# Gene regulatory networks associated with lateral root and nodule development in soybean

**Shuchi Smita[1,4], Jason Kiehne[2], Sajag Adhikari[1,5], Erliang Zeng[3,6], Qin Ma[1,7,*], and Senthil Subramanian[1,*]**

[1] Department of Agronomy, Horticulture and Plant Science and Department of Biology and Microbiology, South Dakota State University, Brookings, SD, USA

[2] Simpson College, Indianola, IA.

[3] Department of Biology and Department of Computer Science, University of South Dakota, Vermillion, USA

[4] Present Address: Department of Computational and Systems Biology and Department of Immunology, School of Medicine, University of Pittsburgh, Pittsburgh, PA 15261, USA

[5] Present Address: Celerion Inc., Lincoln, NE 68502.

[6] Present Address: Division of Biostatistics and Computational Biology, College of Dentistry, University of Iowa, Iowa City, USA.

[7] Present Address: Department of Biomedical Informatics, Ohio State University, Columbus, OH 43210, USA.

* Corresponding Authors Email Id: Qin.Ma@osumc.edu or Senthil.Subramanian@sdstate.edu

22  **Running Title: Soybean root lateral organ gene regulatory networks**

23

24  **Abstract**

25  Legume plants such as soybean produce two major types of root lateral organs, lateral roots and

26  root nodules. A robust computational framework was developed to predict potential gene

27  regulatory networks (GRNs) associated with root lateral organ development in soybean. A

28  genome-scale expression dataset was obtained from soybean root nodules and lateral roots and

29  subjected to biclustering using QUBIC. Biclusters (BCs) and transcription factor (TF) genes with

30  enriched expression in lateral root tissues were converged using different network inference

31  algorithms to predict high confident regulatory modules that are repeatedly retrieved in different

32  methods. The ranked combination of results from all different network inference algorithms into

33  one ensemble solution identified 21 GRN modules of 182 co-regulated genes networks

34  potentially involved in root lateral organ development stages in soybean. The pipeline correctly

35  predicted previously known nodule- and LR-associated TFs including the expected hierarchical

36  relationships. The results revealed high scorer AP2, GRF5, and C3H co-regulated GRN modules

37  during early nodule development; and GRAS, LBD41, and ARR18 co-regulated GRN modules

38  late during nodule maturation. Knowledge from this work supported by experimental validation

39  in the future is expected to help determine key gene targets for biotechnological strategies to

40  optimize nodule formation and enhance nitrogen fixation.

41

## Introduction

42
43

44       Gene regulation is a fundamental process that controls spatial and temporal patterns of
45  gene expression. Transcription factors (TFs) are central to gene regulation as their activities
46  determine the expression patterns and function of multiple genes (1). A TF is a functional protein
47  that binds to short sequences (called TF binding site; TFBS or *cis*-regulatory elements) on the
48  upstream promoter region of genes to regulate their transcription. One TF can regulate multiple
49  genes including other TFs in a signaling, developmental, or metabolic pathway and so act as
50  master regulators of the pathways. The nested group of all different TF regulators and their
51  downstream target genes form gene regulatory networks (GRNs) (2).  Identification of gene
52  regulatory networks and key TFs that are part of these networks is an effective approach to
53  answer multiple biological questions on genotype to phenotype relationships. For example,
54  potential TFs, their co-regulators, and downstream signaling pathways, and target genes
55  associated with specific biological processes can be predicted by constructing GRNs.

56

57       Clustering of large-scale datasets such as global gene expression profiles obtained by
58  RNA-sequencing to identify co-regulated TFs and the targeting genes is a promising approach to
59  model and infer the GRNs at a systems level (3, 4). Briefly, genes/TFs with similar expression
60  patterns (i.e. co-expressed genes) with a tendency to co-activate across a group of samples might
61  give insight on TFs regulated gene network and related biological process. In fact, multiple
62  levels of gene regulation affect transcriptional regulatory capabilities (5). Recruitment and
63  binding of other protein such as "co-factors" in complexes and other small protein molecules to
64  target DNA sequences is one of the major mechanisms (6). Often, this interactions between
65  different TFs and co-factor partners are studied using protein-protein interaction (PPI) assays
66  which provide immediate insights into their potential biological function (7, 8). GRNs can be
67  validated by PPI data, as PPIs can reveal signaling, regulatory and/or biochemical roles of
68  proteins based on their interactomes (9).

69       The combined use of high-throughput data and mathematical models to build gene co-
70  expression and regulatory networks is the core principle of systems biology approaches (10).
71  However, these large-scale datasets are likely to be noisy, and GRN predictions using these big
72  datasets may contain many false positives. Additionally, GRN inference is a computationally

73    intensive job; so filtered datasets consisting of well-defined/accurate datasets (such as

74    significantly co-expressed genes set) might dramatically reduce the computational complexity

75    and time. Most importantly, it would reduce the true search space for the prediction of regulators

76    (TFs) and their potential target genes. In order to obtain significantly co-expressed genes,

77    "biclustering" is a desirable method as it allows two-way clustering of genes as well as samples

78    i.e. a similar expression pattern (co-expressed genes) under a subset of all samples.

79    Subsequently, this sorted biclustering-filtered data fed into GRN inference algorithms might

80    improve and accurate predictions of a regulator and their target genes. We applied this approach

81    to determine gene regulatory networks associated with root lateral organ development in

82    soybean.

83

84    Plants produce lateral organs such as leaves, flowers, and axillary branches in the shoot,

85    and lateral roots in the roots. Pools of stem cells present in the growing tip of the shoot (the shoot

86    apical meristem) contribute to the formation of aerial/shoot lateral organs. Lateral organs in the

87    root are unique in that they are derived via "*de novo*" differentiation of mature cells in the root.

88    Lateral roots are present in all vascular plants, but a group of *Fabids* clade plants is capable of

89    producing another root lateral organ, called root nodules. These arise from specific and

90    coordinated interactions with a set of nitrogen-fixing bacteria collectively called rhizobia. For

91    example, the interaction of soybeans with *Bradyrhizobium diazoefficiens* results in root nodules.

92    Biological nitrogen fixation in root nodules helps reduce the need for chemical nitrogen

93    fertilizers, which are expensive and cause environmental pollution. Similarly, proper patterns of

94    lateral root formation (root branching) are crucial for plants to access water and other nutrients in

95    the soil. Therefore, these two root lateral organs play important roles in the development of

96    soybeans, a major crop in the United States as well as in other countries. Many functional

97    genomics studies have identified genes expressed during nodule development in soybean and

98    other legumes, but gene expression profiles during lateral root formation have not been evaluated

99    in legumes (11, 12).

100    Recently, we obtained transcriptomes of emerging nodules, mature nodules, emerging

101    lateral roots, and young lateral roots in soybean (13), we present a robust computational

102    framework, which we applied to predict TFs and their target GRNs associated with soybean root

103    nodule development. This approach consists of the following steps (Figure 1): (i) preparing a

104    compendium of soybean lateral organ transcriptome data and cataloging TFs enriched in root

105    nodules; (ii) Initial biclustering of transcriptome data using QUBIC (14, 15, 16) to identify all

106    (nodule development stage-specific) co-expressed gene modules; (iii) GRN construction and

107    inference based on identified gene modules and reliable network construction programs, Lemon-

108    Tree (17) and Inferelator (18); (iv) Augmentation of GRNs with evidence from physical or direct

109    and indirect regulatory interaction information from PPI and *cis*-regulatory element enrichment

110    analysis; and  (v) building a consensus from different modes of GRN inference for potential

111    regulators and their predicted GRNs. We ran two modes of Lemon-Tree, one with default mode,

112    where Lemon-Tree itself produce the co-expressed clusters and the other mode with reinforced

113    bicluster (BC) information from QUBIC. This study provides a template framework for GRN

114    construction and augmentation by exploiting big data sets, which are increasingly generated,

115    deposited and available (making use of available data) in public domain.

116

117    **Material and methods**

118

119    **RNA-seq dataset for root lateral organ development in soybean**

120        We utilized the genome-wide soybean transcriptome dataset generated for root lateral

121    organs (13). This dataset contains the transcriptomes of two different developmental stages of

122    two root lateral organs collected in three biological replicates: emerging nodules (EN), mature

123    nodules (MN), emerging lateral roots (ELR) and young lateral roots (YLR). Adjacent root

124    sections above and below these organs devoid of any lateral organs (designated as ABEN,

125    ABMN, ABELR, and ABYLR respectively) were used to construct respective age- and

126    inoculation-status appropriate control tissue libraries. Comparison of gene expression profiles

127    between each lateral organ tissue type and the corresponding control tissue type (e.g., EN vs.

128    ABEN, ELR vs. ABELR and so on) helped identify organ-specific/enriched genes. In total, 24

129    RNA-seq libraries (four target tissue types, four control tissue-types, three biological replicates

130    each) were prepared, sequenced, and analyzed. Expression patterns of preciously known marker

131    genes, consistency between replicates, high sequence quality of this dataset indicated that it was

132    of very high quality and well-suited for global gene expression analysis (13). A total of 113,210

133    gene transcripts (FPKM threshold $\geq$ 1 in at least one sample) with their normalized expression

134    values in 24 different tissues from the above dataset were utilized here.

