



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 With the advent of high-throughput phenotyping platforms, plant breeders have a means  
25 to assess many traits for large breeding populations. However, understanding the genetic  
26 interdependencies among high-dimensional traits in a statistically robust manner remains  
27 a major challenge. Since multiple phenotypes likely share mutual relationships, elucidating  
28 the interdependencies among economically important traits can better inform breeding de-  
29 cisions and accelerate the genetic improvement of plants. The objective of this study was to  
30 leverage confirmatory factor analysis and graphical modeling to elucidate the genetic interde-  
31 pendencies among a diverse agronomic traits in rice. We used a Bayesian network to depict  
32 conditional dependencies among phenotypes, which can not be obtained by standard multi-  
33 trait analysis. We utilized Bayesian confirmatory factor analysis which hypothesized that 48  
34 observed phenotypes resulted from six latent variables including grain morphology, morphol-  
35 ogy, flowering time, physiology, yield, and morphological salt response. This was followed  
36 by studying the genetics of each latent variable, which is also known as factor, using single  
37 nucleotide polymorphisms. Bayesian network structures involving the genomic component  
38 of six latent variables were established by fitting four algorithms (i.e., Hill Climbing, Tabu,  
39 Max-Min Hill Climbing, and General 2-Phase Restricted Maximization algorithms). Phys-  
40 iological components influenced the flowering time and grain morphology, and morphology  
41 and grain morphology influenced yield. In summary, we show the Bayesian network coupled  
42 with factor analysis can provide an effective approach to understand the interdependence  
43 patterns among phenotypes and to predict the potential influence of external interventions  
44 or selection related to target traits in the interrelated complex traits systems.

## 45 Introduction

46 A primary objective in plant breeding is the develop high yielding varieties with specific  
47 grain qualities, resilience to pests and abiotic stresses, and superior adaption to the target  
48 environment. As a result, plant breeders devote considerable resources to extensive pheno-  
49 typic evaluation of germplasm and select on multiple traits. These traits are often correlated  
50 at a genetic level through common genetic effects (e.g. pleiotropy) or linkage disequilibrium  
51 between quantitative trait locus (QTL). Since multiple phenotypes may exhibit mutual re-  
52 lationships, knowledge of the interdependence among agronomically important traits can  
53 improve the efficacy of selection and rate of genetic improvement in systems with complex  
54 traits.

55 In a standard quantitative genetic analysis, multivariate phenotypes can be modeled  
56 through multi-trait models (MTM) of Henderson and Quaas (1976) or some genomic coun-  
57 terparts (e.g., Calus and Veerkamp, 2011; Jia and Jannink, 2012) by leveraging genetic or  
58 environmental correlations among traits. In particular, MTM has been useful in deriving  
59 genetic correlations and enhancing the prediction accuracy of breeding values for traits with  
60 low heritability or scarce records via joint modeling with one or more genetically correlated,  
61 highly heritable traits (Mrode, 2014). Conventional MTM strategies may provide impor-  
62 tant insight into the genetic relations between agronomically important traits, but they fail  
63 to explain how these traits are related. For instance, consider a case where we have three  
64 genetically correlated traits:  $y_1$ ,  $y_2$ , and  $y_3$ . With MTM, we cannot address whether the  
65 relationship between  $y_1$  and  $y_3$  is due to direct effects, or if the relationship is driven by  
66 indirect effects mediated by  $y_2$ . Bayesian Networks (BN) offer an effective approach to elu-  
67 cidate the underlying network structure in multivariate data and infer network relationships  
68 between correlated variables. A BN is a probabilistic graphical model that represents condi-  
69 tional dependencies among a set of variables via a directed acyclic graph (DAG) (Neapolitan  
70 et al., 2004). In the DAG, the variables are represented by nodes, while their conditional

71 dependencies between nodes are indicated with directed edges. In the context of plant  
72 breeding, BN can be used to elucidate the interdependencies among traits and inform selection  
73 decisions for simultaneously improving multiple traits. For instance in the latter case above  
74 ( $y_1 \rightarrow y_2 \rightarrow y_3$ ), selection directly on  $y_2$  will affect the quantity of  $y_3$  without an effect on  $y_1$ .

75 With the advent of high throughput phenotyping platforms, plant breeders have been pro-  
76 vided with a suite of tools for phenotypic evaluation of large populations (Shakoor et al.,  
77 2017). These platforms leverage robotics, precise environmental control, and remote sensing  
78 techniques to provide accurate, repeatable and high resolution phenotypes for large breed-  
79 ing populations throughout the growing season (Araus and Cairns, 2014; Shakoor et al.,  
80 2017; Araus et al., 2018). These data can be used to redefine characteristics underlying  
81 superior agronomic performance by quantifying secondary traits associated with seedling  
82 vigor, plant architecture, photosynthesis, transpiration, disease resistance, and stress toler-  
83 ance (Cabrera-Bosquet et al., 2016; Sun et al., 2017; Crain et al., 2018). However given these  
84 new approaches, breeders are faced with the new challenge of efficiently utilizing these large  
85 multidimensional data sets to improve selection efficiency. The primary challenges associated  
86 with multivariate analysis and BN approaches using HTP data is that robust parameter  
87 estimates can be untenable because the number of estimated parameters within the model  
88 increases with the increasing number of phenotypes. Moreover even in cases where MTM or  
89 BN can be applied, interpreting of interrelationships among a large number of phenotypes  
90 can be difficult.

91 One approach to characterize high-dimensional phenotypes is by using factor analysis  
92 (FA). The central idea of FA approaches is to reduce the dimensions of multivariate data  
93 sets by constructing unobserved, latent factors, or modules, from correlated phenotypes  
94 (de los Campos and Gianola, 2007). The biological importance of these latent factors can be  
95 interpreted by inspecting the phenotypes that contribute to each factor. Thus, the advantage  
96 of FA for large, multivariate data sets is two fold. First, FA provides a means to reduce  
97 the dimensions of multivariate data sets thereby providing statistically sound parameter

98 estimates, and easing visualization and interpretation. Secondly, the latent variables/factors  
99 themselves may be representative of underlying biological processes that cannot be observed  
100 or measured in the population. For instance, several studies have highlighted the effects  
101 of plant hormones such as GA on multiple morphological attributes (Wang and Li, 2006;  
102 Lo et al., 2008; Umehara et al., 2008; Bhattacharya et al., 2010; Brewer et al., 2013; Zhou  
103 et al., 2013). Thus, a latent factor constructed from these morphological traits may provide  
104 information on the biosynthesis or sensitivity of these hormones for individuals within the  
105 population. If a certain amount of knowledge regarding the biological role of the variables is  
106 already known, a variant of FA, confirmatory factor analysis (CFA), can be used to estimate  
107 latent variables based on predetermined biological classes of observed traits (Jöreskog, 1969).  
108 These latent variables underlie observed phenotypes and can be evaluated for how well the  
109 data support the hypothesis. For instance, Peñagaricano et al. (2015) performed CFA in  
110 swine to derive five latent variables from 19 phenotypic traits and inferred BN structures  
111 among those latent variables, thereby demonstrating the potential of this approach.

112 This study aimed to leverage CFA and graphical modeling to elucidate the genetic inter-  
113 dependencies among traits typically recorded in breeding programs (e.g., yield, plant mor-  
114 phology, phenology, and stress resilience). First, we constructed latent variables, using prior  
115 biological knowledge obtained from the literature. Then we connected the observed high-  
116 dimensional phenotypes with these to establish latent variables via Bayesian confirmatory  
117 factor analysis (BCFA) to reduce the dimensions of the dataset. Further, factor scores com-  
118 puted from BCFA were considered new phenotypes for a Bayesian multivariate analysis to  
119 separate breeding values from noise. This was followed by adjustment of breeding values via  
120 Cholesky decomposition to eliminate the dependencies introduced by genomic relationships.  
121 Finally, the adjusted breeding values were considered inputs to assess the causal network  
122 structure between latent variables by conducting a Gaussian BN analysis. This study is the  
123 first, to our knowledge, in rice to characterize various phenotypes with graphical modeling  
124 such as BCFA and BN.

