

1. For each node in the network, create a Node Vicinity Network (NVN), by taking n steps out from the query node.

2. For each node in a NVN, iteratively remove nodes that show higher number of edges to the outside of NVN.

3. Each NVN is transformed into a Stable Putative Cluster (SPC).

4. Greedy selection of clusters from the population of SPCs. This is done by ranking the SPCs according to in/out edge ratio for all nodes in a SPC. Starting from the SPC with highest ratio in the population, a next non-overlapping SPC is accepted as cluster, until SPCs are depleted. This step permits selection of properties of accepted clusters, such as cluster size.

5. Remove clustered nodes found in step 4 from the network. Go to step 1.

**Supplemental Figure 1. Principles of Heuristic ClusterChiseling Algorithm (HCCA).** HCCA consists of five steps; 1: generation of NVN for each node in the network by taking in the n-step neighborhood (3-step neighborhood used in this study), 2: removal of nodes more connected to the nodes outside of NVN which results in, 3: generation of SPCs, 4: selection of clusters from the population of SPCs (preferred cluster size range 40-200 was selected for this study), 5: removal of clustered nodes from the network and return to step 1. A more detailed explanation and comparison of HCCA to other clustering algorithms is available in Mutwil et al. 2010.

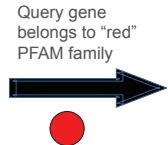
1) Enter gene of interest

2) Generate NVN for each member of the "red" PFAM gene family

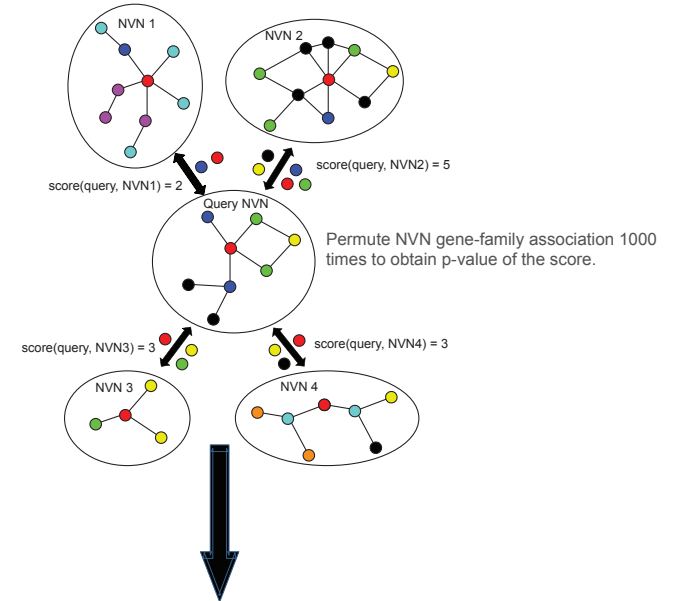
3) NetworkComparator - comparison of query NVN to NVNs centered around red PFAM

Computation of similarity of query NVN to every NVN centered around a gene belonging to the red pfam family for all species in the database. Score of 1 is given for each gene family to query NVN and target NVN have in common. Colors of nodes depict different PFAM families.

Probe set ID  
Locus ID  
Keyword

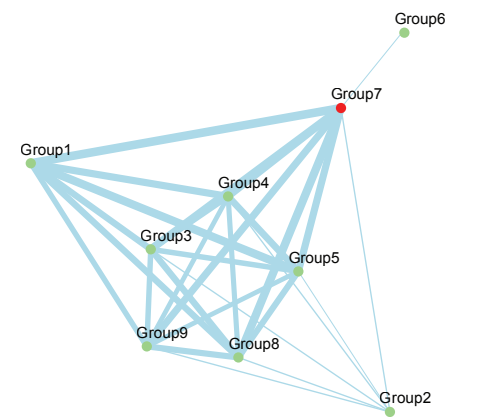


NVNs centered around genes belonging to red PFAM families are generated for all species in database.



4) Selection step: Analysis of co-expression group similarities and selection of NVNs for further analysis

The analysis returns three elements:



The table below shows score of your query gene to all genes found in the family: PsdD

| Group   | Network         | Probeset ID         | Transcript ID                        | Score | p-value |
|---------|-----------------|---------------------|--------------------------------------|-------|---------|
| Group 1 | ExpMatBar.lhr   | Contg4015_at        | p36213                               | 23    | 0.0     |
| Group 2 | ExpMatMed.hlr   | Mtr_34726.1.S1_at   | medtr_lg132590.1                     | 4     | 0.0     |
| Group 3 | ExpMatMed.hlr   | Mtr_19267.1.S1_at   | medtr5g06220.1                       | 21    | 0.0     |
| Group 4 | ExpMatRice.lhr  | Oa_7974.1.S1_at     | LOC_Os03g09220.3                     | 20    | 0.0     |
| Group 5 | ExpMatWheat.lhr | Ta_24304.2.S1_a_at  | UniRef90_P36213                      | 20    | 0.0     |
| Group 6 | ExpMatWheat.lhr | TaAdf_58068.1.S1_at | UniRef90_P36213                      | 1     | 0.033   |
| Group 7 | ExpMatAra.lhr   | atg03130            | atg02770                             | 9     | 0.0     |
| Group 8 | ExpMatPop.lhr   | Pip_5240.1.S1_s_at  | gippoptr_1_11564827 eugene3.00081422 | 31    | 0.0     |
| Group 9 | ExpMatSoy.lhr   | Gma_10852.1.S1_s_at | gippoptr_1_11566261 eugene3.00100819 | 8     | 0.0     |
|         |                 | Gma_10852.2.S1_at   | Glyma10g39460.1                      | 19    | 0.0     |
|         |                 |                     | Glyma20g28300.1                      | 16    | 0.0     |

-A table showing score of the query gene NVN to all NVNs belonging to the same PFAM family. Red arrows indicate NVNs selected for further analysis.

The table below shows which pfam families are enriched in vicinity networks of PsdD family in analyzed species.

| Family                      | Description  | Times found in analyzed vicinity networks |
|-----------------------------|--|---|
| PsdD                        | This family consists of PsdD from plants and cyanobacteria   | 12  |
| Chloroa_b_bind              | Chloroal A-B binding protein   | 11  |
| PsdB                        | This family consists of the 23 kDa subunit of oxygen evolving system of photosystem II or PsdB from various plants (where it is encoded by the nuclear genome) and Cyanobacteria | 10  |
| PSI_PsaH                    | Photosystem I reaction centre subunit VI   | 9   |
| MSP                         | Manganese stabilizing protein; photosystem II polypeptide  | 8   |
| PSI_PSAK                    | Photosystem I psaf6_psaK   | 8   |
| AAA                         | Oxidoreductase NAD-binding domain; Similar photosynthetic, but also bind FAD/NAD, have essentially no similarity   | 8   |
| NAD_binding_1               | ATPase family associated with various cellular activities (AAA)  | 8   |
| PSII_PsbL                   | Photosystem II protein L (PsbL)  | 7   |
| PSII_PsbM                   | Photosystem II reaction centre S protein (PsbM)  | 7   |
| PSI_PsaI                    | Photosystem I reaction centre subunit IV - PsaI  | 7   |
| Fructose_1,6_bisphosphatase | Fructose 1,6-bisphosphatase  | 7   |
| Copper_binding              | Copper binding proteins; photosystem antenna family  | 7   |
| PsbQ                        | Oxygen evolving enhancer protein I (PsbQ)  | 6   |
| PsdB                        | Photosystem I reaction centre subunit VI (PsbQ or PSI-N)   | 6   |
| PsaI                        | Photosystem I reaction centre subunit XI   | 6   |
| PSI_PsaF                    | Photosystem I reaction centre subunit III  | 6   |

-A table reporting the frequency of specific PFAM families being present in NVNs of the red PFAM family.