135

136       Further, for expression comparisons at different steps during our analysis, we utilized the

137    public datasets, soybean gene atlas encompassing RNA-seq data from 14 different soybean

138    tissues  (19) and Soybean eFP browser http://bar.utoronto.ca/efpsoybean/cgi-bin/efpWeb.cgi

139    comprising RNA-seq data from soybean root hair and other tissues (20, 21). Soybean genome

140    sequence     assembly     version     7.0     (Gmax_109_gene.gff3.gz";     ftp://ftp.jgi-

141    psf.org/pub/compgen/phytozome/v9.0/Gmax/annotation/) was used for gene annotation and

142    Arabidopsis orthologs information.

143

144    **Cataloging TFs enriched in root lateral organ development stages in soybean**

145       To achieve our objective of identifying regulator TFs and prediction of GRNs associated

146    with root nodules, we used soybean transcription factor annotations from the Plant transcription

147    factor database (PlantTFDB v3.0; http://planttfdb.cbi.pku.edu.cn/) (22) as a starting point.

148    Among 58 TF families annotated in soybean, 48 TF families had at least one member

149    differentially expressed in at least one of the four organ tissue types. For each TF family, we

150    summed the unique transcripts that were enriched in EN and/or MN to calculate the total number

151    of family members enriched in nodule tissues. Similarly, we calculated the number of TFs

152    enriched in lateral root tissues. By comparing the number of family members enriched in nodule

153    vs. lateral root tissues, we identified nodule-specific or -enriched, lateral root-specific or -

154    enriched, and lateral organ non-specific (equal number of transcripts in lateral root and nodules)

155    TF families (Figure 1; Table S1). Statistical analysis (Fisher's Exact test, P<0.05) of nodule- vs.

156    lateral root- specific enrichment showed that TALE, MYB-related, MIKC, C2H2, bZIP, G2-like,

157    WRKY, and NFYB were either nodule-specific or significantly enriched in nodules (Figure 2).

158    Overall, very distinct families of TFs appear to be active in nodule and lateral roots despite

159    reported morphological similarities between these organs.

160

161       We selected a set of 294 TFs, which were differentially expressed and specifically

162    enriched in EN, and MN tissues in our dataset as possible regulators (see Results, Supplementary

163    Table S1). This approach led us to focus on regulators and their GRNs acting specifically during

164    nodule development. We also included 22 previously characterized TFs/ regulator genes reported

165    elsewhere in literature for their respective role in root lateral organ development in model crop

166    plants as positive control marker genes for validation and relevancy of parameters

167    (Supplementary Table S2). For example, ENOD40, FWL1, LBC_A, LBC_C1, LBC_C2, and

168    LBC_C3 genes were used as marker genes, and NIN1 and NSP1 were used as marker regulators

169    for nodule development. ARF5, CRF2, GATA23, LRP1, and TMO7 genes were used as marker

170    regulators for lateral root development. Together, we used 316 TFs of interest as a starting point

171    for the identification of GRNs.

172

173    **Initial biclustering of transcriptome data**

174    We utilized normalized expression values of all the 113,210 gene transcripts in 24

175    libraries for initial biclustering, rather than only significantly differentially expressed gene

176    (DEGs) transcripts. We reasoned that irrespective of enrichment, the TFs and their target gene

177    clusters tend to have similar expression patterns in the root lateral organs, making this an

178    unbiased approach. We chose biclustering (two-way clustering), over traditional clustering for

179    simultaneously clustering using QUBIC (QUalitative Biclustering) (15) to identify all the

180    statistically significant biclusters (BCs) of target genes with TFs, if any as well as samples from

181    the above transcriptome data. Different combinations of QUBIC's parameters were tuned to

182    optimize biclustering to retain the majority of TFs while keeping the total number of transcripts

183    to the minimum. The program first discretizes the data using the parameters $q$ and $r$ and then a

184    heuristic algorithm applied to identify biclusters, where $q$ is the proportion of affected expression

185    data under all conditions for each gene; and $r$ represents the rank of the regulating conditions

186    detected by the parameter $q$. It is suggested to select a smaller $q$ to focus on a local regulator

187    (15). Parameter $f$ controls the overlap between different BCs, and $k$ controls the minimum

188    number of samples in BCs. Another important parameter $c;$ which controls the level of

189    consistency in BCs, was tested to balance the number of TFs and a total number of genes

190    covered in BCs. We obtained 219 BCs that contained 240 of the 316 TFs (76%) and 30, 639 out

191    of 113,210 transcripts (~27%; See Results for details). This "filtered" dataset was used for

192    regulator and GRN prediction. All programs were tested and implemented on a Linux server

193    with Intel x86-64 processor and 32 cores with 1TB RAM configuration.

194

195    **Prediction of potential TF regulators and their GRN inference**

196    To improve the confidence of regulator and GRN prediction, we utilized two module-

197    based GRN inference methods: Lemon-Tree (v.3.0) (17) and Inferelator (v.2015.08.05) (23). We

198    compared and scored the regulatory prediction made by both methods to select high confidence

199    regulators and their target genes in GRN.

200

201    **Lemon-Tree**

202    Lemon-Tree has the option to integrate cluster information; hence, we ran it in two

203    modes: (i) where clusters were generated by Lemon-Tree from the "filtered" dataset (mode I)

204    and (ii) where BC information from QUBIC was fed to Lemon-Tree as co-expressed gene

205    modules for GRN inference (mode II). For mode I, we ran ten independent Gibbs sampler runs

206    of Lemon-Tree (with default parameters) to identify the most confident regulatory modules and

207    TF regulators. The results were used to extract representative module solution (tight clusters)

208    from an ensemble of all possible statistical models using the Gibbs sampler method. Lemon-Tree

209    modules are clustered (hierarchical tree) based on samples with similar mean and standard

210    deviation. This tight cluster corresponds to sets of genes, frequently associated across all

211    clustering solutions. For mode II, we prepared this tight cluster dataset using BCs information

212    from QUBIC, but otherwise used the same settings used for mode I.

213

214    In the next step, the Lemon-Tree algorithm provides a list of weighted TFs with a ranked

215    probability score, and the top 1% were selected as true regulators for each cluster of co-

216    expressed genes. A global score reflecting the statistical confidence of the regulator assigned to

217    each node in a hierarchal tree manner for each set of co-expressed genes modules. The regulator

218    score takes into account the number of trees a regulator is assigned to, with what score (posterior

219    probability), and at which level of the tree (24). An empirical distribution of scores for randomly

220    assigned regulators-to-module is also provided to assess significance (17). In this dataset, the

221    lowest score of a regulator in the top 1% list was at least 3x higher than that of the highest score

222    for a randomly assigned regulator (See Result section for details). Therefore, either the top 1% or

223    at least a 3-fold higher score than randomly assigned regulators appears to be a good threshold to

224    determine true regulators.

225

226    **Inferelator**

227    Inferelator (20 bootstraps) with default settings was utilized to build regulatory networks.

228    Similar to Lemon-Tree, it also uses the gene expression matrix to predict the regulator TFs and

229    their target genes. However, unlike Lemon-Tree, Inferelator does not take cluster information as

230    input, but generates its own clusters. The program generated a ranked list of target genes for each

231    regulator TF utilizing the gene expression matrix and the TFs of our interest. Unlike Lemon-

232    Tree, there is no "score-based" selection of TFs in Inferelator, while there are score-based

233    regulatory interactions between TF and their target genes. Inferelator-generated scores (s) for TF

234    ($x$) regulating gene ($g$) using input gene expression matrix ($RNA$-$seq$) as:

$$s(x \to g|RNAseq) = Inferelator(x \to g|RNAseq) \, X \, sign(cor(x,g))$$

235    where a regulatory interaction confidence score is multiplied by the sign of the correlation

236    coefficient between the TF and the putative target gene to differentiate putative activating from

237    repressing interactions (positive and negative scores, respectively) (18, 25).

238

239    **Combined scoring of regulatory predictions for consensus GRN**

240        By taking advantage of the top regulator prediction feature of Lemon-Tree and top-

241    ranked regulatory target prediction of Inferelator, we compared and combined TF and targeted

242    module genes from all three-inference solutions: Lemon-Tree mode I, II, and Inferelator

243    (described above). The regulatory TFs and corresponding target genes common among all three-

244    inference solutions using Linux "*comm*" command, were rated as potential consensus regulators

245    and their targeted GRN interactions. Ranked score function for every predicted regulatory

246    interaction was calculated by normalizing scores produced by each inference solution (score

247    divided by the highest score in each inference solution) and then averaging normalized score

248    calculated from all three-inference solutions. These ranked scores were used to select high

249    confidence candidate TF-target interactions. These were showed in edges in the GRN modules,

250    visualized and analyzed using Cytoscape (version 3.3.0) (26).

$$Average \, score, As = \frac{\sum Ns(L - mode \, I), Ns(L - mode \, II), Ns(Inferelator)}{3}$$

$$where, Ns = \frac{x}{X}$$

251    Ns=normalized score

252    x = probabilistic score from each mode

253    X= maximum score in each mode

254    L-mode I = Lemon-Tree mode I

255    L-mode II = Lemon-Tree mode II

256

**Protein-protein interaction (PPI) network evidence for physical interaction**

258    Most eukaryotic TFs recruit various co-factors for their DNA-binding specificities and

259    regulatory activities through PPIs. To evaluate potential PPIs that are part of the predicted GRNs,

260    a total of 31,932,066 predicted/experimentally validated soybean protein interactions (NCBI

261    taxon-Id:3847) were obtained from the STRING database (version 10.0) (search tool for the PPI

262    network) (27). This database provides information on both experimental and predicted

263    interactions from varied sources based on co-expression, experiments and literature mining, etc.