## 125 **Materials and Methods**

### 126 **Phenotypic and genotypic data**

127 The rice dataset comprised  $n = 374$  accessions sampled from six subpopulations: temperate  
128 japonica (92), tropical japonica (85), indica (77), aus (52), aromatic (12), and admixture  
129 of japonica and indica (56) (Zhao et al., 2011). The improvement status of each accession  
130 was obtained from the USDA-ARS Germplasm Resources Information Network. We used  
131  $t = 48$  phenotypes and data regarding 44,000 single-nucleotide polymorphisms (SNP). After  
132 removing SNP markers with minor allele frequency less than 0.05, 374 accessions and 33,584  
133 markers were used for further analysis. Of those, 27 phenotypes were reported in Zhao et al.  
134 (2011) and McCouch et al. (2016). These phenotypes can be classified into four categories:  
135 flowering time (flowering time at three locations, photoperiod sensitivity), grain morphology  
136 (seed length, seed width, seed surface area, seed length to width ratio, seed volume), plant  
137 morphology (culm habit/angle, flag leaf length and width, plant height at maturity), and  
138 yield traits (panicle fertility, seed number per panicle, number of primary branches on the  
139 main panicle, panicle length, and the number of panicles on each plant). Zhao et al. (2011)  
140 evaluated flowering time-related traits using data from three locations, while the remaining  
141 traits were evaluated at one location (Arkansas). The remaining phenotypes were assessed  
142 from the salinity stress experiments conducted in Campbell et al. (2017a). These traits were  
143 classified into three categories: morphological salt response, ionic components of salt stress,  
144 and plant morphology. The class morphological salt response represents how plant growth is  
145 affected by salinity stress and is composed of the ratio of shoot biomass of salt stressed plants  
146 to control, the ratio of root biomass of salt stressed plants to control, the ratio of the number  
147 of tillers for salt stressed plants to control, and two metrics that represent the ratio of shoot  
148 height of salt stressed plants to control. Ionic components of salt stress is composed of traits  
149 that quantify ions important for salinity tolerance ( $\text{Na}^+$  and  $\text{K}^+$ ) in both root and shoot

150 tissues. Morphology traits are those that describe the growth of the plant in both control  
151 and saline conditions (e.g. shoot biomass, root biomass, shoot height, and tiller number).  
152 The data used from Campbell et al. (2017a) were derived from three to six independent  
153 greenhouse experiments performed between July and October 2013. Information for all  
154 experiments were combined and best linear unbiased estimators were calculated for each line  
155 as described in Campbell et al. (2017a). The detailed descriptions of the phenotypes are  
156 summarized in Supplementary Table S1.

## 157 **Bayesian confirmatory factor analysis**

A CFA 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, ionic components of salt stress, yield, and morphological 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

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

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.

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



158 where  $\Phi$  is the variance of latent variables, and  $\Psi$  is the variance of specific effects (Brown,  
159 2014). Six latent variables were assumed to account for the covariance in the observed  
160 phenotypes. Moreover, latent variables were assumed to be correlated with each other. Prior  
161 distributions were assigned to all unknown parameters. The non-zero coefficients within  
162 factor loading matrix  $\Lambda$  were assumed to follow a Gaussian distribution with mean of 0  
163 and variance of 0.01. The variance-covariance matrix  $\Phi$  was assigned an inverse Wishart  
164 distribution with a  $6 \times 6$  identity scale matrix  $\mathbf{I}_{66}$  and a degree freedom of 7,  $\Phi \sim \mathcal{W}^{-1}(\mathbf{I}_{66}, 7)$   
165 and an inverse Gamma distribution with scale parameter 1 and shape parameter 0.5 was  
166 assigned to  $\Psi \sim \Gamma^{-1}(1, 0.5)$ .

167 We employed the blavaan R package (Merkle and Rosseel, 2018) jointly with JAGS  
168 (Hornik et al., 2003) to fit the above BCFA. The blavaan runs the runjags R package (Den-  
169 wood, 2016) to summarize the Markov chain Monte Carlo (MCMC) and samples unknown  
170 parameters from the posterior distributions. Three MCMC chains, each of 5,000 samples  
171 with 2,000 burn-in, were used to infer the unknown model parameters. The convergence of  
172 the parameters was investigated with trace plots and potential scale reduction factor (PSRF)  
173 less than 1.2 (Brooks and Gelman, 1998). The PSRF computes the difference between es-  
174 timated variances among multiple Markov chains and estimated variances within the chain.  
175 A large difference indicates non-convergence and may require additional Gibbs sampling.

176 Subsequently, the posterior means of factor scores ( $\mathbf{F}$ ), which reflect the contribution of  
177 latent variables to each accession were estimated. Within each draw of Gibbs sampling,  $\mathbf{F}$   
178 was sampled from the conditional distribution of  $p(\mathbf{F}|\boldsymbol{\theta}, \mathbf{T})$ , where  $\boldsymbol{\theta}$  refers to the unknown  
179 parameters in  $\Lambda$ ,  $\Phi$ , and  $\Psi$ . This conditional distribution was derived with data augmenta-  
180 tion (Tanner and Wong, 1987) assuming  $\mathbf{F}$  as missing data (Lee and Song, 2012).

## 181 Multivariate genomic best linear unbiased prediction

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

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

182 where  $\boldsymbol{\mu}$  is the vector of intercept,  $\mathbf{X}$  is the incidence matrix of covariates,  $\mathbf{b}$  is the vector of  
183 covariate effects,  $\mathbf{Z}$  is the incidence matrix relating accessions with additive genetic effects,  $\mathbf{u}$   
184 is the vector of additive genetic effects, and  $\boldsymbol{\epsilon}$  is the vector of residuals. The incident matrix  
185  $\mathbf{X}$  included subpopulation information (temperate japonica, tropical japonica, indica, aus,  
186 aromatic, and admixture), as the rice diversity panel used herein shows a clear substructure  
187 (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

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

188 where  $\mathbf{G}$  represents the second genomic relationship matrix of VanRaden (2008),  $\mathbf{I}$  is the  
189 identity matrix,  $\boldsymbol{\Sigma}_{\mathbf{u}}$  and  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}$  refer to  $6 \times 6$  dimensional genetic and residual variance-covariance  
190 matrices, respectively. An inverse Wishart distribution with a  $6 \times 6$  identity scale matrix  
191 of  $\mathbf{I}_{66}$  and a degree of freedom 6 was assigned as prior for  $\boldsymbol{\Sigma}_{\mathbf{u}}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}} \sim \mathcal{W}^{-1}(\mathbf{I}_{66}, 6)$ . These  
192 parameters were selected so that relatively uninformative priors were used. The Bayesian  
193 multivariate genomic best linear unbiased prediction model was implemented using the MTM  
194 R package (<https://github.com/QuantGen/MTM>). Posterior mean estimates of genomic cor-  
195 relation between latent variables and predicted breeding values ( $\hat{\mathbf{u}}$ ) were then obtained. The  
196 convergence of the estimated parameters was verified by trace plots.

