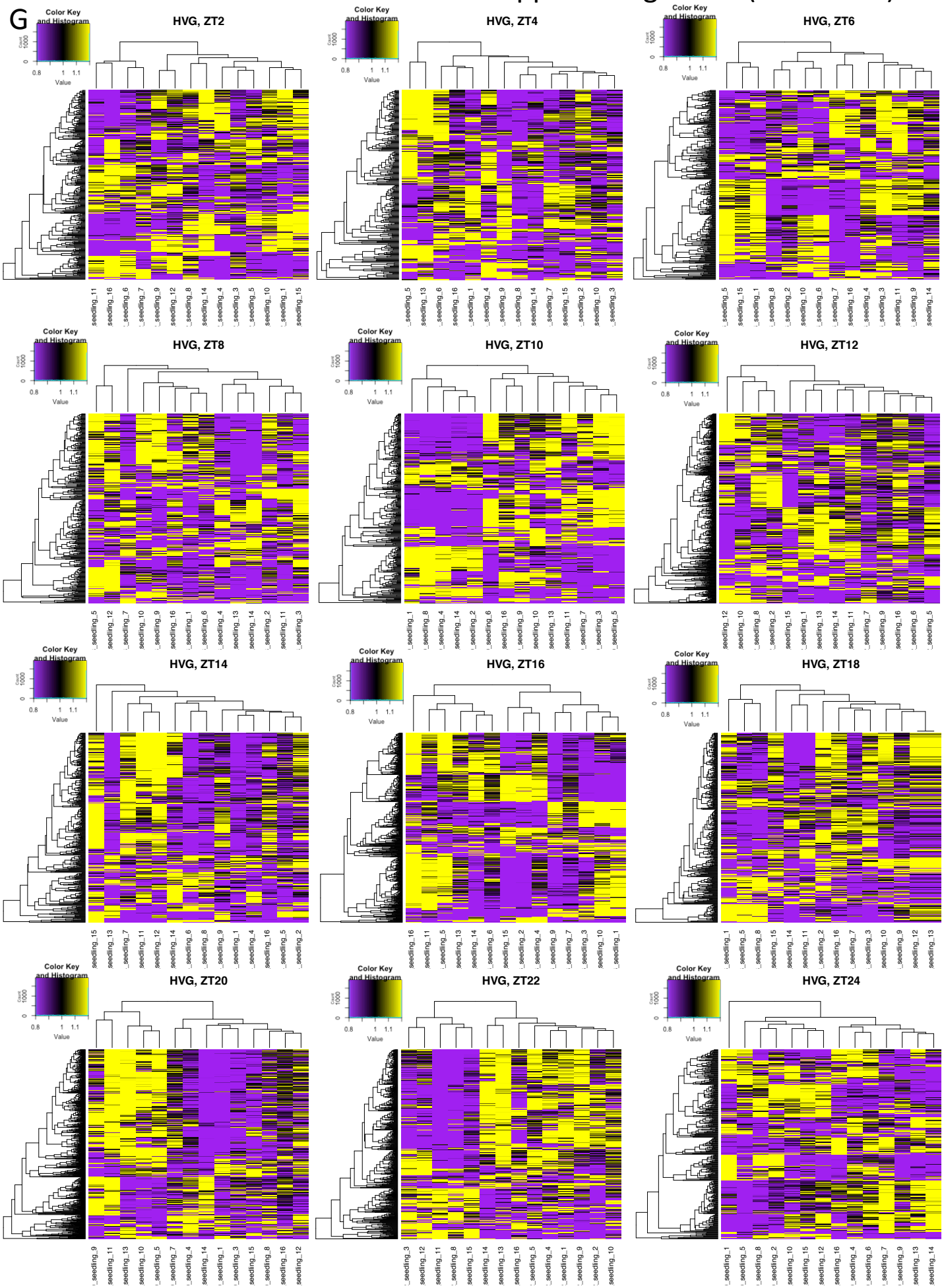
B



C

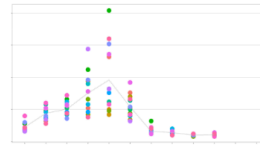
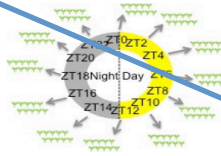






H

About page of AraNoisy



What is AraNoisy

AraNoisy is a web-based tool for accessing inter-individual transcriptional variability in *Arabidopsis thaliana*, through a 24h diurnal cycle. It is based on single seedlings RNA-seq from Cortijo et al 2018 performed on 14 individual seedlings every 2 hours during a 24 hours time course. ZT2 (2 hours after dawn) to ZT12 have been collected during the day, and ZT14 to ZT24 during the night. ZT12 was collected just a few minutes before dusk, and ZT24 just a few minutes before dawn.

The global trend of the square coefficient of variation ( $CV^2$ ) relative to expression level was measured for each time point, and used to calculate a corrected variability:  $\log_2(CV^2 / \text{trend at the same expression level})$ . Highly variable genes were detected for each time point, and are represented by blue dots in the plot of the corrected variability. Genes with expression level lower than 5 TPM and/or with a size lower than 150bp are not analysed. If data for a few time-points are missing for a gene, this means its expression is less than 5 TPM at these time-points.

To see the level of variability for a gene, go to the 'Single seedling profiles' tab, enter the name of your gene of interest and click 'Draw plots'. The results can be accessed graphically online and downloaded to your computer for further use.

How to cite us

Cortijo, S., Aydin, Z., Ahnert, S., & Locke, J. (2018). Widespread inter-individual gene expression variability in *Arabidopsis thaliana*. *bioRxiv*.

How to access the raw data

Include GEO link

Report a bug

If you have any problems with ShinyNoisy or would like to make a suggestion, send us an email at [sandra.cortijo@slcu.cam.ac.uk](mailto:sandra.cortijo@slcu.cam.ac.uk)