264    We evaluated and compared if the predicted TFs and targets from the different inference

265    solutions (Lemon-Tree mode I, II and Inferelator) were potential PPI partners using all the

266    31,932,066 STRING PPI interactions in soybean. Non-redundant dataset, ignoring the transcript

267    numbers of TFs, targets (from TF-target interactions) predicted by three individual inference

268    solutions and PPI from STRING were compared using the Linux "*comm*" command to identify

269    TF-target pair common in STRING dataset and their PPI scores.

270

***Cis*-regulatory motif and functional enrichment analysis evidence for direct regulation**

272    *Cis-regulatory* motif enrichment was carried out using potential promoter sequences of

273    target genes for all potential regulator TFs predicted by all three inference solutions (Lemon-

274    Tree mode I, II and Inferelator). Motif enrichment and Gene Ontology were performed by

275    ShinyGO (http://www.ge-lab.org:3838/go/) using p-value cutoff (FDR) < 0.05 to determine

276    regulation and function.

277

278    **Results**

279

**Optimization of QUBIC parameters for initial biclustering**

281    The primary goal for biclustering in our analysis was to optimize the total number of

282    significant BCs; where the majority of the TFs (out of TFs of interest and marker TFs) are

283    retained while keeping the total number of genes to a minimum for true GRN prediction. In order

284    to evaluate this condition, we iterated various runs in several steps to empirically optimize key

285    QUBIC parameters. For example - *q* to focus on a local regulator, and as regulatory networks are

286  quite small networks, we chose smaller $q$ values. To control the overlap by checking the
287  overlapping genes and the number of TFs in between produced BCs, we iterated the run with $f =$
288  0.5 to 0.65 (by 0.1). We used k = 6 presumably to retain at least three replicates each from either
289  early or late developmental stages or from lateral root or nodule tissue types in one BC.
290  Importantly, the consistency level of BCs was tested using parameter "$c$" iterated from $c = 0.5$ to
291  1 (by 0.1) to balance the number of TFs and a total number of genes covered in BCs. We noticed
292  that the lower consistency level "$c$" values led to the increased size of BCs. We evaluated the
293  produced BCs to determine the "$c$" value at which we covered the greatest number of TFs in
294  comparison to a total number of genes without losing much consistency ($c$). At $c = 0.98$, 76% of
295  the TFs of interest were retained with just 27% of the genes covered in BCs (Figure 3).
296  Interestingly, the maximum number of marker TFs (18 out of 22) cataloged for root lateral
297  organs were covered at $c = 0.98$. On the other hand, at the best consistency level (c=1), only
298  three marker TFs were covered in BCs (not shown). Overall, based on results from several
299  iterations and optimizing for the inclusion of greater number of TFs in BCs, we finalized the
300  following parameters: $r = 1$, $q= 0.2$, $c = 0.98$, $o = 500$, $f = 0.25$, $k = 6$; which produced 219
301  statistically significant BCs (Supplementary Table S3). These 219 BCs comprised ~27% (30,
302  639 out of 113,210) of total gene transcripts. Notably, ~76% (240 out of 316 TFs of our interest)
303  of the TFs of interest and marker TFs were retained in 141 of the 219 total BCs produced. The
304  first cluster was the largest cluster with a total of 446 genes. We conclude that the empirical
305  determination of biclustering parameters depending on the biological question and the associated
306  experimental objective is crucial for useful outcomes.
307
308  **Evaluation of QUBIC biclusters using characterized TFs and co-expressed genes from**
309  **public lateral root organ-related datasets**
310  We observed organ-specific bicluster each for lateral root (both ELR and LR; BC001)
311  and nodule (both EN and MN; BC013) tissues that included all three biological replicate samples
312  in one bicluster, suggesting that these are likely to be highly consistent and reproducible. Four
313  BCs each were specific to all three replicates of ELR (BC015, 019, 033 and 101) and MN (044,
314  048, 152 and 155) tissue types (Supplementary Table S3). To test the rationality of BCs, we
315  compared the expression patterns of co-expressed genes with marker TFs in publicly available
316  transcriptome data (19). The transcription factor "*NSP*1 (*Glyma16g01020*)" crucial for nodule

11

317    development was present in BC037 and BC045 (Supplementary Table S3B). BC037 was specific

318    to nodule tissues and comprised of 367 co-expressed genes. Among these, 52% had more than

319    two-fold up-regulation in EN and MN tissues in our RNA-seq data. A marker gene highly

320    enriched in nodule tissues, *ENOD*40 (*Glyma02g04180*), was found in five BCs (BC013, 22, 40,

321    45, 53 and 95) with different combinations of nodule samples clustered together in each BC.  All

322    genes in BC013 that showed specificity for nodule tissue samples with all three replicates in EN

323    and MN in our study. Also, 50% of the genes from this BC showed greater expression in nodule

324    tissue relative to other tissues types in the soybean gene expression atlas (19) (Supplementary

325    Table S4). Gene Ontology (GO) enrichment analysis for this BC showed enrichment of nucleic

326    acid metabolic process GO term with a significant p-value (FDR; 0.02) and molecular function

327    GO term "Purine ribonucleoside triphosphate binding (FDR; 0.05); both of which are associated

328    with biological nitrogen fixation, a process specific to nodule tissues. For example, soybean

329    nodules export nitrogen in the form of ureides (purines) (28). The above observations indicate

330    the appropriate clustering of relevant transcripts and validate the parameters used for clustering.

331    Notably, we observed few novel transcripts and genes with unknown function, co-expressed in

332    the nodule-specific biclusters (Supplementary Table S4). This observation suggests a potential

333    role for these genes in nodule development and offers candidate genes for functional

334    characterization.

335

336         Further, we took advantage of the time course data for IAA-induced lateral root

337    development in *Arabidopsis* (29), to select and evaluate marker genes present in LR-related BCs

338    in soybean. For example, the LR marker TF, *GmTMO7* (*Glyma04g34080*), a potential ortholog

339    of Arabidopsis *TMO7* identified in the above study, was present in BCs 110, 120 and 173

340    (Supplementary Table S3). Of the 113 genes present in BC120, 96 showed coordinated up-

341    regulation with *TMO*7 in LR tissues, whereas 17 showed negative co-expression. Upon

342    comparison with the Arabidopsis LR induction time course dataset (29), we found 15 co-

343    expressed soybean orthologs (13 positively co-expressed and 2 negatively co-expressed). Where

344    seven (out of 13) from positively co-expressed gene orthologs set were mostly induced in the

345    later stage of lateral root development, one (out of two) from negatively co-expressed had down-

346    regulation in a later stage of lateral root development (see marked blue and red box in

347    Supplementary Table S5). The other lateral root marker *LRP*1 was in BC019 that comprised of

12

348    845 genes. Among these genes, 746 were positively and 99 were negatively co-expressed with

349    LRP1 in all three replicates of ELR. Interestingly, 30 (out of 46 matched genes) of the positively

350    co-expressed genes were potential orthologs of Arabidopsis genes that also showed induction

351    during a similar stage of lateral root development (Supplementary Table S5) in the LR induction

352    time course dataset (29). These comparisons enabled us to evaluate the ability of biclustering

353    parameters and GRN algorithms to appropriately identify regulators and regulatory relationships

354    of target genes during root lateral organ development.

355

356    **Regulatory TF and their Gene Regulatory Networks (GRN) related to root lateral organ**

357    **development in soybean**

358    For the prediction of regulators and inference of corresponding GRNs, we utilized only

359    those 141 BCs that contained our TFs of interest and marker TFs (240 TFs) which comprised

360    25.8% (29,270 out of 113,210) of expressed gene transcripts. This approach potentially reduced

361    the computational complexity and time required for modeling GRNs relevant to our study. This

362    sum expression matrix of 29,270 genes and 240 TF genes (Supplementary Table S6) was used as

363    input for GRN inference by Lemon-Tree mode I, mode II and Inferelator.

364

365    Lemon-Tree produced 828 tight clusters in step 1 from the input expression matrix. A

366    higher number of clusters (828 vs. 141 BCs from QUBIC) suggested that Lemon-Tree clusters

367    were relatively more discrete/smaller in comparison to QUBIC BCs. In step 2, two separate

368    options/modes were utilized (See methods and Figure 1). In mode I, we utilized the 828 tight-

369    clustered modules generated by Lemon-Tree (mode I) and in mode II, the 141 BCs produced by

370    QUBIC (mode II). In mode I, 176 TFs were ranked as the top 1% regulators, whereas in mode II,

371    92 TFs were ranked as top 1% regulators (Supplementary Table S7). Score evaluation was

372    performed for top 1% and randomly predicted regulators from both modes. In both the cases, the

373    minimum score for a top regulator (14.22; mode I and 12.13 mode II) was ~3 times higher than

374    the maximal score (4.99; mode I and 4.23; mode II) for a randomly assigned regulator (Figure 4).

