

1 Ranking genome-wide correlation measurements 2 improves microarray and RNA-seq based global and 3 targeted co-expression networks

4 Franziska Liesecke¹, Dimitri Daudu¹, Rodolphe Dugé de Bernonville¹, Sébastien
5 Besseau¹, Marc Clastre¹, Vincent Courdavault¹, Johan-Owen de Craene¹, Joel Crèche¹,
6 Nathalie Giglioli-Guivarc'h¹, Gaëlle Glévarec¹, Olivier Pichon¹, and Thomas Dugé de
7 Bernonville^{1,*}

8 ¹Université de Tours, EA2106 Biomolécules et Biotechnologies végétales, Tours, 37200, France

9 *corresponding author: thomas.duge@univ-tours.fr

10 ABSTRACT

Co-expression networks are essential tools to infer biological associations between gene products and predict gene annotation. Global networks can be analyzed at the transcriptome wide scale or after querying them with a set of guide genes to capture the transcriptional landscape of a given pathway in a process named Pathway Level Correlation (PLC). A critical step in network construction remains the definition of gene co-expression. In the present work, we compared how Pearson Correlation Coefficient (PCC), Spearman Correlation Coefficient (SCC), their respective ranked values (Highest Reciprocal Rank (HRR)), Mutual Information (MI) and Partial Correlations (PC) performed on global networks and PLCs. This evaluation was conducted on the model plant *Arabidopsis thaliana* using microarray and differently pre-processed RNA-seq datasets. We particularly evaluated how dataset x distance measurement combinations performed in 5 PLCs corresponding to 4 well described plant metabolic pathways (phenylpropanoid, carbohydrate, fatty acid and terpene metabolisms) and the cytokinin signaling pathway. Our present work highlights how PCC ranked with HRR is better suited for global network construction and PLC with microarray and RNA-seq data than other distance methods, especially to cluster genes in partitions similar to biological subpathways.

12 Introduction

13 Constructing global gene co-expression networks is a popular approach to highlight transcriptional relationships (edges)
14 between genes (vertices). The ‘Guilt-by-Association’ (GBA) principle supposes that genes sharing similar functions are
15 preferentially connected and aims at predicting new functions for proteins by determining how their respective encoding
16 genes are co-expressed with others using a reference dataset containing known gene functions such as the Gene Ontology
17 (GO)¹. Defining edges connecting genes remains a critical step in global co-expression network construction. Expression
18 data (microarray or RNA-seq) are used to construct expression matrices (genes x samples) and to calculate a distance or
19 a similarity for each possible gene pair. The resulting pairwise distance matrix is then thresholded to obtain an adjacency
20 matrix that discriminates relevant edges. Only edges with a distance below (or a similarity above) the set threshold are
21 considered significant and retained for network construction. The procedure is expected to remove non biologically relevant
22 gene associations while retaining the relevant ones and can be assessed with any reference dataset. Alternatively, guide gene sets
23 may be used to extract more human-readable information from large networks in a process named Pathway-Level Correlation
24 (PLC)²⁻⁶. This approach aims at capturing the best transcriptional associations of a gene set and at highlighting functional gene
25 groups such as known subpathways in this set. There are two types of approaches to determine transcriptional associations
26 of genes: those that are supervised and those that are unsupervised. Supervised approaches such as regression and machine
27 learning based methods require a prior knowledge which is used as a training dataset to recover biologically relevant gene
28 associations. The superiority of supervised methods in extracting potential physical regulatory interactions between genes has
29 been demonstrated using simulated and real *E. coli* and *S. cerevisiae* subnetworks⁷. This study has revealed that prediction
30 accuracy is higher with smaller networks and concluded that inferring genome-scale networks remains elusive unless performing
31 a feature selection step to reduce inference problem size (because of the under determined nature of current expression datasets).
32 Among the unsupervised methods, four are commonly used and have been thoroughly tested. The first approach is Mutual
33 Information (MI) which measures a statistical dependence between two variables⁷. It is based on density function estimates
34 and has been shown to perform well with non linear relationships⁸. The second approach which relies on integrating multiple
35 transcriptional associations is Partial Correlation (PC). PCs are generally calculated from multiple linear regression and include

a variable selection step⁹. PCs aim at explaining a gene's expression profile by a small number of strongly correlated genes after eliminating those less correlated that do not significantly explain this gene's expression profile. The two last methods are Correlation Coefficients (CCs), either Pearson CC (PCC) or Spearman CC (SCC), which are the classical estimators of linear transcriptional relationship among genes^{9,10}. CCs are 2-dimensional distance measurements because a CC between two genes does not take into account the expression of the remaining transcripts in the whole transcriptome. To compensate for this lack, these approaches have been improved by using ranked CCs instead of raw values. Ranking CC implies that for every gene, all CCs calculated with the N-1 remaining genes (where N is the number of genes) are ranked from 1 to N. Within a pair of genes A and B, rank(A to B) differs from rank(B to A) because the two genes display different expression profiles and different relationships with the remaining transcripts in the transcriptome. Two related ranking methods have been developed. One is mutual ranking (MR, geometric mean of the two ranks) which has been shown to improve GO term recovery with PCC using large microarray data from Arabidopsis, Human, mouse and rat¹¹. MR has been successfully used in multispecies analysis of co-expression modules¹². Another is Highest Reciprocal Ranking (HRR, maximum value of the two ranks)¹³. MR and HRR are thought to be more integrative than unranked CCs because they depend on other CC values around that of a gene pair. Although not as robust as supervised methods, unsupervised methods can efficiently capture relevant gene associations as previously shown⁸. These authors have shown that non parametric CC and MI calculations were more efficient than PCC on a small dataset. Among other unsupervised methods, SCC calculations have been similarly shown to outperform other distance measurements in Human expression data¹⁴. In this case, SCC were calculated from RNA-seq or microarray data in order to construct several smaller networks subsequently aggregated to yield the final network. We firmly believe that genome-scale networks inferred with CCs, especially when combined with a ranking procedure, are helpful to find new associations between genes. Although CCs are not efficient in detecting non linear associations⁸, gene-to-gene relationships have been predicted to be essentially linear¹⁵ suggesting that CCs are valuable distance measurements. To date, there is no clear evaluation of how ranked CCs affect genome-scale network reconstruction with RNA-seq data in comparison with other unsupervised methods. We evaluated ranked CC, raw CC, MI, and PC performance in global and targeted network construction using Arabidopsis microarray and differentially processed RNA-seq expression data (Figure 1). Performance was measured as network ability to capture biologically relevant gene associations found in a Gene Ontology (GO) annotation reference set but also to correctly cluster guide genes in PLC. Global network quality was first evaluated according to the different dataset x distance measurement combinations. The resulting global networks were next interrogated in PLC analyses with five different guide gene sets corresponding to four different metabolic pathways and one signaling pathway. Whereas metabolic pathways have relatively clearly defined and partially linear partitions, signaling pathways usually involve post transcriptional regulations and a more intricate organization, which might render gene transcriptional associations less evident. We looked at the dataset x distance measurement combinations optimizing pathway reconstruction and maximizing co-occurrence quality between microarray and RNA-seq networks. Our results show that, of the six methods evaluated, PCC ranked with HRR generated the best biologically relevant networks according to initial guide gene representation and clustering in distinct modules. In addition, it offers the possibility to merge subgraphs obtained by microarrays and RNA-seq to generate high confidence networks.

Results

Inferring global co-expression networks and comparing correlation measurements

Large co-expression networks were obtained by varying the confidence threshold (correlation value above or rank below) within lists containing the 10 million best gene pairs from eight different datasets and six data measurement combinations (Figure 1). Each of the 10 million best pair lists was filtered at different confidence thresholds (1, 5, 10, 20, 40, 60 or 80% best pairs from these lists) to evaluate the effect of network size on performance. Expression datasets included a microarray-based expression matrix and seven RNA-seq based expression matrices normalized with different methods to evaluate their effect on network inference: transcript per Million (TPM), log₂ TPM, sample scaled (ss) TPM, ss log₂ TPM, raw counts, variance stabilized transformed (VST) raw counts and VST-TPM. The six distance measurements were: raw PCC, raw SCC, PCC-HRR, SCC-HRR, PC and MI. Each network performance was considered as a network ability to capture edges corresponding to functional associations found in the GO reference dataset and was evaluated in 4 different ways (Figure 2): GO term enrichment (GO terms that are significantly enriched with gene pairs from the co-expression network), a ROC curve constructed with TPR and FPR calculated for each confidence threshold and two ROC analyses based on the GBA concept, an average 3-fold cross validated neighbor voting (NV) AUROC and a global AUROC. AUROCs correspond to Area Under Receiver Operating Characteristic curves calculated for every network either from each GO (with three test sets obtained after hiding part of the gene labels, NV AUROC corresponding to the average of AUROCs for all GO terms) or the whole annotation dataset (global AUROC). AUROCs are used as global indicators of a dataset performance, a value of 0.5 indicating a random attribution of labels in the network and a value of 1 indicating a perfect match with the reference dataset. AUROC>0.6 may be considered as moderate¹⁴. In global TPR vs FPR curves, the line extending from (0,0) to (1,1) has an AUROC=0.5 and points above this line indicate more predictive networks than a random selection (Figure 2). The GO annotation table was filtered to perform these

90 analyses by removing weakly represented or non-specific GO terms (>5 or <100 genes).

91 Figure 3 displays TPM network evaluation at different confidence thresholds and Figure 4 shows networks having 1 million
92 of edges across all dataset x distance combinations. Metrics for all other dataset x distance measurement combinations are
93 presented in Supplementary Figure 1 online. All networks combined, pairwise correlations between enriched GO counts, global
94 and NV AUROC performance metrics were moderate (Spearman's $\rho > 0.4$) but significant ($p < 0.001$) indicating these three
95 performance metrics evaluated networks in different ways. The highest correlation was observed between NV AUROC and
96 enriched GO counts ($\rho = 0.70$, $p < 0.001$) showing their consistency. The NV AUROC was the most positively correlated with
97 edge number ($\rho = 0.55$, $p < 0.001$) suggesting that decreasing the confidence threshold and adding more edges in networks did
98 not result in a significant increase in false positives. This was confirmed by the partial ROC curves (obtained for a maximum
99 FPR at 10 million edges) drawn from the TPR and FPR (Figure 3, Supplementary Figure 1 online), where up to 10 million best
100 pairs, TPR increased faster than FPR. Although counts of significantly enriched GO terms were positively correlated to NV
101 AUROC, we observed a slight decline in the largest networks which might reveal a saturation in these enriched GO terms. It is
102 possible that with the hypergeometric testing, some GO classes are fully enriched in smaller networks leading to a decrease
103 in their significance as network size increases. The global AUROC displayed a very low variation (min=0.55, average=0.61,
104 max=0.68) and was significantly correlated to vertex number ($\rho = 0.43$, $p < 0.001$) only. This observation suggests that the
105 global AUROC is not an appropriated measure in our case.

