

1 **Genomic Bayesian confirmatory factor**  
2 **analysis and Bayesian network to characterize**  
3 **a wide spectrum of rice phenotypes**

4 Haipeng Yu<sup>1</sup>, Malachy Campbell<sup>2</sup>, Qi Zhang<sup>3</sup>, Harkamal Walia<sup>2</sup>, and Gota  
5 Morota<sup>1</sup>

6 <sup>1</sup>Department of Animal and Poultry Sciences, Virginia Polytechnic Institute  
7 and State University, Blacksburg, VA 24061

8 <sup>2</sup>Department of Agronomy and Horticulture, University of  
9 Nebraska-Lincoln, Lincoln, NE 68583

10 <sup>3</sup>Department of Statistics, University of Nebraska-Lincoln, Lincoln, NE  
11 68583

12 Keywords: Bayesian network, factor analysis, multi-phenotypes, rice

13

14 Running title: Network analysis in rice

15

16 Corresponding author:

17 Gota Morota

18 Department of Animal and Poultry Sciences

19 Virginia Polytechnic Institute and State University

20 Blacksburg, VA 24061, USA.

21 E-mail: morota@vt.edu

22

## 23 Abstract

24 Drawing biological inferences from large data generated to dissect the genetic basis of com-  
25 plex traits remains a challenge. Since multiple phenotypes likely share mutual relationships,  
26 elucidating the interdependencies among economically important traits can accelerate the  
27 genetic improvement of plants and animals. A Bayesian network depicts a probabilistic di-  
28 rected acyclic graph representing conditional dependencies among variables. This study aims  
29 to characterize various phenotypes in rice (*Oryza sativa*) via confirmatory factor analysis and  
30 Bayesian network. Confirmatory factor analysis under the Bayesian treatment hypothesized  
31 that 48 observed phenotypes resulted from six latent variables including grain morphology,  
32 morphology, flowering time, physiology (e.g., ion content), yield, and morphological salt re-  
33 sponse. This was followed by studying the genetics of each latent variable. Bayesian network  
34 structures involving the genomic component of six latent variables were established by fitting  
35 four different algorithms. Negative genomic correlations were obtained between salt response  
36 and yield, salt response and grain morphology, salt response and physiology, and morphology  
37 and yield, whereas a positive correlation was obtained between yield and grain morphology.  
38 There were four common directed edges across the different Bayesian networks. Physiolog-  
39 ical components influenced the flowering time and grain morphology, and morphology and  
40 grain morphology influenced yield. This work suggests that the Bayesian network coupled  
41 with factor analysis can provide an effective approach to understand the interdependence  
42 patterns among phenotypes and to predict the potential influence of external interventions or  
43 selection related to target traits in the high-dimensional interrelated complex traits systems.

## 44 Introduction

45 Genetic correlation constitutes a major aspect of quantitative genetics (Lush, 1948; Falconer  
46 and Mackay, 1996). In its simplest form, a single gene or mutation may affect several bio-  
47 logical pathways leading to correlated phenotypes. This phenomenon, known as pleiotropy,  
48 induces genetic correlations among multiple traits at the population level. In plant and  
49 animal breeding, more than one phenotype is generally assessed to account for the overall  
50 performance of individuals. Because multiple phenotypes may exhibit mutual relationships,  
51 knowledge of the interdependence among economically important traits can bring more ef-  
52 fective selection and genetic improvement in systems with complex traits. In a standard  
53 quantitative genetic analysis, multivariate phenotypes can be modeled through multi-trait  
54 models (MTM) of Henderson and Quaas (1976) or some genomic counterparts (e.g., Calus  
55 and Veerkamp, 2011; Jia and Jannink, 2012) by leveraging genetic or environmental corre-  
56 lations among traits. In particular, MTM has been useful in deriving genetic correlations  
57 and enhancing the prediction accuracy of breeding values for traits with low heritability via  
58 joint modeling with one or more genetically correlated, highly heritable traits (Mrode, 2014).  
59 However, genetic selection for breeding requires causal assumptions, as the effects of exter-  
60 nal interventions on interrelated complex traits cannot be predicted on the basis of these  
61 associations (Pearl, 2009). This modeling step is essential to verify that predictors consid-  
62 ered for selection accurately reflect genetic causal effects (Valente et al., 2015). Although  
63 Bayesian network (BN) analysis or causal structure inference from observational data has  
64 been an active research area in plant and animal breeding (Valente et al., 2010; Töpner et al.,  
65 2017), the primary challenge associated with multivariate analysis is that computation can be  
66 untenable. This is because the number of estimated parameters within the model increases  
67 with the increasing number of phenotypes and the difficulty of interpreting interrelationships  
68 among multiple phenotypes. This is a particularly persistent challenge in plant breeding,  
69 owing to the availability of high-dimensional and diverse phenotypes currently being gener-

70 ated via high-throughput, image-based phenomics platforms in addition to the conventional,  
71 non-image phenotypes (Awada et al., 2018).

72 One approach to characterize high-dimensional phenotypes is by using factor analysis,  
73 which facilitates modeling correlated responses through underlying unobserved latent vari-  
74 ables, which are also known as factors or modules (de los Campos and Gianola, 2007).  
75 Confirmatory factor analysis, a variant of factor analysis, hypothesizes that observed pheno-  
76 types result from lower-dimensional latent variables specified by prior biological knowledge  
77 (Jöreskog, 1969). These latent variables underlie observed phenotypes and can be evalu-  
78 ated for how well the data support the hypothesis. For instance, Peñagaricano et al. (2015)  
79 performed confirmatory factor analysis in swine to derive five latent variables from 19 pheno-  
80 typic traits and inferred BN structures among those latent variables, thereby demonstrating  
81 the potential of this approach.

82 This study aimed to obtain a first glimpse of the utility of graphical modeling to char-  
83 acterize a wide range of phenotypes in rice by studying the genetics of each latent variable.  
84 First, we constructed latent variables, using prior biological knowledge obtained from the  
85 literature. Then we connected the observed high-dimensional phenotypes with these to  
86 establish latent variables via Bayesian confirmatory factor analysis (BCFA) to reduce the  
87 dimensions of the dataset. Further, factor scores computed from BCFA were considered  
88 new phenotypes for a Bayesian multivariate analysis to separate breeding values from noise.  
89 This was followed by adjustment of breeding values via Cholesky decomposition to eliminate  
90 the dependencies introduced by genomic relationships. Finally, the adjusted breeding values  
91 were considered inputs to assess the causal network structure between latent variables by  
92 conducting a Gaussian BN analysis. This study is the first, to our knowledge, in rice to  
93 characterize various phenotypes with graphical modeling such as BCFA and BN.

## 94 **Materials and Methods**

### 95 **Sources of phenotypic and genotypic data**

96 The rice dataset comprised  $n = 413$  accessions sampled from six subpopulations: temperate  
97 japonica (92), tropical japonica (85), indica (77), aus (52), aromatic (12), and admixture  
98 of japonica and indica (56). We used  $t = 48$  phenotypes and data regarding 44,000 single-  
99 nucleotide polymorphisms (SNP). Of those, 34 phenotypic records were reported in Zhao  
100 et al. (2011). The remaining phenotypes were assessed from the abiotic stress experiments  
101 conducted in Campbell et al. (2017a). The detailed descriptions of the phenotypes used  
102 can be found in Zhao et al. (2011) and Campbell et al. (2017a), and are summarized in  
103 Supplementary Table S1. After removing SNP markers with minor allele frequency less than  
104 0.05, 374 accessions and 33,584 markers were used for further analysis.

### 105 **Bayesian confirmatory factor analysis**

A confirmatory factor analysis under the Bayesian framework was performed to model 48 phenotypes. The number of factors and the pattern of phenotype-factor relationships need to be specified in BCFA prior to model fitting. We constructed six latent variables ( $q = 6$ ) from previous reports (Acquaah, 2009; Zhao et al., 2011; Campbell et al., 2017a). The six latent variables derived from our analysis represent the grain morphology, morphology, flowering time, physiology, yield, and salt response (Table S1). Each latent variable captures common signals spanning genetic and environmental effects across all its phenotypes. The latent variables, which determine the observed phenotypes can be modeled as

$$\mathbf{T} = \mathbf{\Lambda F} + \mathbf{s},$$

where  $\mathbf{T}$  is the  $t \times n$  matrix of observed phenotypes,  $\mathbf{\Lambda}$  is the  $t \times q$  factor loading matrix,  $\mathbf{F}$  is the  $q \times n$  latent variables matrix, and  $\mathbf{s}$  is the  $t \times n$  matrix of specific effects. Here,  $\mathbf{\Lambda}$  maps latent variables to the observed variables and can be interpreted as the extent of contribution each latent variable to phenotype. This can be derived by solving the following variance-covariance model.

$$\text{var}(\mathbf{T}) = \mathbf{\Lambda}\mathbf{\Phi}\mathbf{\Lambda}' + \mathbf{\Psi},$$

106 where  $\mathbf{\Phi}$  is the variance of latent variables, and  $\mathbf{\Psi}$  is the variance of specific effects (Brown,  
107 2014). Six latent variables were assumed to account for the covariance in the observed  
108 phenotypes. Moreover, latent variables were assumed to be correlated with each other. Prior  
109 distributions were assigned to all unknown parameters. The non-zero coefficient within factor  
110 loading matrix  $\mathbf{\Lambda}$  was assumed to follow a Gaussian distribution with mean of 0 and variance  
111 of 0.01. The variance-covariance matrix  $\mathbf{\Phi}$  was assigned an inverse Wishart distribution  
112 with a  $6 \times 6$  identity scale matrix  $\mathbf{I}_{66}$  and a degree freedom of 7,  $\mathbf{\Phi} \sim \mathcal{W}^{-1}(\mathbf{I}_{66}, 7)$  and an  
113 inverse Gamma distribution with scale parameter 1 and shape parameter 0.5 was assigned  
114 to  $\mathbf{\Psi} \sim \Gamma^{-1}(1, 0.5)$ .