375    This suggested that the scores for top regulators are greater than what could be expected by

376    chance. Inferelator algorithm predicted 132 TFs as potential regulators and five predicted groups

377    (Supplementary Table S7). Comparison of 176, 92, and 132 TFs predicted as regulators

378    respectively, by Lemon-Tree mode I, mode II, and Inferelator, revealed that 56 TFs (~27%) were

379   predicted by all three different modes (Figure 5A). We ranked these common 56 TFs as high

380   confidence TF regulators. In addition, ~62% of the TFs predicted as regulator by Lemon-Tree

381   mode I were also identified as regulators by Lemon-tree mode II and/or Inferelator

382   (Supplementary Figure S1A).

383

384        Furthermore, a total of 113,668 non-redundant TF-target regulatory interactions were

385   predicted by all three modes (Lemon-Tree mode I – 26,012, mode II – 95,845 and Inferelator -

386   3,287) (Supplementary Table S8). A higher number of regulatory interactions in Lemon-Tree

387   mode II is likely due to larger BCs produced by QUBIC. There was relatively smaller overlap

388   among the three modes (Supplementary Figure S1B). We evaluated if the known LR and nodule

389   marker TFs were predicted as regulators as a measure of successful TF prediction by the three

390   different modes. Soybean orthologs of lateral root marker TFs, LRP1 (Glyma14g03900), ARF5

391   (Glyma14g40540), CRF2 (Glyma08g02460), and *TMO7* (Glyma04g34080 and

392   Glyma06g20400) were predicted as regulators by all three inference modes. Additional orthologs

393   of ARF5 (Glyma17g37580) and CRF2 (Glyma05g37120) were predicted as regulators by

394   Lemon-Tree mode I and II. However, orthologs of GATA23 (Glyma03g39220,

395   Glyma19g41780), and LRP1 (Glyma02g44860, Glyma07g35780) were not identified as

396   regulators by any of the modes. These four genes were not enriched in LR tissues

397   (Supplementary Table S2) potentially why they were not predicted as a regulator in this dataset.

398   Successful prediction of four of the five LR-associated markers correctly as regulators by all

399   three modes suggested that the pipeline was reliable and would be used in predicting previously

400   unknown regulators of nodule development.

401

402        A number of TFs were demonstrated to play a crucial role in nodule development through

403   genetic evidence from model legumes (4, 30). These include *NODULE INCEPTION* (*NIN*)

404   (RWP-RK family; (31), NODULATION SIGNALING PATHWAY1 and 2 (NSP1 and NSP2;

405   GRAS domain proteins), Nuclear Factor Y (NF-YA1; (32)), Ethylene Response Factors

406   Required for Nodulation (ERN1 and ERN2; AP2/ERF family; (33)), and CYCLOPS (coiled-coil

407   domain protein) (34–37). In addition, a MYB TF that interacts with NSP2, an ARID domain

408   protein that interacts with SymRK, a bHLH and a set of HD-ZIP IIIs involved in nodule vascular

409   development, and a C2H2 Zn finger TF involved in bacteroid development are also known (38).

14

410 A potential soybean ortholog of NIN, Glyma02g48080 (34), belonging to orthogroup OGEF1237

411 was predicted as a regulator by Lemon-Tree mode I. Only one other NIN-like gene in this

412 orthogroup (Glyma04g00210) was included in our list of input TFs based on expression

413 enrichment in nodules, but was not predicted as a regulator by any mode. Two other NIN-like

414 genes outside of this orthogroup (Glyma12g05390 and Glyma01g36360) were predicted to be

415 regulators by Lemon-Tree modes I and II. Nodule-enriched NFY-As (Glyma02g35190 and

416 Glyma10g10240) were identified as regulators by Lemon-Tree mode I and Inferelator. In *Lotus*

417 *japonicus*, two Nuclear Factor-Y (NF-Y) subunit genes, *LjNF-YA1* and *LjNF-YB1*, were

418 identified as transcriptional targets of NIN (39). In agreement, our analysis predicted that one of

419 the soybean NIN-like genes, Glyma12g05390, regulates NF-YA1 (Glyma10g10240; Lemon-

420 Tree mode II) and the other NIN-like gene, Glyma01g36360, regulates NF-YA2

421 (Glyma02g35190; Lemon-Tree mode I; Supplementary Table S7).

422 Two potential orthologs of LjERN1 (Glyma02g08020 and Glyma19g29000) were

423 predicted as regulators by Lemon-Tree modes I and II. Among the major nodulation TFs, only

424 NSP1 was not predicted to be a regulator by our GRN pipeline. In summary, the pipeline

425 correctly predicted known nodulation and LR TFs including the expected relationships between

426 NIN, NF-YA, and ERN1.

427

428 **Putative protein-protein interactions (PPI) identified in root lateral organ-related GRNs**

429 Co-expressed and co-regulated genes have a higher likelihood of having an indirect

430 functional interaction or direct physical interaction (40). Many TFs form a complex with other

431 proteins for proper molecular and cellular activity. PPIs are the physical interactions between

432 two or more proteins which form the crux of a functional protein complex formation (41). To

433 evaluate if potential regulators identified by us undergo PPIs with other co-regulated proteins, we

434 compared all 113,668 unique TF-target predicted regulatory interactions from three modes of

435 GRN inference method against experimentally verified and/or predicted PPIs based on

436 experimental data reported in the STRING database (see methods for details). We identified, 843

437 potential interactions among 69 TFs with PPI confidence scores ranging from 150 to 995

438 (Supplementary Figure S2, Supplementary Table S9). The high scorer (>800) PPIs were

439 observed from Lemon-Tree mode II run. It was previously suggested that a score < 800 were

440 probably false positives that originated from prediction methods (42). Also, the maximum

441  number (~64%) of PPI interactions were identified by Lemon-Tree mode II, while only four PPI
442  were predicted by all three modes (Supplementary Figure S1C). A likely explanation is the
443  comparatively bigger BCs in this mode generated by QUBIC. While overall, in comparison to all
444  predicted interactions by each mode independently, Inferelator had a greater frequency (2%) of
445  interactions in PPI, i.e., out of total predicted 3288, 61 were observed in PPI, followed by
446  Lemon-Tree Mode I (1%) and then mode II (0.65%). Two ARF5 lateral root markers
447  Glyma14g40540 and Glyma17g37580 were predicted to interact with Glyma13g43050 (PPI
448  score 980) and Glyma15g13640 (PPI score 530) present in GRNs predicted by Lemon-Tree
449  mode I and Inferelator respectively. Glyma13g43050 is an ortholog of Arabidopsis IAA28 which
450  has been demonstrated to interact with AtARF5 (43), and this regulatory module plays a key role
451  in lateral root development (44).

452

453  **High confidence TF regulators and their GRNs associated with root lateral organ**
454  **development in soybean**

455      To determine high-confident regulatory interactions and build a consensus GRN, we
456  evaluated if interactions were conserved across all three modes of GRN prediction (Lemon-Tree
457  modes I, II and Inferelator). Results showed that 182 co-regulatory interactions (for 21 TFs) were
458  commonly predicted by all three modes (Figure 4B, Supplementary Table S10). Therefore, for
459  38% of the TFs predicted as a regulator (21 of 56), have also predicted common target genes
460  independently by all three modes. These 21 TFs made independent GRN with their co-regulated
461  target genes (Figure 6). We ranked the consensus interactions by computing the average of the
462  normalized score given by all three GRN inference modes (ranged from min = 0.19, max=0.88)
463  (See materials and methods for full detail). Table 1 shows the score for 21 commons TFs and
464  their common regulatory interaction predicted from different methods (Lemon-Tree mode I,
465  Lemon-Tree mode II and Inferelator). The complete list of modules together with their high-
466  scorer regulators for this study is available in the Supplementary Table S10. Based on the
467  expression of the TF regulator and their predicted target (Figure 7), we categorized GRN
468  enriched in specific lateral organ tissues.