106 At equivalent edge numbers, different distance measurements generated networks varying considerably in vertex number
107 (Figure 4A, Supplementary Figure 1 online). Considering all datasets and distance measurements, raw PCC, raw SCC and
108 MI resulted on average in fewer vertices and higher node degree (vertex number/node degree: 13,164/511, 9,986/465 and
109 14,074/468 respectively) than PCC-HRR, SCC-HRR or PC (26,645/116, 24,731/124 and 23,966/166 respectively). This
110 trend was clearly observed when setting an edge number to 1 million (Figure 4A). Expression networks constructed from
111 microarrays, TPM, TPM log2, and counts displayed very similar ROC curves: PC based networks followed random predictions
112 (NV AUROC=0.5) and the other distance measurements were above the random prediction with similar AUC (Supplementary
113 Figure 1 online). This was confirmed for PC by NV AUROC and enriched GO term counts. Performance of the other distance
114 measurements in the global TPR/FPR curves did not exactly match that measured with AUROCs. Taking the TPM dataset as
115 an illustration (Figure 3), the MI ROC curve was above the others while NV AUROC for similar edge numbers was slightly
116 below that measured for SCC. This was probably due to differences in network topologies (see above) and the procedures
117 underpinning the two evaluations. The global TPR/FPR curve does not measure a network predictability *per se* as NV AUROC
118 does and considering any gene pair sharing a same GO term as valid could have overestimated TP (Figure 2). As a general trend,
119 raw PCC and raw SCC generated smaller networks than PCC-HRR and SCC-HRR but displayed similar TPR/FPR curves, *i.e.*
120 for a similar performance, HRR-ranked CC networks had more vertices and fewer edges than raw CC based networks (Figure 3).
121 CC ranked with HRR always generated relevant networks for TPM ss, TPM log2 ss, TPM VST and counts VST, which was not
122 the case for raw CC (Supplementary Figure 1 online). These normalizations induced strong biases in CC distribution as revealed
123 by thresholds used to obtain the 10 million best pairs (Supplementary Table 1 online) but these biases were compensated by
124 HRR. Taken together, these results revealed that HRR CCs are able to generate complete genome-wide networks with good
125 performances similar to other classical measures such a MI and PC. Node degree AUROC measures whether genes are more
126 likely associated according to their number of connections rather than to their function. A positive correlation was found
127 between NV AUROC and degree AUROC ($\rho = 0.47$, $p < 2e-16$) indicating that highly predictive networks (NV AUROC > 0.7)
128 also had a higher node degree AUROC. Node degree AUROC was generally under 0.55. We therefore considered that in our
129 conditions, this bias was only limited. Concerning edge co-occurrence between the different dataset x distance combinations,
130 the lowest conservation was observed with raw (MI, PCC and SCC) RNA-seq datasets and PC networks and microarrays
131 networks (Figure 4B, area 1). More co-occurring edges were found when microarray networks were compared to RNA-seq
132 networks obtained with CC-HRR (mean of 97,646 vs 25,277; Figure 4B, area2). This indicated that microarrays and RNA-seq
133 networks were more comparable when obtained with HRR, reinforcing their validity. The previous section focused on global
134 network properties. Community detection procedures can be applied to such global networks to cluster tightly connected genes
135 into modules. In our case, we rather used a knowledge-driven approach known as Pathway-Level Correlation (PLC) to extract
136 gene pairs associated within a given pathway (Supplementary Figure 2 online). PLC are particularly interesting in plants for
137 example to decipher incomplete specialized metabolic pathways. It aims at capturing a transcriptional landscape for genes
138 known to be involved in a given pathway, in order to highlight their organization as well as finding new genes (transporters,
139 transcription factors,...) associated with the process. In the next part, we evaluated the ability of all previous networks to capture
140 relevant information associated with four metabolic and one signaling pathways. We selected two primary metabolic pathways
141 (carbohydrate and fatty acid metabolisms), two specialized (secondary) pathways (phenylpropanoid and terpenoid metabolisms)
142 and the cytokinin signaling pathway.

Assessing PLC quality: trade-off between GO term representation and guide genes

The PLC procedure is expected to cluster together guide genes with many co-expressed genes ('associated genes') and to reflect the subpathway organization (Figure 5A). For PLC, we systematically removed all genes showing a degree value of 1 (*i.e.*, those connected to only one guide gene). However we included edges between associated genes if they were found among edges retained at the selected threshold. Using five pathways (Table 1, Figure 5B, Supplementary Table 2 and Supplementary Figure 3 online), we extracted five PLC from the global networks generated above to determine the best suitable dataset x distance measurement combinations. All pathways have modular structures with gene sets forming specific sub-pathways (also called partitions or modules). We expected that PLC would be able to reconstruct such a partitioning, by connecting guide genes with associated genes. The phenylpropanoid pathway contains a core module composed of 3 genes leading to a precursor used by 3 other distinct subpathways¹⁶⁻¹⁹ (Figure 5B, Table 1, Supplementary Table 2 online). The three other metabolic pathways, carbohydrates, fatty acids and terpenoids, were structured in modules as described on the KEGG database²⁰ (Table 1, Supplementary Table 2 online). The fatty acid pathway contains 97 genes divided into 6 modules. The central carbohydrate metabolism contains 202 genes partitioned in 8 modules. Finally, the terpene pathway has 64 genes partitioned into 6 modules. Pathway organizations were used as indicated in the KEGG database (apart from phenylpropanoid pathway which was manually curated from our previous work) and compared to PLC subnetworks. The plant cytokinin (CK) pathway is known to regulate many processes in plant physiology and is hierarchically organized in three levels: a histidine kinase receptor, a transducer (histidine phosphotransfer proteins) and a response regulator (type A/B/C) which may act as a transcription factor²¹ (Table 1, Supplementary Table 2 online). Although CK pathway members are relatively well known, each level is represented by several members which may have specific roles and it is still unclear how they biologically interact with each other to drive a specific physiological response. We expected that PLC would group some of these actors according to specific physiological responses. CK pathway includes both transcription activating and repressing activities (via response regulators) and post-transcriptional (phosphorylations) and would therefore be an excellent test of PLC applicability on associations expected to be more complex than in metabolic pathways. In addition, we included other histidine kinases integrating other signals and known to crosstalk with the CK pathway²². We therefore included 2 ethylene receptors, ETR1 and ERS1 to determine whether they could be clustered with CK histidine kinase. The initial pathway was not partitioned into sub-pathways but rather into 5 levels (receptor, transducer, type A/B/C response regulator) because interactions between specific actors of each level are not completely understood.

Pathway	Genes	Number of subpathways	Subpathway names (KEGG module accession)
Phenylpropanoids	43	4	core phenylpropanoid (PP), flavonoids, monolignols, phenolamides
Fatty acid	97	6	fatty acid biosynthesis (initiation (M00082), elongation (M00083), its ER-localized part (M00415)), jasmonic acid phytohormone biosynthesis (M00113) and β -oxidation (M00086 and M00087)
Carbohydrate	202	8	glycolysis (Embden-Meyerhof pathway (M00001) and the core module involving three-carbon compounds (M00002)), neoglucogenesis (M00003), pyruvate oxidation (M00307), citrate cycle (M00010), pentose phosphate pathway (M00004, M00006 and M00007)
Terpenes	64	6	mevalonate (M00095), methylerythritol (M00096), C10-C20 isoprenoid (M00366), beta-carotene (M00097), abscisic acid hormone (M00372) and phytosterol (M00371) biosynthetic blocks
Cytokinin signaling	37	?	?

Table 1. Pathway description. ? Indicates that partition in sub-pathway is not known.

Subgraphs of global networks were constructed for each pathway by retrieving edges involving at least one guide gene and were partitioned into communities with a fast greedy algorithm designed to maximize network modularity and which has been shown to extract relevant communities from large networks²³. We compared guide gene distribution in these communities to target subpathways using a normalized χ^2 test which values range from 0 to 1, 1 being the expected partition and 0 a random partition of guide genes or very few guide genes (Figure 5A). All networks having a χ^2 p -value > 0.05 were considered to have a χ^2 statistic equal to 0. PLC performance in recovering GO terms was evaluated by counting significantly enriched GO terms and by calculating a NV AUROC for each network. A good PLC was expected to contain a large number of guide genes and to have both a good score in grouping them into expected partitions (high normalized χ^2 value) and a good score in overall biologically relevant edge recovery (NV AUROC > 0.6). We first analyzed correlations between all these metrics (NV AUROC, number of guide genes and χ^2 statistic) together with two topological metrics (mean node degree and

modularity), for each pathway separately (Figure 5B). Strongest correlations were observed between NV AUROC and mean node degree ($\rho > 0.5$, $p < 0.001$) and between modularity and normalized χ^2 ($\rho > 0.59$, $p < 0.001$). We found that PLC performance (NV AUROC) was almost negatively correlated with normalized χ^2 ($\rho < -0.2$) indicating that guide genes were clustered correctly at the expense of capturing GO associated gene pairs. Given the CK pathway structure, partitioning based on protein functions (receptor, transducer or response regulator) did not result in high χ^2 values, suggesting that partitions in the co-expression networks contained guide genes from different levels, reinforcing the existence of specific sub-pathways. These results indicated a trade-off in PLC between edge quality and guide gene partitioning. A visual examination of PLC with either lower modularity and higher NV AUROC (Figure 5D) or higher modularity and lower NV AUROC (Figure 5E) revealed that PLC with higher modularity as well as higher χ^2 values displayed a biologically relevant organization. Such subgraphs had generally a lower average node degree and a higher representation of guide genes rendering their analysis more convenient. Taking the phenylpropanoid pathway as an example, the PCC-HRR based TPM network (Figure 5E, with a higher modularity) correctly clustered genes from the core phenylpropanoid (PP) and the flavonoid modules while the raw PCC network did not (Figure 5D, with a higher NV AUROC). Similar results were observed with the four other pathways with either microarray or RNA-seq datasets (Supplementary Figure 3 online). Modularity and normalized χ^2 could therefore be considered as consistent quality metrics for PLC. NV AUROC should also be considered to ensure that subgraphs had a minimum predictability (> 0.55).

196 HRR-CCs optimize recovery and clustering of guide genes in PLC

197 The best performing dataset x distance measurement combinations were searched by analyzing NV AUROC, modularity and
198 normalized χ^2 among networks with a χ^2 $p < 0.05$. Statistical effects of dataset, distance, ranking and their interactions
199 on subgraph characteristics were analyzed by ANOVA for each pathway. Ranking and distance measurements had generally
200 the strongest effects on modularity and normalized χ^2 ($p < 2e-5$) (Figure 6A & B). Ranking had a significant effect on NV
201 AUROC ($p < 0.01$) but was weaker than distance measurement ($p < 1e-4$). Datasets only had a significant effect on modularity
202 ($p < 0.002$). Significant interactions were rarely observed between these three factors (*i.e.* in few pathways and with a weak
203 effect). This revealed that the different RNA-seq normalizations had only minor effects on these PLCs. Taken as a whole,
204 networks obtained with raw datasets had a significant higher NV AUROC (t-test, mean in raw=0.58, mean in HRR=0.57,
205 $p < 0.01$) but significant lower modularity (mean in raw=0.35, mean in HRR=0.68, $p < 2.2e-16$) and lower normalized χ^2 value
206 (mean in raw=0.20, mean in HRR=0.35, $p < 2.2e-16$) (Figure 6A). It therefore appeared that clustering guide genes correctly
207 was improved with CC ranked with HRR at the expense of performance. NV AUROCs in HRR-based networks were generally
208 higher than 0.55, indicating an average low performance in GO capture (Figure 6A). In non-ranked distances, PC resulted in the
209 weakest NV AUROC, while MI and raw SCC based networks displayed the highest NV AUROC (Figure 6B). This weakness in
210 PC based networks was compensated neither by a higher modularity nor by a higher normalized χ^2 statistic.