5) Generation of ancestral network and revealing transcript identities

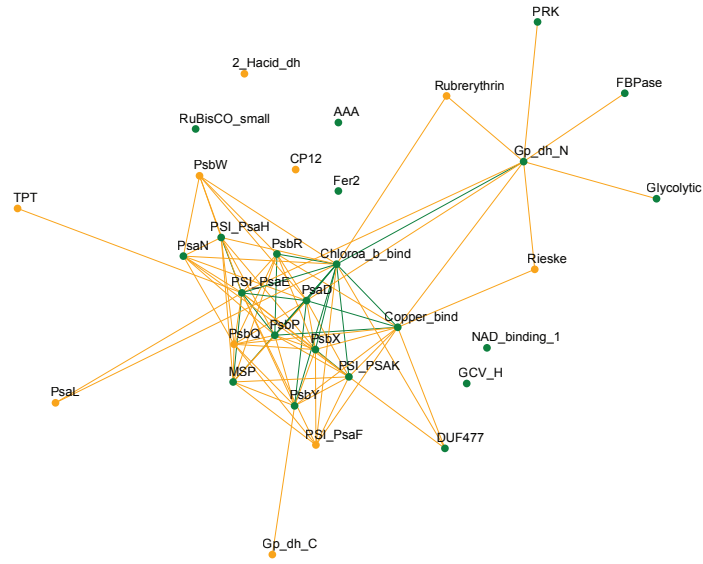
The program displays the ancestral network and table that displays the identity of potential functional homology.



6) Result

The analysis returns two elements:

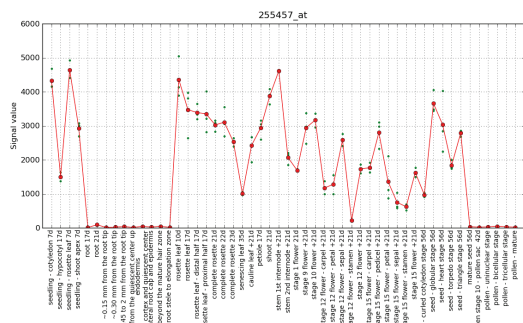
-Ancestral network depicting conserved PFAM families and co-expression relationships between them. Green, orange and red nodes and edges represent PFAMs and co-expression relationships present in ≥75%, ≥50% and ≥25% of NVNs selected in step 4.



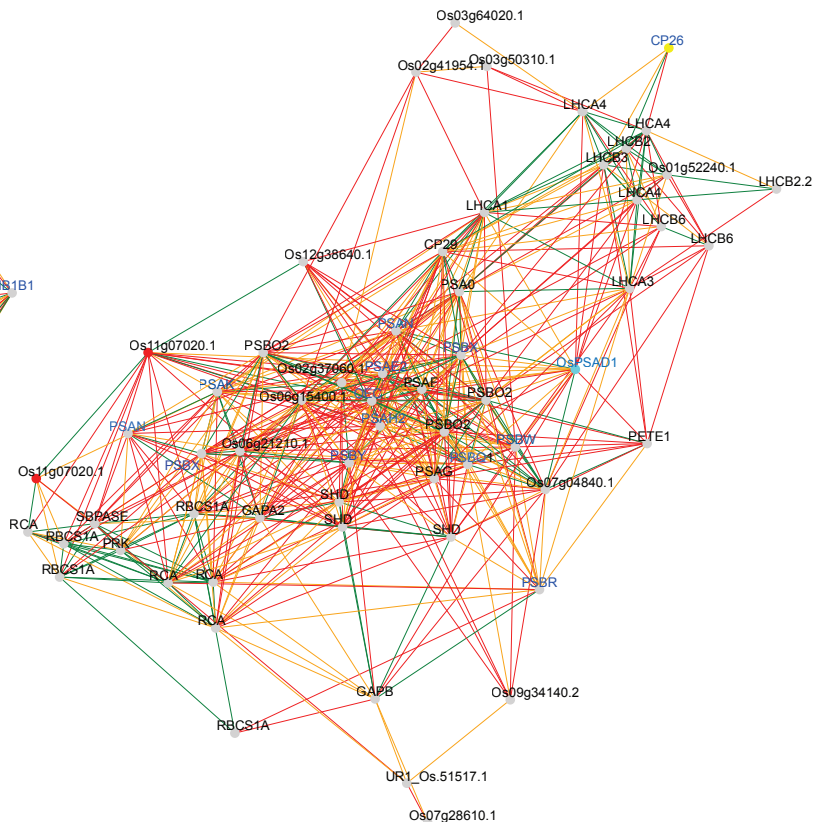
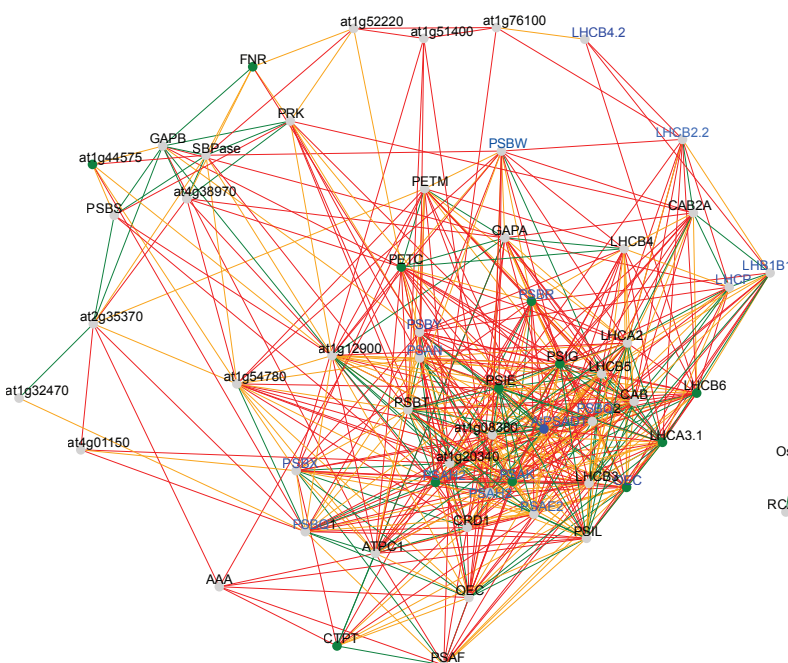
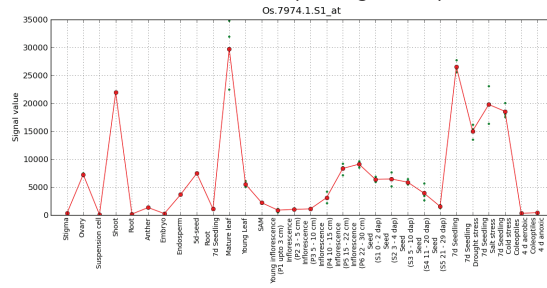
-A table revealing the probesets (genes) constituting the enriched PFAMs.