115 We employed the blavaan R package (Merkle and Rosseel, 2018) jointly with JAGS  
116 (Hornik et al., 2003) to fit the above BCFA. The blavaan runs the runjags R package (Den-  
117 wood, 2016) to summarize the Markov chain Monte Carlo (MCMC) and samples unknown  
118 parameters from the posterior distributions. Three MCMC chains, each of 5,000 samples  
119 with 2,000 burn-in, were used to infer the unknown model parameters. The convergence of  
120 the parameters was investigated with trace plots and potential scale reduction factor (PSRF;  
121 Gelman and Rubin, 1992). The PSRF computes the difference between estimated variances  
122 among multiple Markov chains and estimated variances within the chain. A large difference  
123 indicates non-convergence and may require additional Gibbs sampling.

124 Subsequently, the posterior means of factor scores ( $\mathbf{F}$ ), which reflect the contribution of

125 latent variables to each accession were estimated. Within each draw of Gibbs sampling,  $\mathbf{F}$   
126 was sampled from the conditional distribution of  $p(\mathbf{F}|\boldsymbol{\theta}, \mathbf{T})$ , where  $\boldsymbol{\theta}$  refers to the unknown  
127 parameters in  $\boldsymbol{\Lambda}$ ,  $\boldsymbol{\Phi}$ , and  $\boldsymbol{\Psi}$ . This conditional distribution was derived with data augmenta-  
128 tion (Tanner and Wong, 1987) assuming  $\mathbf{F}$  as missing data (Lee and Song, 2012).

## 129 Multivariate genomic best linear unbiased prediction

We fitted a Bayesian multivariate genomic best linear unbiased prediction to separate breed-  
ing values from population structure and noise in the six factor scores computed previously.

$$\mathbf{F} = \boldsymbol{\mu} + \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon},$$

130 where  $\boldsymbol{\mu}$  is the vector of intercept,  $\mathbf{X}$  is the incidence matrix of covariates,  $\mathbf{b}$  is the vector of  
131 covariate effects,  $\mathbf{Z}$  is the incidence matrix relating accessions with additive genetic effects,  $\mathbf{u}$   
132 is the vector of additive genetic effects, and  $\boldsymbol{\epsilon}$  is the vector of residuals. The incident matrix  
133  $\mathbf{X}$  included subpopulation information (temperate japonica, tropical japonica, indica, aus,  
134 aromatic, and admixture), as the rice diversity panel used herein shows a clear substructure  
135 (Zhao et al., 2011).

A flat prior was assigned to  $\boldsymbol{\mu}$  and  $\mathbf{b}$ , and the joint distribution of  $\mathbf{u}$  and  $\boldsymbol{\epsilon}$  follows  
multivariate normal

$$\begin{pmatrix} \mathbf{u} \\ \boldsymbol{\epsilon} \end{pmatrix} \sim N \left[ \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_u \otimes \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_\epsilon \otimes \mathbf{I} \end{pmatrix} \right],$$

136 where  $\mathbf{G}$  represents the second genomic relationship matrix of VanRaden (2008),  $\mathbf{I}$  is the  
137 identity matrix,  $\boldsymbol{\Sigma}_u$  and  $\boldsymbol{\Sigma}_\epsilon$  refer to  $6 \times 6$  dimensional genetic and residual variance-covariance  
138 matrices, respectively. An inverse Wishart distribution with a  $6 \times 6$  identity scale matrix of  $\mathbf{I}_{66}$   
139 and a degree of freedom 6 was assigned as prior for  $\boldsymbol{\Sigma}_u, \boldsymbol{\Sigma}_\epsilon \sim \mathcal{W}^{-1}(\mathbf{I}_{66}, 6)$ . These parameters  
140 were selected so that relatively uninformative priors were used. The Bayesian multivariate



141 genomic best linear unbiased prediction model was implemented using the MTM R pack-  
142 age (<https://github.com/QuantGen/MTM>). Posterior mean estimates of genetic correlation  
143 between latent variables and predicted breeding values ( $\hat{\mathbf{u}}$ ) were then obtained.

## 144 Sample independence in the Bayesian network

Theoretically, BN learning algorithms assume sample independence. In the multivariate genomic best linear unbiased prediction, the residuals between phenotypes were assumed independent through  $\mathbf{I}_{374 \times 374}$ . However, phenotypic dependencies were introduced by the  $\mathbf{G}$  matrix for the additive genetic effects, thereby potentially serving as a confounder. Thus, a transformation of  $\hat{\mathbf{u}}$  was carried out to derive an adjusted  $\hat{\mathbf{u}}^*$  by eliminating the dependencies in  $\mathbf{G}$ . For a single trait model, the adjusted  $\hat{\mathbf{u}}^*$  can be computed by premultiplying  $\hat{\mathbf{u}}$  by  $\mathbf{L}^{-1}$ , where  $\mathbf{L}$  is a lower triangular matrix derived from the Choleskey decomposition of  $\mathbf{G}$  matrix ( $\mathbf{G} = \mathbf{L}\mathbf{L}'$ ). Since  $\mathbf{u} \sim \mathcal{N}(0, \mathbf{G}\sigma_u^2)$ , the distribution of  $\hat{\mathbf{u}}^*$  follows  $\mathcal{N}(0, \mathbf{I}\sigma_u^2)$  (Vazquez et al., 2010)

$$\begin{aligned} \text{Var}(\mathbf{u}^*) &= \text{Var}(\mathbf{L}^{-1}\mathbf{u}) \\ &= \mathbf{L}^{-1}\text{Var}(\mathbf{u})(\mathbf{L}^{-1})' \\ &= \mathbf{L}^{-1}\mathbf{G}(\mathbf{L}^{-1})'\sigma_u^2 \\ &= \mathbf{L}^{-1}\mathbf{L}\mathbf{L}'(\mathbf{L}')^{-1}\sigma_u^2 \\ &= \mathbf{I}\sigma_u^2. \end{aligned}$$

145 This transformation can be extended to a multi-traits model by defining  $\mathbf{u}^* = \mathbf{M}^{-1}\mathbf{u}$ , where  
146  $\mathbf{M}^{-1} = \mathbf{I}_{\mathbf{qq}} \otimes \mathbf{L}^{-1}$  (Töpner et al., 2017). Under the multivariate framework,  $\mathbf{u}$  follows  
147  $\mathcal{N}(0, \mathbf{\Sigma}_{\mathbf{u}} \otimes \mathbf{G})$  and the variance of  $\mathbf{u}^*$  is

$$\begin{aligned}
 \text{Var}(\mathbf{u}^*) &= \text{Var}(\mathbf{M}^{-1}\mathbf{u}) \\
 &= (\mathbf{I}_{\text{tt}} \otimes \mathbf{L}^{-1})(\boldsymbol{\Sigma}_{\mathbf{u}} \otimes \mathbf{G})(\mathbf{I}_{\text{qq}} \otimes \mathbf{L}^{-1})' \\
 &= (\mathbf{I}_{\text{qq}} \otimes \mathbf{L}^{-1})(\boldsymbol{\Sigma}_{\mathbf{u}} \otimes \mathbf{L}\mathbf{L}')(\mathbf{I}_{\text{qq}} \otimes \mathbf{L}^{-1})' \\
 &= \boldsymbol{\Sigma}_{\mathbf{u}} \otimes \mathbf{I}_{\text{nn}},
 \end{aligned}$$

148 where  $\mathbf{L}^{-1}\mathbf{L}\mathbf{L}'(\mathbf{L}^{-1})' = \mathbf{I}_{\text{nn}}$ . This adjusted  $\hat{\mathbf{u}}^*$  was used to learn BN structures between  
 149 predicted breeding values.

## 150 Bayesian network

A BN depicts the joint distribution of random variables regarding their probabilistic conditional dependencies (Scutari and Denis, 2014)

$$\mathcal{BN} = (\mathcal{G}, X_V),$$

where  $\mathcal{G}$  represents a directed acyclic graph (DAG) =  $(V, E)$  with nodes  $(V)$  connected by one or more edges  $(E)$  conveying the probabilistic relationships and the random vector  $X_V = (X_1, \dots, X_K)$  is  $K$  random variables. The joint probability distribution can be factorized as

$$P(X_V) = P(X_1, \dots, X_K) = \prod_{v=1}^K P(X_v | Pa(X_v)),$$

151 where  $Pa(X_v)$  denotes a set of parent nodes of child node  $X_v$ . The DAG and joint prob-  
 152 ability distribution are governed by the Markov condition, which states that every random  
 153 variable is independent of its non-descendants conditioned on its parents. A BN is known  
 154 as a Gaussian BN, when all variables or phenotypes are defined as marginal or conditional  
 155 Gaussian distribution as in the present study.