211 A more detailed examination of best PLC subgraphs maximizing either modularity or NV AUROC, revealed that each of
212 the five pathways involved specific dataset x distance measurement combinations. PCC-HRR based networks were always
213 found to maximize modularity (Figure 6C) and normalized χ^2 (Figure 6D) with almost all datasets. Raw distance based PLCs
214 had a higher NV AUROC and some of them also had a good modularity but they also had a lower normalized χ^2 statistic
215 indicating they contained fewer guide genes (*e.g.* raw RNA-seq counts with raw SCC in the terpene PLC). The results suggest
216 that PCC-HRR could be used as a reliable distance measurement whatever the dataset. Careful analysis of PLC obtained from
217 PCC-HRR revealed the presence of relevant associations in each PLC (Supplementary Figure 3 and Supplementary Table 3
218 online). For example, community 12 from the phenylpropanoid PLC obtained with microarray data processes with PCC-HRR
219 (Figure 5D) contained AT1G06000 encoding a Flavonol 7-O-rhamnosyltransferase and was clearly associated with other genes
220 from the flavonoid sub-pathway. This gene was not detected in the raw PCC PLC (Figure 5C). Other examples are highlighted
221 in yellow in Supplementary Table 3 online.

222 Vertex and edge co-occurrence in microarray and RNA-seq based PLC subgraphs

223 Edge co-occurrence in networks constructed from expression datasets obtained by different technologies may be considered as
224 a further validation. Quantifying gene expression with microarrays relies on probe hybridization by sequence complementary
225 while with RNA-seq, short reads are mapped back *in silico* to the reference transcriptome. The two main differences between
226 these technologies are (i) the number of quantified transcripts (due to the completion of genome annotation) and (ii) the dynamic
227 range (fluorescent probe intensities for microarrays, *in silico* read counts for RNA-seq). Because microarrays and RNA-seq
228 technologies differ, edges co-occurring in networks obtained from these two technologies are probably more relevant. In
229 Figure 4C, we analyzed co-occurrence in global networks and found that HRR ranked CCs apparently increased the number
230 of co-occurring edges between microarrays and RNA-seq. To get more insights into co-occurrence in PLCs, common edges
231 and vertices were counted in pairwise intersections of networks (RNA-seq vs microarrays) obtained with the six distance
232 measurements and set at a 1,000 vertices. The resulting intersection networks were further characterized by the number of
233 represented guide genes, their normalized χ^2 statistic, modularity and NV AUROC. This evaluation was performed with the

234 RNA-seq dataset expressed as TPM only because we showed in the previous section that normalization methods had a minor
235 impact on PLC. In addition, TPM networks with raw distance methods had enough vertices to correctly extract PLC (it was not
236 the case with raw distances, e.g. for TPM normalized with VST as revealed by their very low normalized χ^2 statistics; Figure
237 6A).

238 Many more co-occurring edges were generally recovered when raw CC and MI networks were compared (e.g. 18,334
239 averaged over the five pathways with MI networks vs 550 with PCC-HRR networks; Supplementary Figure 4 online). At a
240 1,000 vertices, all raw networks but PC contained more edges (221,297 and 85,059 in average for microarrays and TPM) than
241 HRR-CCs networks (12,431 and 12,877). This might have resulted in more co-occurrences between MI networks. PC networks
242 had the lowest number of co-occurring vertices (94 in average) but intersections from MI and/or raw CC had comparable vertex
243 number (268) to intersection networks from CC-HRR (267 in average) (Supplementary Figure 4 online). These results suggest
244 that HRR-based networks have strong overlaps. Intersections of PCC-HRR subgraphs were able to maximize the % of guide
245 genes (mean of 75% over the 5 PLC), modularity (0.78) and normalized χ^2 statistic (0.70) (Figure 7). Detailed characteristics
246 for each PLC are presented in Supplementary Figure 4 online. Modularity was generally high in the intersection between
247 CC-HRR networks (>0.70) but intersections with SCC-HRR displayed lower normalized χ^2 values (<0.6). Intersection
248 network performance in recovering GO terms was globally low (Figure 7D). The highest NV AUROCs were observed in
249 intersections between MI networks (0.52), MI (microarrays) – raw SCC (TPM)(0.54) and raw PCC (microarrays) – raw SCC
250 (TPM)(0.52) (Figure 7D). Intersection networks and their contents are available in Supplementary Figure 5 and Supplementary
251 Table 4 online. Again, we found candidate genes not included in the guide gene sets that were correctly associated with other
252 guide genes (highlighted in yellow in Supplementary Table 4 online). Taking the phenylpropanoid pathway as an example,
253 Figure 7E shows edge and vertex co-occurrence between MI networks and Figure 7F between PCC-HRR networks. The
254 co-occurrence network obtained from MI contained fewer guide genes (26 vs 41) and displayed lower modularity (0.49 vs
255 0.67) and normalized χ^2 statistic (0.39 vs 0.66). Although it had a higher NV AUROC (0.6 vs 0.46), its structure did not
256 reflect that of the expected pathway (Figure 5C). For example, phenolamide related genes were not represented. Average guide
257 gene degree (33) was below the average degree of the remaining nodes (100) indicating that guide genes were only slightly
258 connected to other genes in this co-occurrence network from MI PLC. By contrast, guide gene degree (11.4) was very similar
259 to the other node degree (11.1) revealing a uniform integration of guide genes with other genes in the co-occurrence network
260 of PCC-HRR PLCs. As observed in co-occurrence in large networks (Figure 4B), RNA-seq TPM normalized with VST had
261 slightly more edges in common with microarray networks. We therefore compared PCC-HRR PLC between microarrays
262 and RNA-seq TPM normalized with VST. Intersection networks had very similar characteristics to that observed between
263 microarrays and RNA-seq TPM. Although it contained slightly more co-occurring vertices and edges in average (360 and 1,252
264 respectively with TPM VST vs 240 and 550 with TPM), it displayed fewer guide genes (54 vs 57). TPM normalized with
265 VST could therefore be an interesting alternative to TPM. PLC intersection networks and their description are available in
266 Supplementary Figure 6 and Supplementary Table 5 online.

267 Discussion

268 Pathway Level-Correlation (PLC) is an interesting approach to capture biologically relevant transcriptional relationships using
269 guide genes (e.g. genes involved in a same metabolic pathway) from transcriptome-wide co-expression networks. Our present
270 work highlights that distances between genes calculated with highest reciprocally ranked PCC (PCC-HRR) improve PLC. The
271 main improvement was guide gene representation. PCC-HRR based PLCs contained more guide genes than observed with
272 other distances and they were generally more correctly partitioned into expected sub-pathways in the co-expression network.
273 This was associated with a lower mean node degree and a higher modularity but also with a slightly weaker performance in
274 GO term recovery. Our results propose that modularity and normalized χ^2 values could be used as reliable indicators of
275 PLC quality. We also observed that edge and vertex co-occurrences in PLCs obtained with PCC-HRR and microarray and
276 RNA-seq TPM data can be used to construct relevant networks. A surprising observation was that in our conditions, for most
277 combinations tested, true positive rates remained higher than false positive rates in spite of increasing network sizes. A similar
278 trend using small *E. coli* and *S. cerevisiae* networks (<110 nodes) has been previously observed with CCs⁷. This suggests that
279 co-expression studies should test different confidence thresholds to efficiently capture gene associations. Evaluating network
280 quality was done in respect of the Arabidopsis reference GO annotation set. We found that the NV AUROC¹⁴ evaluates
281 networks efficiently and was generally in accordance with significantly enriched GO term counts and TPR vs FPR curves. NV
282 AUROC has the advantage of being a more global measure of predictability (values above 0.6 can be considered as moderate).
283 Different distance measurements displayed different efficiencies according to the dataset but as a general trend, performance of
284 the different combinations were similar (e.g. between microarrays and RNA-seq TPM in Figure 3B). The same performance
285 was obtained for different topologies: high node degree (more edges and fewer vertices) for MI and raw CC networks vs
286 lower node degree (fewer edges and more vertices) for CC-HRR networks. PC networks displayed a high performance with
287 microarray data only, suggesting that PCs calculated with ‘corpcor’ R package may not be recommended for RNA-seq data. A

288 recent study has focused on metabolic pathways in plants using mutual ranks, another CC ranking method¹². Complementary
289 to this previous work, we found that ranking CCs increases vertex number without penalizing absolute network performance.
290 Contrastingly, an opposite trend was observed in another study²⁴, where larger networks displayed a lower Matthew Coefficient
291 when compared to protein-protein interactions or regulatory networks. This indicates that different absolute performance
292 measurements lead to different results and interpretations but this might also be due to our datasets which were larger than
293 theirs. Another advantage of CC-HRR was that it clearly homogenized network characteristics from differently normalized
294 RNA-seq datasets in addition to increase the number of co-occurring edges between microarrays and RNA-seq.