1. To go to AraNoisy tool, click 'Single seedling profiles'

AraNoisy tool page



Gene Variability Profile

Input a gene name (ex: AT5G52310) and click on the 'Draw plot' button

Draw plots

Download top plot

Download bottom plot

Download expression table for this gene

Download variability table for this gene

2. Enter a gene id, using the following format: AT5G52310

3. Click 'Draw plot' You need to click 'Draw plot' each time you enter a new gene id



Gene Variability Profile

Input a gene name (ex: AT5G52310) and click on the 'Draw plot' button

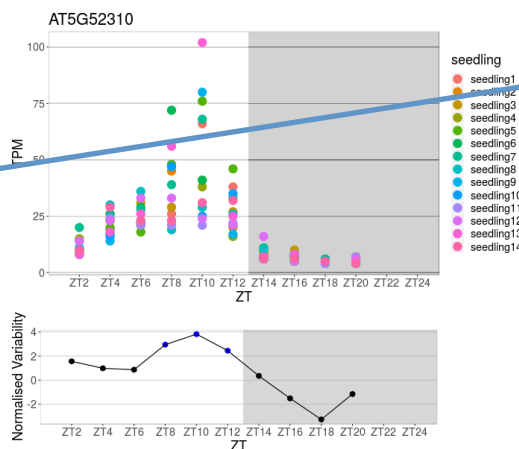
Draw plots

Download top plot

Download bottom plot

Download expression table for this gene

Download variability table for this gene

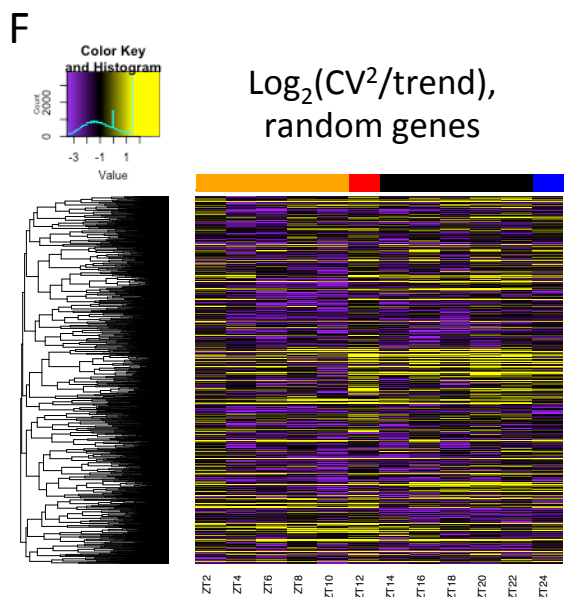
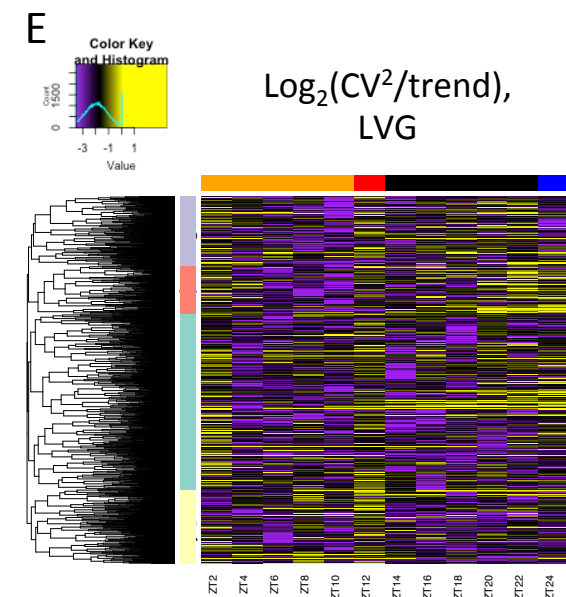
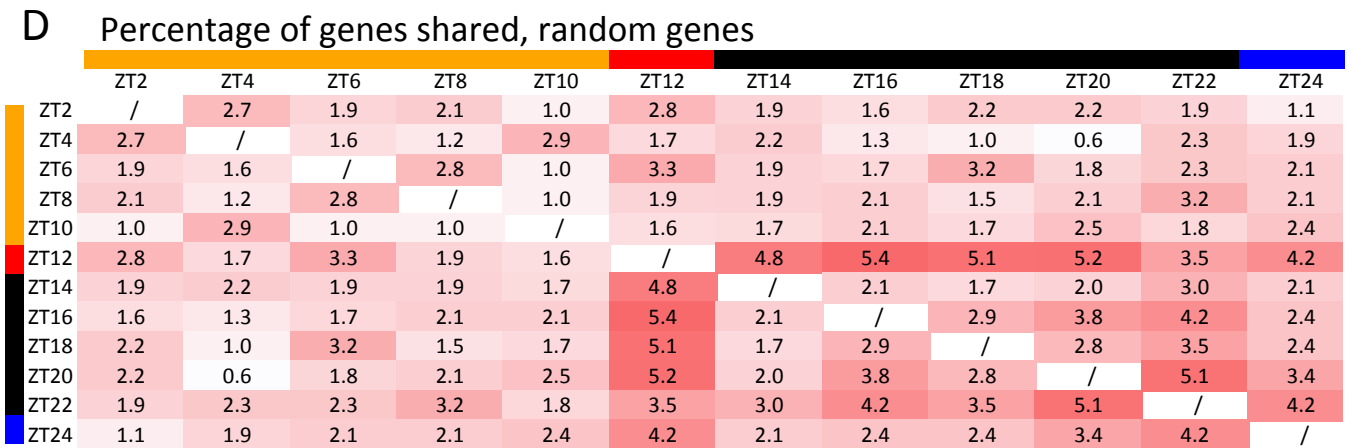
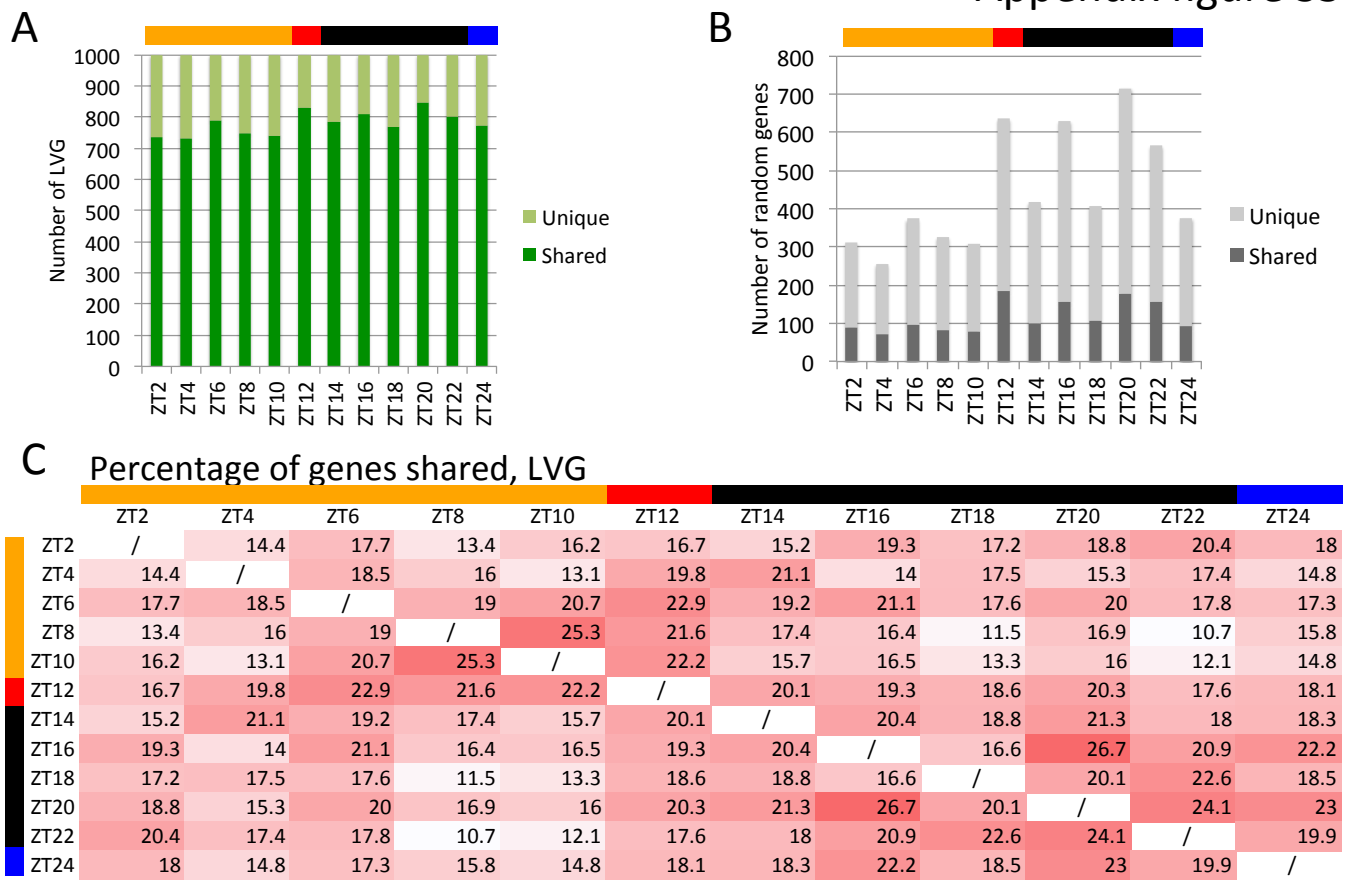


4. Click to download the top plot (expression levels during the time course in the single seedlings) or the bottom plot (variability level during the time course)



## Appendix Figure S2:

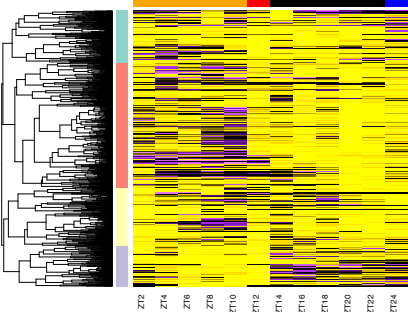
- A.** Trend for the global  $CV^2$  profile for each of the 12 time-points.
- B.** Identification of HVGs (blue), LVGs (green) and of one set of random genes (grey) for each of the 12 time-points.
- C.** Distribution of the corrected  $CV^2$  for LVGs (green) HVGs (blue) at each time-point. The distribution of average (black) and 95% confidence interval (dotted grey) for the thousand sets of random genes (of same size as HVGs) is also represented at each time-point.
- D.** Expression profiles (top) in the 14 seedlings over a 24 hour time course (with 12 time-points) for two non-variable genes. Each dot is expression level (TPM) for a single seedling. Variability profiles (bottom) of the  $\log_2(CV^2 / \text{trend})$  for the same genes are also shown.
- E.** Expression profiles (top) in the 14 seedlings over a 24 hour time course (with 12 time-points) for two highly variable genes. Each dot is the expression level (TPM) for a single seedling. Variability profiles (bottom) of the  $\log_2(CV^2 / \text{trend})$  for the same genes are also shown. Blue dots indicate time-points for which the gene is identified as being variable.
- F.** Expression profiles (top) in the 14 seedlings over a 24 hour time course (with 12 time-points) for two genes with the level of variability changing across the 24 hour. Each dot is the expression level (TPM) for a single seedling. Variability profiles (bottom) of the  $\log_2(CV^2 / \text{trend})$  for the same genes are also shown. Blue dots indicate time-points for which the gene is identified as being variable.
- G.** Hierarchical clustering for each time-point of HVGs and of individual seedling using the mean normalised expression levels. The result is represented as a heatmap where yellow indicates a high mean normalised expression.
- H.** Step by step guide showing how to use the AraNoisy web-application (<https://jlggroup.shinyapps.io/AraNoisy/>). The top panel shows the about page, which contains a description of the tool. The page with the AraNoisy tool is accessed by clicking on the 'Single seedlings profiles' tab. This page and explanation of how to generate graphs, and export them, for variability profiles and expression profiles in the 14 seedlings are detailed in the middle and bottom panels.



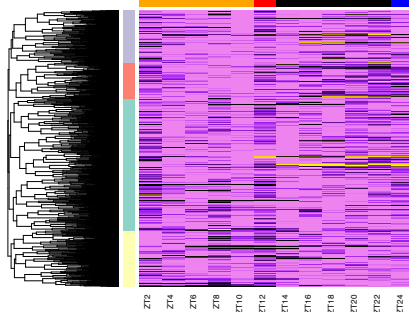
G



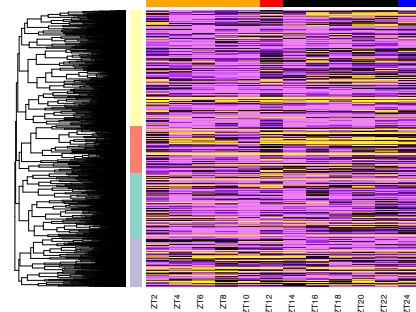
$\text{Log}_2(\text{CV}^2/\text{trend})$ ,  
HVG



$\text{Log}_2(\text{CV}^2/\text{trend})$ ,  
LVG

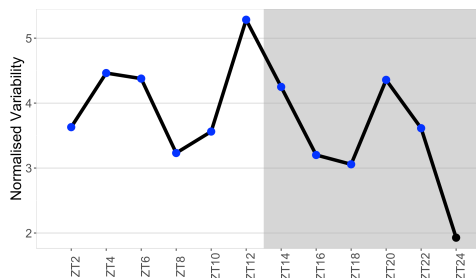
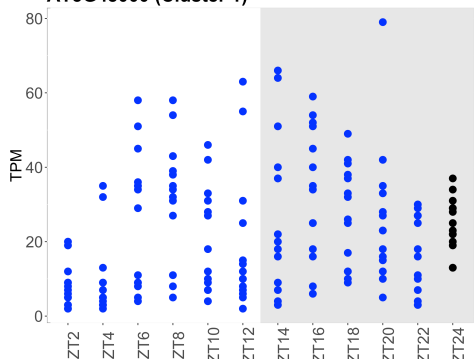


$\text{Log}_2(\text{CV}^2/\text{trend})$ ,  
Random genes

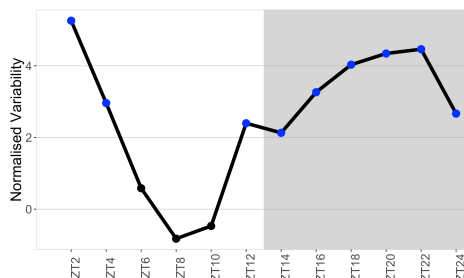
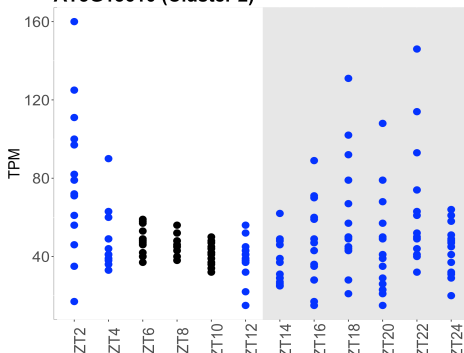