156 The adjusted breeding values  $\hat{\mathbf{u}}^*$  were used to infer a genomic network structure among  
157 the aforementioned six latent variables. There are three types of structure-learning algo-  
158 rithms for BN: constraint-based algorithms, score-based algorithms, and a hybrid of these  
159 two (Scutari and Denis, 2014). The constraint-based algorithms can be originally traced  
160 to the inductive causation algorithm (Verma and Pearl, 1991), which uses conditional in-  
161 dependence tests for network inference. Briefly, the first step is to identify a d-separation  
162 set for each pair of nodes and confer an undirected edge between the two if they are not  
163 d-separated. The second step is to identify a v-structure for each pair of non-adjacent nodes,  
164 where a common neighbor is the outcome of two non-adjacent nodes. In the last step, com-  
165 pelled edges were identified and oriented, where neither cyclic graph nor new v-structures  
166 are permitted. The score-based algorithms are based on heuristic approaches, which first  
167 assign a goodness-of-fit score for an initial graph structure and then maximize this score by  
168 updating the structure (i.e., add, delete, or reverse the edges of initial graph). The hybrid  
169 algorithm includes two steps, restrict and maximize, which harness both constrain-based and  
170 score-based algorithms to construct a reliable network. In this study, the two score-based  
171 (Hill Climbing and Tabu) and two hybrid algorithms (Max-Min Hill Climbing and General  
172 2-Phase Restricted Maximization) were used to perform structure learning.

173 We quantified the strength of edges and uncertainty regarding the direction of networks,  
174 using 500 bootstrapping replicates with a size equal to the number of accessions and per-  
175 formed structure learning for each replicate in accordance with Scutari and Denis (2014).  
176 Non-parametric bootstrap resampling aimed at reducing the impact of the local optimal  
177 structures by computing the probability of the arcs and directions. Subsequently, 500 learned  
178 structures were averaged with a strength threshold of 85% or higher to produce a more robust  
179 network structure. This process, known as model averaging, returns the final network with  
180 arcs present in at least 85% among all 500 networks. Candidate networks were compared  
181 on the basis of the Bayesian information criterion (BIC) and Bayesian Gaussian equivalent  
182 score (BGe). The BIC accounts for the goodness-of-fit and model complexity, and BGe aims

183 at maximizing the posterior probability of networks per the data. All BN were learned via  
184 the bnlearn R package (Scutari, 2010). In bnlearn, the BIC score is rescaled by -2, which  
185 indicates that the larger BIC refers to a preferred model.

## 186 **Data availability**

187 Genotypic data regarding the rice accessions can be downloaded from the rice diversity panel  
188 website (<http://www.ricediversity.org/>). Phenotypic data used herein are available in  
189 Zhao et al. (2011) and Campbell et al. (2017b).

## 190 Results

### 191 Latent variable modeling

192 The BCFA model grouped the observed phenotypes into the underlying latent variables  
193 on the basis of prior biological knowledge, assuming these latent variables determine the  
194 observed phenotypes. This allowed us to study the genetics of each latent variable. A  
195 measurement model derived from BCFA evaluating the six latent variables is shown in Figure  
196 1. Forty-eight observed phenotypes were hypothesized to result from the six latent variables:  
197 7 for flowering time, 14 for morphology, 5 for yield, 11 for grain morphology, 6 for physiology,  
198 and 5 for salt response. The convergence of the parameters was confirmed graphically with  
199 the trace plots and a PSRF value less than 1.2 (Merkle and Rosseel, 2018).

200 The six latent factors showed strong contributions to the 48 observed phenotypes, with  
201 standardized regression coefficients ranging from -0.668 to 0.980 for flowering time, -0.112 to  
202 0.903 for morphology, -0.113 to 0.977 for yield, -0.501 to 0.986 for grain morphology, -0.016  
203 to 0.829 for physiology, and 0.011 to 0.929 for salt response. The latent factor flowering time  
204 showed a strong positive contribution to flowering time in Arkansas (Fla) and Flowering  
205 time in Arkansas in 2007 (Fla7), indicating that larger values for the latent factor can be  
206 interpreted as a greater number of days from sowing to emergence of the inflorescence.  
207 The latent factor morphology showed the largest positive contributions to traits describing  
208 height during the vegetative stage (e.g. height to newest ligule in salt (Hls), height to  
209 newest ligule in control (Hlc), height to the tip of first fully expanded leaf in salt (Hfs), and  
210 height to tip of first fully expanded leaf in control (Hfc)), suggesting that this latent factor  
211 is an overall representation of plant size. Yield showed large positive contributions to the  
212 observed phenotypes primary panicle branch number (Ppn) and seed number per panicle  
213 (Snpp), suggesting that larger values for yield indicate a higher degree of branching and seed  
214 number. Observed phenotypes describing seed size (e.g. seed volume (Sv) and brown rice

215 volume (Bvl)) were most strongly associated with grain morphology. The latent factor ionic  
216 components of salt stress showed strong positive contributions to two observed phenotypes  
217 that quantify the ionic components of salt stress (shoot  $\text{Na}^+:\text{K}^+$  (Kslm) and shoot  $\text{Na}^+$   
218 (Nas)), indicating that higher values for the latent factor result in greater shoot  $\text{Na}^+$  and  
219  $\text{Na}^+:\text{K}^+$ . Finally, the latent factor describing morphological salt response showed strong  
220 positive contributions to the observed phenotype describing the effect of salt treatment on  
221 plant height (ratio of height to tip of newest fully expanded leaf in salt to that of control  
222 plants (Hfr)), thus larger values for the latent factor may indicate a more tolerant growth  
223 response to salinity.

## 224 **Genomic correlation among latent variables**

225 To understand the genetic relationships between latent variables, genomic correlation analy-  
226 sis was performed. Genomic correlation is due to pleiotropy or linkage disequilibrium between  
227 quantitative trait locus (QTL). The genomic correlations among latent variables are shown  
228 in Figure 2. Negative correlations were observed between salt response (Slr) and all other  
229 five latent variables. In particular, flowering time (-0.5), yield (-0.54), and grain morphology  
230 (-0.74) were moderately correlated with morphological salt response. These results suggest  
231 that accessions that harbor alleles for more tolerant morphological salt responses may also  
232 have alleles associated with longer flowering times, smaller seeds, and low yield. Similarly,  
233 a moderate negative correlation was observed between morphology and yield (-0.56) and  
234 between morphology and grain morphology (-0.31). Thus, accessions with alleles associated  
235 with large plant size may also have alleles that result in low yield, small grain volume, and  
236 lower shoot  $\text{Na}^+$  and  $\text{Na}^+:\text{K}^+$ . In contrast, a positive moderate correlation was observed  
237 between grain morphology and yield (0.49) and between grain morphology and ionic com-  
238 ponents of salt stress (0.4). Thus, selection for large grain may result in improved yield, and  
239 higher shoot  $\text{Na}^+$  and  $\text{Na}^+:\text{K}^+$ .

## 240 Bayesian network

241 To infer the possible causal structure between latent variables, BN was performed. Prior  
242 to BN, the normality of latent variables was assessed using histogram plots combined with  
243 density curves as shown in Figure S1. Overall, all the six latent variables approximately  
244 followed a Gaussian distribution.

245 The Bayesian networks learned with the score-based and hybrid algorithms are shown  
246 in Figures 3, 4, 5, and 6. The structures of BN were refined by model averaging with 500  
247 networks from bootstrap resampling to reduce the impact of local optimal structures. The  
248 labels of the arcs measure the uncertainty of the arcs, corresponding to strength and direc-  
249 tion (in parenthesis). The former measures the frequency of the arc presented among all 500  
250 networks from the bootstrapping replicates and the latter is the frequency of the direction  
251 shown conditional on the presence of the arc. We observed minor differences in the structures  
252 presented within and across the two types of algorithms used. In general, small differences  
253 were observed within algorithm types compared to those across algorithms. The two score-  
254 based algorithms produced a greater number of edges than two hybrid algorithms. In Figure  
255 3, the Hill Climbing algorithm produced seven directed connections among the six latent  
256 variables. Three connections were indicated towards flowering time from morphological salt  
257 response, ionic components of salt stress, and morphology, and two edges to yield from mor-  
258 phology and from grain morphology. Other two edges were observed from ionic components  
259 of salt stress to grain morphology and from grain morphology to morphological salt response.  
260 A similar structure was generated by the Tabu algorithm, except that the connection be-  
261 tween salt response and grain morphology presented an opposite direction (Figure 4). The  
262 Max-Min Hill Climbing hybrid algorithm yielded six directed edges from morphological salt  
263 response to grain morphology, from ionic components of salt stress to grain morphology, from  
264 ionic components of salt stress to flowering time, from flowering time to morphology, from  
265 morphology to yield, and from grain morphology to yield (Figure 5). An analogous structure

266 with the only difference observed in the directed edge from morphology to flowering time was  
267 inferred with the General 2-Phase Restricted Maximization algorithm as shown in Figure 6.  
268 Across all four algorithms, there were four common directed edges: from ionic components  
269 of salt stress to flowering time and to grain morphology, and from morphology and grain  
270 morphology to yield. The most favorable network was considered the one from the Tabu  
271 algorithm, which returned the largest network score in terms of BIC (1086.61) and BGe  
272 (1080.88). Collectively, these results suggest that there may be a direct genetic influence of  
273 morphology and grain morphology on yield, and physiological components of salt tolerance  
274 on grain morphology and flowering time.



## 275 Discussion

276 This study is based on the premise that most phenotypes interact to greater or lesser de-  
277 grees with each other through underlying physiological and molecular pathways. While these  
278 physiological pathways are important for the development of agronomically important char-  
279 acteristics, they are often unknown or difficult to assess in large populations. The approach  
280 utilized here leverages phenotypes that can be readily assessed in large populations to quan-  
281 tify these underlying unobserved phenotypes, and elucidates the relationships between these  
282 variables.