295 As revealed recently²⁵, highlighting correlations between genes may require specific data processing or distance algorithms
296 best suited to their query pathway. We also found that each of the five PLCs performed best with specific RNA-seq normalizations
297 (Figure 6C & D) but RNA-seq TPM processed with PCC-HRR always provided informative networks which can be used as
298 reliable starting point because they matched well expected pathway structure. In our case, the different data normalizations
299 had a relatively weak effect on PLC characteristics especially when CCs were used with HRR. In a comparative analysis²⁴,
300 the authors have shown that PCC networks from VST normalized counts were more comparable to those from microarrays.
301 In our case, VST normalization slightly improved the overlap between RNA-seq TPM and microarrays both at the global
302 and targeted levels. This normalization can thus be further considered for co-expression studies. A fast greedy approach
303 maximizing modularity was used to detect communities within PLC subgraphs. Guide gene partitioning in these communities
304 was compared to expected partitions in subpathways with a normalized χ^2 test (Figure 5A). We found that correct guide
305 gene partitioning was negatively correlated with NV AUROCs but positively with modularity. Subnetworks with highest
306 NV AUROCs but lower modularity such as those obtained with MI represented fewer guide genes and displayed large edge
307 numbers. In these networks, guide genes formed inappropriate structures (Supplementary Figure 5 online). We applied PLC to
308 five pathways varying in size and nature. For the four metabolic pathways, PLC extracted from PCC-HRR based networks
309 were able to cluster guide genes in the proper subpathways (Supplementary Figure 5 online). Guide genes were associated
310 in communities resembling subpathways and containing genes not included in the query gene set but known to be involved
311 with the given pathway or being good candidates to be functionally validated (Supplementary Table 4 online). A similar PLC
312 approach has been recently performed¹² using the Arabidopsis aliphatic glucosinolate pathway. The authors have successfully
313 reconstructed this pathway and identified a new candidate glucosyltransferase that could be part of it. This demonstrated
314 again that PLC is a powerful approach to complete biological pathways. When tested with a signaling pathway, we found that
315 PLCs also displayed meaningful communities. For example, the CK signaling pathway is physiologically well known but its
316 organization at the molecular level is far from being understood²¹. In particular, it is unclear how multi-family members of
317 each signaling level (receptor, transducer and response regulator) interact with each other to drive a specific physiological
318 response. In the PLC dedicated to the CK signaling pathway, PCC-HRR with microarrays suggested preferential transcriptional
319 associations that have been described in the literature [36]. For example AHP2, AHP3 and AHP5 were grouped in the same
320 module (module 7 Supplementary Figure 5 and Supplementary Table 4 online). These three AHPs have been reported to
321 negatively regulate tolerance to abiotic stress [40]. The same community also contained AHK3, ARR1 and ARR2. Those
322 three members are known to regulate primary root meristem activity and senescence²⁶. AHK4, AHK2 and ARR14 which
323 have been shown to regulate shoot apical meristem activity were grouped in the same community²⁷. In addition, we saw clear
324 associations between ETR1 and AHK3 in individual PLC subgraphs. Such association highlights crosstalk already known
325 between CK and ethylene signaling pathways²². The co-occurrence pathway was relatively sparse in contrast to the metabolic
326 pathways (Supplementary Figure 5 online). It is possible that vertex number for this analysis (1,000) might have been too
327 small to capture complex associations within this signaling pathway. Using VST normalized TPM increased edge and vertex
328 number in the co-occurrence network (Supplementary Figure 6 and Supplementary Table 5 online). The above-described
329 associations were also found in this co-occurrence network. While effective in revealing strong gene associations, merging PLC
330 from microarray and RNA-seq data could miss other relevant associations. First, experimental conditions represented by each
331 starting dataset are not completely overlapping. Together with inherent differences due to dynamic range, this leads to networks
332 with very different edge compositions and node degrees¹⁴, explaining the relative weak overlap between networks. Second,
333 RNA-seq expression data include genes that are not included in the GPL198 microarray. As an example, some important genes
334 in aliphatic glucosinolate biosynthesis were not represented in a previous Arabidopsis microarray dataset but found in RNA-seq
335 expression matrices from other related species¹².

336 To capture transcriptional environment of a query gene list, distance calculations have to be performed on the whole tran-
337 scriptome. Calculating partial correlations was particularly challenging but using a covariance shrinkage estimator worked well
338 in terms of computing performance. It took less than 2h for RNA-seq expression matrices but more than 12h for the microarray
339 dataset. By contrast, our program which is freely available at (<https://github.com/EA2106-Universite-Francois-Rabelais-Expression-network-analysis>)
340 was able to calculate PCC-HRR in less than 3h for both datasets. As PCC-HRR
341 resulted in relevant networks, this tool can be useful for further studies requiring many computations such as analyzing sample
342 size impact on PLC or testing other normalization methods.

343 The present work demonstrates that Pearson's Correlation Coefficients (PCC) on which highest reciprocal ranking (HRR)
344 was applied can be used to construct reliable global and targeted networks. When considering Pathway Level Correlation
345 (PLC) with a set of guide genes, three reliable measures can be used for evaluation: NV AUROC as a global indicator of GO
346 recovery (expecting values >0.5), modularity (between 0 and 1, 1 being the best network partition) and normalized Chi statistic
347 (between 0 and 1, 1 indicating a perfect match with an expected partition). Clustering guide genes correctly was at the expense
348 of capturing GO terms and dataset x distance measurement combination should be carefully selected to construct reliable
349 PLC. Although specific RNA-seq data normalizations may be adapted to each pathway of interest, using TPM with PCC-HRR
350 generated accurate and safe PLC. Using PCC with HRR also increased the quality of co-occurrence networks between RNA-seq
351 and microarrays.

352 **Methods**

353 **Microarray data preparation**

354 Experiment accessions (GSE) for GPL198 (*Arabidopsis* ATH1, 22,746 genes) were retrieved from ArrayExpress (Supplementary
355 Table 6 online). Signal intensities per probe were generated with R [16] using the 'arrayexpress' package²⁸. The function
356 'getAE' was used to convert the raw signal CEL files. Array normalization was performed per GSE using the 'justRMA'
357 function of the 'affy' package. This procedure applies a background correction together with a quantile normalization to correct
358 for biases within arrays and finally returns log₂-transformed corrected signal intensities. All 10,095 arrays were combined into
359 a single file and subjected to a quality control based on upper quartile dispersion (75%) and Kolmogorov-Smirnov statistical
360 testing for outliers using an empirical cumulative distribution function as described previously²⁹. A total of 142 arrays were
361 considered outliers in the two tests and discarded from the final matrix. Each array was finally centered and scaled individually.

362 **RNA-seq data preparation**

363 2,549 RNA-seq accessions obtained for *A. thaliana* were retrieved from ArrayExpress. Fastq files were obtained from the SRA
364 after converting .sra files with the SRA ToolKit function 'fastq-dump' with the -split-files option for paired-end sequencing runs.
365 Reads were systematically trimmed with Trimmomatic using adapter files according to the Illumina platform used for the runs
366 (ref). Trimmed reads were pseudo-aligned to predicted transcripts from the representative gene models of *Arabidopsis* TAIR
367 genome v10 (33,604 transcripts) with Salmon v0.7.2 using the variational Bayesian EM algorithm mode to improve abundance
368 estimation³⁰. Only samples displaying a mapping rate of reads >30% were kept, resulting in a final matrix containing 1,676
369 samples (Supplementary Table 6 online). RNA-seq counts were used as non-normalized raw counts or expressed as Transcript
370 per Million to correct for sequencing depth. Normalization by Variance Stabilizing Transformation (VST) was performed with
371 the DESeq2 R package. This normalization method aims at limiting the variance dependence to the mean³¹.

372 **Distance calculations**

373 Before calculations, zero-variance genes were discarded. CCs (Pearson or Spearman) are computationally intensive particularly
374 in the case of large matrices. Highest Reciprocal Ranking (HRR) of CCs for genes A and B is calculated as $\max(\text{rank}(\text{CC}(A,B)),$
375 $\text{rank}(\text{CC}(B,A)))$. For each gene, all CC values are first transformed as ranks, with 0 corresponding to the gene rank against
376 itself. Ranks are subsequently compared and the highest value is retained for each gene pair. We developed a tool written
377 in C allowing the easy parallelization of these computations. Briefly, for a given initial matrix containing n genes and p
378 samples, the number of cores c allocated is used to split the dataset into n/c submatrices. In case of non-integer value,
379 the last line of the matrix is replicated (without incidence on PCC or rank values) so that n/c is an integer. PCC or HRR
380 are then calculated for each gene pair using communication between CPUs with Message Passing Interface. The program
381 delivers c files containing $n/c \times n$ values corresponding to PCC or HRR. This program is freely available on Github (<https://github.com/EA2106-Universite-Francois-Rabelais/Expression-network-analysis>). To calcu-
382 late SCCs, expression values were first ranked in R. Mutual information (MI) which is reported to better capture non linear
383 relationships⁸ were calculated with the 'knni.all' function of the Parmigene R package³². This function estimates MI using a
384 k -nearest neighbor. Partial correlations were challenging to compute on genome scale expression matrices. Partial correlations
385 are usually calculated from multiple linear regressions or by inverting the correlation matrix and used in Graphical Gaussian
386 Models³³. Our expression matrices had many more variables (genes) than samples therefore regression methods would have
387 required a Lasso or Ridge penalization to estimate coefficients. However, this procedure generally leads to memory errors when
388 considering more than 30,000 variables. We found that the most computationally appropriate method in our case was to estimate
389 shrinkages of partial correlations with the R package 'corpcor' (<http://strimmerlab.org/software/corpcor/>).
390 This package is maintained by Korbinian Strimmer's team^{34,35}. We used 'pcor.shrink' function which relies on the inversion of
391 the shrunken estimated covariance matrix to estimate partial correlations and which is suited for matrices with more genes than
392 samples.
393

394 Reference dataset

395 We used the Arabidopsis Gene Ontology (GO) standard dataset to assess network quality. The annotation file provided by the
396 AGRIGO database³⁶ and was filtered out to remove all terms with a IEA evidence code and keep only functionally attributed
397 terms. We also removed GO terms represented by less than 5 genes or more than 100 to remove non-specific terms.

398 Global Network analysis

399 Construction: For each dataset x distance combination, we dynamically set a threshold to obtain arbitrary lists of 10 million
400 best gene pairs (with CC above or HRR below that threshold), *i.e.* less than 2% of the total possible edges. Networks were then
401 constructed with the 1, 5, 10, 20, 40, 60 or 80% best pairs from these lists. Thresholds used to get the 10 million gene pairs are
402 reported in S2 Table. Global networks were analyzed as adjacency matrices in R. Network characteristics: besides classical
403 topological characteristics such as vertex and edge numbers and mean node degree (the average number of connections for
404 each vertex), we evaluated network quality by comparison with the reference dataset (Figure 2). In a first approach, we built
405 a confusion matrix by classifying edges as false or true positives, considering edges as valid if both genes were annotated
406 with at least one same GO term. In this confusion matrix, true positives (TP) corresponded to gene pairs also found in the GO
407 annotation, false positive (FP) to genes associated in the network but not in the GO annotation, false negatives (FN) to pairs in
408 the GO annotation not predicted in the network and finally true negatives (TN) genes pairs not predicted in the network and the
409 annotation table. This confusion matrix was used to calculate True Positive Rates (TPR) and False Positive Rates (FPR). TPR
410 and FPR were obtained at various confidence thresholds (*i.e.* for networks differing in sizes) and used to draw a TPR vs FPR
411 curve as described elsewhere¹¹. These curves were only partial because we included only the first 10 million best pairs. This
412 was useful to pinpoint the importance of low FPR³⁷. In the second and third approaches, we relied on the guilt-by-association
413 principle to estimate network predictability. In the second method, we used the ‘predictions’ function of EGAD R package³⁸.
414 For each gene, this function counts the number of connected genes annotated with an identical GO term and divides this count
415 by the gene’s degree. These scores are next ordered decreasingly to construct a TPR vs FPR curve for each network. It differs
416 from the first approach described above because here TPR and FPR are not obtained from different confidence thresholds (and
417 from different networks) but from all possible true positive and false positive edges in the current network. A global Area
418 Under Receiver Operating Characteristic (global AUROC) was calculated from each of these TPR/FPR curves. In the third
419 method, predictability was evaluated using a neighbor voting (NV) algorithm. In this case, an AUROC is calculated for each
420 GO term from the ability of genes to predict the GO annotation of their direct neighbors in a 3-fold cross-validation^{14,39}. A
421 mean NV AUROC was calculated for each network. In addition to ROC analysis, we counted GO terms that were significantly
422 enriched with gene pairs using a hypergeometric test with R.