H

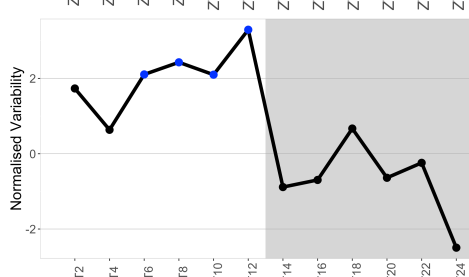
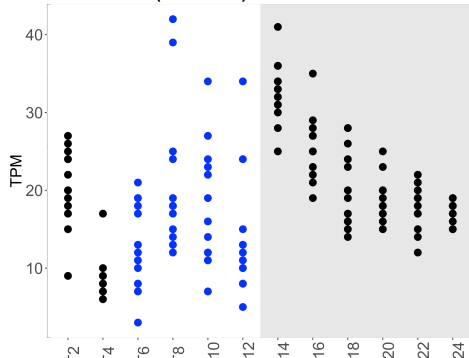
AT5G48000 (Cluster 1)



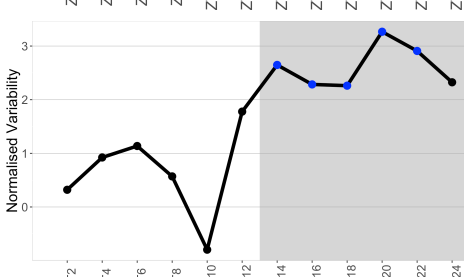
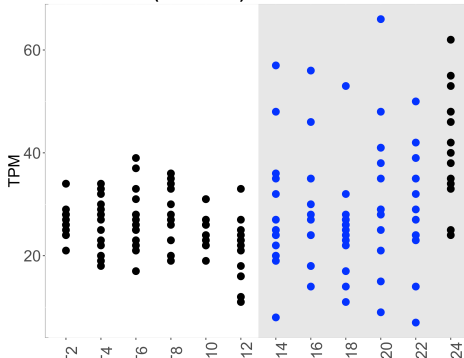
AT3G13610 (Cluster 2)



AT1G22400 (Cluster 3)



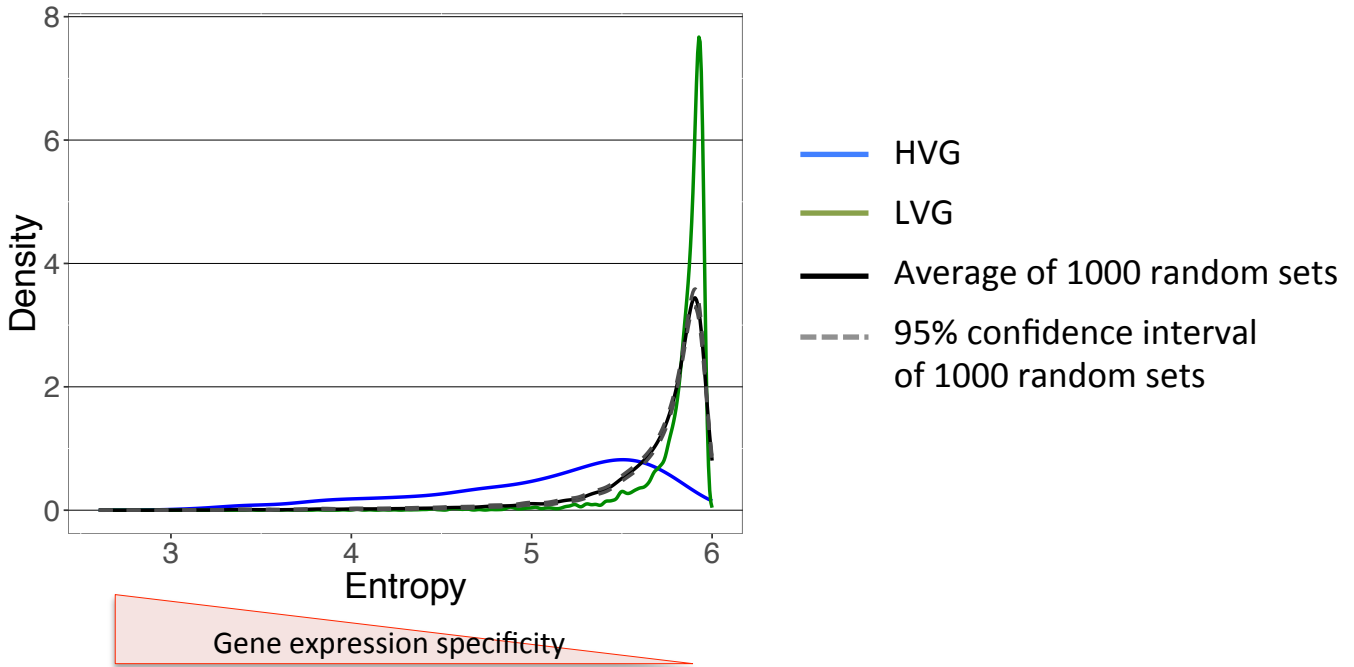
AT3G13435 (Cluster 4)



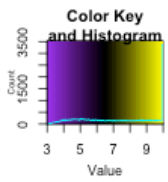
### Appendix Figure S3:

- A.** Number of LVGs for each time-point. These genes are separated between those that are also LVGs in at least one other time-point ('shared', dark green) or selected in only one time-point ('unique', light green). The top bar indicates time-points harvested during the day (orange), just before dusk (red), during the night (black) and just before dawn (blue).
- B.** Example for one set of random genes selected for each time-point, corresponding to the number of HVGs at the corresponding time-point. These genes are separated between those that are also selected in at least one other time-point (dark grey) and those selected in only one time-point (light grey). The top bar indicates time-points harvested during the day (orange), just before dusk (red), during the night (black) and just before dawn (blue).
- C.** Heatmap of the percentage of LVGs shared between time-points. Red indicates a high percentage in common between two time-points. The top and side bars indicate time-points harvested during the day (orange), just before dusk (red), during the night (black) and just before dawn (blue).
- D.** Heatmap of the percentage of random genes shared between time-points. Red indicates a high percentage in common between two time-points. The top and side bars indicate time-points harvested during the day (orange), just before dusk (red), during the night (black) and just before dawn (blue).
- E.** Hierarchical clustering of LVGs based on the  $\log_2(CV^2/\text{trend})$  at each time point. The result is represented as a heatmap where yellow indicates a high  $\log_2(CV^2/\text{trend})$ . The genes were separated into four clusters, indicated by the side colored bar. The top bar indicates time-points harvested during the day (orange), just before dusk (red), during the night (black) and just before dawn (blue). See Appendix Fig S3G for heatmaps with the same colour cutoffs for HVGs, LVGs and random genes.
- F.** Hierarchical clustering of one set of random genes based on the  $\log_2(CV^2/\text{trend})$  at each time point. The result is represented as a heatmap where yellow indicates a high  $\log_2(CV^2/\text{trend})$ . The genes were separated into four clusters, indicated by the side colored bar. The top bar indicates time-points harvested during the day (orange), just before dusk (red), during the night (black) and just before dawn (blue). See Appendix Fig S3G for heatmaps with the same colour cutoffs for HVGs, LVGs and random genes.
- G.** Hierarchical clustering of HVG (left), LVG (middle) and one set of random genes (right) based on the  $\log_2(CV^2/\text{trend})$  at each time point. The result is represented as a heatmap where yellow indicates a high  $\log_2(CV^2/\text{trend})$ . The same color cutoff is used in all three heatmaps.
- H.** Example of expression levels (TPM) in the 14 seedlings across the 12 time-points (top panel) for 4 genes that are representative of each of the 4 clusters identified in Fig 2D for HVGs based on the  $\log_2(CV^2/\text{trend})$ . Each dot is the expression level (TPM) for a single seedling from the RNAseq. AT5G48000, AT3G13610, AT1G22400 and AT3G13435 are representative genes of the clusters 1, 2, 3 and 4 respectively. The normalised  $CV^2$  over the time course (bottom panel) is also shown. Blue dots indicate time-points for which the gene is identified as being highly variable.