283 Understanding the behaviors among phenotypes in the complex traits is critical for genetic  
284 improvement of agricultural species (Hickey et al., 2017). Graphical modeling offers an av-  
285 enue to decipher bi-directional associations or probabilistic dependencies among variables of  
286 interest in plant and animal breeding. For instance, BN and L1-regularized undirected net-  
287 work can be used to model interrelationships of linkage disequilibrium (LD) (Morota et al.,  
288 2012; Morota and Gianola, 2013) or phenotypic, genetic, and environmental interactions  
289 (Xavier et al., 2017) in a systematic manner. Importantly, MTM elucidates both direct and  
290 indirect relationships among phenotypes. Inaccurate interpretation of these relationships  
291 may substantially bias selection decisions (Valente et al., 2015; Gianola et al., 2015). Thus,  
292 we applied BCFA to reduce the dimension of the responses by hypothesizing 48 manifest  
293 phenotypes originated from the underlying six constructed latent variables as shown in Fig-  
294 ure 1 assuming that these latent traits are most important, followed by application of BN to  
295 infer the structures among the six biologically relevant latent variables (Figures 3,4, 5, and  
296 6). The BN represents the conditional dependencies between variables. Care must be taken  
297 in interpreting these relationships as a causal effect. Although a good BN is expected to  
298 describe the underlying causal structure per the data, when the structure is learned solely  
299 on the basis of the observed data, it may return multiple equivalent networks that describe  
300 the data well. In practice, searching such a causal structure with observed data needs three

301 additional assumptions (Scutari and Denis, 2014): 1) each variable is independent of its  
302 non-effects (i.e., direct and indirect) conditioned on its direct causes, 2) the probability dis-  
303 tribution of variables is supported by a DAG, where the d-separation in DAG provides all  
304 dependencies in the probability distribution, and 3) no additional variables influence the  
305 variables within the network. Although it may be difficult to meet these assumptions in the  
306 observed data, a BN is equipped with suggesting potential causal relationships among la-  
307 tent variables, which can assist in exploring data, making breeding decisions, and improving  
308 management strategies in breeding programs (Rosa et al., 2011).

## 309 **Biological meaning of latent variables and their relation-** 310 **ships**

311 We performed BCFA to summarize the original 48 phenotypes with the six latent variables.  
312 The number of latent variables and which latent variables load onto phenotypes were deter-  
313 mined from the literature. The latent variable morphological salt response (Slr) contributed  
314 strongly to salt indices for shoot biomass, root biomass, and two indices for plant height.  
315 Thus, morphological salt response can be interpreted as the morphological responses to  
316 salinity stress, with higher values indicating a more tolerant growth response. The latent  
317 variable yield is a representation of overall grain productivity, and contributed strongly to  
318 the observed phenotypes primary panicle branch number, seed number per panicle, and pan-  
319 icle length. The positive loading scores on these observable phenotypes indicates that more  
320 highly branched, productive panicles will have higher values for yield. Seed width, seed vol-  
321 ume, and seed surface area contributed significantly to the latent variable grain morphology  
322 (Grm). Therefore, these results indicate that the grain morphology is a summary of the  
323 overall shape of the grain, where high values represent large, round grains, while low values  
324 represent small, slender grains. Considering the grain characteristics of rice subpopulations,  
325 temperate japonica accessions are expected to have high values for grain morphology, while

326 indica accessions have lower values for grain morphology. Latent variable morphology (Mrp)  
327 is a representation of plant biomass during the vegetative stage (28-day-old plants). Shoot  
328 biomass, root biomass, and two metrics for plant height contributed largely to morphol-  
329 ogy, suggesting that accessions with high values for morphology are tall plants with a large  
330 biomass.

331 Genomic correlation analysis among the six latent variables showed moderate correlations  
332 among several pairs. These genetic correlations can either be caused by linkage or pleitropy.  
333 The former is likely to prevail in species with high LD, which is the case in rice where  
334 LD ranges from 100 to 200kb (Huang et al., 2010). A strong negative relationship was  
335 observed between morphological salt response and three other latent variables. For instance,  
336 a negative correlation between morphological salt response and yield indicates that accessions  
337 of samples harboring alleles for superior morphological salt responses (e.g. those that are  
338 more tolerant) tend to also harbor alleles for poor yield. The rice diversity panel we used  
339 is a representative sample of the total genetic diversity within cultivated rice and contains  
340 many unimproved traditional varieties and modern breeding lines (Eizenga et al., 2014).  
341 While traditional varieties exhibit superior adaptation to abiotic stresses, they often have  
342 very poor agronomic characteristics including low yield, late flowering, and high photoperiod  
343 sensitivity (Thomson et al., 2009, 2010). Moreover, the indica and japonica subspecies have  
344 contrasting salt responses and very different grain morphology. Japonica accessions tend to  
345 have short, round seeds and are more sensitive to salt stress, while indica accessions have  
346 long, slender grains and often are more salt tolerant (Zhao et al., 2011; Campbell et al.,  
347 2017a). The negative relationship observed between salt response and grain morphology  
348 suggests that lines that harbor alleles for high grain morphology (e.g., large, round grains)  
349 tend to also harbor alleles for a tolerant growth response to salt stress. However, no studies  
350 have yet reported an association between alleles for grain morphology and morphological  
351 salt response. Therefore, it remains to be addressed whether this relationship is due to LD  
352 or pleitropy.

353 Genetic correlations observed between other latent variables may suggest a pleiotropic  
354 effect among loci. For instance, a moderate negative relationship was observed between  
355 morphological salt response and ionic components of salt stress, indicating that accessions  
356 harboring alleles associated with superior morphological salt response also tend to harbor  
357 alleles for reduced ion content under salt stress. The relationship between salt tolerance,  
358 measured in terms of growth or yield, and  $\text{Na}^+$  and  $\text{Na}^+:\text{K}^+$  has been documented for  
359 decades (reviewed by Munns and Tester (2008)). Moreover, natural variation for  $\text{Na}^+$  trans-  
360 porters has been utilized to improve growth and yield under saline conditions in rice and  
361 other cereals (Ren et al., 2005; Byrt et al., 2007; Horie et al., 2009; Munns et al., 2012;  
362 Campbell et al., 2017a). Therefore, the negative genetic relationships observed between  
363 morphological salt response and ion content may be due to the pleiotropic effects of some  
364 loci.

365 The genomic relationships among latent variables including morphology, yield, and grain  
366 morphology may have resulted from the selection of alleles associated with good agronomic  
367 characteristics. A moderate positive relationship was observed between yield and grain mor-  
368 phology, suggesting that alleles that positively contribute to productive panicles also may  
369 contribute to large, round grains. Furthermore, the negative genomic correlation observed be-  
370 tween morphology and yield indicates that alleles negatively influencing total plant biomass  
371 also have a positive contribution to traits for productive panicles. This genomic relationship  
372 may reflect the genetics of harvest index, which is defined as the ratio of grain yield to total  
373 biomass. Over the past 50 years, rice breeders have selected high harvest index, resulting  
374 in plants with short compact morphology and many highly productive panicles (Hay, 1995;  
375 Peng et al., 2008).

376 Although BCFA may yield biologically meaningful results, a potential limitation of BCFA  
377 is that we assumed each phenotype does not measure more than one latent variable. This  
378 assumption may not always strictly concur with the observational data. Therefore, further  
379 studies are required to allow each phenotype to potentially load onto multiple factors in the

380 BCFA framework. An alternative approach is to derive the number of latent variables and  
381 determine which latent variables load onto phenotypes directly from observed data, using  
382 exploratory factor analysis. This approach was not pursued here because accurate estimation  
383 of unknown parameters in the exploratory factor analysis requires a large sample size, which  
384 was not the case herein (Brown, 2014).

## 385 **Bayesian network of latent variables**

386 The BN is a probabilistic DAG, which represents the conditional dependencies among phe-  
387 notypes. The genomic correlation among latent variables described in Figure 2 does not  
388 inform the flow of genetic signals nor distinguish direct and indirect associations, whereas  
389 BN displays directions between latent variables and separate direct and indirect associations.  
390 Therefore, the BN describes the possibility that other phenotypes will change if one pheno-  
391 type is intervened (i.e., selection). However, caution is required to interpret this network as  
392 a causal effect, as the causal BN requires more assumptions, which are usually difficult to  
393 meet in observational data (Pearl, 2009).

394 Four common edges or consensus subnetworks across the four BN may be the most  
395 reliable substructure of latent variables and may describe the dependence between agronomic  
396 traits (Figures 3, 4, 5, and 6). For example, edges from grain morphology to yield and  
397 morphology to yield can be interpreted as final grain productivity is dependant on specific  
398 vegetative characteristics as well grain traits. This is because yield, which represents the  
399 overall grain productivity of a plant, depends on morphological characteristics such as the  
400 degree of tillering, an architecture that allows the plant to efficiently capture light and  
401 carbon, and a stature that is resistant to lodging, the degree of panicle branching, as well  
402 as specific grain characteristics such as seed volume and shape. Moreover, there is a direct  
403 biological linkage between specific vegetative architectural traits such as tillering and plant  
404 height, and yield related traits such as panicle branching and number of seeds per panicle.  
405 The degree of branching during both vegetative and reproductive development is dependant

406 on the development and initiation of auxiliary meristems. Several genes have been identified  
407 in this pathway and have shown to have pleiotropic effects on tillering and panicle branching  
408 (reviewed by Liang et al. (2014)). For instance, *OsSPL14* has been shown to be an important  
409 regulator of auxiliary branching in both vegetative and reproductive stages in rice (Jiao  
410 et al., 2010; Miura et al., 2010). Moreover, other genes such as *OsGhd8* have been reported  
411 to regulate other morphological traits such as plant height and yield through increase panicle  
412 branching (Yan et al., 2011). The biological importance of these dependencies can also be  
413 illustrated by viewing them in the context of genetic improvement, as selection for specific  
414 architectural traits (represented by the latent variable morphology) and grain characteristics  
415 have traditionally been used as traits to improve rice productivity in many conventional  
416 breeding programs (Redona and Mackill, 1998; Huang et al., 2013).