423 Pathway Level Correlation

424 Construction: In PLC, subnetworks were constructed from global networks (see above) by keeping edges connecting at least
425 one guide gene. Guide gene lists are indicated in Supplementary Table 2 online. The R package ‘igraph’⁴⁰ v1.0.1 was used to
426 construct and visualize these targeted networks with a force-directed layout (Fruchterman-Reingold). Community Detection:
427 Modules containing densely connected vertices were estimated within each network by using a fast greedy approach which
428 aims at maximizing modularity of the detected communities²³. Modularity measures how good a network partition is by
429 calculating for each gene the number of edges within its community against its total node degree. The fast greedy approach
430 optimizes modularity over all possible divisions of the network and has been shown to perform well on large networks. Guide
431 genes clustering within the communities was compared to expected partitions in sub-pathway with a Pearson’s Chi² test and
432 Monte-Carlo simulated *p*-values with 2,000 replicates. This test was based on a contingency table with dimensions $n \times m$ (n ,
433 sub-pathway number, m community number in the co-expression network) and each entry corresponding to the number of
434 genes being in communities n_i and m_j , with $i=1$ to n and $j=1$ to m . Because Chi² statistic depends on sample number, values
435 were normalized by dividing them to the maximal expected value (the ideal partition) of each pathway. This resulted in a score
436 ranging from 0 to 1, 0 being a random distribution of guide genes in the network and 1 to the exact partitioning.

437 Acknowledgements

438 We deeply acknowledge the Fédération CaSciModOT (CCSC Orléans-Tours, France), Jean-Louis Rouet and Laurent Catherine
439 for help and access to the Région Centre computing grid. We also thanks Yann Jullian for access and help on University
440 computer resources. This study was supported by the Région Centre-Val de Loire, France (SiSCyLi grant). Doctoral Fellow
441 attributed to F.L. and D.D. was jointly funded by the Région Centre-Val de Loire, France and the Ministère de l’Enseignement
442 Supérieur et de la Recherche, France.

443 Author contributions statement

444 F.L., O.P., J.C., N.G. and T.D.D.B. conceived the experiment(s), F.L., D.D., O.P., M.C., S.B., V.C. and R.D.D.B conducted the
445 experiment(s), F.L., S.B., G.G., J.C., J.O.C. and T.D.D.B. analyzed the results. All authors reviewed the manuscript.

446 Competing interests

447 The authors declare no competing interests.

448 Data availability

449 All datasets generated during and/or analyzed during the current study are available from the corresponding author on reasonable
450 request.

451 References

- 452 1. Oliver, S. Proteomics: guilt-by-association goes global. *Nat.* **403**, 601–603 (2000).
- 453 2. Lisso, J., Steinhauser, D., Altmann, T., Kopka, J. & Müssig, C. Identification of brassinosteroid-related genes by means of
454 transcript co-response analyses. *Nucleic Acids Res.* **33**, 2685–2696 (2005).
- 455 3. Wei, H. *et al.* Transcriptional coordination of the metabolic network in arabidopsis. *Plant physiology* **142**, 762–774 (2006).
- 456 4. Ruiz-Sola, M. Á. *et al.* Arabidopsis geranylgeranyl diphosphate synthase 11 is a hub isozyme required for the production
457 of most photosynthesis-related isoprenoids. *New Phytol.* **209**, 252–264 (2016).
- 458 5. Guerin, C. *et al.* Gene coexpression network analysis of oil biosynthesis in an interspecific backcross of oil palm. *The*
459 *Plant J.* **87**, 423–441 (2016).
- 460 6. Coman, D., Rütimann, P. & Gruissem, W. A flexible protocol for targeted gene co-expression network analysis. *Plant*
461 *Isoprenoids: Methods Protoc.* 285–299 (2014).
- 462 7. Maetschke, S. R., Madhamshettiwar, P. B., Davis, M. J. & Ragan, M. A. Supervised, semi-supervised and unsupervised
463 inference of gene regulatory networks. *Briefings bioinformatics* **15**, 195–211 (2013).
- 464 8. de Siqueira Santos, S., Takahashi, D. Y., Nakata, A. & Fujita, A. A comparative study of statistical methods used to identify
465 dependencies between gene expression signals. *Briefings bioinformatics* **15**, 906–918 (2013).
- 466 9. Li, Y., Pearl, S. A. & Jackson, S. A. Gene networks in plant biology: approaches in reconstruction and analysis. *Trends*
467 *plant science* **20**, 664–675 (2015).
- 468 10. Serin, E. A., Nijveen, H., Hilhorst, H. W. & Ligterink, W. Learning from co-expression networks: possibilities and
469 challenges. *Front. plant science* **7** (2016).
- 470 11. Obayashi, T. & Kinoshita, K. Rank of correlation coefficient as a comparable measure for biological significance of gene
471 coexpression. *DNA research* **16**, 249–260 (2009).
- 472 12. Wisecaver, J. H. *et al.* A global co-expression network approach for connecting genes to specialized metabolic pathways in
473 plants. *The Plant Cell Online* tpc-00009 (2017).
- 474 13. Mutwil, M. *et al.* Assembly of an interactive correlation network for the arabidopsis genome using a novel heuristic
475 clustering algorithm. *Plant Physiol.* **152**, 29–43 (2010).
- 476 14. Ballouz, S., Verleyen, W. & Gillis, J. Guidance for rna-seq co-expression network construction and analysis: safety in
477 numbers. *Bioinforma.* **31**, 2123–2130 (2015).
- 478 15. Song, L., Langfelder, P. & Horvath, S. Comparison of co-expression measures: mutual information, correlation, and model
479 based indices. *BMC bioinformatics* **13**, 328 (2012).
- 480 16. Besseau, S. *et al.* Flavonoid accumulation in arabidopsis repressed in lignin synthesis affects auxin transport and plant
481 growth. *The Plant Cell* **19**, 148–162 (2007).
- 482 17. Zhang, Y. *et al.* Phenolic compositions and antioxidant capacities of chinese wild mandarin (*Citrus reticulata* blanco) fruits.
483 *Food chemistry* **145**, 674–680 (2014).
- 484 18. Winkel-Shirley, B. Flavonoid biosynthesis. a colorful model for genetics, biochemistry, cell biology, and biotechnology.
485 *Plant physiology* **126**, 485–493 (2001).

- 486 **19.** Elejalde-Palmett, C. *et al.* Characterization of a spermidine hydroxycinnamoyltransferase in *malus domestica* highlights
487 the evolutionary conservation of trihydroxycinnamoyl spermidines in pollen coat of core eudicotyledons. *J. experimental*
488 *botany* **66**, 7271–7285 (2015).
- 489 **20.** Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. Kegg as a reference resource for gene and protein
490 annotation. *Nucleic acids research* **44**, D457–D462 (2016).
- 491 **21.** Hwang, I., Sheen, J. & Müller, B. Cytokinin signaling networks. *Annu. review plant biology* **63**, 353–380 (2012).
- 492 **22.** Zdarska, M. *et al.* Illuminating light, cytokinin, and ethylene signalling crosstalk in plant development. *J. experimental*
493 *botany* **66**, 4913–4931 (2015).
- 494 **23.** Clauset, A., Newman, M. E. & Moore, C. Finding community structure in very large networks. *Phys. review E* **70**, 066111
495 (2004).
- 496 **24.** Giorgi, F. M., Del Fabbro, C. & Licausi, F. Comparative study of rna-seq-and microarray-derived coexpression networks in
497 *arabidopsis thaliana*. *Bioinforma.* **29**, 717–724 (2013).
- 498 **25.** Uygun, S., Peng, C., Lehti-Shiu, M. D., Last, R. L. & Shiu, S.-H. Utility and limitations of using gene expression data to
499 identify functional associations. *PLoS computational biology* **12**, e1005244 (2016).
- 500 **26.** Jiang, L. *et al.* Strigolactones spatially influence lateral root development through the cytokinin signaling network. *J.*
501 *experimental botany* **67**, 379–389 (2015).
- 502 **27.** Wang, L. & Chong, K. The essential role of cytokinin signaling in root apical meristem formation during somatic
503 embryogenesis. *Front. plant science* **6** (2015).
- 504 **28.** Kauffmann, A. *et al.* Importing arrayexpress datasets into r/bioconductor. *Bioinforma.* **25**, 2092–2094 (2009).
- 505 **29.** Feltus, F. A., Ficklin, S. P., Gibson, S. M. & Smith, M. C. Maximizing capture of gene co-expression relationships through
506 pre-clustering of input expression samples: an arabidopsis case study. *BMC systems biology* **7**, 44 (2013).
- 507 **30.** Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-aware quantification of
508 transcript expression. *Nat. Methods* **14**, 417–419 (2017).
- 509 **31.** Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for rna-seq data with *deseq2*.
510 *Genome biology* **15**, 550 (2014).
- 511 **32.** Sales, G. & Romualdi, C. *parmigene*—a parallel r package for mutual information estimation and gene network reconstruc-
512 tion. *Bioinforma.* **27**, 1876–1877 (2011).
- 513 **33.** López-Kleine, L., Leal, L. & López, C. Biostatistical approaches for the reconstruction of gene co-expression networks
514 based on transcriptomic data. *Briefings functional genomics* **12**, 457–467 (2013).
- 515 **34.** Schäfer, J. & Strimmer, K. Learning large-scale graphical gaussian models from genomic data. In *AIP Conference*
516 *Proceedings*, vol. 776, 263–276 (AIP, 2005).
- 517 **35.** Schaefer, J., Opgen-Rhein, R. & Strimmer, K. *corpcor*: efficient estimation of covariance and (partial) correlation. r
518 package version 1.4. 7 (2007).
- 519 **36.** Du, Z., Zhou, X., Ling, Y., Zhang, Z. & Su, Z. *agrigo*: a go analysis toolkit for the agricultural community. *Nucleic acids*
520 *research* **38**, W64–W70 (2010).
- 521 **37.** Schrynemackers, M., Küffner, R. & Geurts, P. On protocols and measures for the validation of supervised methods for the
522 inference of biological networks. *Front. genetics* **4** (2013).
- 523 **38.** Ballouz, S., Weber, M., Pavlidis, P. & Gillis, J. *Egad*: ultra-fast functional analysis of gene networks. *Bioinforma.* **33**,
524 612–614 (2016).
- 525 **39.** Gillis, J. & Pavlidis, P. The impact of multifunctional genes on “guilt by association” analysis. *PloS one* **6**, e17258 (2011).
- 526 **40.** Csardi, G. & Nepusz, T. The *igraph* software package for complex network research. *InterJournal, Complex Syst.* **1695**,
527 1–9 (2006).

528 **Supplementary Information**

529 **Supplementary Figure 1: Network properties in dataset-distance measurement combinations.** Global network charac-
530 teristics (Number of significantly enriched GO terms, global and NV AUROCs) were expressed in function of vertex or edge
531 number. The horizontal dashed line indicates a 0.6 AUROC value taken as an arbitrary threshold separating good and poor

532 network predictability. For each dataset, TPR=f(FPR) curves are also presented with dashed line corresponding to a random
533 selection (with AUROC <0.5). These curves are partial and the max FPR values were obtained for 10 million gene pairs.