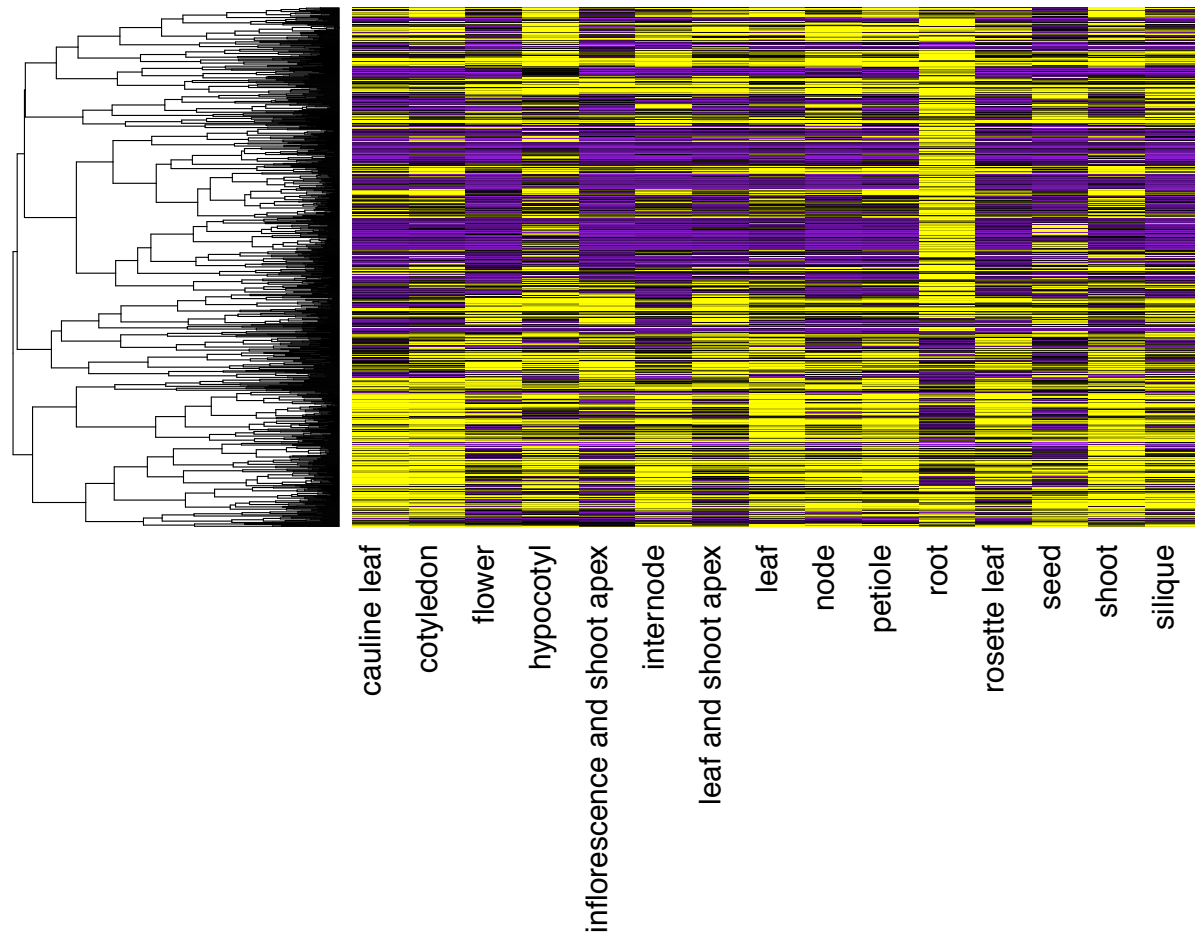
A



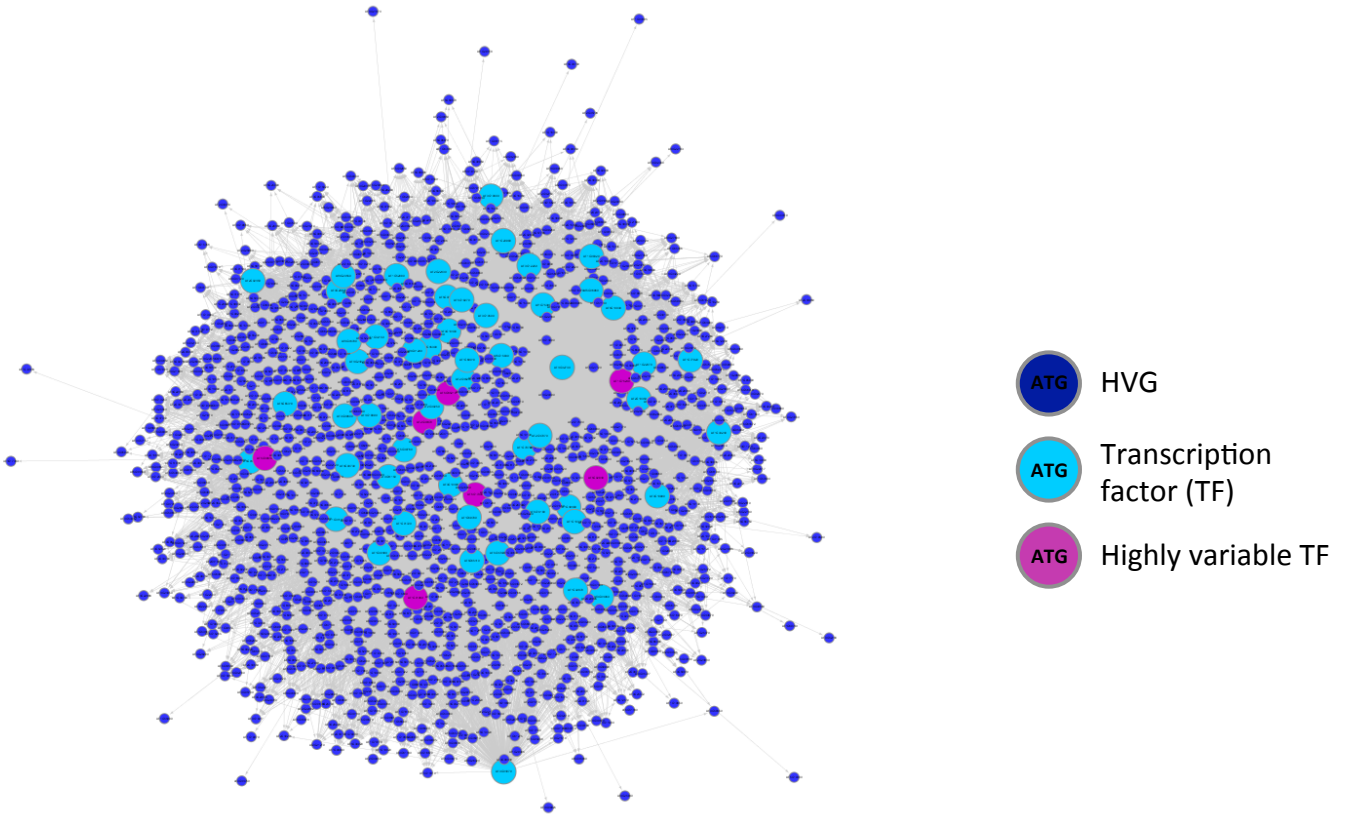
B



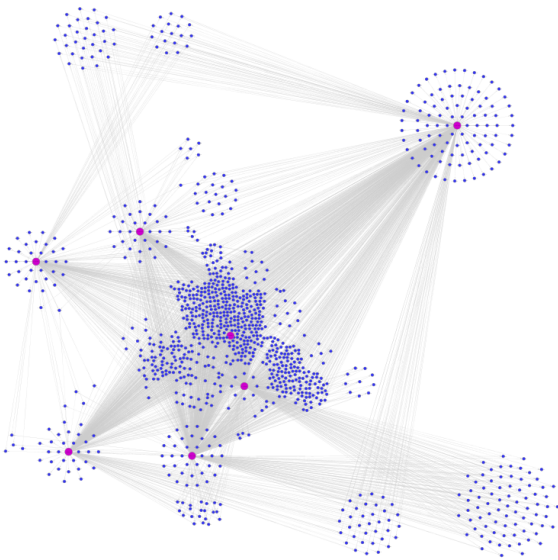
HVG, Tissue specificity expression



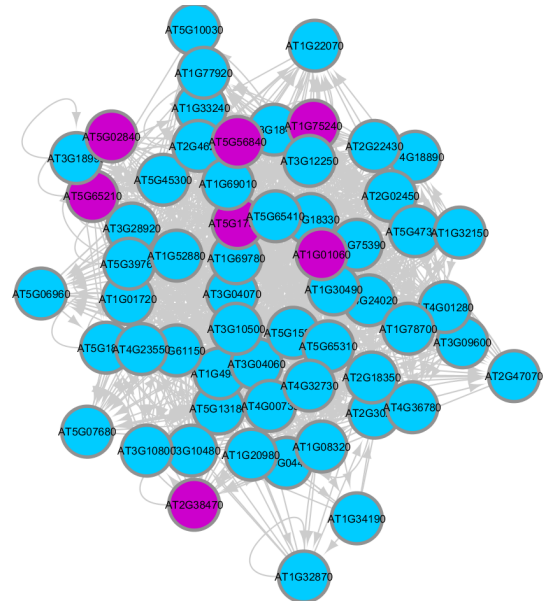
C



D



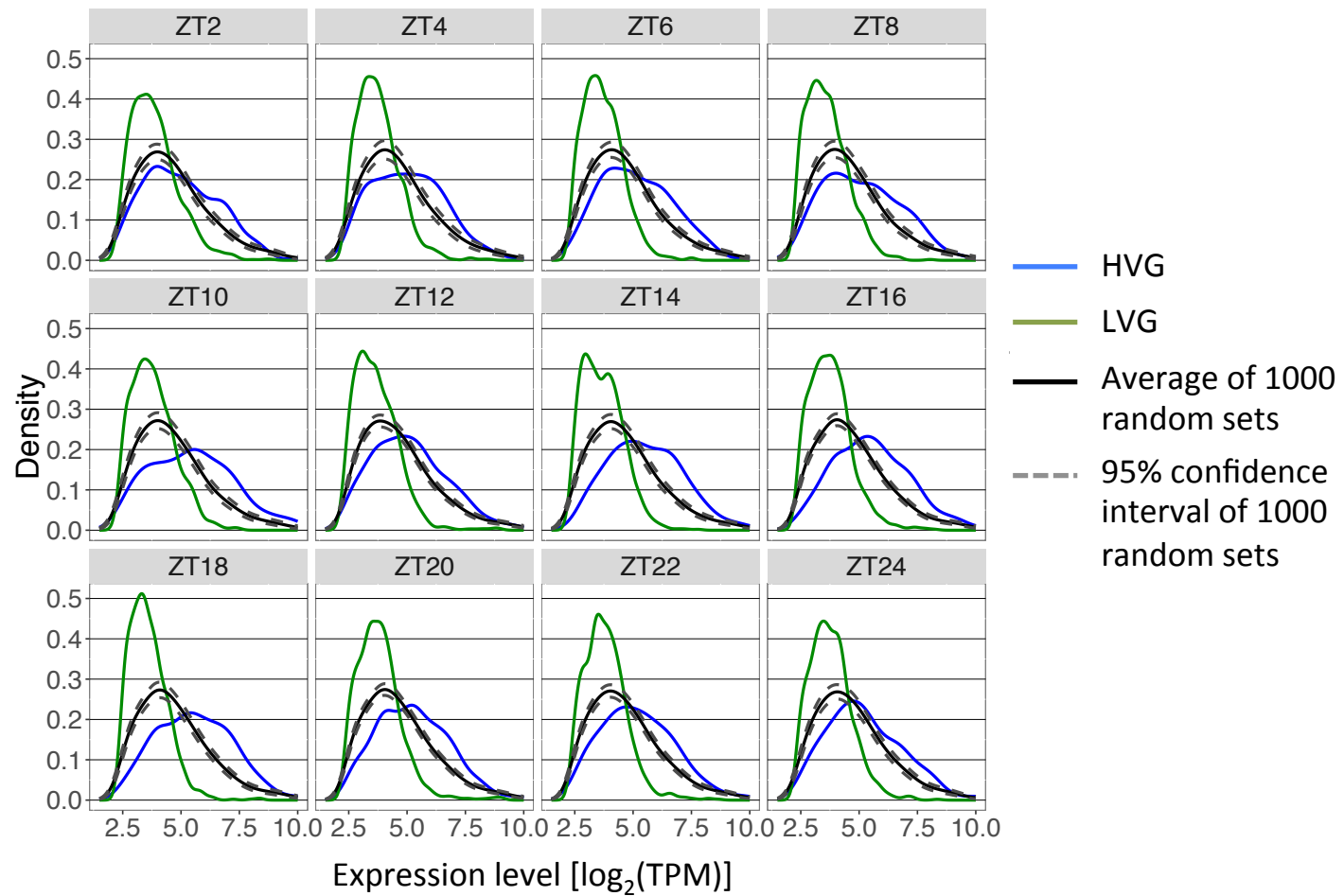
E



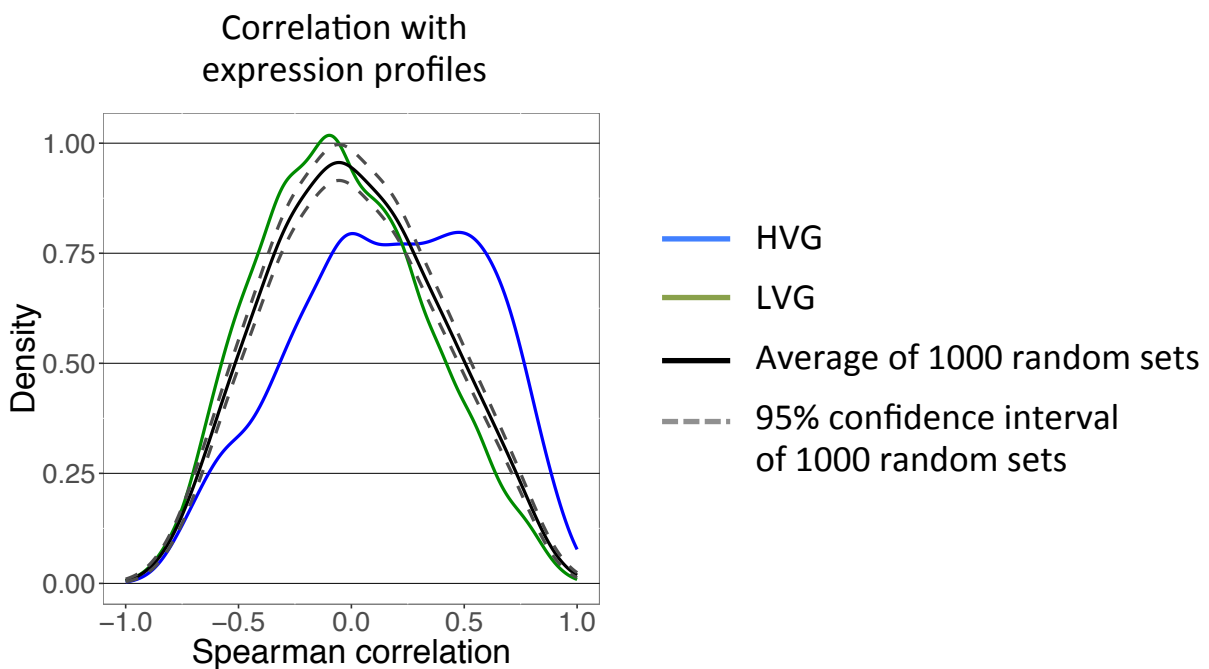
### **Appendix Figure S4:**

- A.** Distribution of the tissue specificity of HVGs (blue) and LVGs (green) as estimated by Shannon entropy calculation. The distribution of average (black) and 95% confidence interval (dotted grey) for the thousand random sets is also represented. Low entropy values indicate high tissue specificity.
- B.** Hierarchical clustering of HVGs based on the expression level in different tissues. The result is represented as a heatmap where yellow indicates a high expression level.