469

470      TF regulators AP2; ANT (AINTEGUMENTA), transcriptional factor B3 family protein,
471  AtGRF5 (Growth-Regulating Factor 5), C3H, AtbZIP52 (*Arabidopsis thaliana* basic leucine

16

472     zipper 52), PC-MYB1, and SHR (Short Root) appear to co-regulate GRN modules during early

473     nodule (EN) development. TF regulators GRAS; scarecrow-like transcription factor 6 (SCL6),

474     LBD41 (LOB Domain-Containing Protein 41), AP2 domain-containing transcription factor

475     TINY, NUC (nutcracker); nucleic acid binding, AtbZIP5 (Arabidopsis thaliana basic leucine-

476     zipper 5), FRU (FER-Like Regulator Of Iron Uptake), ARR18 (Arabidopsis Response Regulator

477     18) and two unknown TF proteins appear to co-regulate GRN modules late during nodule (MN)

478     development. Interestingly, four PPI interaction (out of total 843 PPI network) were also

479     commonly predicted by all three GRN inference networks in our study for LBD41 and FRU in

480     mature nodules (Supplementary Table S10, Supplementary Figure S2B). ARF16 and AUX/IAA-

481     ARF complex were observed for ELR development, whereas TMO7 and ARF10 (Auxin

482     Response Factor 10) co-regulated GRN for YLR development in soybean.

483

484     **Discussion**

485

486         In spite of the economic and environmental importance of biological nitrogen fixation in

487     nodule in soybean, there is still an unanswered question of what key TFs regulate the underlying

488     GRNs in nodules and lateral roots (4). We developed a robust computational framework for

489     GRN construction using genome-scale gene expression data. Specifically, this framework

490     integrates genomic and transcriptomic data to infer the key regulators and GRN associated with

491     nodule development in soybean. The predicted networks consistently included experimentally

492     verified genes, demonstrating the ability of our framework to reveal significant, potentially

493     important GRNs. With a broader impact, the framework can be used as a template for

494     constructing GRNs to address any biological question of interest in any species.

495

496         To reduce the computational complexity and make the predicted regulator TFs and GRNs

497     relevant to our biological question, a biclustering method and a regulatory network inference tool

498     were used, where their parameters were optimized via several iterations for data analysis and

499     modeling. Among existing GRN inference algorithms, Lemon-Tree and Inferelator were

500     successfully applied in different biological questions due to their valued feature i.e. top regulator

501     and top-ranked regulatory target prediction (45–48). Lemon-Tree detects regulatory modules and

502     regulators from gene expression data using probabilistic graphical models (17). Whereas,

17

503     Inferelator learns a system of ordinary differential equations using the Bayesian Best Subset

504     Regression that describes the rate of change in transcription of each gene or gene-cluster, as a

505     function of TFs. It has been shown that predictions made by the Inferelator are highly accurate

506     for top ranking predictions. Stochastic Lemon-Tree and Inferelator perform better if the

507     transcriptional program can be inferred from a pre-specified list of regulators rather than from a

508     full gene list, because erroneous interactions with non-regulators will be eliminated a priori (49).

509     So, we took the differentially expressed TFs and predefined marker TFs with a known role in

510     nodule and LRs to infer GRN.

511

**Novel regulators of nodule development**

513         We distinguished organ (lateral root/ nodule) and/or developmental stage-specific

514     (early/mature) consensus GRNs based on organ-specific enrichment of the TFs, their differential

515     expression and expression pattern of their co-regulated genes in our transcriptome data. In

516     addition, we also employed comparative genomics and information from public tissue atlas and

517     transcriptome data. The analysis correctly predicted four of the five LR regulators with high

518     confidence and known nodulation TFs including the expected relationships between them. For

519     example, the phylogenetic analysis suggested that ERN2 may not be present in legumes that

520     form determinate nodules such as soybean, *L. japonicus*, or common bean (50). The expression

521     of *ERN1* and *ERN2* are under the control of NIN and NF-YA in *Medicago*, a legume that forms

522     indeterminate nodules. In fact, NF-YA binds the promoter of *ERN1* directly regulating its

523     expression in *Medicago*. However, *ERN1* expression does not appear to be regulated by NIN or

524     NF-YA in *L. japonicus* as its expression is not altered in *nin* or *nf-ya* loss of function mutants.

525     Our GRN prediction also did not identify *ERN1* as a target of NF-YA or NIN in soybean. *ERN1*

526     is directly regulated by CYCLOPS in *L. japonicus*. NSP2 and CYCLOPS were not included in

527     the input TF list due to no nodule-specific enrichment and/or incorrect annotation. The inclusion

528     of CYCLOPS in future analyses might reveal regulatory relationships between ERN1 and

529     CYCLOPS in soybean. It remains to be seen if this is conserved among other determinate nodule

530     forming legumes including soybean. Given the reliability of the pipeline in accurately predicting

531     known TFs, we discuss previously unknown regulators of nodule development predicted by the

532     pipeline.

533

534 An identified EN-GRN was enriched with cell division and cycle functions. Three TFs
535 were predicted to drive GRNs specifically associated with emerging nodules, which are soybean
536 orthologues of Arabidopsis ANT (AINTEGUMENTA; At4g37750), AP2/B3 domain
537 transcriptional factor (At5g58280), and AtGRF5 (Growth-Regulating Factor 5). All the three
538 genes are associated with sites of cell proliferation in Arabidopsis. While GRF5 plays a role in
539 cell proliferation during leaf primordia formation and leaf development, ANT is crucial for
540 flower development. At5g58280 shows the highest expression level in the shoot apex,
541 particularly in the central zone. Indeed, it is likely that the soybean TFs associated with EN
542 GRNs direct cell proliferation during early nodule development. Seven other TFs belonging to
543 C3H, bZIP, MYB1, NF-YC, and SHR were also predicted to co-regulate GRN modules in both
544 emerging nodules and emerging lateral roots (Table 1). Soybean ANT ortholog was the regulator
545 with the highest score in our analysis (0.8) and was predicted to co-regulate ten target genes
546 specifically in emerging nodules. Its targets included *ATCSLA*09, *ALDH*2C4, *GCL*1 (*GCR*2-
547 LIKE 1), *AAP*6, and auxin-responsive protein. A maximum of 51 co-regulated target genes were
548 predicted for a C3H TF regulator (enriched in both EN and ELR) by all three modes. Most of the
549 target genes such as glycosyl hydrolase family protein, *CYCA*1;1 (Cyclin A1;1), zinc finger
550 (C3HC4-type RING finger), *CDKB*1, *CMT*3 (chromomethylase 3); DNA (cytosine-5-)-
551 methyltransferase, calmodulin-binding protein-related, *CYC1BAT*; cyclin-dependent protein
552 kinase regulator, mitotic spindle checkpoint protein, putative (*MAD*2), *ATARP*7 (Actin-Related
553 Protein 7); structural constituent of cytoskeleton, kinesin motor protein-related, and *CDC*20.1;
554 signal transducer, were high scoring target genes.

555

556 GO enrichment analysis of genes involved in EN and EN-ELR GRNs showed significant
557 enrichment of regulation of a cell cycle, movement of a cell or subcellular component,
558 microtubule-based movement, cell division, and cell cycle biological process. (Supplementary
559 Table S10). This is consistent with biological processes known to occur early during lateral
560 organ development. *Cis-regulatory* motif GACCGTTA was enriched in the EN related GRN
561 regulated by a Myb/SANT TF (Supplementary Table S10).

562

563 Similarly, MN-GRN involved in mature nodule development was enriched with meristem
564 initiation and growth. Nine TF regulators belonging to GRAS (scarecrow-like transcription

19

565    factor 6, SCL6), LBD41 (LOB Domain-Containing Protein 41), AP2 domain-containing

566    transcription factor TINY, NUC (nutcracker); nucleic acid binding, bZIP5 (Arabidopsis thaliana

567    basic leucine-zipper 5), FRU (FER-Like Regulator Of Iron Uptake), RR18 (Arabidopsis

568    Response Regulator 18), a Myb/SANT-like DNA binding protein, and a SCREAM-like protein

569    appear to co-regulate GRN modules late during nodule (MN) development. Among these TFs,

570    LBD41 had the highest score (0.77). LBD41 was predicted to co-regulate 38 target genes, among

571    which *PDC*2 (pyruvate decarboxylase-2) had the highest normalized score (0.7). Other targets

572    included *PSAT*, *SRO*2 (similar to rcd one 2), MEE14 (maternal effect embryo arrest 14), zinc

573    finger (AN1-like), SNF2, trehalose-6-phosphate phosphatase, hypoxia-responsive family protein,

574    bHLH, wound-responsive family protein, and ASP1 (Aspartate Aminotransferase 1) with

575    normalized score > 0.5 (Figure 7). Arabidopsis LBD41 is associated with hypoxia response and

576    multiple targets predicted for the soybean ortholog of LBD41 in MN were also associated with

577    hypoxia (51). Nodule oxygen concentrations are highly regulated to enable the proper

578    functioning of the oxygen-sensitive nitrogenase enzyme complex. It is tempting to suggest that

579    soybean LBD41 might play a role in regulating response to hypoxia in MN. The Arabidopsis

580    orthologs of SCL-6 a key regulator in MN, play a role in shoot branching by regulating axillary

581    bud development (52). We had previously suggested that nodules and shoot axillary meristems

582    require a similar hormone balance during development. It is possible that some developmental

583    pathways such as those regulated by SCL6 are shared between these organs. Similarly, the role

584    of Arabidopsis NUTCRACKER protein required in periclinal cell divisions (53), that of FRU in

585    uptake of iron (54), and RR18 in positive regulating cytokinin activity (55) are all consistent with

586    biological processes observed in MN tissues (56, 57). GO enrichment analysis for MN-GRN

587    genes showed enrichment of specification of axis polarity, adaxial/abaxial axis specification,

588    meristem initiation, meristem growth and regulation of meristem growth (Supplementary Table

589    S10). While these processes are known to occur in mature nodules, TFs associated with these

590    processes had not been identified previously. Genes involved in MN-GRN had significant

591    enrichment (P-value ≤ 0.05 FDR) for *cis*-regulatory motifs GGGCCCAC, ACCG and TGTCGG

592    in their upstream regulatory regions. These are likely to be regulated by TCP, AP2 and B3 TFs

593    respectively (Supplementary Table S10). The study has revealed potential TFs associated with

594    different functions in nodule development.