534 **Supplementary Table 1: Threshold values to get 10 million best gene pairs.**

535 **Supplementary Figure 2: Workflow for Pathway-Level Correlation. Lists of best co-expressed genes are established
536 for each guide (or bait) gene.** Redundancies among these lists (associated genes) connect guide genes to construct the PLC
537 network. Terms 'guide gene' and 'associated genes' have been introduced by Lisso et al².

538 **Supplementary Table 2: Guide gene accessions.**

539 **Supplementary Figure 3. PLC subgraphs for the carbohydrate (A), fatty acid (B), terpene (C) and cytokinin (D)
540 pathways.** For each PLC, the expected partitioning of guide genes is indicated in the left panel and is compared to PLC
541 subgraphs with higher predictability and lower modularity (center; calculated with MI) or PLC subgraphs with lower pre-
542 dictability and higher modularity (right; calculated with PCC-HRR). Colored vertices correspond to genes encoding enzymes
543 catalyzing steps of similar color in the expected pathway. A and B were drawn from RNA-seq TPM networks while C and D
544 from microarray networks. Community numbers in PCC-HRR networks are indicated in deep blue and can be used to access
545 Supplementary Table 3 online. Polygons surrounding vertices delimit communities.

546 **Supplementary Table 3: Gene lists from PLC obtained with PCC-HRR Genes highlighted in yellow correspond to
547 non-guide genes but known to be involved in the pathway.**

548 **Supplementary Figure 4: PLC based on microarray and TPM data.** Subgraphs were constructed with the 6 distance
549 measurements (MI, PC, raw PCC, raw SCC, PCC-HRR and SCC-HRR) and aligned to find co-occurring edges and vertices.
550 (A) Number of co-occurring vertices and edges. The first distance in each label correspond to microarrays and the second to
551 TPM. Points are half-colored according to the ranking applied to the initial distance. For each intersection graph, % of guide
552 genes (B), normalized Chi² statistic (agreement with expected guide gene partitioning, C), modularity (D) and NV AUROC
553 (GO recovery performance, E) were calculated.

554 **Supplementary Figure 5: Co-occurrence networks from PCC-HRR PLC constructed with microarrays and RNA-
555 seq TPM.** Community numbers are indicated in deep blue and can be used to access Supplementary Table 4 online.

556 **Supplementary Table 4: Gene lists from co-occurrence networks between PCC-HRR PLC obtained with microar-
557 rays and RNA-seq TPM.** Genes highlighted in yellow correspond to non-guide genes but known to be involved in the
558 pathway.

559 **Supplementary Figure 6: Co-occurrence networks from PCC-HRR PLC constructed with microarrays and RNA-
560 seq TPM normalized with VST.** Community numbers are indicated in deep blue and can be used to access Supplementary
561 Table 5 online.

562 **Supplementary Table 5: Gene lists from co-occurrence networks between PCC-HRR PLC obtained with microar-
563 rays and RNA-seq TPM normalized with VST.**

564 **Supplementary Table 6: Microarray and RNA-seq accessions used in this study.**

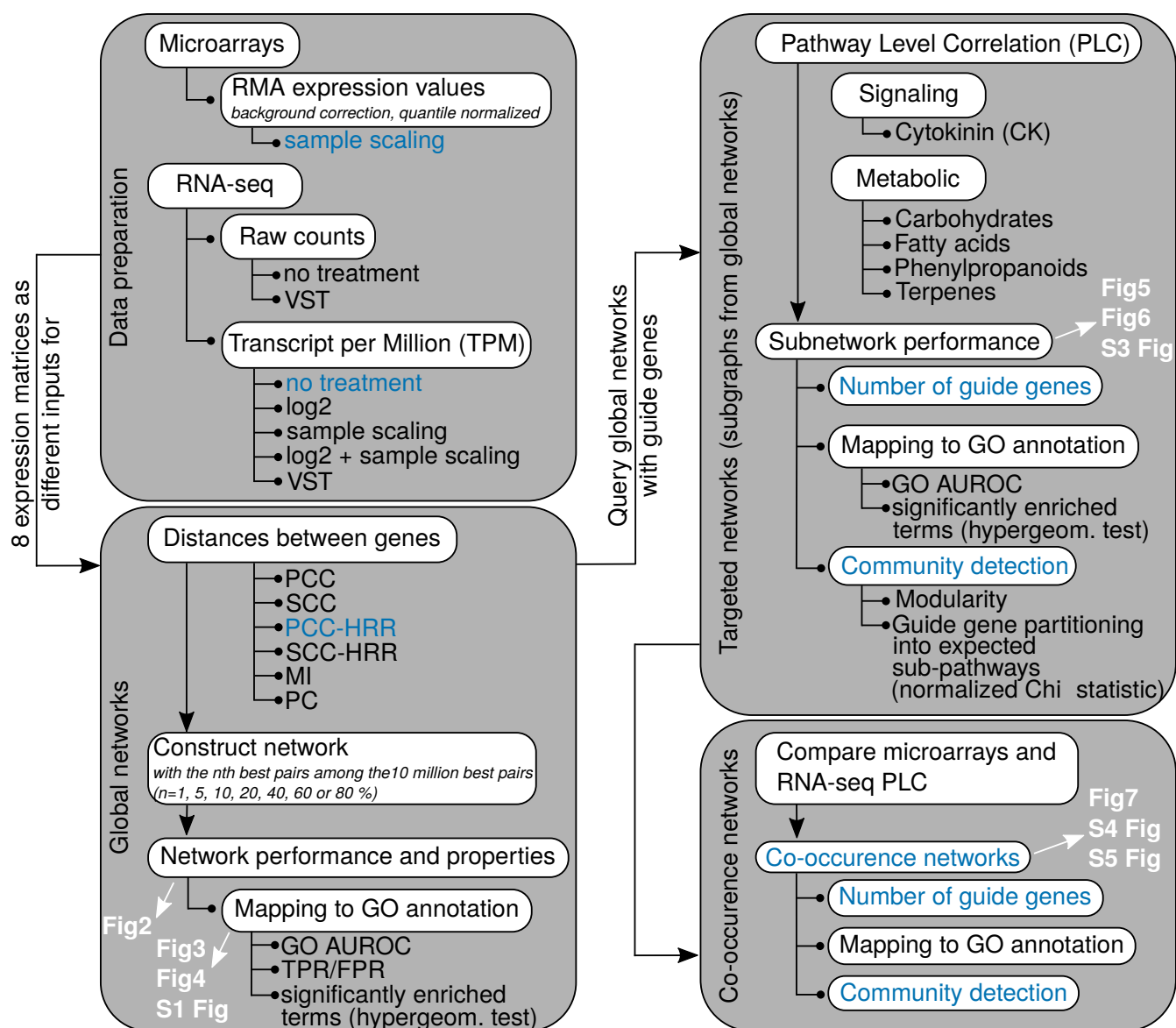


Figure 1. Workflow for global and targeted network analyses. One microarray dataset and a RNA-seq dataset prepared according to 7 normalization procedures were used to generate eight expression matrices analyzed with six different distance measurements (Pearson’s or Spearman’s Correlation Coefficient, unranked or ranked with HRR, Mutual Information (MI) or Partial Correlations (PC)) to obtain 48 distance matrices. Each of these matrices was thresholded to obtain global networks at different confidence thresholds. Global networks were evaluated and also queried with specific guide gene sets reflecting 5 different pathways in a process named Pathway Level Correlation (PLC). The resulting subnetworks were evaluated and used to construct co-occurrence networks between microarray and RNA-seq datasets. In white are indicated the figures corresponding to the different steps analyzed. Dataset x distance combinations are indicated in blue and characteristics that are improved by these combinations.

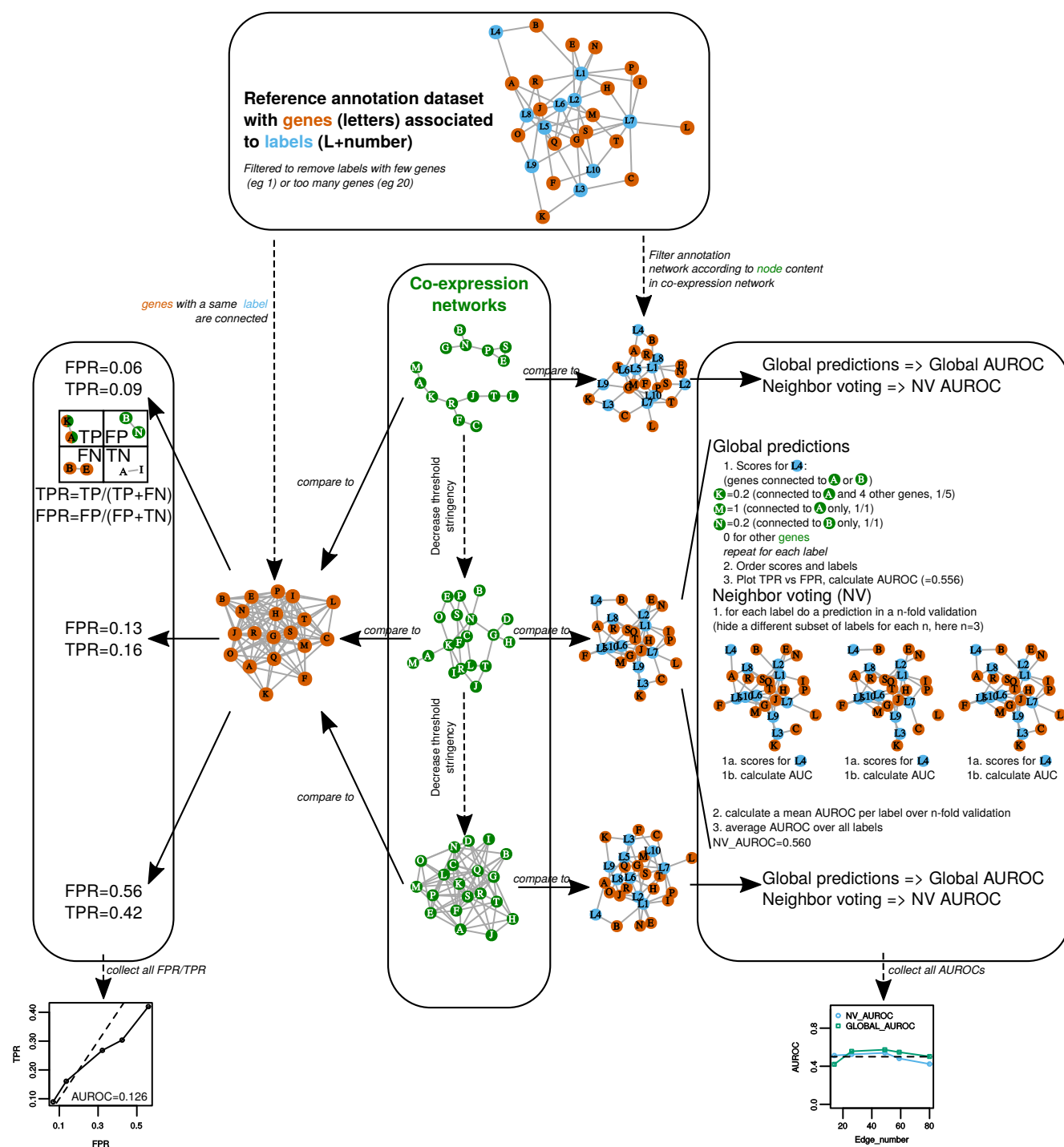


Figure 2. Network performance. This small example describe strategies to evaluate networks according to a reference functional annotation. Co-expression networks were obtained for each dataset x distance measurement combination (Figure 1) at different confidence thresholds, resulting in networks increasing in size with lower stringency. A total evaluation was made with True Positive Rate (TPR) vs False Positive Rate (FPR) analysis (left panel) by classifying edges as True positives (TP), False Positives (FP), False Negatives (FN) or True Negatives (TN). Single network evaluation was performed by calculating AUROCs with the EGAD R package, either as a global prediction or using a neighbor voting (NV) algorithm with a 3-fold cross validation (right panel). All indicated values are in accordance with the small networks in this example. In addition to these 3 evaluations (FPR vs TPR, global AUROC and NV AUROC), GO term significant enrichment was statistically tested with a hypergeometric distribution (not shown in this example).