- C.** Gene regulatory network derived from the DAP-seq data for HVGs and the 60 TFs with targets enriched in HVGs. Dark blue nodes are HVGs, light blue nodes are TFs and purple nodes are highly variable TFs.
- D.** Gene regulatory network derived from the DAP-seq data for HVGs and the 7 highly variable TFs with targets enriched in HVGs. Dark blue nodes are HVGs and purple nodes are highly variable TFs.
- E.** Gene regulatory network derived from the DAP-seq data for 60 TFs with targets enriched in HVGs. Light blue nodes are TFs and purple nodes are highly variable TFs.

A

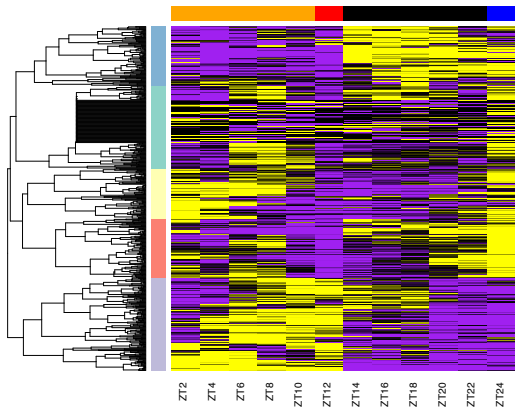


B

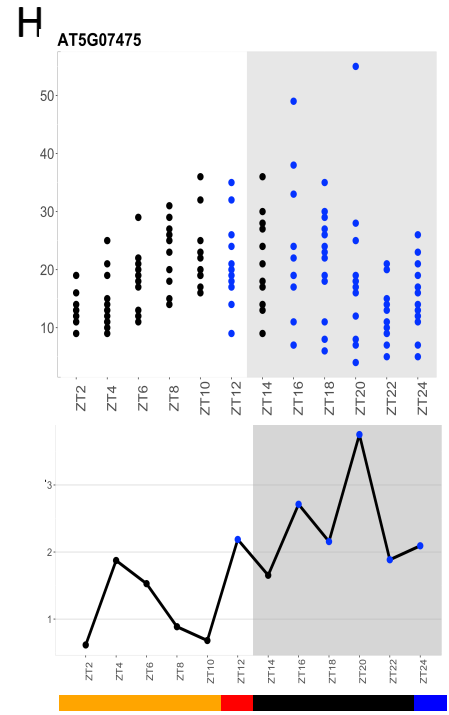
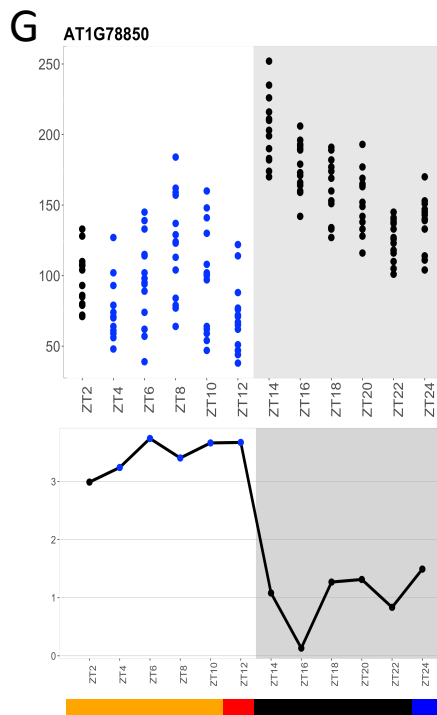
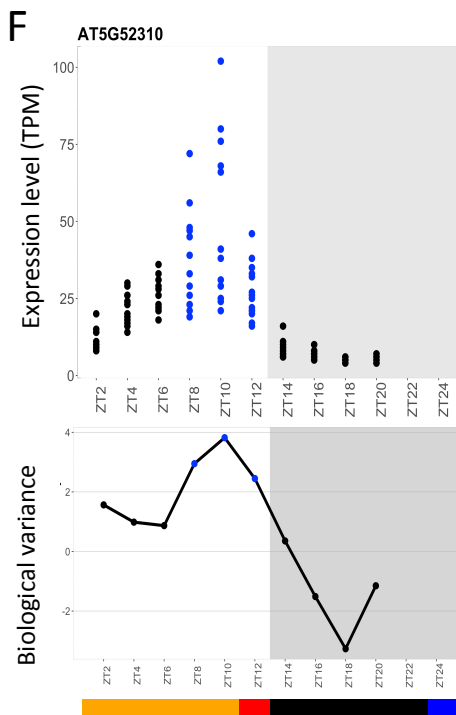
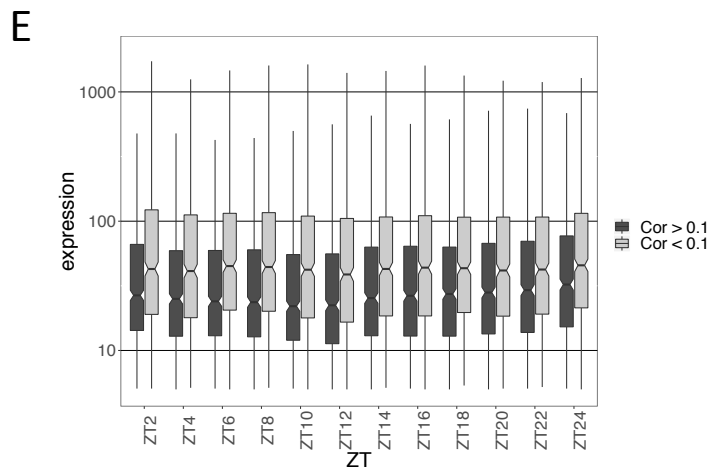
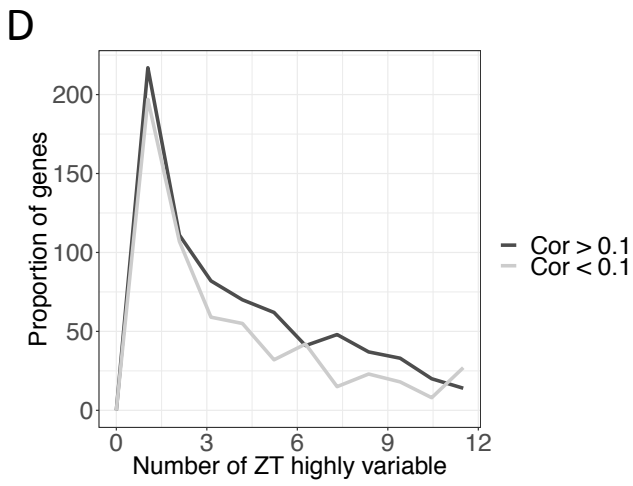
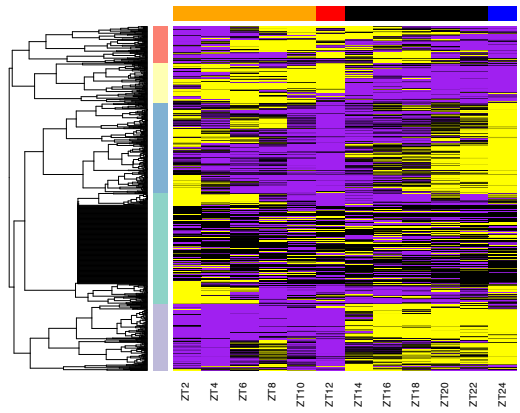