417 While the above example provides a plausible network structure between latent variables,  
418 edges from ionic components of salt stress to flowering time and to grain morphology are an  
419 example of instances where caution should be used to infer causation. As mentioned above,  
420 there is an inherent difference in salt tolerance and grain morphological traits between the  
421 indica and japonica subspecies. The edges observed for these two latent variables (ionic  
422 components of salt stress and grain morphology) in BN may be driven by LD between alleles  
423 associated with grain morphology and alleles for salt tolerance rather than pleiotropy. Thus,  
424 given the current data set, genetic effects for grain morphology may still be conditionally  
425 dependant on ionic components of salt stress and the BN may be true, even if there is no  
426 direct overlap in the genetic mechanisms for the two traits.

427 We found that there are some uncertain edges among BN. For instance, direction from  
428 salt response to grain morphology is supported by 65% (Figure 4), 58% (Figure 5), and 58%  
429 (Figure 6) bootstrap sampling, whereas the opposite direction is supported by 56% bootstrap  
430 sampling (Figure 3). An analogous uncertainty was also observed between morphology and  
431 flowering time, i.e., the path from morphology to flowering time was supported 60% (Figure  
432 3), 51% (Figure 4), and 52% (Figure 6), while the reverse direction was supported 51%

433 (Figure 6) upon bootstrapping. In addition, the two score-based algorithms captured edges  
434 between morphological salt response and flowering time with 70% and 76% bootstrapping  
435 evidence. However, this connection was not detected in the two hybrid algorithms. In  
436 general, inferring the direction of edges was harder than inferring the presence or absence of  
437 undirected edges. Finally, the whole structures of BN were evaluated in terms of the BIC  
438 score and BGe. Ranking of the networks was consistent across BIC and BGe and the two  
439 score-based algorithms produced networks with greater goodness-of-fit than the two hybrid  
440 algorithms. The optimal network was produced by the Tabu algorithm. This is consistent  
441 with the previous study reporting that the score-based algorithm produced a better fit of  
442 networks in data on maize (Töpner et al., 2017).

443 In conclusion, the present results show the utility of factor analysis and network analysis  
444 to characterize various phenotypes in rice. We showed that the joint use of BCFA and  
445 BN can be applied to predict the potential influence of external interventions or selection  
446 associated with target traits such as yield in the high-dimensional interrelated complex traits  
447 system. We contend that the approaches used herein provide greater insights than pairwise-  
448 association measures of multiple phenotypes and can be used to analyze the massive amount  
449 of diverse image-based phenomics dataset being generated by the automated plant phenomics  
450 platforms (e.g., Furbank and Tester, 2011). With a large volume of complex traits being  
451 collected through phenomics, numerous opportunities to forge new research directions are  
452 generated by using network analysis for the growing number of phenotypes.

## References

- 453
- 454 Acquaah, G. (2009). *Principles of plant genetics and breeding*. John Wiley & Sons.
- 455 Awada, L., Phillips, P. W., and Smyth, S. J. (2018). The adoption of automated phenotyping  
456 by plant breeders. *Euphytica*, 214(8):148.
- 457 Brown, T. A. (2014). *Confirmatory factor analysis for applied research*. Guilford Publica-  
458 tions.
- 459 Byrt, C. S., Platten, J. D., Spielmeier, W., James, R. A., Lagudah, E. S., Dennis, E. S.,  
460 Tester, M., and Munns, R. (2007). Hkt1; 5-like cation transporters linked to na<sup>+</sup> exclusion  
461 loci in wheat, nax2 and kna1. *Plant Physiology*, 143(4):1918–1928.
- 462 Calus, M. P. and Veerkamp, R. F. (2011). Accuracy of multi-trait genomic selection using  
463 different methods. *Genetics Selection Evolution*, 43(1):26.
- 464 Campbell, M. T., Bandillo, N., Al Shiblawi, F. R. A., Sharma, S., Liu, K., Du, Q., Schmitz,  
465 A. J., Zhang, C., Véry, A.-A., Lorenz, A. J., et al. (2017a). Allelic variants of oshkt1; 1  
466 underlie the divergence between indica and japonica subspecies of rice (*oryza sativa*) for  
467 root sodium content. *PLoS Genetics*, 13(6):e1006823.
- 468 Campbell, M. T., Du, Q., Liu, K., Brien, C. J., Berger, B., Zhang, C., and Walia, H. (2017b).  
469 A comprehensive image-based phenomic analysis reveals the complex genetic architecture  
470 of shoot growth dynamics in rice (*oryza sativa*). *The Plant Genome*, 10(2).
- 471 de los Campos, G. and Gianola, D. (2007). Factor analysis models for structuring covari-  
472 ance matrices of additive genetic effects: a bayesian implementation. *Genetics Selection  
473 Evolution*, 39(5):481.
- 474 Denwood, M. (2016). runjags: An r package providing interface utilities, model templates,



- 475 parallel computing methods and additional distributions for mcmc models in jags. *Journal*  
476 *of Statistical Software, Articles*, 71(9):1–25.
- 477 Eizenga, G. C., Ali, M., Bryant, R. J., Yeater, K. M., McClung, A. M., McCouch, S. R.,  
478 et al. (2014). Registration of the rice diversity panel 1 for genomewide association studies.  
479 *Journal of Plant Registrations*, 8(1):109–116.
- 480 Falconer, D. and Mackay, T. (1996). *Introduction to Quantitative Genetics*. Pearson.
- 481 Furbank, R. T. and Tester, M. (2011). Phenomics-technologies to relieve the phenotyping  
482 bottleneck. *Trends Plant Sci.*, 16:635–644.
- 483 Gelman, A. and Rubin, D. B. (1992). Inference from iterative simulation using multiple  
484 sequences. *Statistical Science*, pages 457–472.
- 485 Gianola, D., de los Campos, G., Toro, M. A., Naya, H., Schön, C.-C., and Sorensen, D.  
486 (2015). Do molecular markers inform about pleiotropy? *Genetics*, pages genetics–115.
- 487 Hay, R. (1995). Harvest index: a review of its use in plant breeding and crop physiology.  
488 *Annals of applied biology*, 126(1):197–216.
- 489 Henderson, C. and Quaas, R. (1976). Multiple trait evaluation using relatives' records.  
490 *Journal of Animal Science*, 43(6):1188–1197.
- 491 Hickey, J. M., Chiurugwi, T., Mackay, I., Powell, W., Eggen, A., Kilian, A., Jones, C.,  
492 Canales, C., Grattapaglia, D., Bassi, F., et al. (2017). Genomic prediction unifies animal  
493 and plant breeding programs to form platforms for biological discovery. *Nature genetics*,  
494 49(9):1297.
- 495 Horie, T., Hauser, F., and Schroeder, J. I. (2009). Hkt transporter-mediated salinity re-  
496 sistance mechanisms in arabidopsis and monocot crop plants. *Trends in plant science*,  
497 14(12):660–668.

- 498 Hornik, K., Leisch, F., and Zeileis, A. (2003). Jags: A program for analysis of bayesian  
499 graphical models using gibbs sampling. In *Proceedings of DSC*, volume 2, pages 1–1.
- 500 Huang, R., Jiang, L., Zheng, J., Wang, T., Wang, H., Huang, Y., and Hong, Z. (2013).  
501 Genetic bases of rice grain shape: so many genes, so little known. *Trends in plant science*,  
502 18(4):218–226.
- 503 Huang, X., Sang, T., Zhao, Q., Feng, Q., Zhao, Y., Li, C., Zhu, C., Lu, T., Zhang, Z., Li, M.,  
504 et al. (2010). Genome-wide association studies of 14 agronomic traits in rice landraces.  
505 *Nature Genetics*, 42(11):961.
- 506 Jia, Y. and Jannink, J.-L. (2012). Multiple trait genomic selection methods increase genetic  
507 value prediction accuracy. *Genetics*, pages genetics–112.
- 508 Jiao, Y., Wang, Y., Xue, D., Wang, J., Yan, M., Liu, G., Dong, G., Zeng, D., Lu, Z., Zhu,  
509 X., et al. (2010). Regulation of *osspl14* by *osmir156* defines ideal plant architecture in rice.  
510 *Nature genetics*, 42(6):541.
- 511 Jöreskog, K. G. (1969). A general approach to confirmatory maximum likelihood factor  
512 analysis. *Psychometrika*, 34(2):183–202.
- 513 Lee, S.-Y. and Song, X.-Y. (2012). *Basic and advanced Bayesian structural equation model-*  
514 *ing: With applications in the medical and behavioral sciences*. John Wiley & Sons.
- 515 Liang, W.-h., Shang, F., Lin, Q.-t., Lou, C., and Zhang, J. (2014). Tillering and panicle  
516 branching genes in rice. *Gene*, 537(1):1–5.
- 517 Lush, J. L. (1948). The genetics of populations. Technical Report Mimeo.
- 518 Merkle, E. and Rosseel, Y. (2018). blavaan: Bayesian structural equation models via pa-  
519 rameter expansion. *Journal of Statistical Software, Articles*, 85(4):1–30.