| Gene family                 | Description   | Gene ID          | Accession  | Species              |
|-----------------------------|---|------------------|------------|----------------------|
| PsdD                        | This family consists of PsdD from plants and cyanobacteria. PsdD is essential polypeptide of photosystem II (PSII) and is responsible for assembly of PSII reaction centre and is localized in the photosystem I binding site within the reaction centre. PsdD binds to other proteins which is bound by PsdA to form the reaction centre PSII environment.   | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| Chloroa_b_bind              | This family consists of chloroal A-B binding protein. This family consists of chloroal A-B binding protein. This family consists of chloroal A-B binding protein. This family consists of chloroal A-B binding protein.   | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| PsdB                        | This family consists of the 23 kDa subunit of oxygen-evolving system photosystem II. The 23 kDa subunit protein is responsible for PSII in higher eukaryotes and cyanobacteria. The 23 kDa subunit protein is responsible for PSII in higher eukaryotes and cyanobacteria. The 23 kDa subunit protein is responsible for PSII in higher eukaryotes and cyanobacteria.   | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| PSI_PsaH                    | Photosystem I reaction centre subunit VI. Photosystem I (PSI) is an integral membrane protein complex that uses light energy to reduce electron acceptors.  | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| MEP                         | Manganese stabilizing protein; photosystem II polypeptide. This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.  | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| PSI_PSAK                    | Photosystem I psaf6_psaK. This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.   | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| NAD_binding_1               | Oxidoreductase NAD-binding domain; Similar photosynthetic, but also bind FAD/NAD, have essentially no similarity. This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII. | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| AAA                         | ATPase family associated with various cellular activities (AAA). This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.  | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| PSII_PsbL                   | Photosystem II protein L (PsbL). This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.  | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| PSII_PsbM                   | Photosystem II reaction centre S protein (PsbM). This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.  | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| Fructose_1,6_bisphosphatase | Fructose 1,6-bisphosphatase. This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.  | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| Copper_binding              | Copper binding proteins; photosystem antenna family. This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.  | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| PsbQ                        | Oxygen evolving enhancer protein I (PsbQ). This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.  | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| PsdB                        | Photosystem I reaction centre subunit VI (PsbQ or PSI-N). This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.   | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| PsaI                        | Photosystem I reaction centre subunit XI. This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.   | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |
| PSI_PsaF                    | Photosystem I reaction centre subunit III. This family consists of the 33 kDa, photosystem II stabilizing protein that is encoded by the nuclear genome. The protein is also known as the manganese stabilizing protein as it is associated with the manganese stabilizing protein (MSP) and provides the structural support for PSII.  | trn13997.1.S1_at | trn13997.1 | Arabidopsis thaliana |

Arabidopsis *PSAD1* (*at4g02770*)