**C** HVG with Spearman correlation < 0.1



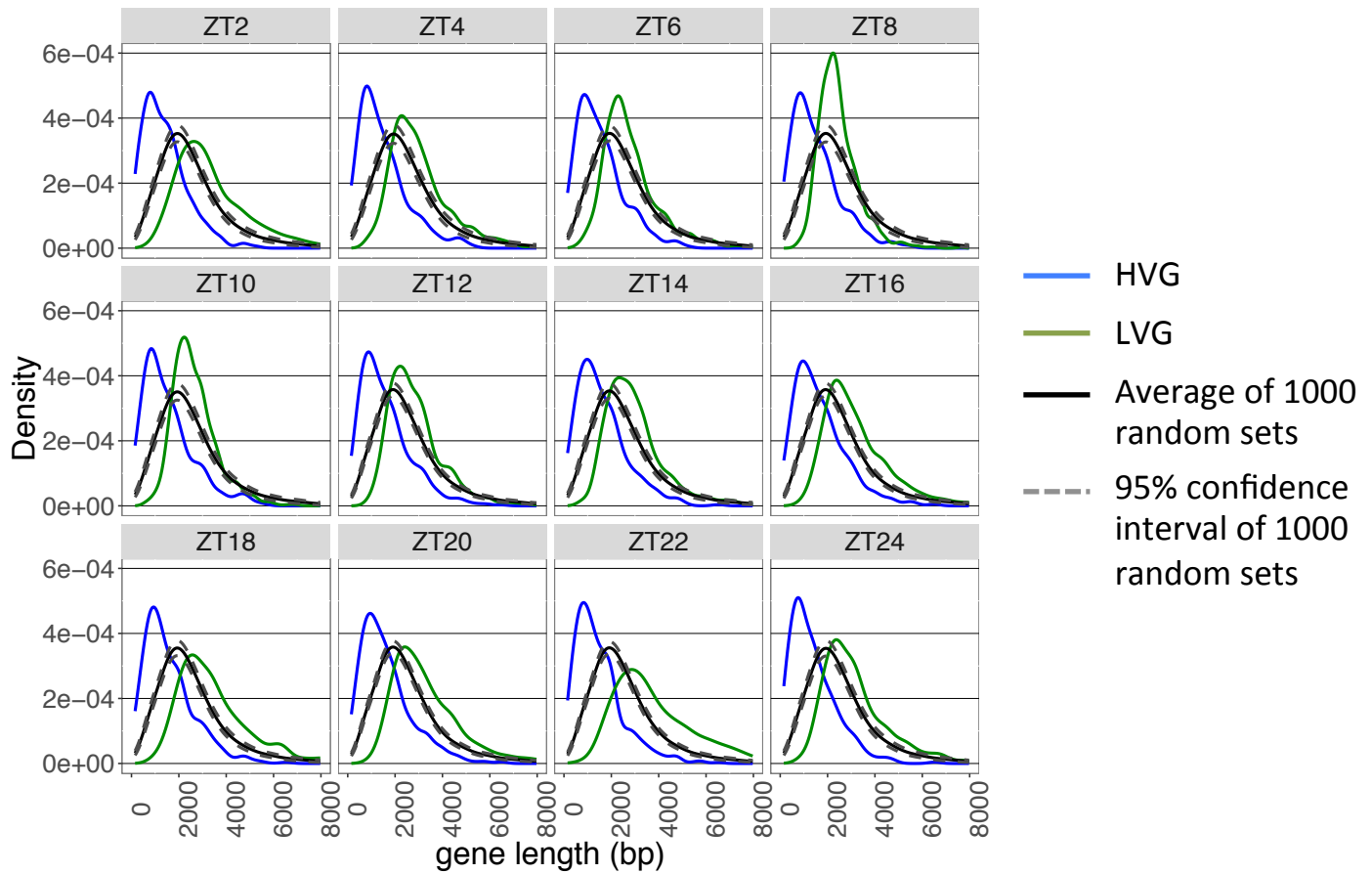
HVG with Spearman correlation < 0.1



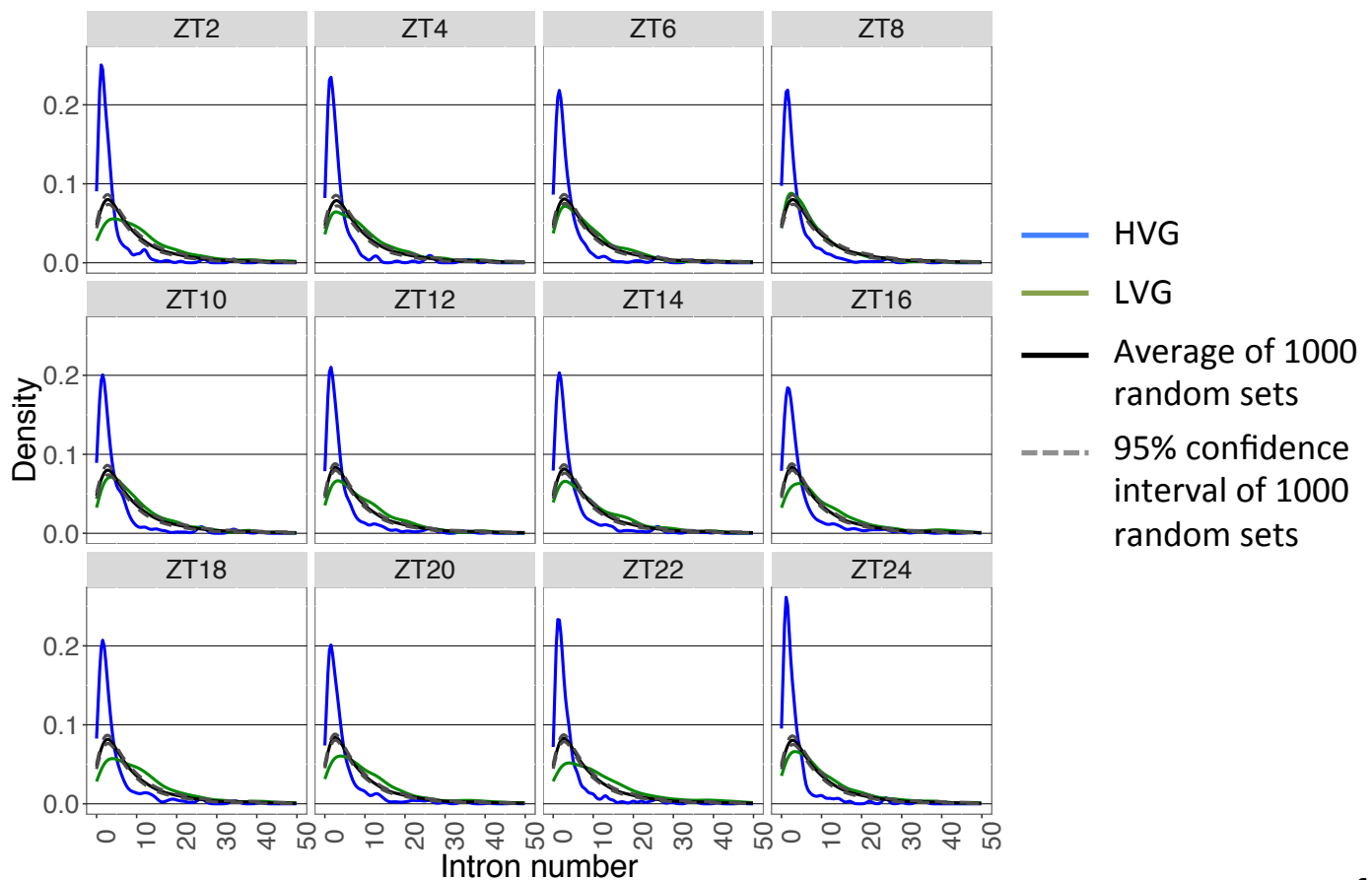
## Appendix Figure S5:

- A.** Distribution of the gene expression levels ( $\log_2(\text{TPM})$ ) for LVGs (green) and HVGs (blue) in each of the 12 time-points. The distribution of average (black) and 95% confidence interval (dotted grey) for the thousand sets of random genes (of same size as HVGs) is also represented at each time-point.
- B.** Distribution of the Spearman correlation between the average expression profile and the  $\log_2(\text{CV}^2/\text{trend})$  profile for LVG (green) and HVG (blue). The distribution of average (black) and 95% confidence interval (dotted grey) for the thousand random sets is also represented.
- C.** Hierarchical clustering of HVGs for which the Spearman correlation between the average expression profile and the  $\log_2(\text{CV}^2/\text{trend})$  is less than 0.1 (left) or more than 0.1 (right) based on mean normalised expression level in the 12 time-points. The result is represented as a heatmap where yellow indicates a high mean normalised expression level. The genes were separated into five clusters, indicated by the side colored bar. The top bar indicates time-points harvested during the day (orange), just before dusk (red), during the night (black) and just before dawn (blue).
- D.** Distribution of the number of time-points at which genes are identified as highly variable, for HVGs for which the Spearman correlation between the average expression profile and the  $\log_2(\text{CV}^2/\text{trend})$  is less than 0.1 (light grey) or more than 0.1 (dark grey).
- E.** Boxplot of the expression level at each time-point for HVGs for which the Pearson correlation between the average expression profile and the  $\log_2(\text{CV}^2/\text{trend})$  is less than 0.1 (light grey) or more than 0.1 (dark grey).
- F.** Example of expression levels (TPM) in the 14 seedlings across the 12 time-points (top panel) for a HVG with a strong positive Spearman correlation (0.86) between the average expression profile and the variability profile. The normalised  $\text{CV}^2$  over the time course (bottom panel) is also shown. Blue dots indicate time-points for which the gene is identified as being highly variable.
- G.** Example of expression levels (TPM) in the 14 seedlings across the 12 time-points (top panel) for a HVG with a strong negative Spearman correlation (-0.84) between the average expression profile and the variability profile. The normalised  $\text{CV}^2$  over the time course (bottom panel) is also shown. Blue dots indicate time-points for which the gene is identified as being highly variable.
- H.** Example of expression levels (TPM) in the 14 seedlings across the 12 time-points (top panel) for a HVG without a strong Spearman correlation (0.06) between the average expression profile and the variability profile. The normalised  $\text{CV}^2$  over the time course (bottom panel) is also shown. Blue dots indicate time-points for which the gene is identified as being highly variable.