595

## Data availability

Gene expression data used to construct gene regulatory networks are available in NCBI Gene Expression Omnibus (GEO), accession number GSE129509. Raw data files are available in NCBI's Sequence Read Archive (SRA) and can be accessed via links available at the GEO record URL: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE129509.

## Conflict of interest

The authors declare no conflict of interest.

## References

1. Eeckhoute,J., Métivier,R. and Salbert,G. (2009) Defining specificity of transcription factor regulatory activities. *J. Cell Sci.*, **122**, 4027–4034.

2. Blais,A. and Dynlacht,B.D. (2005) Constructing transcriptional regulatory networks. *Genes Dev.*, **19**, 1499–1511.

3. Baitaluk,M., Kozhenkov,S. and Ponomarenko,J. (2012) An Integrative Approach to Inferring Gene Regulatory Module Networks. *PLoS ONE*, **7**.

4. Udvardi,M.K., Kakar,K., Wandrey,M., Montanari,O., Murray,J.D., Andriankaja,A., Zhang,J., Benedito,V.A., Hofer,J.M.I. and Chueng,F. (2007) Update on Legume Transcription Factors Legume Transcription Factors⬚: Global Regulators of Plant Development and Response to the Environment 1 [ W ]. In.

629  5. Kaufmann,K., Pajoro,A. and Angenent,G.C. (2010) Regulation of transcription in plants: mechanisms
630      controlling developmental switches. *Nat. Rev. Genet.*, **11**, 830–842.

631  6. Guan,D., Shao,J., Zhao,Z., Wang,P., Qin,J., Deng,Y., Boheler,K.R., Wang,J. and Yan,B. (2014) PTHGRN:
632      unraveling post-translational hierarchical gene regulatory networks using PPI, ChIP-seq and gene
633      expression data. *Nucleic Acids Res.*, **42**, W130–W136.

634  7. Rivas,J.D.L. and Fontanillo,C. (2010) Protein–Protein Interactions Essentials: Key Concepts to Building
635      and Analyzing Interactome Networks. *PLOS Comput. Biol.*, **6**, e1000807.

636  8. Szklarczyk,D., Morris,J.H., Cook,H., Kuhn,M., Wyder,S., Simonovic,M., Santos,A., Doncheva,N.T.,
637      Roth,A., Bork,P., *et al.* (2017) The STRING database in 2017: quality-controlled protein–protein
638      association networks, made broadly accessible. *Nucleic Acids Res.*, **45**, D362–D368.

639  9. Chaturvedi,I., Sakharkar,M.K. and Rajapakse,J.C. (2007) Validation of Gene Regulatory Networks from
640      Protein-Protein Interaction Data: Application to Cell-Cycle Regulation. In *Pattern Recognition in*
641      *Bioinformatics*. Springer, Berlin, Heidelberg, pp. 300–310.

642  10. Sun,N. and Zhao,H. (2009) Reconstructing transcriptional regulatory networks through genomics
643      data. *Stat. Methods Med. Res.*, **18**, 595–617.

644  11. Li,Y. and Jackson,S.A. (2016) Crowdsourcing the nodulation gene network discovery environment.
645      *BMC Bioinformatics*, **17**.

646  12. Zhu,M., Dahmen,J.L., Stacey,G. and Cheng,J. (2013) Predicting gene regulatory networks of soybean
647      nodulation from RNA-Seq transcriptome data. *BMC Bioinformatics*, **14**, 278.

648  13. Adhikari,S., Damodaran,S. and Subramanian,S. (2019) Lateral Root and Nodule Transcriptomes of
649      Soybean. *Data*, **4**, 64.

650  14. Xie,J., Ma,A., Zhang,Y., Liu,B., Cao,S., Wang,C., Xu,J., Zhang,C. and Ma,Q. QUBIC2: a novel and robust
651      biclustering algorithm for analyses and interpretation of large-scale RNA-Seq data.
652      *Bioinformatics*, 10.1093/bioinformatics/btz692.

653  15. Li,G., Ma,Q., Tang,H., Paterson,A.H. and Xu,Y. (2009) QUBIC: a qualitative biclustering algorithm for
654      analyses of gene expression data. *Nucleic Acids Res.*, **37**, e101.

655  16. Zhang,Y., Xie,J., Yang,J., Fennell,A., Zhang,C. and Ma,Q. (2017) QUBIC: a bioconductor package for
656      qualitative biclustering analysis of gene co-expression data. *Bioinformatics*, **33**, 450–452.

657  17. Bonnet,E., Calzone,L. and Michoel,T. (2015) Integrative Multi-omics Module Network Inference with
658      Lemon-Tree. *PLOS Comput. Biol.*, **11**, e1003983.

659  18. Greenfield,A., Madar,A., Ostrer,H. and Bonneau,R. (2010) DREAM4: Combining genetic and dynamic
660      information to identify biological networks and dynamical models. *PloS One*, **5**, e13397.

661  19. Severin,A.J., Woody,J.L., Bolon,Y.-T., Joseph,B., Diers,B.W., Farmer,A.D., Muehlbauer,G.J.,
662      Nelson,R.T., Grant,D., Specht,J.E., *et al.* (2010) RNA-Seq Atlas of Glycine max: A guide to the
663      soybean transcriptome. *BMC Plant Biol.*, **10**, 160.

664   20. Libault,M., Farmer,A., Brechenmacher,L., Drnevich,J., Langley,R.J., Bilgin,D.D., Radwan,O.,
665        Neece,D.J., Clough,S.J., May,G.D., *et al.* (2010) Complete Transcriptome of the Soybean Root
666        Hair Cell, a Single-Cell Model, and Its Alteration in Response to Bradyrhizobium japonicum
667        Infection. *Plant Physiol.*, **152**, 541–552.

668   21. Libault,M., Farmer,A., Joshi,T., Takahashi,K., Langley,R.J., Franklin,L.D., He,J., Xu,D., May,G. and
669        Stacey,G. (2010) An integrated transcriptome atlas of the crop model Glycine max, and its use in
670        comparative analyses in plants. *Plant J. Cell Mol. Biol.*, **63**, 86–99.

671   22. Jin,J., Zhang,H., Kong,L., Gao,G. and Luo,J. (2014) PlantTFDB 3.0: a portal for the functional and
672        evolutionary study of plant transcription factors. *Nucleic Acids Res.*, **42**, D1182–D1187.

673   23. Bonneau,R., Reiss,D.J., Shannon,P., Facciotti,M., Hood,L., Baliga,N.S. and Thorsson,V. (2006) The
674        Inferelator: an algorithm for learning parsimonious regulatory networks from systems-biology
675        data sets de novo. *Genome Biol.*, **7**, R36.

676   24. Joshi,A., De Smet,R., Marchal,K., Van de Peer,Y. and Michoel,T. (2009) Module networks revisited:
677        computational assessment and prioritization of model predictions. *Bioinformatics*, **25**, 490–496.

678   25. Ciofani,M., Madar,A., Galan,C., Sellars,M., Mace,K., Pauli,F., Agarwal,A., Huang,W., Parkurst,C.N.,
679        Muratet,M., *et al.* (2012) A Validated Regulatory Network for Th17 Cell Specification. *Cell*, **151**,
680        289–303.

681   26. Shannon,P., Markiel,A., Ozier,O., Baliga,N.S., Wang,J.T., Ramage,D., Amin,N., Schwikowski,B. and
682        Ideker,T. (2003) Cytoscape: a software environment for integrated models of biomolecular
683        interaction networks. *Genome Res.*, **13**, 2498–2504.

684   27. Szklarczyk,D., Franceschini,A., Wyder,S., Forslund,K., Heller,D., Huerta-Cepas,J., Simonovic,M.,
685        Roth,A., Santos,A., Tsafou,K.P., *et al.* (2015) STRING v10: protein-protein interaction networks,
686        integrated over the tree of life. *Nucleic Acids Res.*, **43**, D447-452.

687   28. Collier,R. and Tegeder,M. (2012) Soybean ureide transporters play a critical role in nodule
688        development, function and nitrogen export. *Plant J.*, **72**, 355–367.

689   29. Lewis,D.R., Olex,A.L., Lundy,S.R., Turkett,W.H., Fetrow,J.S. and Muday,G.K. (2013) A Kinetic Analysis
690        of the Auxin Transcriptome Reveals Cell Wall Remodeling Proteins That Modulate Lateral Root
691        Development in Arabidopsis. *Plant Cell*, **25**, 3329–3346.