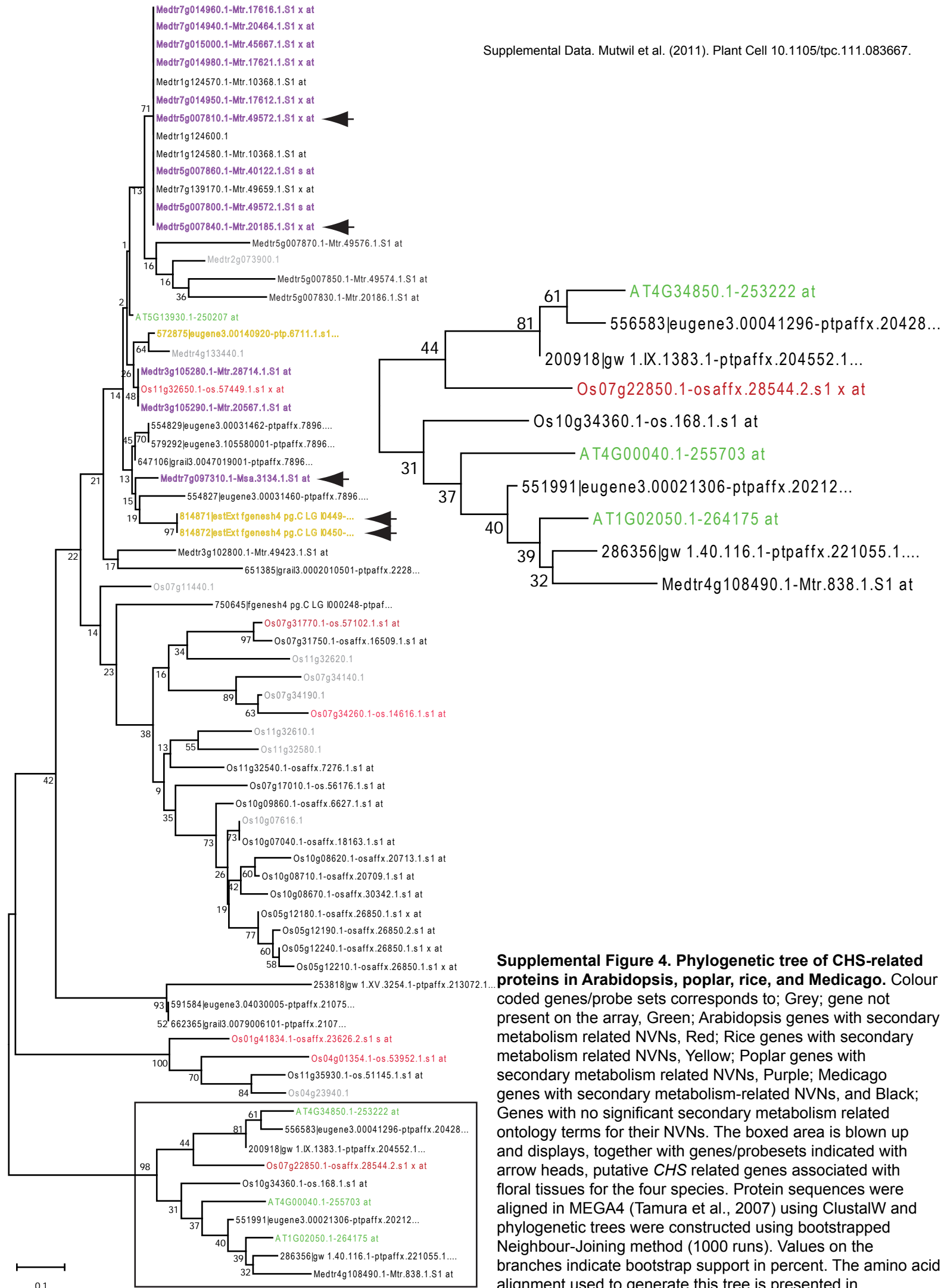


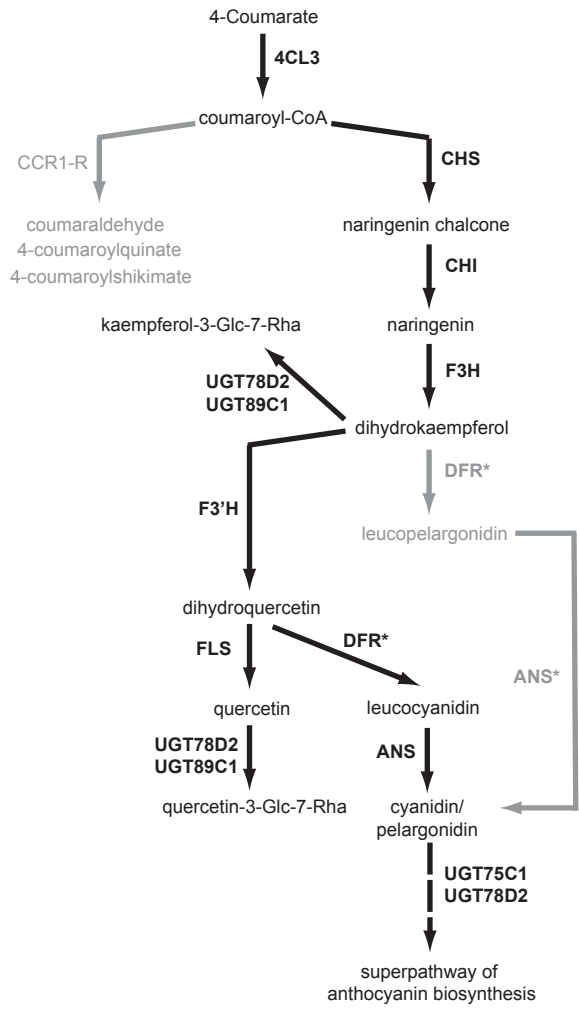
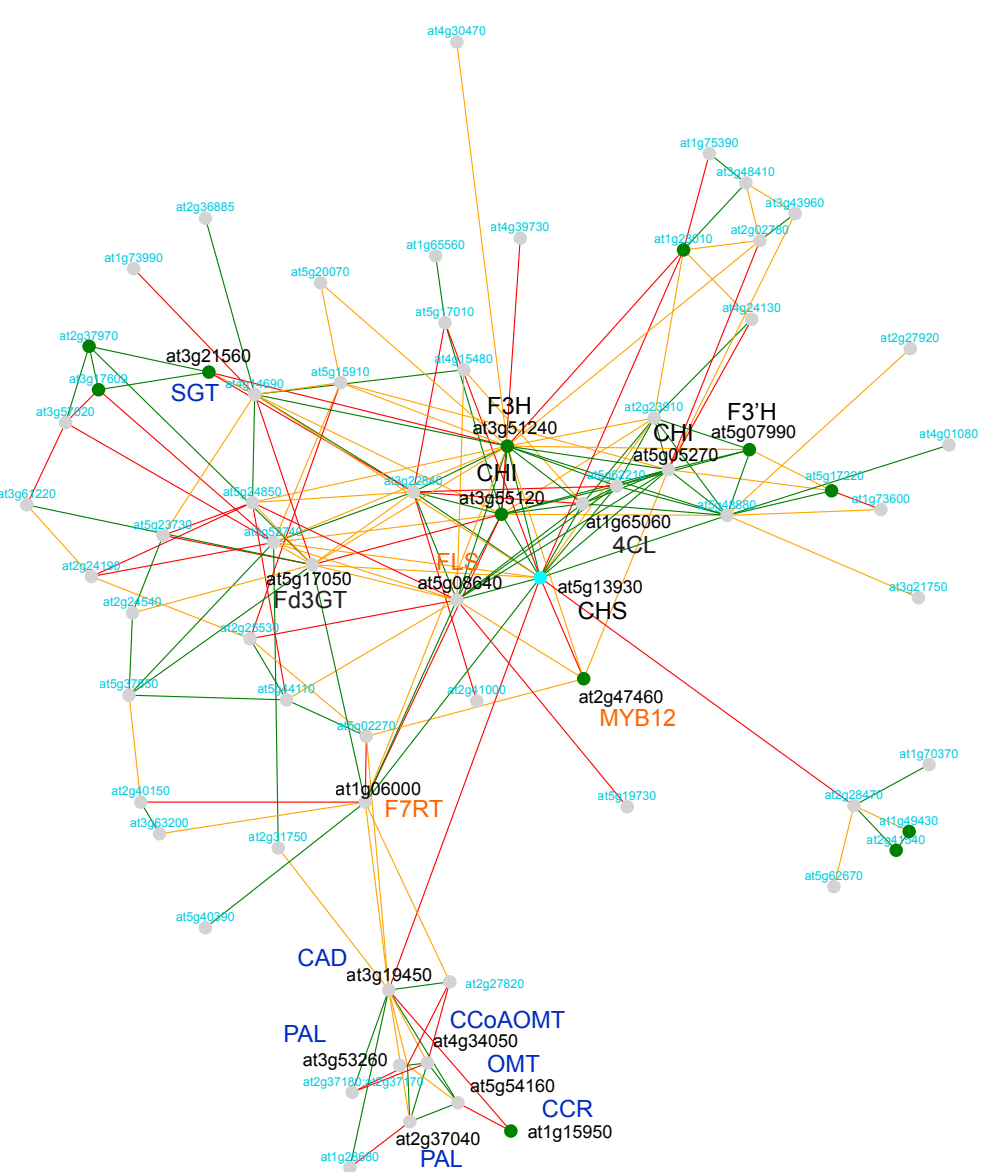
Rice *PSAD* (*Os03g09220*)



**Supplemental Figure 3. Co-expressed gene vicinity networks for *PSA-D* genes in Arabidopsis (left) and Rice (right).** Nodes in the network resemble individual genes and the connecting edges represent co-expressed links. The colouration of nodes and edges are explained in Figure 2. Examples of genes common to the two NVNs are indicated in blue colour with the gene acronyms (see also Supplemental Dataset 2). The expression of the *PSA-D*-related genes across different tissues is displayed above the networks. The interactive networks may be found at <http://aranet.mpimp-golm.mpg.de/aranet/a11983>, and <http://aranet.mpimp-golm.mpg.de/ricenet/r1587>, respectively.