## 197 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 Cholesky 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  $\mathcal{N}(0, \mathbf{I}\sigma_u^2)$  (Callanan and Harville, 1989; Vazquez et al., 2010)

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

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

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

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

## 203 Bayesian network

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

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

where  $\mathcal{G}$  represents a 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

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

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

209 The adjusted breeding values  $\hat{\mathbf{u}}^*$  were used to infer a genomic network structure among  
210 the aforementioned six latent variables. There are three types of structure-learning algo-  
211 rithms for BN: constraint-based algorithms, score-based algorithms, and a hybrid of these  
212 two (Scutari and Denis, 2014). The constraint-based algorithms can be originally traced  
213 to the inductive causation algorithm (Verma and Pearl, 1991), which uses conditional in-  
214 dependence tests for network inference. Briefly, the first step is to identify a d-separation  
215 set for each pair of nodes and confer an undirected edge between the two if they are not

216 d-separated. The second step is to identify a v-structure for each pair of non-adjacent nodes,  
217 where a common neighbor is the outcome of two non-adjacent nodes. In the last step, com-  
218 pelled edges were identified and oriented, where neither cyclic graph nor new v-structures  
219 are permitted. The score-based algorithms are based on heuristic approaches, which first  
220 assign a goodness-of-fit score for an initial graph structure and then maximize this score by  
221 updating the structure (i.e., add, delete, or reverse the edges of initial graph). The hybrid  
222 algorithm includes two steps, restrict and maximize, which harness both constrain-based and  
223 score-based algorithms to construct a reliable network. In this study, the two score-based  
224 (Hill Climbing and Tabu) and two hybrid algorithms (Max-Min Hill Climbing and General  
225 2-Phase Restricted Maximization) were used to perform structure learning.

226 We quantified the strength of edges and uncertainty regarding the direction of networks,  
227 using 500 bootstrapping replicates with a size equal to the number of accessions and per-  
228 formed structure learning for each replicate in accordance with Scutari and Denis (2014).  
229 Non-parametric bootstrap resampling aimed at reducing the impact of the local optimal  
230 structures by computing the probability of the arcs and directions. Subsequently, 500 learned  
231 structures were averaged with a strength threshold of 85% or higher to produce a more robust  
232 network structure. This process, known as model averaging, returns the final network with  
233 arcs present in at least 85% among all 500 networks. Candidate networks were compared  
234 on the basis of the Bayesian information criterion (BIC) and Bayesian Gaussian equivalent  
235 score (BGe). The BIC accounts for the goodness-of-fit and model complexity, and BGe aims  
236 at maximizing the posterior probability of networks per the data. All BN were learned via  
237 the bnlearn R package (Scutari, 2010). In bnlearn, the BIC score is rescaled by -2, which  
238 indicates that the larger BIC refers to a preferred model.

## 239 Data availability

240 Genotypic data regarding the rice accessions can be downloaded from the rice diversity panel  
241 website (<http://www.ricediversity.org/>). Phenotypic data used herein are available in

242 Zhao et al. (2011), Campbell et al. (2017b), and Supplementary File S3.

## 243 Results

244 To elucidate the genetic interdependencies among traits typically recorded in breeding pro-  
245 grams, we utilized a collection of 48 publicly available phenotypes recorded on a panel of  
246 diverse rice accessions (Zhao et al., 2011; Campbell et al., 2017a). The phenotypic data  
247 was derived from two independent studies. The first set of phenotypes was recorded from  
248 materials grown in two field environments in Arkansas and Faridpur Bangladesh, and in  
249 a greenhouse in Aberdeen, UK (Zhao et al., 2011). The 34 phenotypes were recorded at  
250 maturity and were largely associated with yield (panicle characteristics flowering time, plant  
251 morphology (e.g. height and growth habits), and seed morphological traits. The second  
252 study consisted of 14 phenotypes were recorded in a greenhouse environment on plants in  
253 the active tillering stage (e.g. 30 day-old plants) under control and saline (14 days of 9.5  
254 dS m<sup>-2</sup> NaCl stress). The phenotypes from this study can be classified into three cate-  
255 gories: morphological traits (e.g. shoot and root biomass, and plant height), morphological  
256 responses to salinity (e.g. the ratio of morphological traits in saline conditions to control),  
257 and the ionic components of salinity stress (e.g. Na<sup>+</sup>, K<sup>+</sup>, and Na<sup>+</sup>:K<sup>+</sup> in both root and  
258 shoot tissues) (Campbell et al., 2017a). The complete data set provides an in-depth char-  
259 acterization of phenotypic performance at vegetative and reproductive stages in rice using  
260 several classes of traits.

## 261 Latent variable modeling

262 The BCFA model grouped the observed phenotypes into the underlying latent variables  
263 on the basis of prior biological knowledge, assuming these latent variables determine the  
264 observed phenotypes. This allowed us to study the genetics of each latent variable. A  
265 measurement model derived from BCFA evaluating the six latent variables is shown in Figure  
266 1. Forty-eight observed phenotypes were hypothesized to result from the six latent variables:

267 7 for flowering time, 14 for morphology, 5 for yield, 11 for grain morphology, 6 for physiology,  
268 and 5 for salt response. The convergence of the parameters was confirmed graphically with  
269 the trace plots and a PSRF value less than 1.2 (Brooks and Gelman, 1998; Merkle and  
270 Rosseel, 2018).

271 The six latent factors showed strong contributions to the 48 observed phenotypes, with  
272 standardized regression coefficients ranging from -0.668 to 0.980 for flowering time, -0.112 to  
273 0.903 for morphology, -0.113 to 0.977 for yield, -0.501 to 0.986 for grain morphology, -0.016  
274 to 0.829 for physiology, and 0.011 to 0.929 for salt response. The latent factor flowering time  
275 showed a strong positive contribution to flowering time in Arkansas (Fla) and Flowering  
276 time in Arkansas in 2007 (Fla7) (0.99 and 0.926, respectively; Table 1, indicating that larger  
277 values for the latent factor can be interpreted as a greater number of days from sowing to  
278 emergence of the inflorescence. The latent factor morphology showed the largest positive  
279 contributions to traits describing height during the vegetative stage (e.g. height to newest  
280 ligule in salt (Hls), 0.920; height to newest ligule in control (Hlc), 0.899; height to the tip of  
281 first fully expanded leaf in salt (Hfs), 0.907; and height to tip of first fully expanded leaf in  
282 control (Hfc)), 0.925; suggesting that this latent factor is an overall representation of plant  
283 size. Yield showed large positive contributions to the observed phenotypes primary panicle  
284 branch number (Ppn) and seed number per panicle (Snp) (0.790 and 0.780, respectively),  
285 suggesting that larger values for yield indicate a higher degree of branching and seed number.  
286 Observed phenotypes describing seed size (e.g. seed volume (Sv) and brown rice volume  
287 (Bvl) (0.990 and 0.986, respectively)) were most strongly associated with grain morphology.  
288 The latent factor ionic components of salt stress showed strong positive contributions to two  
289 observed phenotypes that quantify the ionic components of salt stress (shoot  $\text{Na}^+:\text{K}^+$  (Kslm)  
290 and shoot  $\text{Na}^+$  (Nas) (0.983 and 0.975, respectively), indicating that higher values for the  
291 latent factor result in greater shoot  $\text{Na}^+$  and  $\text{Na}^+:\text{K}^+$ . Finally, the latent factor describing  
292 morphological salt response showed strong positive contributions to the observed phenotype  
293 describing the effect of salt treatment on plant height (ratio of height to tip of newest fully



294 expanded leaf in salt to that of control plants (Hfr) (0.939), thus larger values for the latent  
295 factor may indicate a more tolerant growth response to salinity.