- 520 Miura, K., Ikeda, M., Matsubara, A., Song, X.-J., Ito, M., Asano, K., Matsuoka, M., Ki-  
521 tano, H., and Ashikari, M. (2010). Ossl14 promotes panicle branching and higher grain  
522 productivity in rice. *Nature genetics*, 42(6):545.
- 523 Morota, G. and Gianola, D. (2013). Evaluation of linkage disequilibrium in wheat with an l1-  
524 regularized sparse markov network. *Theoretical and Applied Genetics*, 126(8):1991–2002.
- 525 Morota, G., Valente, B., Rosa, G., Weigel, K., and Gianola, D. (2012). An assessment  
526 of linkage disequilibrium in holstein cattle using a bayesian network. *Journal of Animal*  
527 *Breeding and Genetics*, 129(6):474–487.
- 528 Mrode, R. A. (2014). *Linear models for the prediction of animal breeding values*. Cabi.
- 529 Munns, R., James, R. A., Xu, B., Athman, A., Conn, S. J., Jordans, C., Byrt, C. S., Hare,  
530 R. A., Tyerman, S. D., Tester, M., et al. (2012). Wheat grain yield on saline soils is  
531 improved by an ancestral na<sup>+</sup> transporter gene. *Nature biotechnology*, 30(4):360.
- 532 Munns, R. and Tester, M. (2008). Mechanisms of salinity tolerance. *Annu. Rev. Plant Biol.*,  
533 59:651–681.
- 534 Pearl, J. (2009). *Causality: Models, Reasoning and Inference*. Cambridge University Press,  
535 New York, NY, USA, 2nd edition.
- 536 Peñagaricano, F., Valente, B., Steibel, J., Bates, R., Ernst, C., Khatib, H., and Rosa,  
537 G. (2015). Searching for causal networks involving latent variables in complex traits:  
538 application to growth, carcass, and meat quality traits in pigs. *Journal of Animal Science*,  
539 93(10):4617–4623.
- 540 Peng, S., Khush, G. S., Virk, P., Tang, Q., and Zou, Y. (2008). Progress in ideotype breeding  
541 to increase rice yield potential. *Field Crops Research*, 108(1):32–38.
- 542 Redona, E. and Mackill, D. (1998). Quantitative trait locus analysis for rice panicle and  
543 grain characteristics. *Theoretical and Applied Genetics*, 96(6-7):957–963.

- 544 Ren, Z.-H., Gao, J.-P., Li, L.-G., Cai, X.-L., Huang, W., Chao, D.-Y., Zhu, M.-Z., Wang,  
545 Z.-Y., Luan, S., and Lin, H.-X. (2005). A rice quantitative trait locus for salt tolerance  
546 encodes a sodium transporter. *Nature genetics*, 37(10):1141.
- 547 Rosa, G. J., Valente, B. D., de los Campos, G., Wu, X.-L., Gianola, D., and Silva, M. A.  
548 (2011). Inferring causal phenotype networks using structural equation models. *Genetics*  
549 *Selection Evolution*, 43(1):6.
- 550 Scutari, M. (2010). Learning bayesian networks with the bnlearn r package. *Journal of*  
551 *Statistical Software, Articles*, 35(3):1–22.
- 552 Scutari, M. and Denis, J.-B. (2014). *Bayesian networks: with examples in R*. Chapman and  
553 Hall/CRC.
- 554 Tanner, M. A. and Wong, W. H. (1987). The calculation of posterior distributions by data  
555 augmentation. *Journal of the American statistical Association*, 82(398):528–540.
- 556 Thomson, M. J., de Ocampo, M., Egdane, J., Rahman, M. A., Sajise, A. G., Adorada, D. L.,  
557 Tumimbang-Raiz, E., Blumwald, E., Seraj, Z. I., Singh, R. K., et al. (2010). Characterizing  
558 the saltol quantitative trait locus for salinity tolerance in rice. *Rice*, 3(2-3):148–160.
- 559 Thomson, M. J., Ismail, A. M., McCouch, S. R., and Mackill, D. J. (2009). Marker assisted  
560 breeding. In *Abiotic Stress Adaptation in Plants*, pages 451–469. Springer.
- 561 Töpner, K., Rosa, G. J., Gianola, D., and Schön, C.-C. (2017). Bayesian networks illustrate  
562 genomic and residual trait connections in maize (*Zea mays* L.). *G3: Genes, Genomes,*  
563 *Genetics*, 7(8):2779–2789.
- 564 Valente, B. D., Morota, G., Peñagaricano, F., Gianola, D., Weigel, K., and Rosa, G. J.  
565 (2015). The causal meaning of genomic predictors and how it affects construction and  
566 comparison of genome-enabled selection models. *Genetics*, 200(2):483–494.

- 567 Valente, B. D., Rosa, G. J., de los Campos, G., Gianola, D., and Silva, M. A. (2010).  
568 Searching for recursive causal structures in multivariate quantitative genetics mixed mod-  
569 els. *Genetics*, 185(2):633–644.
- 570 VanRaden, P. M. (2008). Efficient methods to compute genomic predictions. *J. Dairy Sci.*,  
571 91:4414–4423.
- 572 Vazquez, A., Bates, D., Rosa, G., Gianola, D., and Weigel, K. (2010). An r package for  
573 fitting generalized linear mixed models in animal breeding 1. *Journal of animal science*,  
574 88(2):497–504.
- 575 Verma, T. and Pearl, J. (1991). Equivalence and synthesis of causal models. In *Proceedings*  
576 *of the Sixth Annual Conference on Uncertainty in Artificial Intelligence, UAI '90*, pages  
577 255–270, New York, NY, USA. Elsevier Science Inc.
- 578 Xavier, A., Hall, B., Casteel, S., Muir, W., and Rainey, K. M. (2017). Using unsupervised  
579 learning techniques to assess interactions among complex traits in soybeans. *Euphytica*,  
580 213(8):200.
- 581 Yan, W.-H., Wang, P., Chen, H.-X., Zhou, H.-J., Li, Q.-P., Wang, C.-R., Ding, Z.-H., Zhang,  
582 Y.-S., Yu, S.-B., Xing, Y.-Z., et al. (2011). A major qtl, ghd8, plays pleiotropic roles in  
583 regulating grain productivity, plant height, and heading date in rice. *Molecular plant*,  
584 4(2):319–330.
- 585 Zhao, K., Tung, C.-W., Eizenga, G. C., Wright, M. H., Ali, M. L., Price, A. H., Norton,  
586 G. J., Islam, M. R., Reynolds, A., Mezey, J., et al. (2011). Genome-wide association  
587 mapping reveals a rich genetic architecture of complex traits in oryza sativa. *Nature*  
588 *Communications*, 2:467.

## 589 Tables

Table 1: Standardized factor loadings obtained from the Bayesian confirmatory factor analysis.

Latent variable	Observed phenotype	Loading
Flowering time	Flowering time at Arkansas (Fla)	0.990
Flowering time	Flowering time at Faridpur (Flf)	0.500
Flowering time	Flowering time at Aberdeen (Flb)	0.578
Flowering time	FT ratio of Arkansas/Aberdeen (Flaa)	-0.212
Flowering time	FT ratio of Faridpur/Aberdeen (Flfa)	-0.549
Flowering time	Year07 Flowering time at Arkansas (Fla7)	0.926
Flowering time	Year06 Flowering time at Arkansas (Fla6)	0.886
Morphology	Culm habit (Cuh)	0.227
Morphology	Flag leaf length (FlL)	0.116
Morphology	Flag leaf width (Flw)	-0.044
Morphology	Plant height (Plh)	0.440
Morphology	Shoot BM Control (Sbc)	0.534
Morphology	Shoot BM Salt (Sbs)	0.456
Morphology	Root BM Control (Rbc)	0.418
Morphology	Root BM Salt (Rbs)	0.280
Morphology	Tiller No Salt (Tns)	-0.349
Morphology	Tiller No Control (Tbc)	-0.318
Morphology	Ht Lig Salt (Hls)	0.920
Morphology	Ht Lig Control (Hlc)	0.899
Morphology	Ht FE Salt (Hfs)	0.907
Morphology	Ht FE Control (Hfc)	0.925
Yield	Panicle number per plant (Pnu)	0.190
Yield	Panicle length (Pal)	0.455
Yield	Primary panicle branch number (Ppn)	0.790
Yield	Seed number per panicle (Snpp)	0.780
Yield	Panicle fertility (Paf)	-0.085
Grain Morphology	Seed length (Sl)	0.251
Grain Morphology	Seed width (Sw)	0.876
Grain Morphology	Seed volume (Sv)	0.990
Grain Morphology	Seed surface area (Ssa)	0.901
Grain Morphology	Brown rice seed length (Bsl)	0.158
Grain Morphology	Brown rice seed width (Bsw)	0.837
Grain Morphology	Brown rice surface area (Bsa)	0.902
Grain Morphology	Brown rice volume (Bvl)	0.986
Grain Morphology	Seed length/width ratio (Slwr)	-0.476
Grain Morphology	Brown rice length/width ratio (Blwr)	-0.432
Grain Morphology	Grain length McCouch2016 (Glmc)	0.047
Ionic components of salt stress	Na K Shoot (Ks)	0.983
Ionic components of salt stress	Na Shoot (Nas)	0.975
Ionic components of salt stress	K Shoot Salt (Kss)	-0.265
Ionic components of salt stress	Na K Root (Kr)	0.061
Ionic components of salt stress	Na Root (Nar)	0.000
Ionic components of salt stress	K Root Salt (Krs)	-0.095
Morphological salt response	Shoot BM Ratio (Sbr)	0.410
Morphological salt response	Root BM Ratio (Rbr)	0.395
Morphological salt response	Tiller No Ratio (Tbr)	-0.022
Morphological salt response	Ht Lig Ratio (Hlr)	0.665
Morphological salt response	Ht FE Ratio (Hfr)	0.939