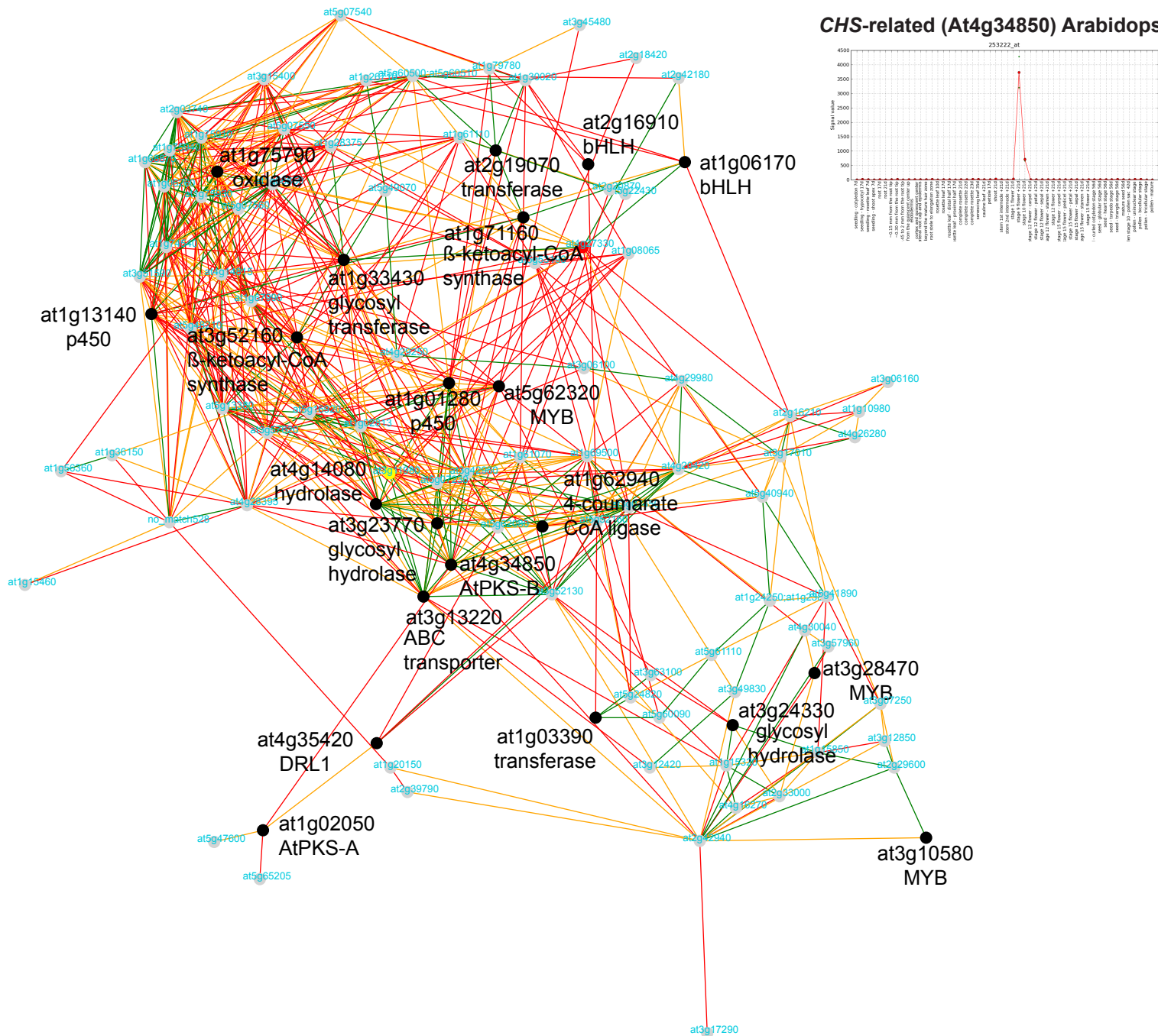






**Supplemental Figure 5. Co-expressed gene vicinity network for Chalcone Synthase (*CHS*; *At5g13930*) in *Arabidopsis*.** Many of the genes co-expressed in this NVN (left) participate in the flavonol (orange font), flavonoid (black font) and phenylpropanoid (blue font) pathway (right) leading to anthocyanin production. \*The alternative pathway in gray is observed in F3'H (TT7, *At5g07990*) knockout mutant. The interactive network may be found at <http://aranet.mpimp-golm.mpg.de/aranet/a17217>.

### CHS-related (*At4g34850*) Arabidopsis



### Supplemental Figure 6. NVN for *CHS*-related gene (*At4g34850*) in Arabidopsis.

Nodes in the network resemble individual genes and the connecting edges represent co-expressed links. The coloration of nodes and edges are explained in Figure 2. Expression of the *CHS*-related gene across different tissues is displayed upper left. The interactive network may be found at <http://aranet.mpimp-golm.mpg.de/aranet/a14232>.









**Supplemental Figure 9. Pseudocode for Column Selection.**

**Input:** matrix  $A_{m \times n}$ , with biological experiments as columns, and a range  $[t_1, t_2]$  for the number of columns to be considered.