## 296 **Genomic correlation among latent variables**

297 To understand the genetic relationships between latent variables, genomic correlation analy-  
298 sis was performed. Genomic correlation is due to pleiotropy or linkage disequilibrium between  
299 QTL. The genomic correlations among latent variables are shown in Figure 2. Negative cor-  
300 relations were observed between salt response (Slr) and all other five latent variables. In  
301 particular, flowering time (-0.5), yield (-0.54), and grain morphology (-0.74) were negatively  
302 correlated with morphological salt response 2. These results suggest that accessions that  
303 harbor alleles for more tolerant morphological salt responses may also have alleles associated  
304 with longer flowering times, smaller seeds, and low yield. Similarly, a negative correlation  
305 was observed between morphology and yield (-0.56) and between morphology and grain mor-  
306 phology (-0.31). Thus, accessions with alleles associated with large plant size may also have  
307 alleles that result in low yield, small grain volume, and lower shoot  $\text{Na}^+$  and  $\text{Na}^+:\text{K}^+$ . In  
308 contrast, a positive correlation was observed between grain morphology and yield (0.49) and  
309 between grain morphology and ionic components of salt stress (0.4). Thus, selection for large  
310 grain may result in improved yield, and higher shoot  $\text{Na}^+$  and  $\text{Na}^+:\text{K}^+$ .

## 311 **Bayesian network**

312 To infer the possible causal structure between latent variables, BN was performed. Prior  
313 to BN, the normality of latent variables was assessed using histogram plots combined with  
314 density curves as shown in Figure S2. Overall, all the six latent variables approximately  
315 followed a Gaussian distribution.

316 The Bayesian networks learned with the score-based and hybrid algorithms are shown  
317 in Figures 3, 4, 5, and 6. The structures of BN were refined by model averaging with 500

318 networks from bootstrap resampling to reduce the impact of local optimal structures. The  
319 labels of the arcs measure the uncertainty of the arcs, corresponding to strength and direc-  
320 tion (in parenthesis). The former measures the frequency of the arc presented among all 500  
321 networks from the bootstrapping replicates and the latter is the frequency of the direction  
322 shown conditional on the presence of the arc. We observed minor differences in the structures  
323 presented within and across the two types of algorithms used. In general, small differences  
324 were observed within algorithm types compared to those across algorithms. The two score-  
325 based algorithms produced a greater number of edges than two hybrid algorithms. In Figure  
326 3, the Hill Climbing algorithm produced seven directed connections among the six latent  
327 variables. Three connections were indicated towards flowering time from morphological salt  
328 response, ionic components of salt stress, and morphology, and two edges to yield from mor-  
329 phology and from grain morphology. Other two edges were observed from ionic components  
330 of salt stress to grain morphology and from grain morphology to morphological salt response.  
331 A similar structure was generated by the Tabu algorithm, except that the connection be-  
332 tween salt response and grain morphology presented an opposite direction (Figure 4). The  
333 Max-Min Hill Climbing hybrid algorithm yielded six directed edges from morphological salt  
334 response to grain morphology, from ionic components of salt stress to grain morphology, from  
335 ionic components of salt stress to flowering time, from flowering time to morphology, from  
336 morphology to yield, and from grain morphology to yield (Figure 5). An analogous structure  
337 with the only difference observed in the directed edge from morphology to flowering time was  
338 inferred with the General 2-Phase Restricted Maximization algorithm as shown in Figure 6.  
339 Across all four algorithms, there were four common directed edges: from ionic components  
340 of salt stress to flowering time and to grain morphology, and from morphology and grain  
341 morphology to yield. The most favorable network was considered the one from the Tabu  
342 algorithm, which returned the largest network score in terms of BIC (1086.61) and BGe  
343 (1080.88). Collectively, these results suggest that there may be a direct genetic influence of  
344 morphology and grain morphology on yield, and physiological components of salt tolerance

345 on grain morphology and flowering time.

## 346 Discussion

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

354 Understanding the behaviors among phenotypes in the complex traits is critical for genetic  
355 improvement of agricultural species (Hickey et al., 2017). Graphical modeling offers an av-  
356 enue to decipher bi-directional associations or probabilistic dependencies among variables of  
357 interest in plant and animal breeding. For instance, BN and L1-regularized undirected net-  
358 work can be used to model interrelationships of linkage disequilibrium (LD) (Morota et al.,  
359 2012; Morota and Gianola, 2013) or phenotypic, genetic, and environmental interactions  
360 (Xavier et al., 2017) in a systematic manner. Importantly, MTM elucidates both direct and  
361 indirect relationships among phenotypes. Inaccurate interpretation of these relationships  
362 may substantially bias selection decisions (Valente et al., 2015; Gianola et al., 2015). Thus,  
363 we applied BCFA to reduce the dimension of the responses by hypothesizing 48 manifest  
364 phenotypes originated from the underlying six constructed latent variables as shown in Fig-  
365 ure 1 assuming that these latent traits are most important, followed by application of BN to  
366 infer the structures among the six biologically relevant latent variables (Figures 3,4, 5, and  
367 6). The BN represents the conditional dependencies between variables. Care must be taken  
368 in interpreting these relationships as a causal effect. Although a good BN is expected to  
369 describe the underlying causal structure per the data, when the structure is learned solely  
370 on the basis of the observed data, it may return multiple equivalent networks that describe  
371 the data well. In practice, searching such a causal structure with observed data needs three

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

## 380 **Biological meaning of latent variables and their relation-** 381 **ships**

382 We performed BCFA to summarize the original 48 phenotypes with the six latent variables.  
383 The number of latent variables and which latent variables load onto phenotypes were deter-  
384 mined from the literature. The latent variable morphological salt response (Slr) contributed  
385 strongly to salt indices for shoot biomass, root biomass, and two indices for plant height (Ta-  
386 ble 1). Thus, morphological salt response can be interpreted as the morphological responses  
387 to salinity stress, with higher values indicating a more tolerant growth response. The la-  
388 tent variable yield is a representation of overall grain productivity, and contributed strongly  
389 to the observed phenotypes primary panicle branch number, seed number per panicle, and  
390 panicle length. The positive loading scores on these observable phenotypes indicates that  
391 more highly branched, productive panicles will have higher values for yield (Table 1). Seed  
392 width, seed volume, and seed surface area contributed significantly to the latent variable  
393 grain morphology (Grm) (Table 1). Therefore, these results indicate that the grain mor-  
394 phology is a summary of the overall shape of the grain, where high values represent large,  
395 round grains, while low values represent small, slender grains. Considering the grain char-  
396 acteristics of rice subpopulations, temperate japonica accessions are expected to have high

397 values for grain morphology, while indica accessions have lower values for grain morphology.  
398 Latent variable morphology (Mrp) is a representation of plant biomass during the vegetative  
399 stage (28-day-old plants) (Table 1). Shoot biomass, root biomass, and two metrics for plant  
400 height contributed largely to morphology, suggesting that accessions with high values for  
401 morphology are tall plants with a large biomass.