692   30. Magne,K., Couzigou,J.-M., Schiessl,K., Liu,S., George,J., Zhukov,V., Sahl,L., Boyer,F., Iantcheva,A.,
693        Mysore,K.S., *et al.* (2018) MtNODULE ROOT1 and MtNODULE ROOT2 Are Essential for
694        Indeterminate Nodule Identity. *Plant Physiol.*, **178**, 295–316.

695   31. Schauser,L., Roussis,A., Stiller,J. and Stougaard,J. (1999) A plant regulator controlling development of
696        symbiotic root nodules. *Nature*, **402**, 191.

697   32. Battaglia,M., Rípodas,C., Clúa,J., Baudin,M., Aguilar,O.M., Niebel,A., Zanetti,M.E. and Blanco,F.A.
698        (2014) A Nuclear Factor Y Interacting Protein of the GRAS Family Is Required for Nodule
699        Organogenesis, Infection Thread Progression, and Lateral Root Growth1[C][W][OPEN]. *Plant*
700        *Physiol.*, **164**, 1430–1442.

701   33. Baudin,M., Laloum,T., Lepage,A., Rípodas,C., Ariel,F., Frances,L., Crespi,M., Gamas,P., Blanco,F.A.,
702         Zanetti,M.E., *et al.* (2015) A Phylogenetically Conserved Group of Nuclear Factor-Y Transcription
703         Factors Interact to Control Nodulation in Legumes1[OPEN]. *Plant Physiol.*, **169**, 2761–2773.

704   34. Hayashi,S., Reid,D.E., Lorenc,M.T., Stiller,J., Edwards,D., Gresshoff,P.M. and Ferguson,B.J. (2012)
705         Transient Nod factor-dependent gene expression in the nodulation-competent zone of soybean
706         (Glycine max [L.] Merr.) roots. *Plant Biotechnol. J.*, **10**, 995–1010.

707   35. Heckmann,A.B., Lombardo,F., Miwa,H., Perry,J.A., Bunnewell,S., Parniske,M., Wang,T.L. and
708         Downie,J.A. (2006) Lotus japonicus nodulation requires two GRAS domain regulators, one of
709         which is functionally conserved in a non-legume. *Plant Physiol.*, **142**, 1739–1750.

710   36. Heckmann,A.B., Sandal,N., Bek,A.S., Madsen,L.H., Jurkiewicz,A., Nielsen,M.W., Tirichine,L. and
711         Stougaard,J. (2011) Cytokinin Induction of Root Nodule Primordia in Lotus japonicus Is
712         Regulated by a Mechanism Operating in the Root Cortex. *Mol. Plant. Microbe Interact.*, **24**,
713         1385–1395.

714   37. Singh,S., Katzer,K., Lambert,J., Cerri,M. and Parniske,M. (2014) CYCLOPS, A DNA-Binding
715         Transcriptional Activator, Orchestrates Symbiotic Root Nodule Development. *Cell Host Microbe*,
716         **15**, 139–152.

717   38. Zhu,H., Chen,T., Zhu,M., Fang,Q., Kang,H., Hong,Z. and Zhang,Z. (2008) A Novel ARID DNA-Binding
718         Protein Interacts with SymRK and Is Expressed during Early Nodule Development in Lotus
719         japonicus. *Plant Physiol.*, **148**, 337–347.

720   39. Soyano,T., Kouchi,H., Hirota,A. and Hayashi,M. (2013) NODULE INCEPTION Directly Targets NF-Y
721         Subunit Genes to Regulate Essential Processes of Root Nodule Development in Lotus japonicus.
722         *PLOS Genet.*, **9**, e1003352.

723   40. Xulvi-Brunet,R. and Li,H. (2010) Co-expression networks: graph properties and topological
724         comparisons. *Bioinformatics*, **26**, 205–214.

725   41. Barabasi,A.-L. and Oltvai,Z.N. (2004) Network biology: understanding the cell's functional
726         organization. *Nat. Rev. Genet.*, **5**, 101–113.

727   42. Isik,Z., Baldow,C., Cannistraci,C.V. and Schroeder,M. (2015) Drug target prioritization by perturbed
728         gene expression and network information. *Sci. Rep.*, **5**.

729   43. De Rybel,B., Vassileva,V., Parizot,B., Demeulenaere,M., Grunewald,W., Audenaert,D., Van
730         Campenhout,J., Overvoorde,P., Jansen,L., Vanneste,S., *et al.* (2010) A Novel Aux/IAA28 Signaling
731         Cascade Activates GATA23-Dependent Specification of Lateral Root Founder Cell Identity. *Curr.*
732         *Biol.*, **20**, 1697–1706.

733   44. Rogg,L.E., Lasswell,J. and Bartel,B. (2001) A gain-of-function mutation in IAA28 suppresses lateral
734         root development. *Plant Cell*, **13**, 465–480.

735   45. Dolinski,K. and Troyanskaya,O.G. (2015) Implications of Big Data for cell biology. *Mol. Biol. Cell*, **26**,
736         2575–2578.

737   46. Finkle,J.D., Wu,J.J. and Bagheri,N. (2018) Windowed Granger causal inference strategy improves
738        discovery of gene regulatory networks. *Proc. Natl. Acad. Sci. U. S. A.*, **115**, 2252–2257.

739   47. Michoel,T., De Smet,R., Joshi,A., Van de Peer,Y. and Marchal,K. (2009) Comparative analysis of
740        module-based versus direct methods for reverse-engineering transcriptional regulatory
741        networks. *BMC Syst. Biol.*, **3**, 49.

742   48. Vermeirssen,V., Joshi,A., Michoel,T., Bonnet,E., Casneuf,T. and Van de Peer,Y. (2009) Transcription
743        regulatory networks in Caenorhabditis elegans inferred through reverse-engineering of gene
744        expression profiles constitute biological hypotheses for metazoan development. *Mol. Biosyst.*, **5**,
745        1817–1830.

746   49. De Smet,R. and Marchal,K. (2010) Advantages and limitations of current network inference methods.
747        *Nat. Rev. Microbiol.*, **8**, 717–729.

748   50. Kawaharada,Y., James,E.K., Kelly,S., Sandal,N. and Stougaard,J. (2017) The Ethylene Responsive
749        Factor Required for Nodulation 1 (ERN1) Transcription Factor Is Required for Infection-Thread
750        Formation in Lotus japonicus. *Mol. Plant. Microbe Interact.*, **30**, 194–204.

751   51. Gasch,P., Fundinger,M., Müller,J.T., Lee,T., Bailey-Serres,J. and Mustroph,A. (2016) Redundant ERF-
752        VII Transcription Factors Bind to an Evolutionarily Conserved cis-Motif to Regulate Hypoxia-
753        Responsive Gene Expression in Arabidopsis. *Plant Cell*, **28**, 160–180.

754   52. Wang,L., Mai,Y.-X., Zhang,Y.-C., Luo,Q. and Yang,H.-Q. (2010) MicroRNA171c-Targeted SCL6-II, SCL6-
755        III, and SCL6-IV Genes Regulate Shoot Branching in Arabidopsis. *Mol. Plant*, **3**, 794–806.

756   53. Long,Y., Smet,W., Cruz-Ramírez,A., Castelijns,B., de Jonge,W., Mähönen,A.P., Bouchet,B.P.,
757        Perez,G.S., Akhmanova,A., Scheres,B., *et al.* (2015) Arabidopsis BIRD Zinc Finger Proteins Jointly
758        Stabilize Tissue Boundaries by Confining the Cell Fate Regulator SHORT-ROOT and Contributing
759        to Fate Specification. *Plant Cell*, **27**, 1185–1199.

760   54. Jakoby,M., Wang,H.-Y., Reidt,W., Weisshaar,B. and Bauer,P. (2004) FRU (BHLH029) is required for
761        induction of iron mobilization genes in Arabidopsis thaliana. *FEBS Lett.*, **577**, 528–534.

762   55. Veerabagu,M., Elgass,K., Kirchler,T., Huppenberger,P., Harter,K., Chaban,C. and Mira-Rodado,V.
763        (2012) The Arabidopsis B-type response regulator 18 homomerizes and positively regulates
764        cytokinin responses. *Plant J.*, **72**, 721–731.

765   56. Breakspear,A., Liu,C., Roy,S., Stacey,N., Rogers,C., Trick,M., Morieri,G., Mysore,K.S., Wen,J.,
766        Oldroyd,G.E.D., *et al.* (2014) The root hair 'infectome' of Medicago truncatula uncovers changes
767        in cell cycle genes and reveals a requirement for Auxin signaling in rhizobial infection. *Plant Cell*,
768        **26**, 4680–4701.

769   57. Reid,D., Nadzieja,M., Novák,O., Heckmann,A.B., Sandal,N. and Stougaard,J. (2017) Cytokinin
770        Biosynthesis Promotes Cortical Cell Responses during Nodule Development. *Plant Physiol.*, **175**,
771        361–375.

772
773

774  **Figures legends**

775  **Figures**

776

777  **Figure 1.** Schematic representation showing our workflows for prediction of regulator

778  transcription factors (TFs) and their Gene Regulatory Networks (GRNs) for root lateral organ

779  development in soybean.