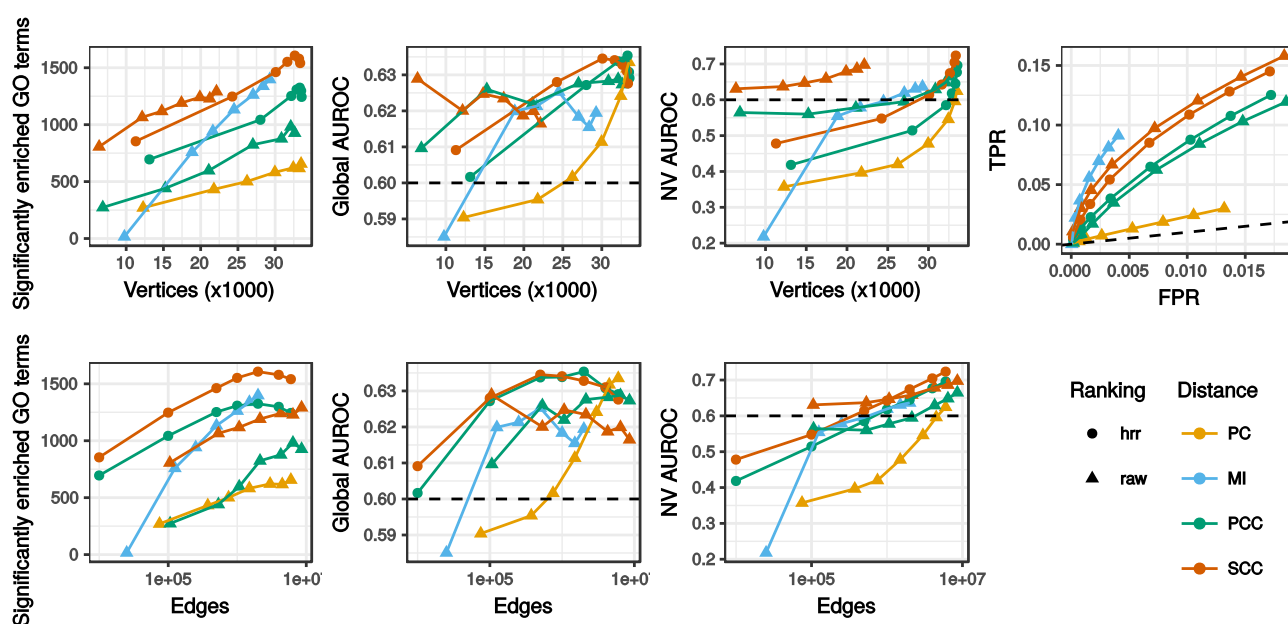


Figure 3. Global network characteristics. Only results for the RNA-seq TPM dataset without further normalization are shown. The horizontal dashed line indicates a 0.6 AUROC value taken as a threshold separating good and poor network predictability. In the $TPR=f(FPR)$ panel, the dashed line corresponds to a random selection (with $AUROC < 0.5$). This panel is partial and the highest FPRs correspond to 10 million gene pairs.

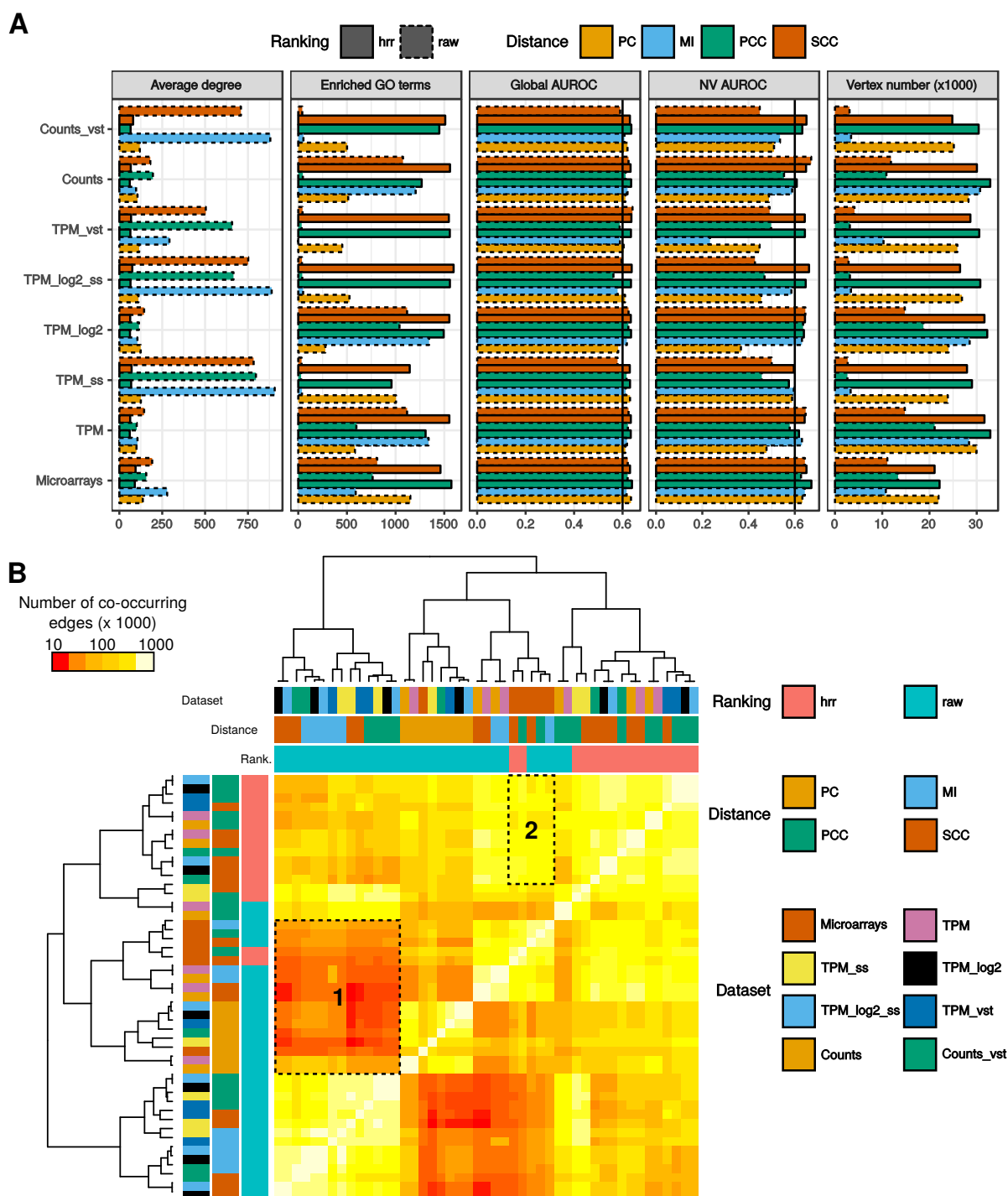


Figure 4. Comparison of dataset x distance measurement combinations for networks with a million gene pairs. Network topology and performance in GO recovery were analyzed (A). Vertical lines at 0.6 indicate AUROCs above which network predictability can be considered as moderate. Co-occurring edges were also counted in every possible comparison between 2 networks (B). Area 1 corresponds to RNA-seq networks having few genes in common with PC networks and microarrays networks and area 2 to combinations maximizing edge co-occurrence between microarray and RNA-seq.