402 Genomic correlation analysis among the six latent variables showed meaningful corre-  
403 lations among several pairs. These genetic correlations can either be caused by linkage or  
404 pleiotropy. The former is likely to prevail in species with high LD, which is the case in  
405 rice where LD ranges from 100 to 200kb (Huang et al., 2010). A negative relationship was  
406 observed between morphological salt response and three other latent variables (Figure 2).  
407 For instance, a negative correlation between morphological salt response and yield indicates  
408 that accessions of samples harboring alleles for superior morphological salt responses (e.g.  
409 those that are more tolerant) tend to also harbor alleles for poor yield (Figure 2). The  
410 rice diversity panel we used is a representative sample of the total genetic diversity within  
411 cultivated rice and contains many unimproved traditional varieties ( $\sim 12\%$  of lines in the  
412 study are landraces and  $\sim 33\%$  classified as cultivars; Supplementary File S2) and modern  
413 breeding lines (Eizenga et al., 2014). While traditional varieties exhibit superior adaptation  
414 to abiotic stresses, they often have very poor agronomic characteristics including low yield,  
415 late flowering, and high photoperiod sensitivity (Thomson et al., 2009, 2010). Moreover,  
416 the indica and japonica subspecies have contrasting salt responses and very different grain  
417 morphology. Japonica accessions tend to have short, round seeds and are more sensitive to  
418 salt stress, while indica accessions have long, slender grains and often are more salt tolerant  
419 (Zhao et al., 2011; Campbell et al., 2017a). The negative relationship observed between salt  
420 response and grain morphology suggests that lines that harbor alleles for high grain mor-  
421 phology (e.g., large, round grains) tend to also harbor alleles for a tolerant growth response  
422 to salt stress. However, no studies have yet reported an association between alleles for grain  
423 morphology and morphological salt response. Therefore, it remains to be addressed whether

424 this relationship is due to LD or pleiotropy.

425 Genetic correlations observed between other latent variables may suggest a pleiotropic  
426 effect among loci. For instance, a negative relationship was observed between morphological  
427 salt response and ionic components of salt stress, indicating that accessions harboring alleles  
428 associated with superior morphological salt response also tend to harbor alleles for reduced  
429 ion content under salt stress (Figure 2). The relationship between salt tolerance, measured in  
430 terms of growth or yield, and  $\text{Na}^+$  and  $\text{Na}^+:\text{K}^+$  has been a documented for decades (reviewed  
431 by Munns and Tester (2008)). Moreover, natural variation for  $\text{Na}^+$  transporters has been  
432 utilized to improve growth and yield under saline conditions in rice and other cereals (Ren  
433 et al., 2005; Byrt et al., 2007; Horie et al., 2009; Munns et al., 2012; Campbell et al., 2017a).  
434 Therefore, the negative genetic relationships observed between morphological salt response  
435 and ion content may be due to the pleiotropic effects of some loci.

436 The genomic relationships among latent variables including morphology, yield, and grain  
437 morphology may have resulted from the selection of alleles associated with good agronomic  
438 characteristics. A positive relationship was observed between yield and grain morphology,  
439 suggesting that alleles that positively contribute to productive panicles also may contribute  
440 to large, round grains. Furthermore, the negative genomic correlation observed between  
441 morphology and yield indicates that alleles negatively influencing total plant biomass also  
442 have a positive contribution to traits for productive panicles. This genomic relationship may  
443 reflect the genetics of harvest index, which is defined as the ratio of grain yield to total  
444 biomass. Over the past 50 years, rice breeders have selected high harvest index, resulting  
445 in plants with short compact morphology and many highly productive panicles (Hay, 1995;  
446 Peng et al., 2008).

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

451 the BCFA framework. An alternative approach is to derive the number of latent variables  
452 and determine which latent variables load onto phenotypes directly from observed data,  
453 using exploratory FA. This approach was not pursued here because accurate estimation of  
454 unknown parameters in the exploratory FA requires a large sample size, which was not the  
455 case herein (Brown, 2014).

## 456 **Bayesian network of latent variables**

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

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



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

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

498 We found that there are some uncertain edges among BN. For instance, direction from  
499 salt response to grain morphology is supported by 65% (Figure 4), 58% (Figure 5), and 58%  
500 (Figure 6) bootstrap sampling, whereas the opposite direction is supported by 56% bootstrap  
501 sampling (Figure 3). An analogous uncertainty was also observed between morphology and  
502 flowering time, i.e., the path from morphology to flowering time was supported 60% (Figure  
503 3), 51% (Figure 4), and 52% (Figure 6), while the reverse direction was supported 51%

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

514 In conclusion, the present results show the utility of CFA and network analysis to char-  
515 acterize various phenotypes in rice. We showed that the joint use of BCFA and BN can be  
516 applied to predict the potential influence of external interventions or selection associated with  
517 target traits such as yield in the high-dimensional interrelated complex traits system. We  
518 contend that the approaches used herein provide greater insights than pairwise-association  
519 measures of multiple phenotypes and can be used to analyze the massive amount of di-  
520 verse image-based phenomics dataset being generated by the automated plant phenomics  
521 platforms (e.g., Furbank and Tester, 2011). With a large volume of complex traits being  
522 collected through phenomics, numerous opportunities to forge new research directions are  
523 generated by using network analysis for the growing number of phenotypes.

## References

- 524
- 525 Acquaah, G. (2009). *Principles of plant genetics and breeding*. John Wiley & Sons.
- 526 Araus, J. L. and Cairns, J. E. (2014). Field high-throughput phenotyping: the new crop  
527 breeding frontier. *Trends in plant science*, 19(1):52–61.
- 528 Araus, J. L., Kefauver, S. C., Zaman-Allah, M., Olsen, M. S., and Cairns, J. E. (2018).  
529 Translating high-throughput phenotyping into genetic gain. *Trends in plant science*.
- 530 Bhattacharya, A., Kourmpetli, S., and Davey, M. R. (2010). Practical applications of manip-  
531 ulating plant architecture by regulating gibberellin metabolism. *Journal of plant growth*  
532 *regulation*, 29(2):249–256.
- 533 Brewer, P. B., Koltai, H., and Beveridge, C. A. (2013). Diverse roles of strigolactones in  
534 plant development. *Molecular plant*, 6(1):18–28.
- 535 Brooks, S. P. and Gelman, A. (1998). General methods for monitoring convergence of  
536 iterative simulations. *Journal of computational and graphical statistics*, 7(4):434–455.
- 537 Brown, T. A. (2014). *Confirmatory factor analysis for applied research*. Guilford Publica-  
538 tions.
- 539 Byrt, C. S., Platten, J. D., Spielmeier, W., James, R. A., Lagudah, E. S., Dennis, E. S.,  
540 Tester, M., and Munns, R. (2007). Hkt1; 5-like cation transporters linked to na<sup>+</sup> exclusion  
541 loci in wheat, nax2 and kna1. *Plant Physiology*, 143(4):1918–1928.
- 542 Cabrera-Bosquet, L., Fournier, C., Bricet, N., Welcker, C., Suard, B., and Tardieu, F.  
543 (2016). High-throughput estimation of incident light, light interception and radiation-use  
544 efficiency of thousands of plants in a phenotyping platform. *New Phytologist*, 212(1):269–  
545 281.

- 546 Callanan, T. P. and Harville, D. A. (1989). *Some new algorithms for computing maximum*  
547 *likelihood estimates of variance components*. Iowa State University. Department of Statis-  
548 tics. Statistical Laboratory.
- 549 Calus, M. P. and Veerkamp, R. F. (2011). Accuracy of multi-trait genomic selection using  
550 different methods. *Genetics Selection Evolution*, 43(1):26.
- 551 Campbell, M. T., Bandillo, N., Al Shiblawi, F. R. A., Sharma, S., Liu, K., Du, Q., Schmitz,  
552 A. J., Zhang, C., Véry, A.-A., Lorenz, A. J., et al. (2017a). Allelic variants of *oshkt1*; 1  
553 underlie the divergence between indica and japonica subspecies of rice (*oryza sativa*) for  
554 root sodium content. *PLoS Genetics*, 13(6):e1006823.
- 555 Campbell, M. T., Du, Q., Liu, K., Brien, C. J., Berger, B., Zhang, C., and Walia, H. (2017b).  
556 A comprehensive image-based phenomic analysis reveals the complex genetic architecture  
557 of shoot growth dynamics in rice (*oryza sativa*). *The Plant Genome*, 10(2).
- 558 Crain, J., Mondal, S., Rutkoski, J., Singh, R. P., and Poland, J. (2018). Combining high-  
559 throughput phenotyping and genomic information to increase prediction and selection  
560 accuracy in wheat breeding. *The plant genome*.
- 561 de los Campos, G. and Gianola, D. (2007). Factor analysis models for structuring covari-  
562 ance matrices of additive genetic effects: a bayesian implementation. *Genetics Selection*  
563 *Evolution*, 39(5):481.
- 564 Denwood, M. (2016). runjags: An r package providing interface utilities, model templates,  
565 parallel computing methods and additional distributions for mcmc models in jags. *Journal*  
566 *of Statistical Software, Articles*, 71(9):1–25.
- 567 Eizenga, G. C., Ali, M., Bryant, R. J., Yeater, K. M., McClung, A. M., McCouch, S. R.,  
568 et al. (2014). Registration of the rice diversity panel 1 for genomewide association studies.  
569 *Journal of Plant Registrations*, 8(1):109–116.