## 590 Figures

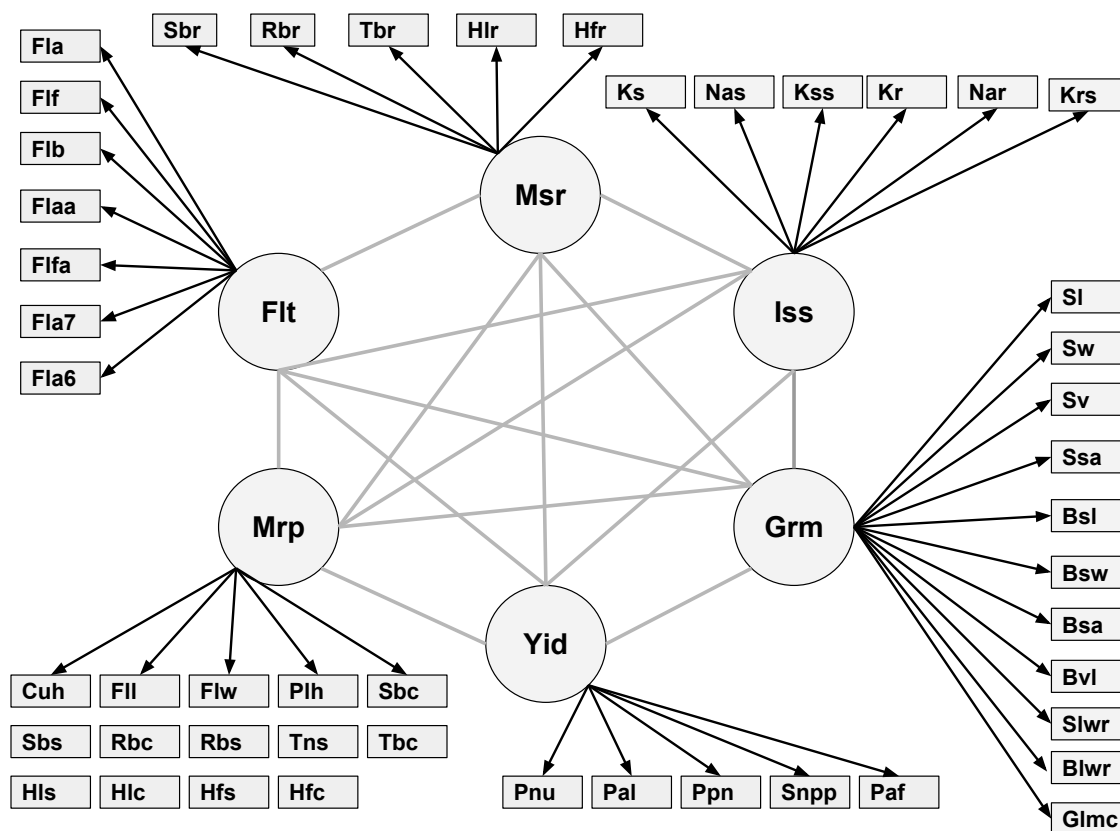


Figure 1: Relationship between six latent variables and observed phenotypes. Msr: morphological salt response; Iss: ionic components of salt stress; Grm: grain morphology; Yid: yield; Mrp: morphology; Flt: flowering time. Abbreviations of observed phenotypes are shown in Table S1.

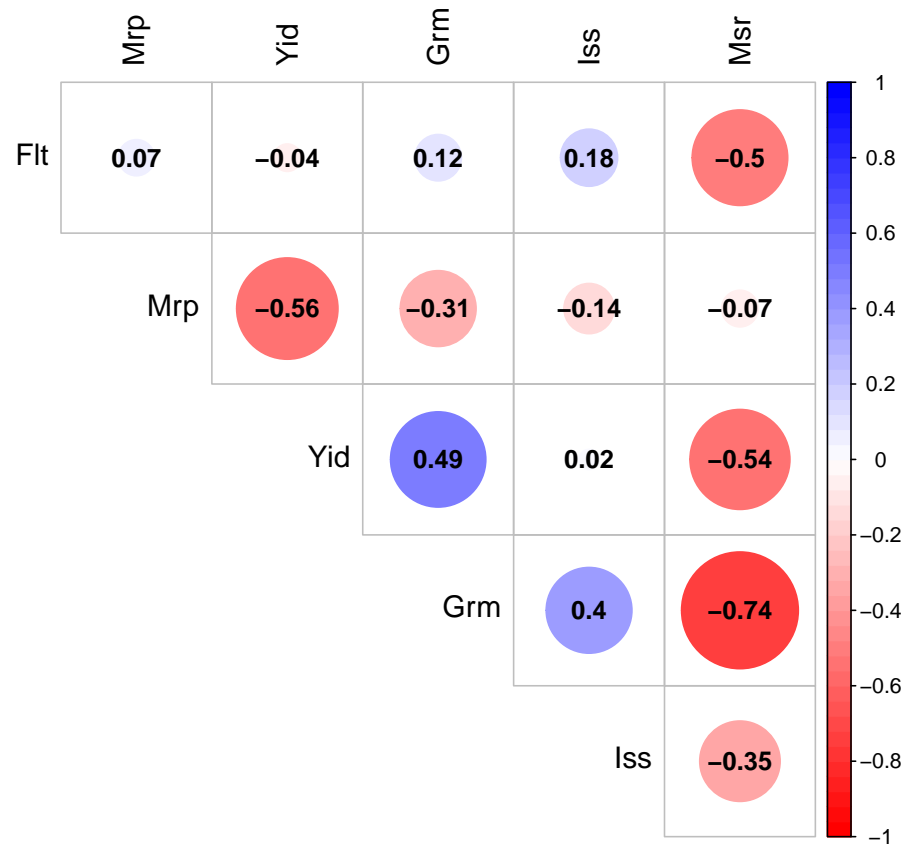


Figure 2: Genomic correlation of six latent variables. The size of each circle, degree of shading, and value reported correspond to the correlation between each pair of latent variables. Msr: morphological salt response; Iss: ionic components of salt stress; Grm: grain morphology; Yid: yield; Mrp: morphology; Flt: flowering time.



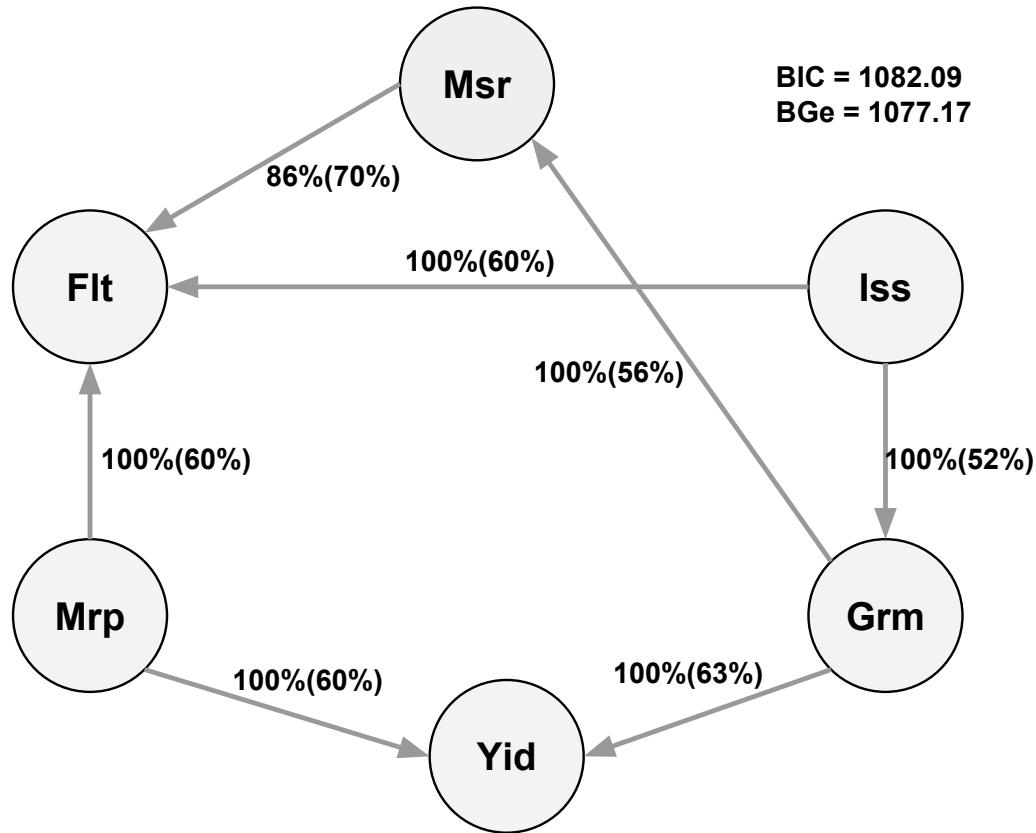


Figure 3: Bayesian network between six latent variables based on the Hill Climbing algorithm. The quality of the structure was evaluated by bootstrap resampling and model averaging across 500 replications. Labels of the edges refer to the strength and direction (parenthesis) which measure the confidence of the directed edge. The strength indicates the frequency of the edge is present and the direction measures the frequency of the direction conditioned on the presence of edge. BIC: Bayesian information criterion score. BGe: Bayesian Gaussian equivalent score. Msr: morphological salt response; Iss: ionic components of salt stress; Grm: grain morphology; Yid: yield; Mrp: morphology; Flt: flowering time.