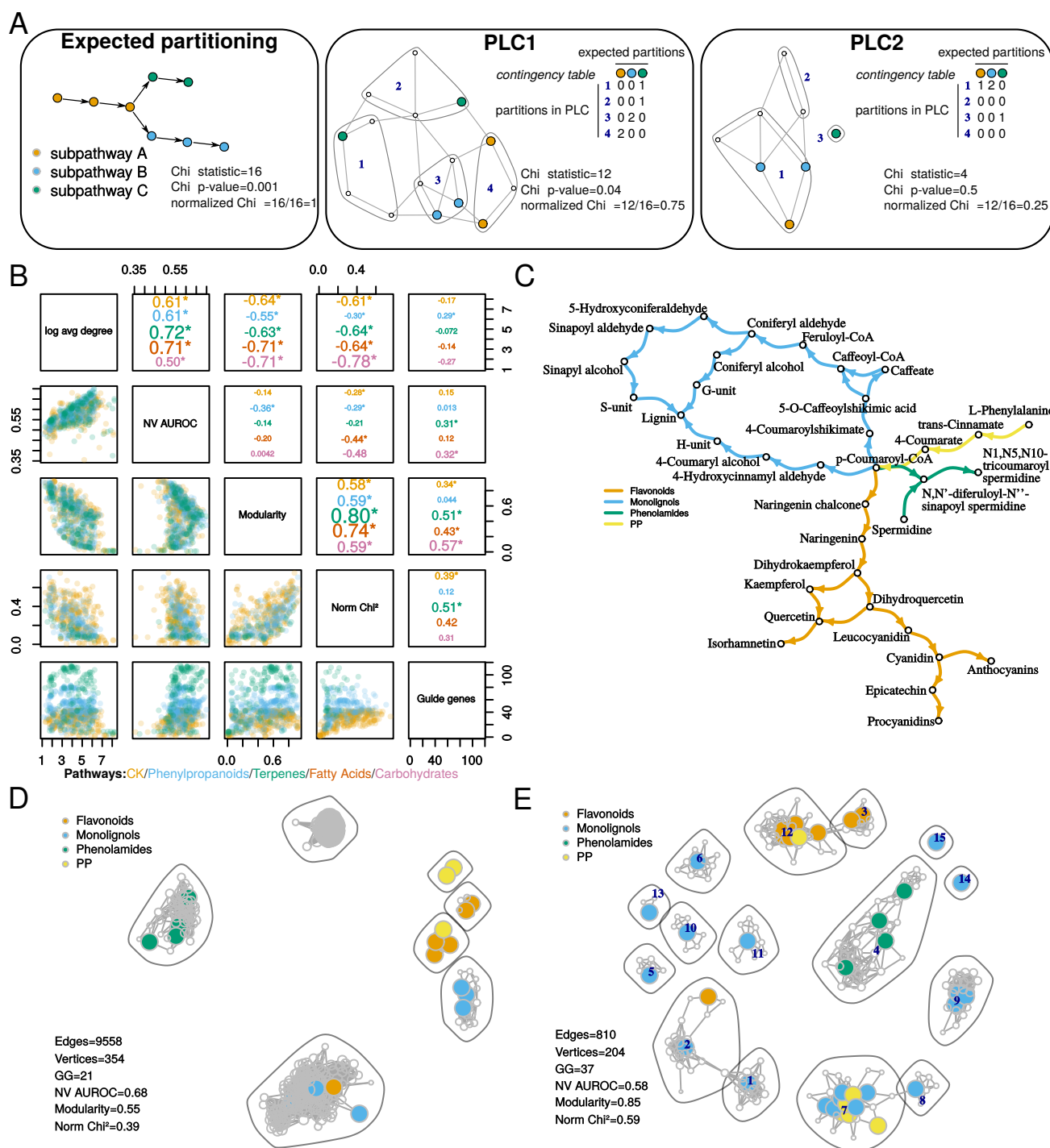


Figure 5. Trade-off in PLC subnetworks between performance in GO term recovery and partitioning guide genes into expected communities. (A) Example showing normalized χ^2 statistic and p -value calculations comparing guide gene distribution into PLC communities (numbers in deep blue within polygons) to the expect partitioning (left; 3 subpathways). Two PLCs (one with a good partitioning (center); one with a weak partitioning (right)) are shown here but the contingency matrix used in χ^2 calculations is described for only one of them (center). (B) Pair plot showing correlations (Spearman's rho, asterisks show significance $p < 0.001$, upper panel) and scatterplots (lower panel) between average network node degree, NV AUROC, normalized χ^2 , modularity and the number of guide genes in the network. Each point in the lower panels (scatterplots) represent one network for which 2 characteristics (eg NV AUROC and modularity) are compared. Data are presented for each pathway separately with a specific color. (C) The expected partitioning of phenylpropanoid related guide genes was compared to two PLC: (D) higher predictability and lower modularity (microarrays raw PCC) and (E) lower predictability and higher modularity (microarrays PCC-HRR). In D and E, colored vertices correspond to genes encoding enzymes catalyzing steps of similar color in C. Community (surrounded by grey polygons) numbers in E are indicated in deep blue and can be used to access Supplementary Table 3 online.

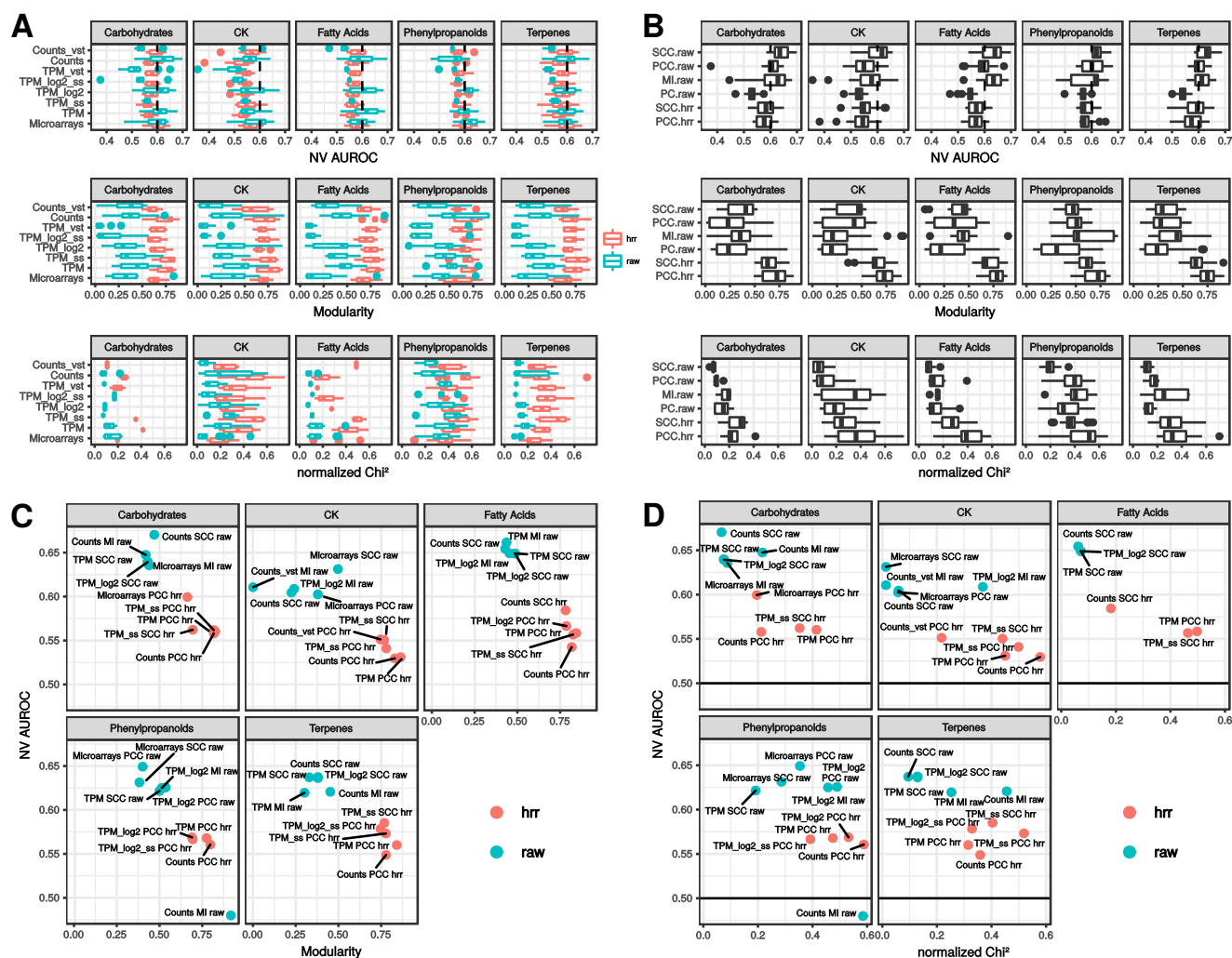


Figure 6. PLC subnetwork performance. Performance in capturing GO terms (NV AUROC), modularity and normalized χ^2 value distribution in interactions between datasets and ranking methods (A) and between distance measurement and ranking methods (B) showing the dominant effect of the ranking procedure (raw vs HRR) on these metrics. (C) Modularity and NV AUROC of the five top NV AUROC networks and 5 top modularity networks. (D) Normalized χ^2 statistic and NV AUROC for the same networks.

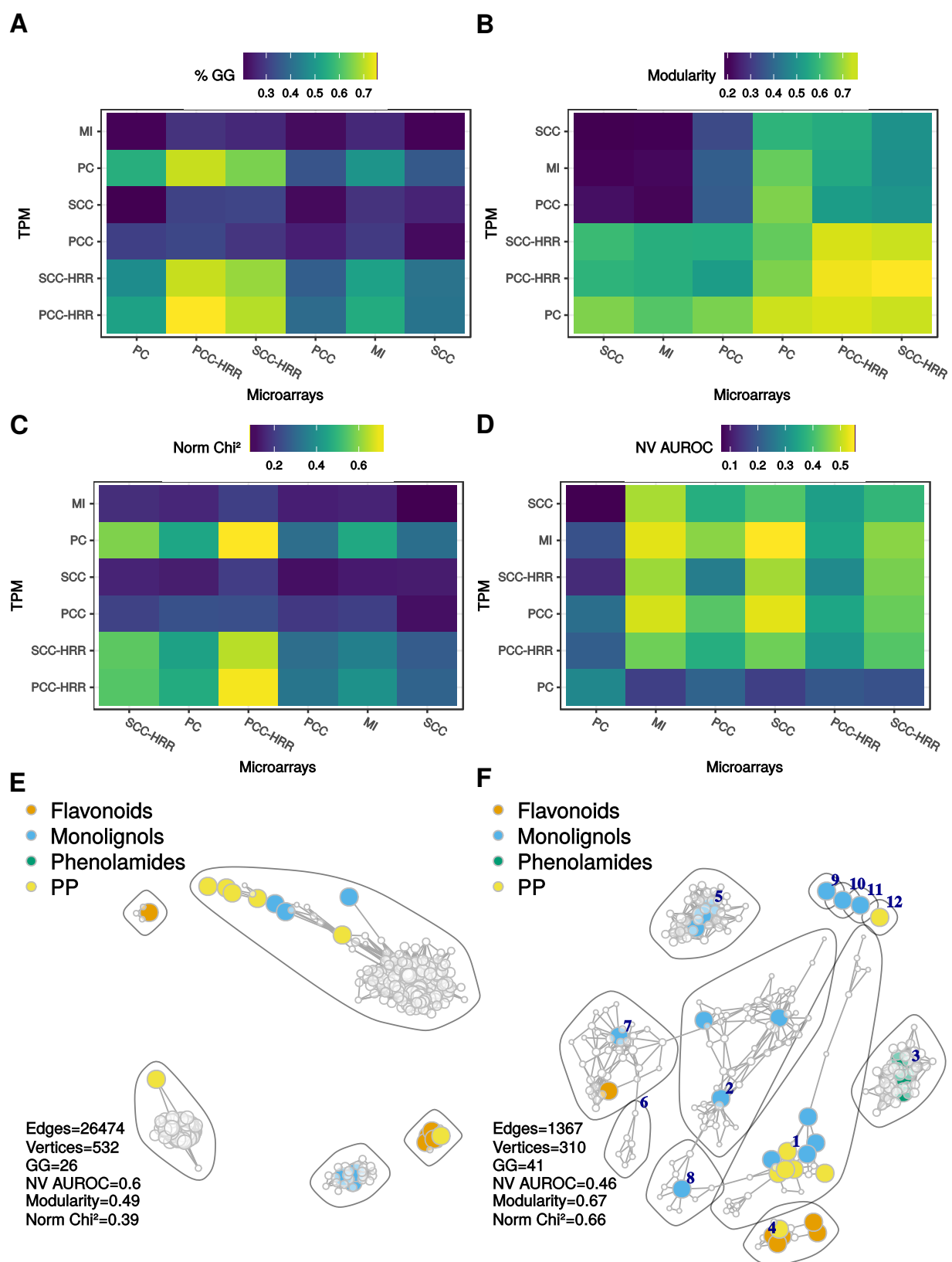


Figure 7. Characteristics of co-occurrence networks between microarrays and RNA-seq TPM. Percentage of guide genes (GG; A), modularity (B), normalized Chi² statistic (agreement with guide gene partitioning, C) and NV AUROC (GO term performance, D) were averaged over the 5 PLCs. Labels are ordered according to a hierarchical clustering. Co-occurrence networks obtained from phenylpropanoid PLC obtained with MI (E) or PCC-HRR (F). GG corresponds to guide gene number in the networks. Community numbers in F are indicated in deep blue and can be used to access Supplementary Table 4 online.