- 570 Furbank, R. T. and Tester, M. (2011). Phenomics-technologies to relieve the phenotyping  
571 bottleneck. *Trends Plant Sci.*, 16:635–644.
- 572 Gianola, D., de los Campos, G., Toro, M. A., Naya, H., Schön, C.-C., and Sorensen, D.  
573 (2015). Do molecular markers inform about pleiotropy? *Genetics*, pages genetics–115.
- 574 Hay, R. (1995). Harvest index: a review of its use in plant breeding and crop physiology.  
575 *Annals of applied biology*, 126(1):197–216.
- 576 Henderson, C. and Quaas, R. (1976). Multiple trait evaluation using relatives' records.  
577 *Journal of Animal Science*, 43(6):1188–1197.
- 578 Hickey, J. M., Chiurugwi, T., Mackay, I., Powell, W., Eggen, A., Kilian, A., Jones, C.,  
579 Canales, C., Grattapaglia, D., Bassi, F., et al. (2017). Genomic prediction unifies animal  
580 and plant breeding programs to form platforms for biological discovery. *Nature genetics*,  
581 49(9):1297.
- 582 Horie, T., Hauser, F., and Schroeder, J. I. (2009). Hkt transporter-mediated salinity re-  
583 sistance mechanisms in arabidopsis and monocot crop plants. *Trends in plant science*,  
584 14(12):660–668.
- 585 Hornik, K., Leisch, F., and Zeileis, A. (2003). Jags: A program for analysis of bayesian  
586 graphical models using gibbs sampling. In *Proceedings of DSC*, volume 2, pages 1–1.
- 587 Huang, R., Jiang, L., Zheng, J., Wang, T., Wang, H., Huang, Y., and Hong, Z. (2013).  
588 Genetic bases of rice grain shape: so many genes, so little known. *Trends in plant science*,  
589 18(4):218–226.
- 590 Huang, X., Sang, T., Zhao, Q., Feng, Q., Zhao, Y., Li, C., Zhu, C., Lu, T., Zhang, Z., Li, M.,  
591 et al. (2010). Genome-wide association studies of 14 agronomic traits in rice landraces.  
592 *Nature Genetics*, 42(11):961.

- 593 Jia, Y. and Jannink, J.-L. (2012). Multiple trait genomic selection methods increase genetic  
594 value prediction accuracy. *Genetics*, pages genetics–112.
- 595 Jiao, Y., Wang, Y., Xue, D., Wang, J., Yan, M., Liu, G., Dong, G., Zeng, D., Lu, Z., Zhu,  
596 X., et al. (2010). Regulation of *osspl14* by *osmir156* defines ideal plant architecture in rice.  
597 *Nature genetics*, 42(6):541.
- 598 Jöreskog, K. G. (1969). A general approach to confirmatory maximum likelihood factor  
599 analysis. *Psychometrika*, 34(2):183–202.
- 600 Lee, S.-Y. and Song, X.-Y. (2012). *Basic and advanced Bayesian structural equation model-*  
601 *ing: With applications in the medical and behavioral sciences*. John Wiley & Sons.
- 602 Liang, W.-h., Shang, F., Lin, Q.-t., Lou, C., and Zhang, J. (2014). Tillerling and panicle  
603 branching genes in rice. *Gene*, 537(1):1–5.
- 604 Lo, S.-F., Yang, S.-Y., Chen, K.-T., Hsing, Y.-I., Zeevaart, J. A., Chen, L.-J., and Yu, S.-M.  
605 (2008). A novel class of gibberellin 2-oxidases control semidwarfism, tillering, and root  
606 development in rice. *The Plant Cell*, 20(10):2603–2618.
- 607 McCouch, S. R., Wright, M. H., Tung, C.-W., Maron, L. G., McNally, K. L., Fitzgerald,  
608 M., Singh, N., DeClerck, G., Agosto-Perez, F., Korniliev, P., et al. (2016). Open access  
609 resources for genome-wide association mapping in rice. *Nature communications*, 7:10532.
- 610 Merkle, E. and Rosseel, Y. (2018). blavaan: Bayesian structural equation models via pa-  
611 rameter expansion. *Journal of Statistical Software, Articles*, 85(4):1–30.
- 612 Miura, K., Ikeda, M., Matsubara, A., Song, X.-J., Ito, M., Asano, K., Matsuoka, M., Ki-  
613 tano, H., and Ashikari, M. (2010). *Osspl14* promotes panicle branching and higher grain  
614 productivity in rice. *Nature genetics*, 42(6):545.
- 615 Morota, G. and Gianola, D. (2013). Evaluation of linkage disequilibrium in wheat with an l1-  
616 regularized sparse markov network. *Theoretical and Applied Genetics*, 126(8):1991–2002.

- 617 Morota, G., Valente, B., Rosa, G., Weigel, K., and Gianola, D. (2012). An assessment  
618 of linkage disequilibrium in holstein cattle using a bayesian network. *Journal of Animal*  
619 *Breeding and Genetics*, 129(6):474–487.
- 620 Mrode, R. A. (2014). *Linear models for the prediction of animal breeding values*. Cabi.
- 621 Munns, R., James, R. A., Xu, B., Athman, A., Conn, S. J., Jordans, C., Byrt, C. S., Hare,  
622 R. A., Tyerman, S. D., Tester, M., et al. (2012). Wheat grain yield on saline soils is  
623 improved by an ancestral na<sup>+</sup> transporter gene. *Nature biotechnology*, 30(4):360.
- 624 Munns, R. and Tester, M. (2008). Mechanisms of salinity tolerance. *Annu. Rev. Plant Biol.*,  
625 59:651–681.
- 626 Neapolitan, R. E. et al. (2004). *Learning bayesian networks*, volume 38. Pearson Prentice  
627 Hall Upper Saddle River, NJ.
- 628 Pearl, J. (2009). *Causality: Models, Reasoning and Inference*. Cambridge University Press,  
629 New York, NY, USA, 2nd edition.
- 630 Peñagaricano, F., Valente, B., Steibel, J., Bates, R., Ernst, C., Khatib, H., and Rosa,  
631 G. (2015). Searching for causal networks involving latent variables in complex traits:  
632 application to growth, carcass, and meat quality traits in pigs. *Journal of Animal Science*,  
633 93(10):4617–4623.
- 634 Peng, S., Khush, G. S., Virk, P., Tang, Q., and Zou, Y. (2008). Progress in ideotype breeding  
635 to increase rice yield potential. *Field Crops Research*, 108(1):32–38.
- 636 Redona, E. and Mackill, D. (1998). Quantitative trait locus analysis for rice panicle and  
637 grain characteristics. *Theoretical and Applied Genetics*, 96(6-7):957–963.
- 638 Ren, Z.-H., Gao, J.-P., Li, L.-G., Cai, X.-L., Huang, W., Chao, D.-Y., Zhu, M.-Z., Wang,  
639 Z.-Y., Luan, S., and Lin, H.-X. (2005). A rice quantitative trait locus for salt tolerance  
640 encodes a sodium transporter. *Nature genetics*, 37(10):1141.