**Output:** matrix  $C_{m \times k_{opt}}$ , with columns of A

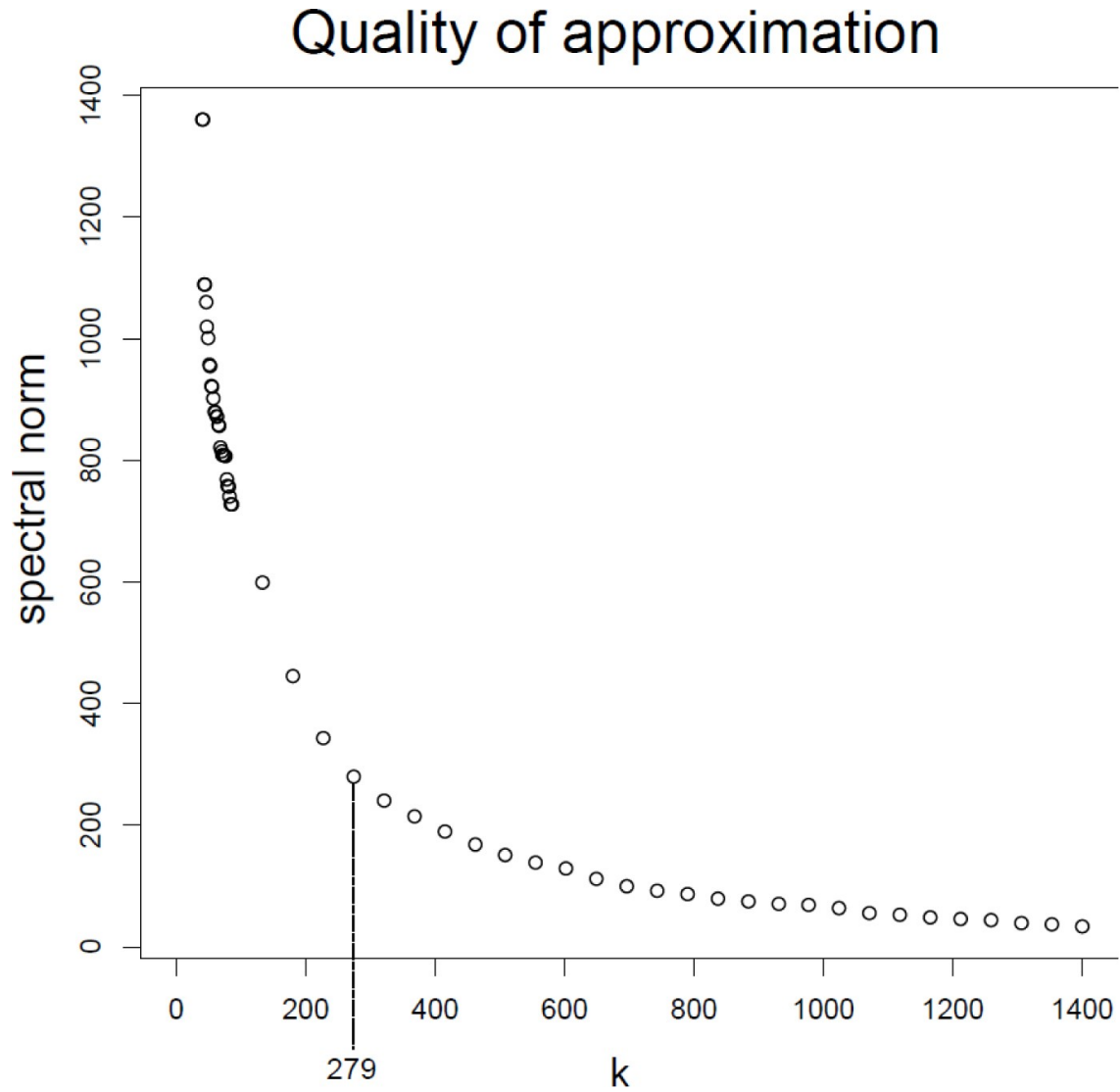
```

1. for k in  $t_1$  to  $t_2$  do
  //Randomized step
  //randomly select  $c = k \cdot \log(k)$  columns of A and form a matrix  $T_{m \times c}$ //
2.   for l in 1 to 40 do
3.     for j in 1 to n do
4.       assign a score to the jth column of A:  $p_j = \|A_i\|_2^2 / \|A\|_F^2$ 
5.       select the jth column with probability  $\min(1, c \cdot p_j)$ 
6.       if the jth column is chosen
7.         keep a rescaling factor  $\sqrt{1/\min(1, c \cdot p_j)}$ 
8.     end for
9.     keep the c columns of minimum spectral norm
  //Deterministic step
  //deterministically select k columns from T
  //form the matrix  $C_m \times k$  using RRQR decomposition with pivoting.
10.   $C(k) = RRQR(T)$ 
11. end for
12. apply elbow criterion on the spectral norms of  $C(k)$  to determine  $k_{opt}$ 
13. return  $C(k_{opt})$ 

```

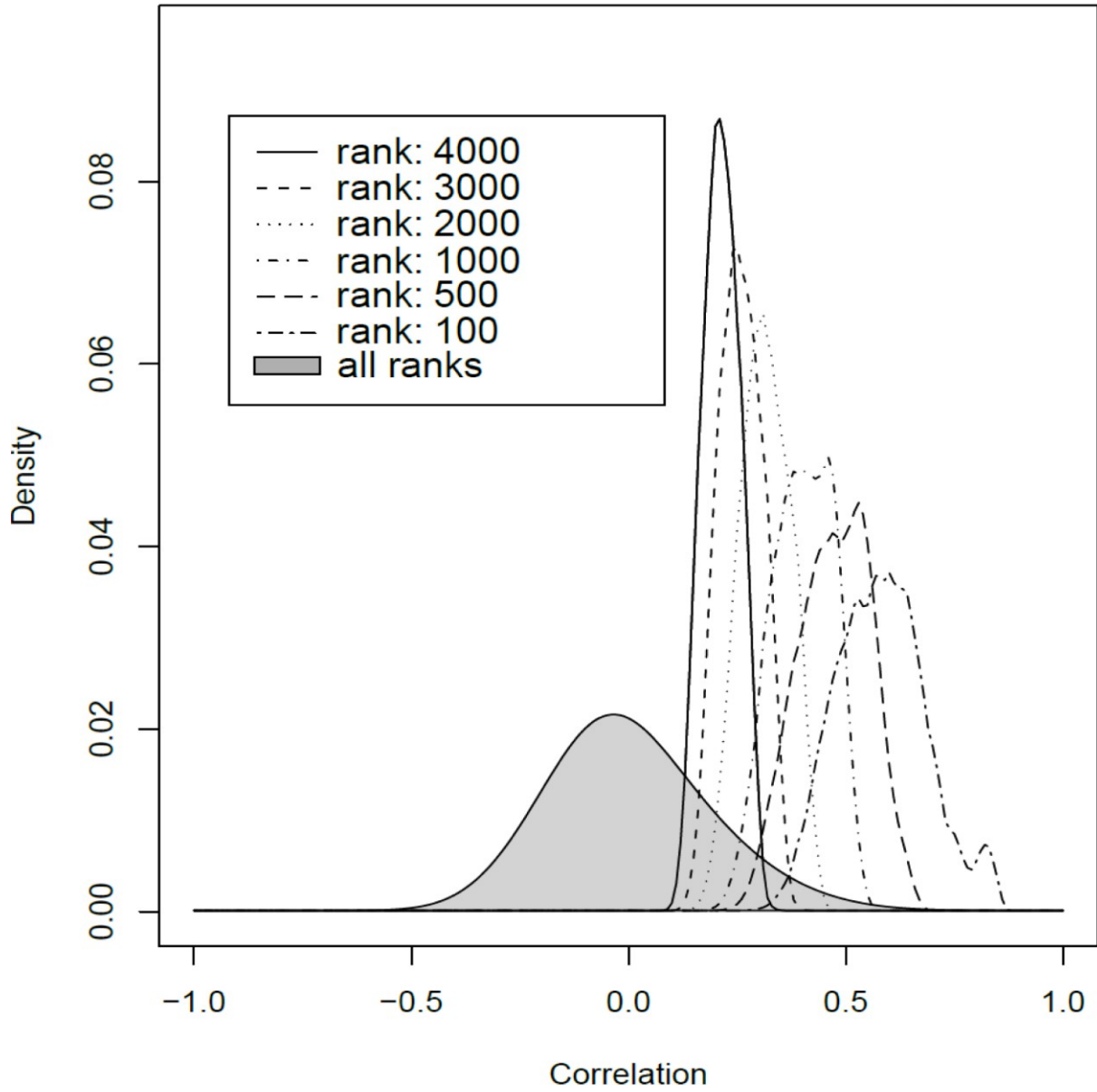
In the first step, for a given matrix A and a parameter k, one performs randomized sampling of size  $k \cdot \log(k)$  columns with probabilities given by  $p_j = \|A_i\|_2^2 / \|A\|_F^2$ . After repeating the randomized sampling 40 times, suggested and tested in (Boutsidis et al. (2008)), one selects those c columns which minimize the spectral norm of the submatrix T. In the second step, we rely on the deterministic rank revealing QR (RRQR) decomposition. Given a matrix A, its QR-decomposition is a representation of A as  $Q \cdot R$ , where R is an upper triangular matrix (all entries below the principle diagonal are zero) and Q is an orthogonal matrix (i.e., a matrix satisfying  $Q \cdot Q^T = I$  with  $Q^T$  being the transpose of Q and I is the identity matrix). We applied the RRQR decomposition to find the k columns, from the c preselected in T, which are most mutually independent. Note that in our computations, we employ the RRQR decomposition with pivoting, as implemented in LAPACK (Anderson et al., 1999). The randomized step can be avoided to arrive at the column selection problem of Golub (1965) with lower bounds on the spectral norm, still acceptable for applications. Due to the large size of the analyzed data compendia, here we rely on the deterministic step in order to avoid calculations of singular value decomposition necessary for establishing the probabilities for selecting columns in the randomized step. Here we provide an extension to this method to allow automatic selection of the number of most mutually independent columns. To this end, we first determine the minimum spectral norm for a number k of columns in the range  $[t_1, t_2]$  from A. thaliana data compendium,