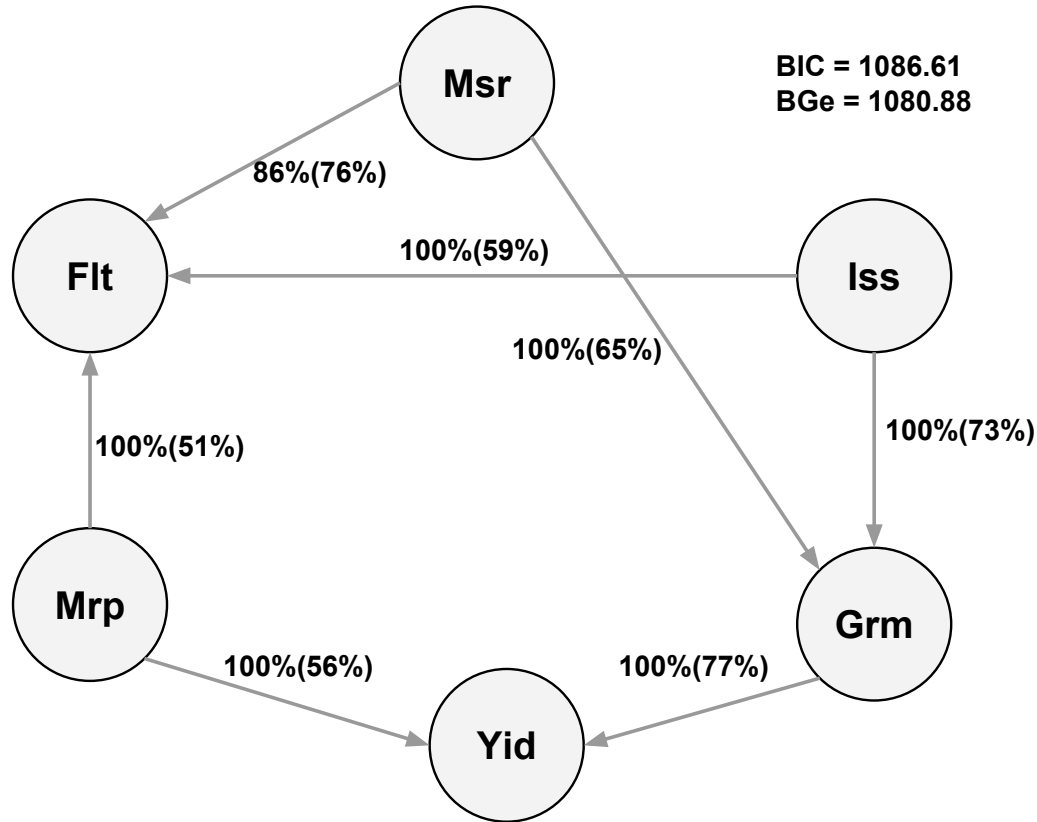


Figure 4: Bayesian network between six latent variables based on the Tabu algorithm. The quality of the structure was evaluated by bootstrap resampling and model averaging across 500 replications. Labels of the edges refer to the strength and direction (parenthesis) which measure the confidence of the directed edge. The strength indicates the frequency of the edge is present and the direction measures the frequency of the direction conditioned on the presence of edge. BIC: Bayesian information criterion score. BGe: Bayesian Gaussian equivalent score. Msr: morphological salt response; Iss: ionic components of salt stress; Grm: grain morphology; Yid: yield; Mrp: morphology; Flt: flowering time.

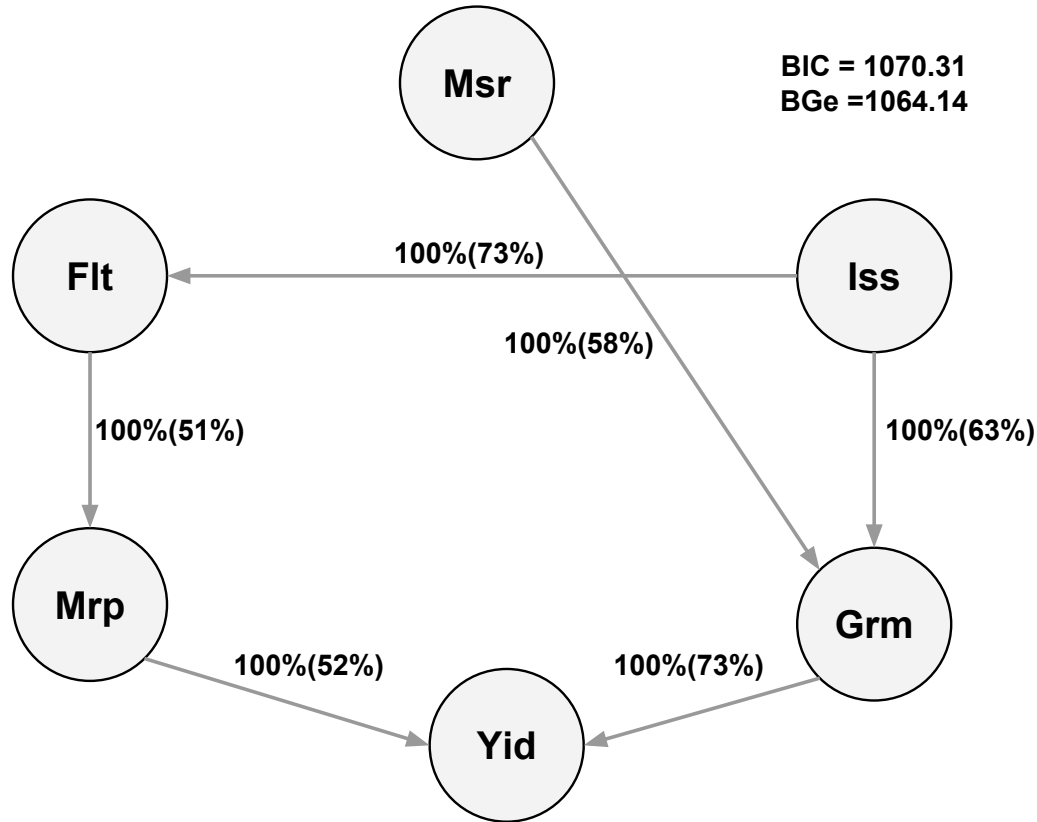


Figure 5: Bayesian network between six latent variables based on the Max-Min Hill Climbing algorithm. The quality of the structure was evaluated by bootstrap resampling and model averaging across 500 replications. Labels of the edges refer to the strength and direction (parenthesis) which measure the confidence of the directed edge. The strength indicates the frequency of the edge is present and the direction measures the frequency of the direction conditioned on the presence of edge. BIC: Bayesian information criterion score. BGe: Bayesian Gaussian equivalent score. Msr: morphological salt response; Iss: ionic components of salt stress; Grm: grain morphology; Yid: yield; Mrp: morphology; Flt: flowering time.

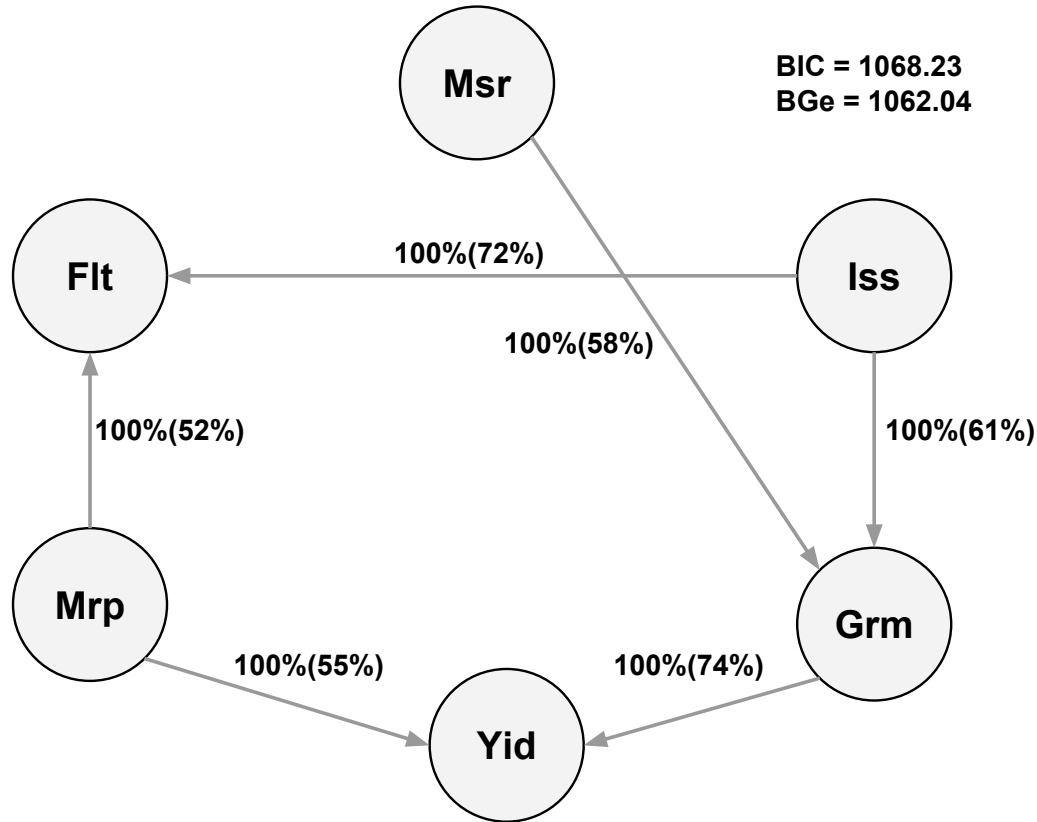


Figure 6: Bayesian network between six latent variables based on the General 2-Phase Restricted Maximization algorithm. The quality of the structure was evaluated by bootstrap resampling and model averaging across 500 replications. Labels of the edges refer to the strength and direction (parenthesis) which measure the confidence of the directed edge. The strength indicates the frequency of the edge is present and the direction measures the frequency of the direction conditioned on the presence of edge. BIC: Bayesian information criterion score. BGe: Bayesian Gaussian equivalent score. Msr: morphological salt response; Iss: ionic components of salt stress; Grm: grain morphology; Yid: yield; Mrp: morphology; Flt: flowering time.