- 641 Rosa, G. J., Valente, B. D., de los Campos, G., Wu, X.-L., Gianola, D., and Silva, M. A.  
642 (2011). Inferring causal phenotype networks using structural equation models. *Genetics*  
643 *Selection Evolution*, 43(1):6.
- 644 Scutari, M. (2010). Learning bayesian networks with the bnlearn r package. *Journal of*  
645 *Statistical Software, Articles*, 35(3):1–22.
- 646 Scutari, M. and Denis, J.-B. (2014). *Bayesian networks: with examples in R*. Chapman and  
647 Hall/CRC.
- 648 Shakoor, N., Lee, S., and Mockler, T. C. (2017). High throughput phenotyping to accelerate  
649 crop breeding and monitoring of diseases in the field. *Current opinion in plant biology*,  
650 38:184–192.
- 651 Sun, J., Rutkoski, J. E., Poland, J. A., Crossa, J., Jannink, J.-L., and Sorrells, M. E.  
652 (2017). Multitrait, random regression, or simple repeatability model in high-throughput  
653 phenotyping data improve genomic prediction for wheat grain yield. *The plant genome*.
- 654 Tanner, M. A. and Wong, W. H. (1987). The calculation of posterior distributions by data  
655 augmentation. *Journal of the American statistical Association*, 82(398):528–540.
- 656 Thomson, M. J., de Ocampo, M., Egdane, J., Rahman, M. A., Sajise, A. G., Adorada, D. L.,  
657 Tumimbang-Raiz, E., Blumwald, E., Seraj, Z. I., Singh, R. K., et al. (2010). Characterizing  
658 the saltol quantitative trait locus for salinity tolerance in rice. *Rice*, 3(2-3):148–160.
- 659 Thomson, M. J., Ismail, A. M., McCouch, S. R., and Mackill, D. J. (2009). Marker assisted  
660 breeding. In *Abiotic Stress Adaptation in Plants*, pages 451–469. Springer.
- 661 Töpner, K., Rosa, G. J., Gianola, D., and Schön, C.-C. (2017). Bayesian networks illustrate  
662 genomic and residual trait connections in maize (*Zea mays* L.). *G3: Genes, Genomes,*  
663 *Genetics*, 7(8):2779–2789.



- 664 Umehara, M., Hanada, A., Yoshida, S., Akiyama, K., Arite, T., Takeda-Kamiya, N.,  
665 Magome, H., Kamiya, Y., Shirasu, K., Yoneyama, K., et al. (2008). Inhibition of shoot  
666 branching by new terpenoid plant hormones. *Nature*, 455(7210):195.
- 667 Valente, B. D., Morota, G., Peñagaricano, F., Gianola, D., Weigel, K., and Rosa, G. J.  
668 (2015). The causal meaning of genomic predictors and how it affects construction and  
669 comparison of genome-enabled selection models. *Genetics*, 200(2):483–494.
- 670 VanRaden, P. M. (2008). Efficient methods to compute genomic predictions. *J. Dairy Sci.*,  
671 91:4414–4423.
- 672 Vazquez, A., Bates, D., Rosa, G., Gianola, D., and Weigel, K. (2010). An r package for  
673 fitting generalized linear mixed models in animal breeding 1. *Journal of animal science*,  
674 88(2):497–504.
- 675 Verma, T. and Pearl, J. (1991). Equivalence and synthesis of causal models. In *Proceedings*  
676 *of the Sixth Annual Conference on Uncertainty in Artificial Intelligence, UAI '90*, pages  
677 255–270, New York, NY, USA. Elsevier Science Inc.
- 678 Wang, Y. and Li, J. (2006). Genes controlling plant architecture. *Current Opinion in*  
679 *Biotechnology*, 17(2):123–129.
- 680 Xavier, A., Hall, B., Casteel, S., Muir, W., and Rainey, K. M. (2017). Using unsupervised  
681 learning techniques to assess interactions among complex traits in soybeans. *Euphytica*,  
682 213(8):200.
- 683 Yan, W.-H., Wang, P., Chen, H.-X., Zhou, H.-J., Li, Q.-P., Wang, C.-R., Ding, Z.-H., Zhang,  
684 Y.-S., Yu, S.-B., Xing, Y.-Z., et al. (2011). A major qtl, *ghd8*, plays pleiotropic roles in  
685 regulating grain productivity, plant height, and heading date in rice. *Molecular plant*,  
686 4(2):319–330.

687 Zhao, K., Tung, C.-W., Eizenga, G. C., Wright, M. H., Ali, M. L., Price, A. H., Norton,  
688 G. J., Islam, M. R., Reynolds, A., Mezey, J., et al. (2011). Genome-wide association  
689 mapping reveals a rich genetic architecture of complex traits in *oryza sativa*. *Nature*  
690 *Communications*, 2:467.

691 Zhou, F., Lin, Q., Zhu, L., Ren, Y., Zhou, K., Shabek, N., Wu, F., Mao, H., Dong, W.,  
692 Gan, L., et al. (2013). D14-scf d3-dependent degradation of d53 regulates strigolactone  
693 signalling. *Nature*, 504(7480):406.

## 694 Tables

Table 1: Standardized factor loadings obtained from the Bayesian confirmatory factor analysis. PSD refers to the posterior standard deviation of standardized factor loadings.

Latent variable	Observed phenotype	Loading	PSD
Flowering time	Flowering time at Arkansas (Fla)	0.990	0.002
Flowering time	Flowering time at Faridpur (Flf)	0.500	0.045
Flowering time	Flowering time at Aberdeen (Flb)	0.578	0.038
Flowering time	FT ratio of Arkansas/Aberdeen (Flaa)	-0.212	0.053
Flowering time	FT ratio of Faridpur/Aberdeen (Flfa)	-0.549	0.041
Flowering time	Year07 Flowering time at Arkansas (Fla7)	0.926	0.008
Flowering time	Year06 Flowering time at Arkansas (Fla6)	0.886	0.013
Morphology	Culm habit (Cuh)	0.227	0.027
Morphology	Flag leaf length (FlL)	0.116	0.057
Morphology	Flag leaf width (Flw)	-0.044	0.058
Morphology	Plant height (Plh)	0.440	0.047
Morphology	Shoot BM Control (Sbc)	0.534	0.042
Morphology	Shoot BM Salt (Sbs)	0.456	0.048
Morphology	Root BM Control (Rbc)	0.418	0.048
Morphology	Root BM Salt (Rbs)	0.280	0.054
Morphology	Tiller No Salt (Tns)	-0.349	0.051
Morphology	Tiller No Control (Tbc)	-0.318	0.052
Morphology	Ht Lig Salt (Hls)	0.920	0.011
Morphology	Ht Lig Control (Hlc)	0.899	0.014
Morphology	Ht FE Salt (Hfs)	0.907	0.013
Morphology	Ht FE Control (Hfc)	0.925	0.011
Yield	Panicle number per plant (Pnu)	0.190	0.020
Yield	Panicle length (Pal)	0.455	0.057
Yield	Primary panicle branch number (Ppn)	0.790	0.041
Yield	Seed number per panicle (Snp)	0.780	0.043
Yield	Panicle fertility (Paf)	-0.085	0.081
Grain Morphology	Seed length (Sl)	0.251	0.029
Grain Morphology	Seed width (Sw)	0.876	0.015
Grain Morphology	Seed volume (Sv)	0.990	0.002
Grain Morphology	Seed surface area (Ssa)	0.901	0.012
Grain Morphology	Brown rice seed length (Bsl)	0.158	0.055
Grain Morphology	Brown rice seed width (Bsw)	0.837	0.019
Grain Morphology	Brown rice surface area (Bsa)	0.902	0.012
Grain Morphology	Brown rice volume (Bvl)	0.986	0.002
Grain Morphology	Seed length/width ratio (Slwr)	-0.476	0.045
Grain Morphology	Brown rice length/width ratio (Blwr)	-0.432	0.047
Grain Morphology	Grain length McCouch2016 (Glmc)	0.047	0.064
Ionic components of salt stress	Na K Shoot (Ks)	0.983	0.003
Ionic components of salt stress	Na Shoot (Nas)	0.975	0.004
Ionic components of salt stress	K Shoot Salt (Kss)	-0.265	0.051
Ionic components of salt stress	Na K Root (Kr)	0.061	0.052
Ionic components of salt stress	Na Root (Nar)	0.000	0.053
Ionic components of salt stress	K Root Salt (Krs)	-0.095	0.052
Morphological salt response	Shoot BM Ratio (Sbr)	0.410	0.047
Morphological salt response	Root BM Ratio (Rbr)	0.395	0.051
Morphological salt response	Tiller No Ratio (Tbr)	-0.022	0.057
Morphological salt response	Ht Lig Ratio (Hlr)	0.665	0.036
Morphological salt response	Ht FE Ratio (Hfr)	0.939	0.019