with  $t_1 = 20$  and  $t_2 = 1400$ . The corresponding matrix will be termed optimal with respect to the spectral norm. We define the optimum  $k$  as that which appears at the elbow of the curve, determined by the number of columns and the spectral norm, above which only a negligible improvement is observed. By inspecting this curve for *A. thaliana* data compendium (Supplemental Figure 10), we observe that  $k_{opt} = 279$ , and its corresponding spectral norm is 15% of the spectral norm of the optimal 20-column submatrix. A computational speed-up can be achieved by testing from right to left and using previous results for  $k$  independent columns. To provide fair comparison with the data compendia from the other species used in the analysis, we rely on the values obtained from applying our method on *A. thaliana* data compendium: We seek that  $k_{opt}$  in another species, whose corresponding spectral norm is also 15% of its corresponding optimal 20-column submatrix. This resulted in the following values for  $k_{opt} = 181, 156, 163, 83, 165, 171$ , corresponding to barley, rice, Medicago, poplar, wheat and soybean.

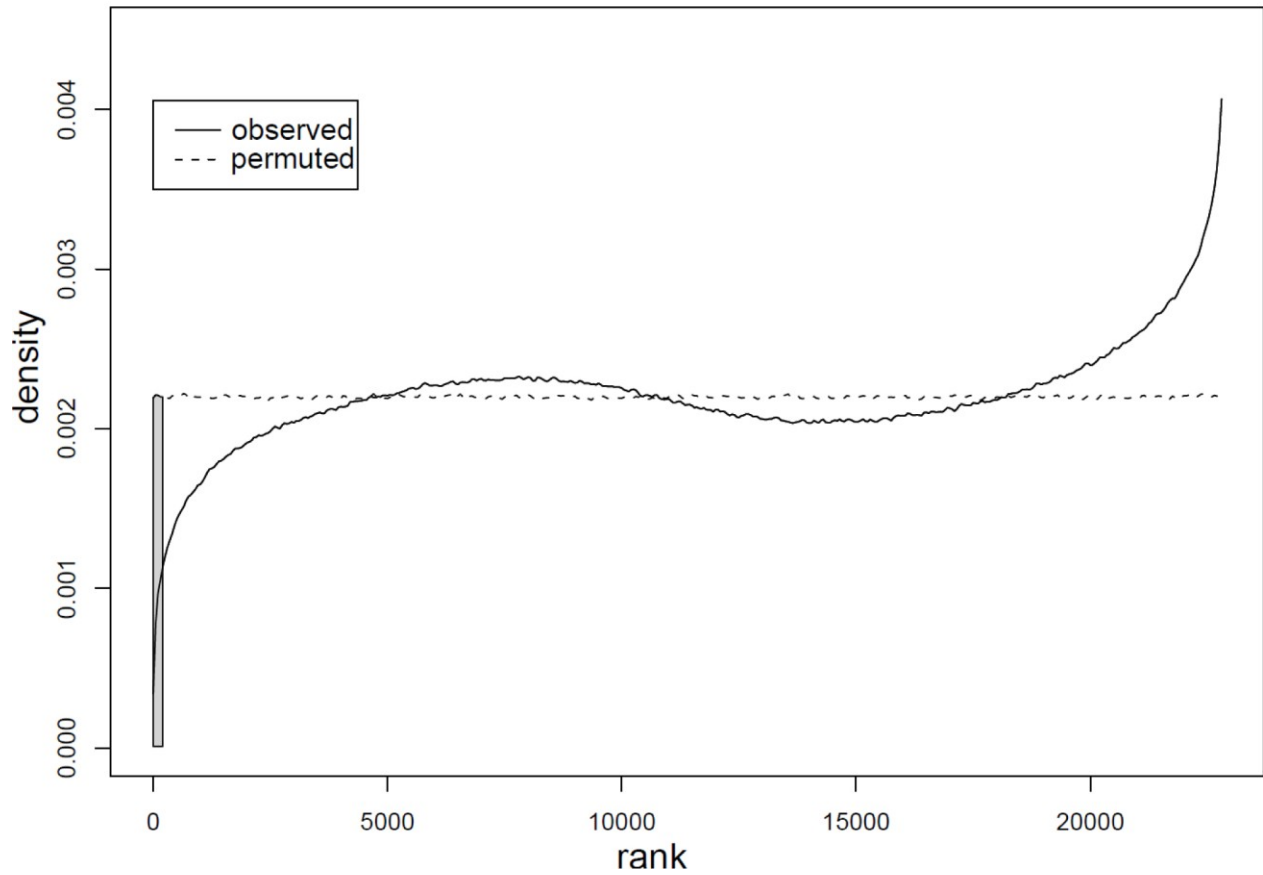