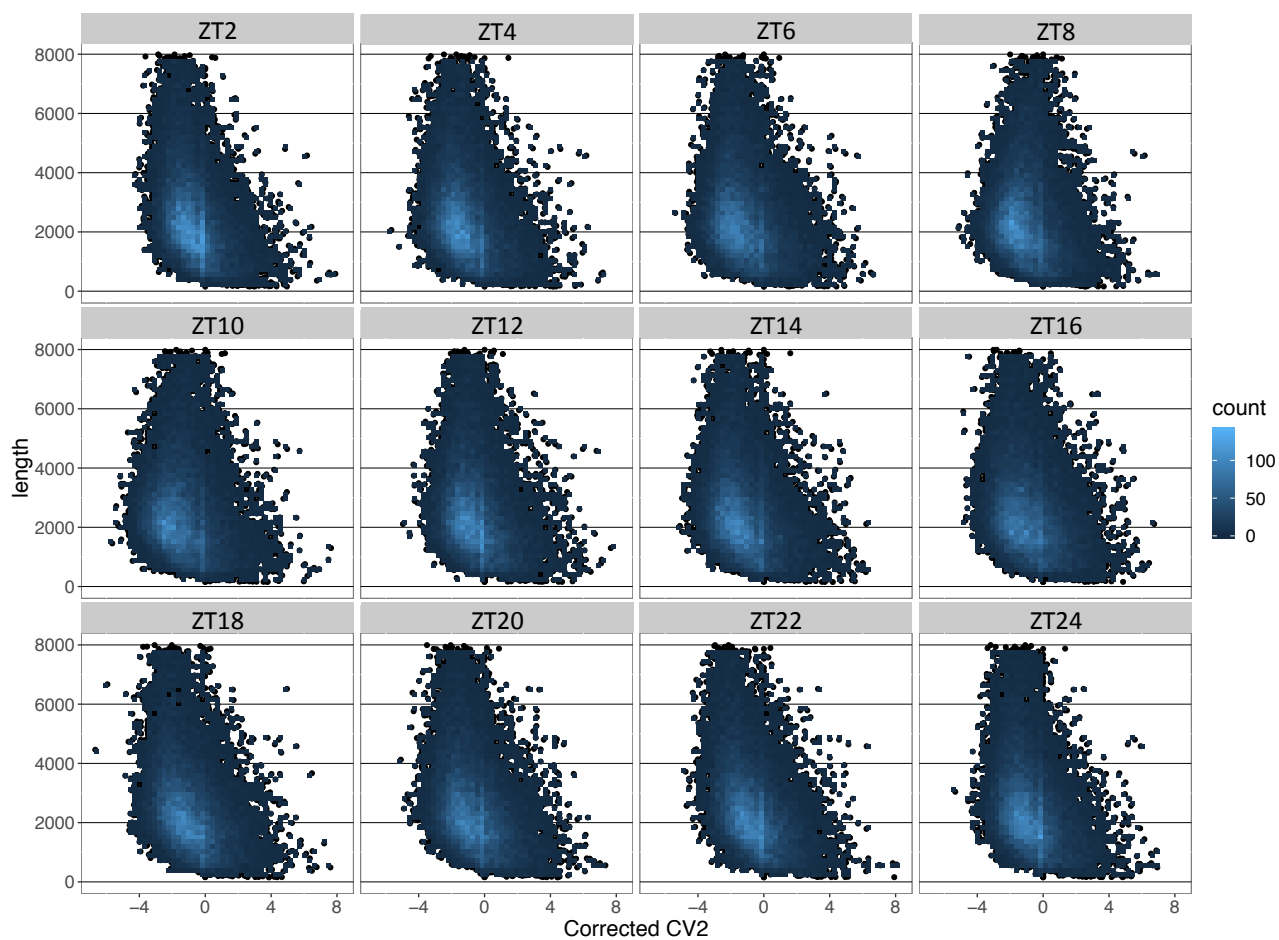
A



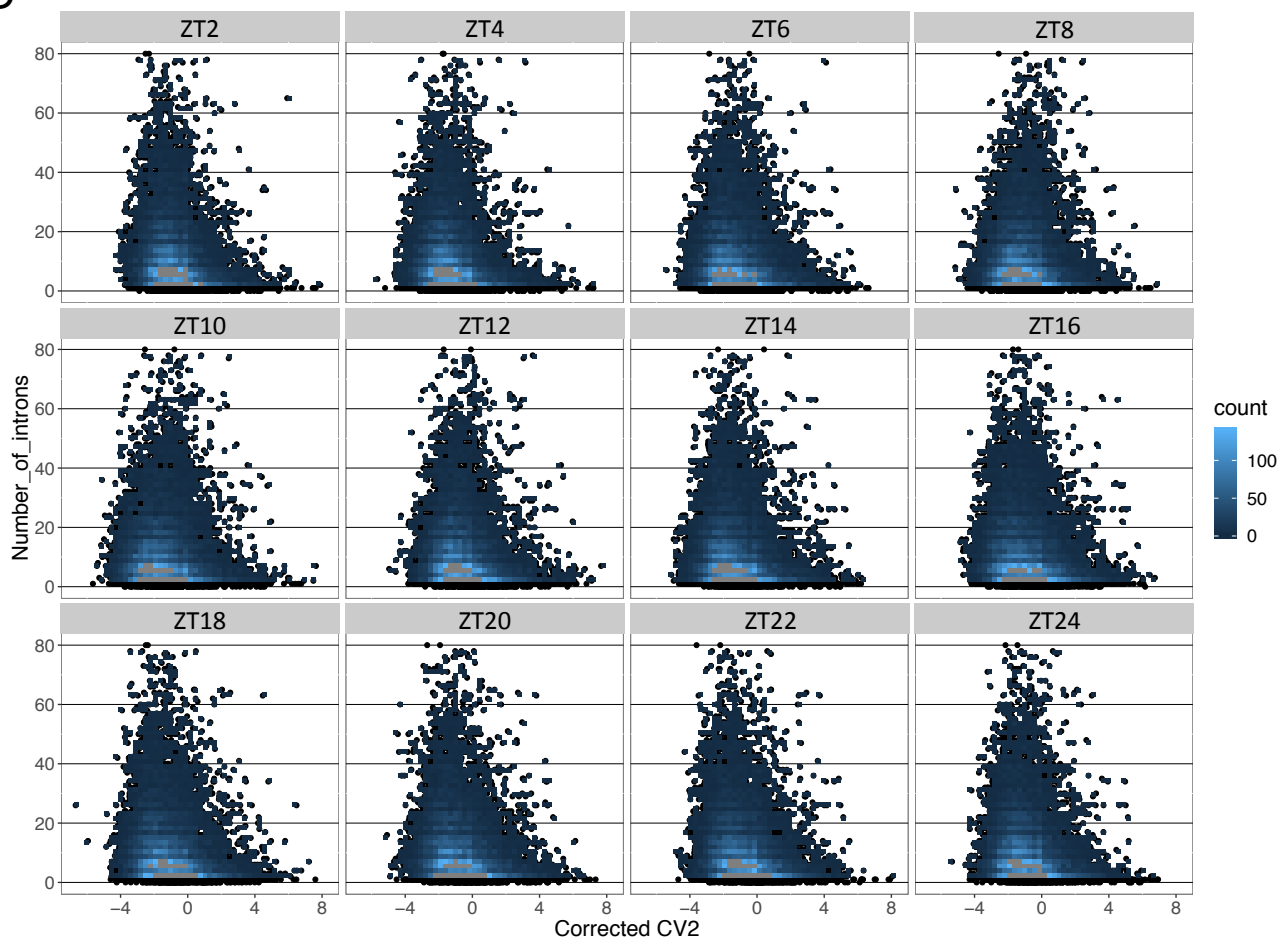
B



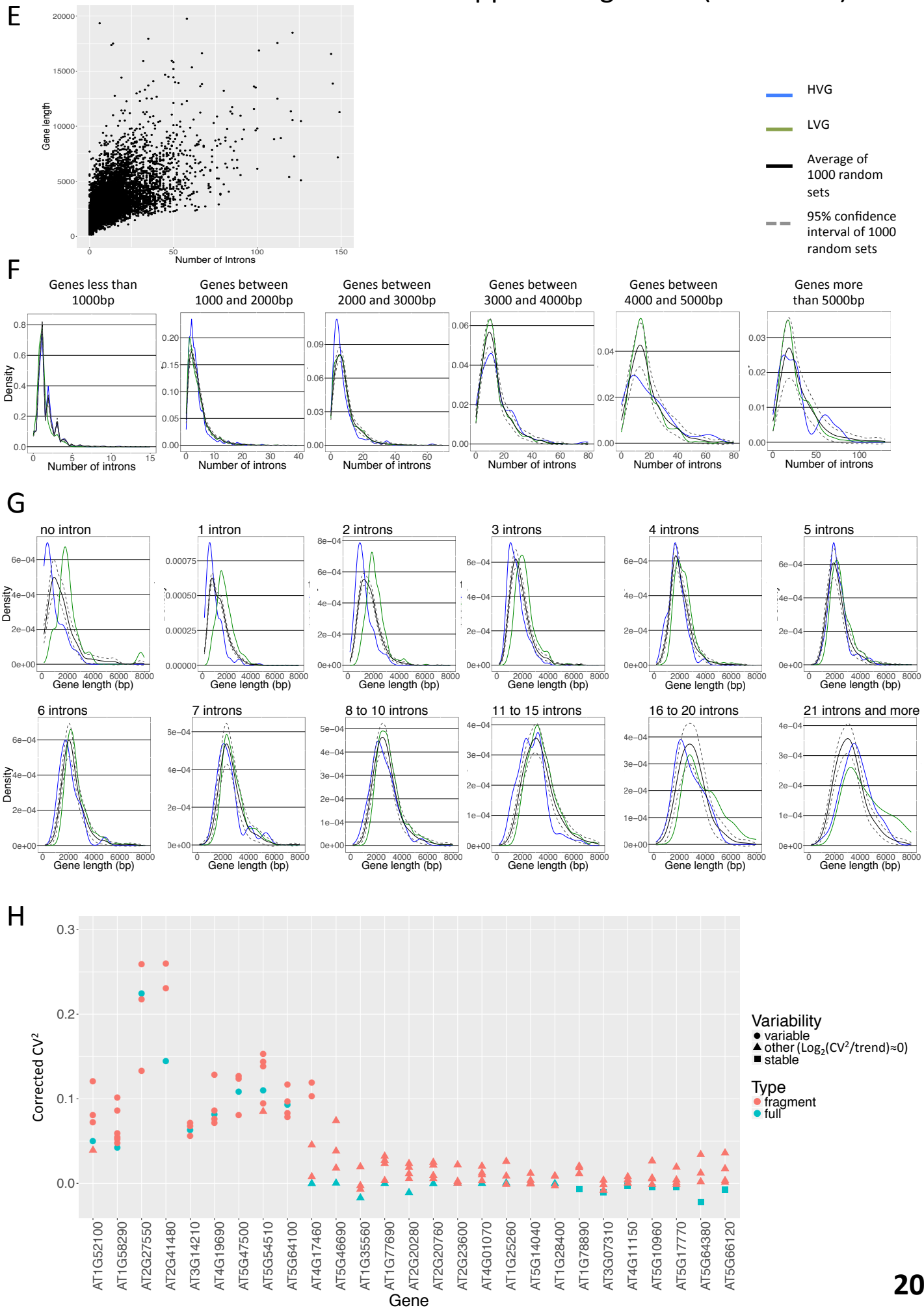
C



D



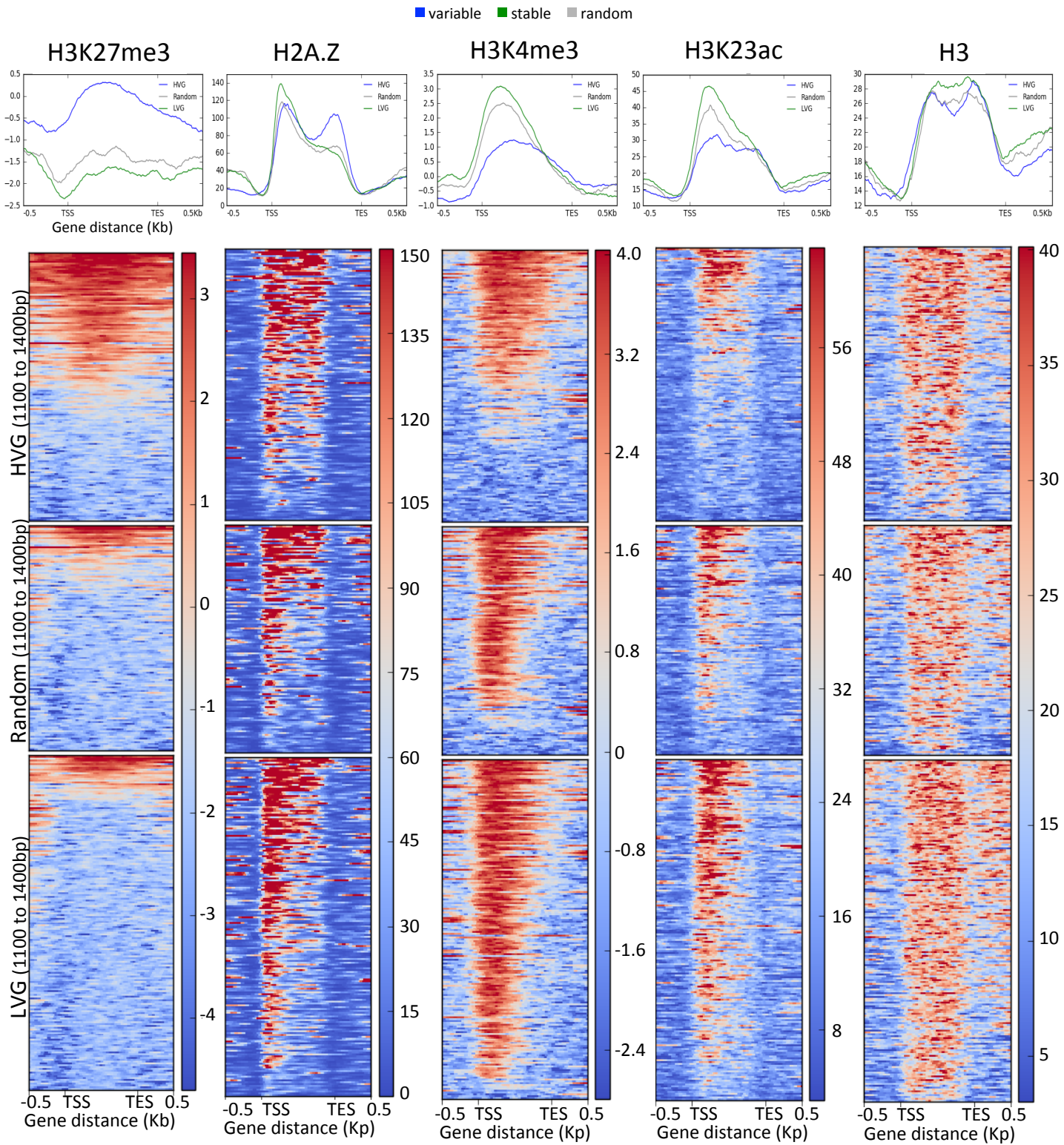
# Appendix figure S6 (continued)



## **Appendix Figure S6:**

- A.** Distribution of the gene length for LVGs (green) and HVGs (blue) in each time-point. The distribution of average (black) and 95% confidence interval (dotted grey) for the thousand sets of random genes (of same size as HVGs) is also represented at each time-point.
- B.** Distribution of the number of introns for LVGs (green) and HVGs (blue) in each time-point. The distribution of average (black) and 95% confidence interval (dotted grey) for the thousand sets of random genes (of same size as HVGs) is also represented at each time-point.
- C.** Plot of the gene length versus the  $\log_2(\text{CV}^2/\text{trend})$  for all genes in each time-point.
- D.** Plot of the number of introns versus the  $\log_2(\text{CV}^2/\text{trend})$  for all genes in each time-point.
- E.** Plot of the gene length versus the number of introns for all genes.
- F.** Distribution of the number of introns for LVGs (green) and HVGs (blue) for genes with several windows of gene length. The distribution of average (black) and 95% confidence interval (dotted grey) for the thousand sets of random genes (of same size as HVGs) is also represented at each time-point.