780

781  **Figure 2.** Transcription factor (TF) families enriched in specific root lateral organs. Bar graphs

782  indicated the number of family members enriched in nodules (blue) or lateral roots (orange). TF

783  annotations are based on Plant Transcription Factor Databases (PlantTFDB). Asterisks indicate

784  TF families that were significantly enriched either in nodule or lateral root (Fisher's exact test; P

785  < 0.05).

786

787  **Figure 3.** Optimization of QUBIC parameter. Relationship between QUBIC's consistency

788  parameter "$c$" (tested from 1 to 0.6) and the number of target transcription factors (TFs) included

789  in bicluster (BC) versus the size of the BC (total number of genes). Each block denotes –c value,

790  TF included in BCs, and total number of genes at that "$c$" value. The optimal "$c$" value selected

791  for final analysis is highlighted.

792

793  **Figure 4.** Distribution of Lemon-Tree scores of true and random regulators for root lateral organ

794  development in soybean. Histogram shows the distribution of score for true and randomly

795  assigned regulator from Lemon-Tree mode I (orange) and mode II (green) produced network.

796  Arrows indicate the minimum and maximum scores from each category with values in

797  parenthesis.

798

799  **Figure 5.** Overlap and differences among outputs from the three different network approaches.

800  (A) Transcription factors (TF) predicted as regulators from three different network approaches.

801  (B) Regulatory interactions predicted by three different network approaches. Numbers in center

802  indicate the number of potential regulators (in A) and interactions (in B) recovered by all the

803  three different approaches playing role in root lateral organ development in soybean.

804

26

805    **Figure 6.** Figure showing consensus 182 co-regulatory interactions predicted and recovered by

806    three different modes chosen in this study. Nodes in diamond denote regulator transcription

807    factors (TFs) and circles denote predicted target genes. Edges denote the normalized score of

808    interaction calculated by all three different modes. Broader the edges, higher the interaction

809    score.

810

811    **Figure 7.** Heat map showing normalized expression from varied samples of root lateral organ

812    development in soybean for regulator transcription factors (TFs) and their co-regulatory target

813    genes in consensus network predicted by three different modes chosen in our study. Row

814    annotation for 21 regulator TFs and their co-regulatory partners are shown in different colors.
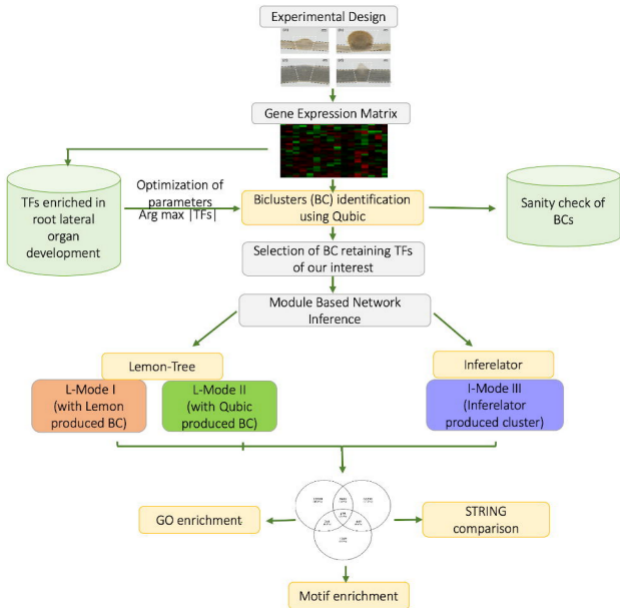
815

816

817 **Table1.** List of transcription factors predicted as regulator by all three workflows used in our
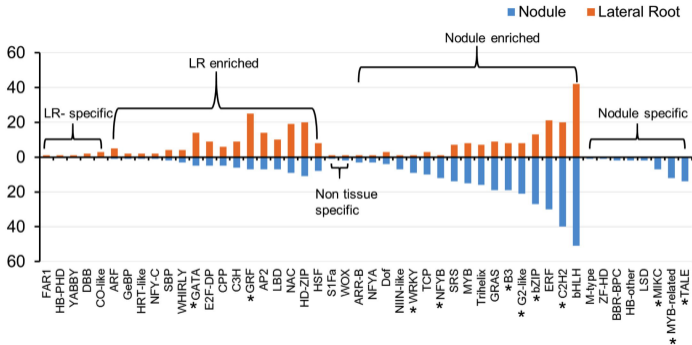
818 study.

| 21 TFs IDs | TF annotation | Enrichment (log$_2$ fold change) in each organ | | | |
|---|---|---|---|---|---|
| | | EN | MN | ELR | YLR |
| Glyma03g27050 | AP2 domain-containing protein (TINY) | | 2.32 | | |
| Glyma17g08380 | ARR18 (RESPONSE REGULATOR 18) | | 2.96 | | |
| Glyma11g04920 | AtbZIP5 (basic leucine-zipper 5) | | | | |
| Glyma13g39650 | FRU (FER-LIKE REGULATOR OF IRON UPTAKE) | | 1.8 | -1.74 | -3.04 |
| Glyma03g03760 | GRAS TF; scarecrow-like 6 (SCL6) | | 2.29 | 1.22 | |
| Glyma19g06280 | LBD41 (LOB DOMAIN-CONTAINING PROTEIN 41) | | 1.19 | | |
| Glyma06g44080 | NUC (nutcracker) | | 1.53 | | |
| Glyma03g34730 | Putative transcription factor | | 2.49 | | |
| Glyma01g32130 | Unknown protein | | 2.45 | -0.87 | |
| Glyma06g05170 | AP2; ANT (AINTEGUMENTA) | 1.42 | -2.23 | | |
| Glyma09g07990 | AtGRF5 (GROWTH-REGULATING FACTOR 5) | 3.34 | | | |
| Glyma02g40400 | Transcriptional factor B3 family protein | 2.76 | | | |
| Glyma14g38460 | AtbZIP52 (basic leucine zipper 52) | 1.51 | | 1.25 | |
| Glyma16g01296 | C3H | 2.01 | | 2.17 | |
| Glyma05g22460 | SHR (SHORT ROOT) | 1.78 | | 1.55 | |
| Glyma06g08660 | PC-MYB1 | 1.4 | | 1.42 | |
| Glyma11g37130 | NFYC | 3.57 | | 1.49 | |
| Glyma11g20490 | ARF10 (AUXIN RESPONSE FACTOR 10) | | | 1.91 | 2.7 |
| Glyma06g20400 | bHLH family protein | | | 2.35 | |
| Glyma10g06080 | ARF16 (AUXIN RESPONSE FACTOR 16) | | | | |
| Glyma19g36571 | AUX/IAA-ARF complex | | | | |

819

820
821

Legend: ■ Nodule  ■ Lateral Root

Figure showing a scatter plot with "Number of transcription factors in biclusters" on the x-axis (ranging from 0 to 350) and "Number of genes in biclusters" on the y-axis (ranging from 0 to 120,000).

Data points labeled:
- -c 0.6, 316, 99898
- -c 0.7, 316, 98562
- -c 0.8, 315, 95602
- -c 0.9, 275, 51546
- -c 0.92, 242, 40304
- -c 0.91, 259, 44243
- -c 0.93, 251, 39168
- -c 0.96, 240, 36639
- -c 0.95, 250, 37783
- -c 0.97, 233, 35040
- -c 0.99, 209, 27398
- -c 0.94, 243, 37975
- -c 0.98, 240, 30639
- -c 1, 39, 5213

(A)

100 ▼ (95.8)
80
60
40
20 ▼ (14.2)
0
True regulators L-mode I

(B)

6
5 ▼ (4.9)
4
3
2
1 ▼ (2.4)
0
Random regulators L-mode I

(C)

100
80 ▼ (85.4)
60
40
20 ▼ (12.1)
0
True regulators L-mode II

(D)

6
5
4 ▼ (4.2)
3
2
1 ▼ (1.8)
0
Random regulators L-mode II

(A) Predicted TF regulators

Inferelator

24

49    3

56

45    26    7

Lemon-Tree mode I    Lemon-Tree mode II

(B) Predicted TF-target gene interaction

Inferelator

2409

255    441

182

21071    4504    84806

Lemon-Tree mode I    Lemon-Tree mode II

**TFmoduleAno**

- AP2_ANT
- AP2_TINY
- ARF10
- ARF16
- ARR18
- AtbZIP5
- AtbZIP52
- AtGRF5
- AUX/IAA-ARF
- C3H
- FRU
- GRAS
- LBD41
- NFYC
- NUC
- PC-MYB1
- SHR
- TMO7
- transcription_factorGT-2
- transcriptional_factorB3
- unknown_protein

TFmoduleAno

ELR_Rep1
ELR_Rep2
ELR_Rep3
EN_Rep3
LR_Rep2
EN_Rep1
EN_Rep2
LR_Rep1
LR_Rep3
MN_Rep1
MN_Rep2
MN_Rep3
ABEN_Rep3
ABEN_Rep2
ABMN_Rep1
ABEN_Rep1
ABMN_Rep2
ABMN_Rep3
ABLR_Rep1
ABLR_Rep3
ABELR_Rep2
ABLR_Rep2
ABELR_Rep1
ABELR_Rep3