**Supplemental Figure 10. Spectral norm of optimal k-column submatrices for Arabidopsis in the range of 20 to 1400.**

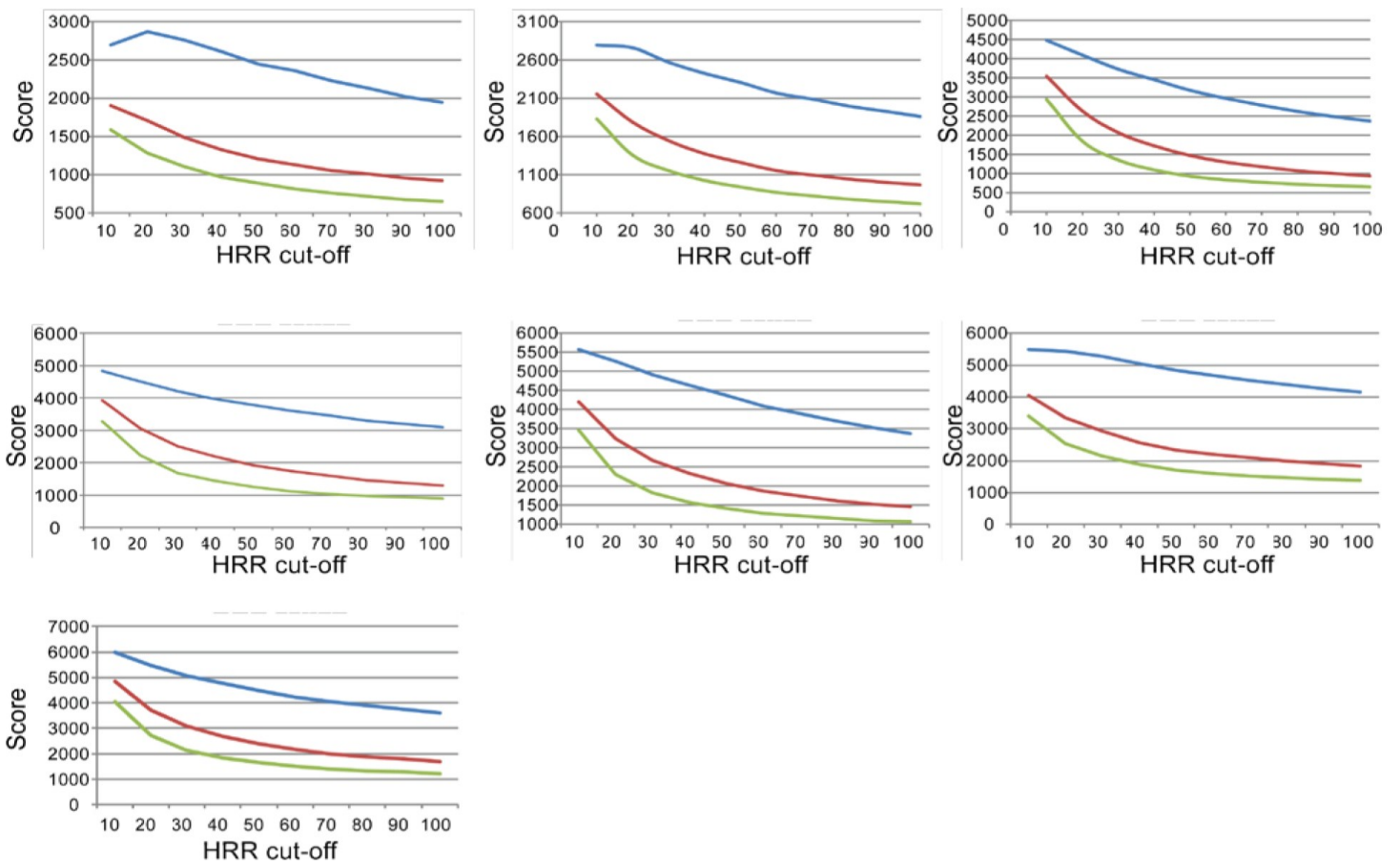




**Supplemental Figure 11. Distribution of correlation coefficients observed for specific ranks on the Arabidopsis data collection.**

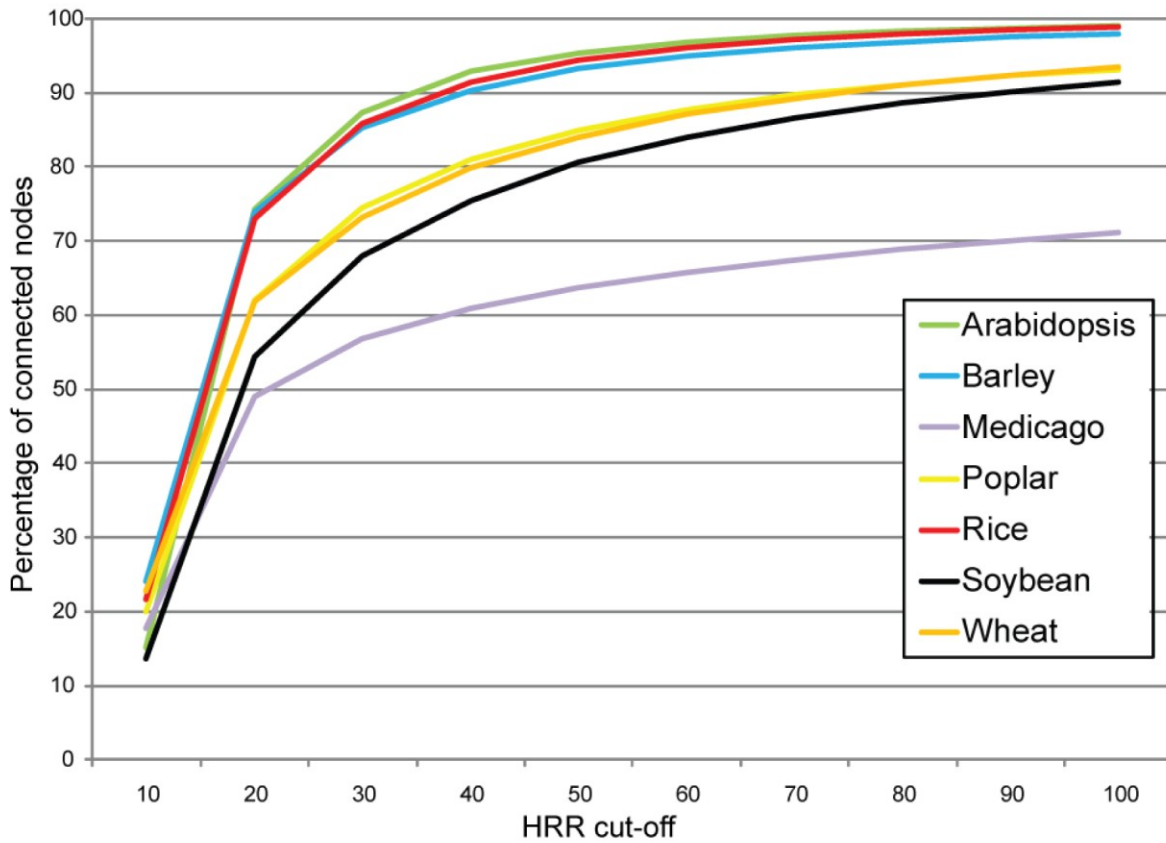


**Supplemental Figure 12. Figure 4: Distribution of reciprocal ranks on the Arabidopsis data collection.** Permuted ranks are distributed uniformly. The area under the curve shaded in grey determines 1% of the overall area of the permuted ranks and coincides with a rank cut-off of 228.



**Supplemental Figure 13. Influence of HRR cut-off and NVN step size on the biological relevance for the seven analyzed species.** The x-axis represents HRR value cut-off which was used to construct a co-expression network, while the y-axis represents the score. Blue, red and green lines represent step sizes of 1,2 and 3 used to generate NVNs.





**Supplemental Figure 14. Influence of HRR cut-off on the percentage of connected nodes in the networks of the seven analyzed species.** The different species are color coded, as shown in the legend.