- G.** Distribution of the gene length for LVGs (green) and HVGs (blue) for genes with fixed number of introns. The distribution of average (black) and 95% confidence interval (dotted grey) for the thousand sets of random genes (of same size as HVGs) is also represented at each time-point.
- H.** Corrected  $\text{CV}^2$  for 9 HVG (round), 7 LVG (square) and 11 other genes (triangle), for the full genes (blue) as well as ~250bp fragments of these genes (red). Variability levels of full genes and their fragments were analysed together.

A

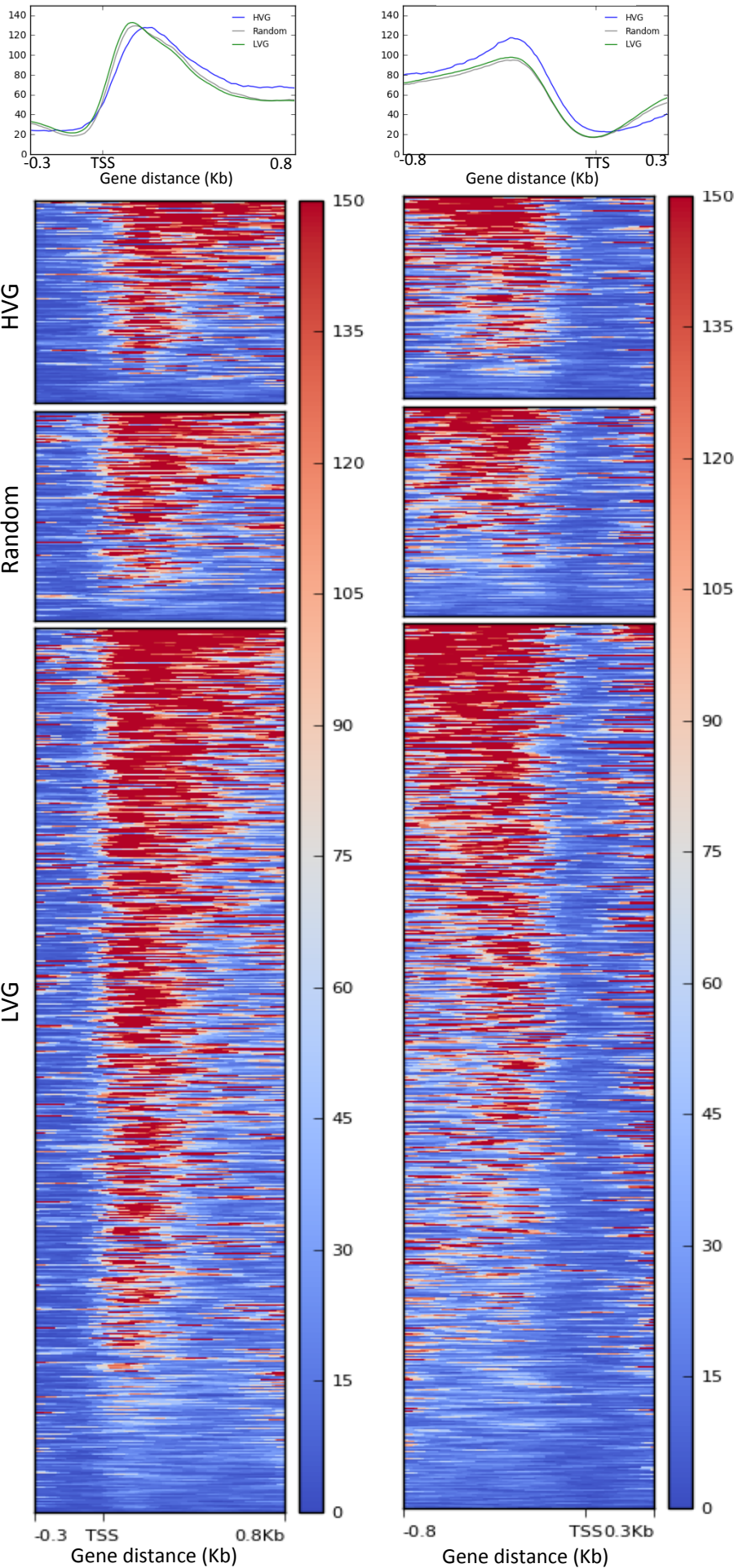


■ variable ■ stable ■ random

B

Mnase at TSS

Mnase at TTS





**Appendix Figure S7:**

**A.** Average profile and heatmap of H3K27me3, H2A.Z, H3K4me3, H3K23ac and H3 marks for HVGs (top), random genes (middle) and LVGs (bottom) that are 1100bp to 1400bp long. Red means a high level and blue means a low level for the chromatin marks.

**B.** Average profile and heatmap of MNase signal for HVGs (top), random genes (middle) and LVGs (bottom) around TSS (left) and TTS (right) of genes. Red means a high level and blue means a low level for the MNase signal.