## 695 Figures

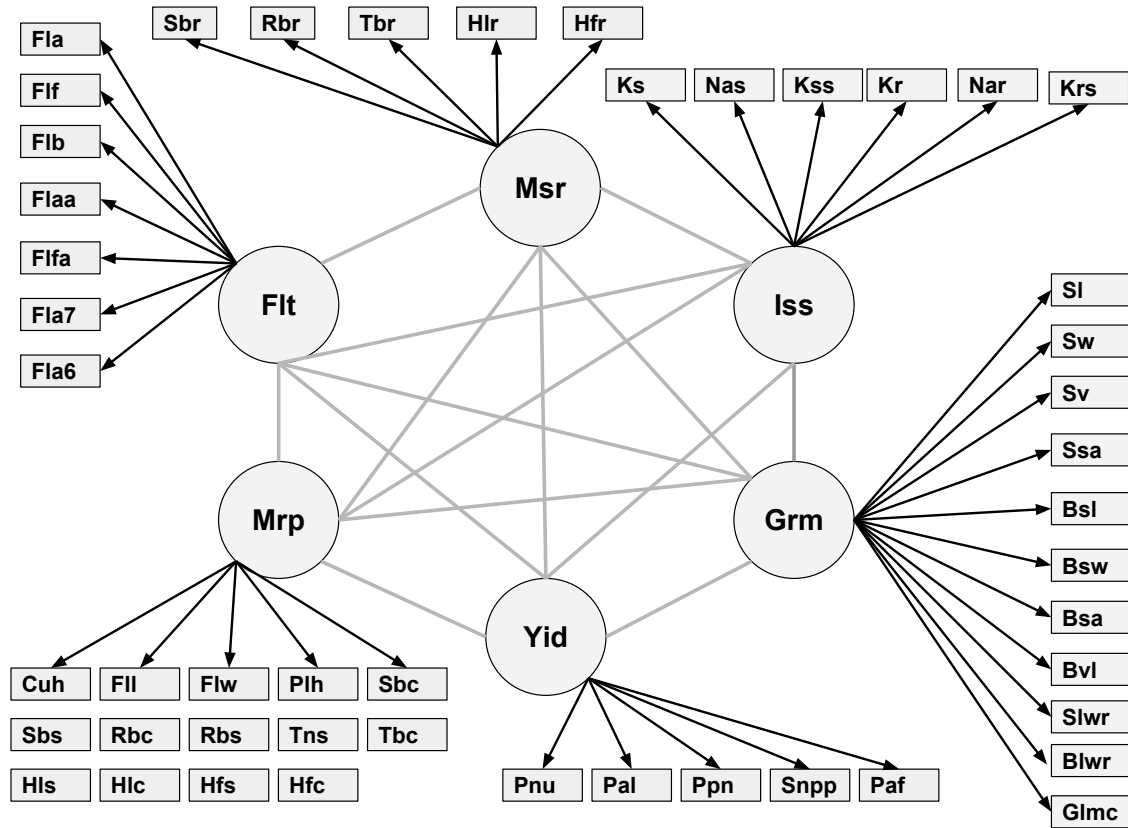


Figure 1: Relationship between six latent variables and observed phenotypes. Msr: morphological salt response; Iss: ionic components of salt stress; Grm: grain morphology; Yid: yield; Mrp: morphology; Flt: flowering time. Abbreviations of observed phenotypes are shown in Table S1.

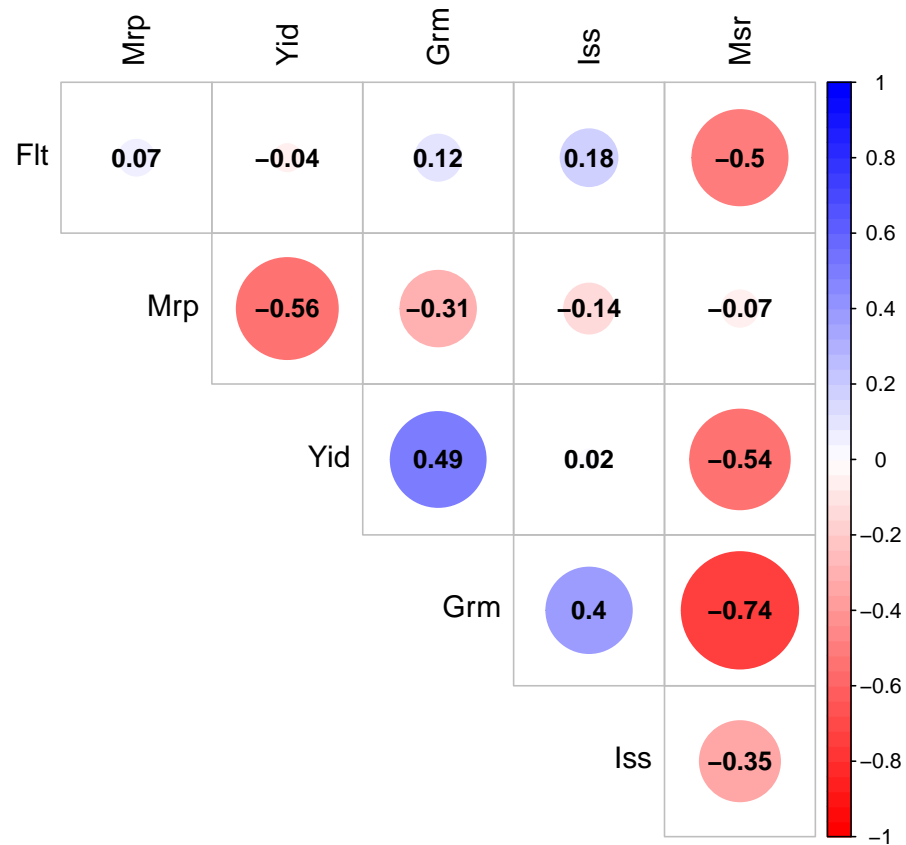


Figure 2: Genomic correlation of six latent variables. The size of each circle, degree of shading, and value reported correspond to the correlation between each pair of latent variables. Msr: morphological salt response; Iss: ionic components of salt stress; Grm: grain morphology; Yid: yield; Mrp: morphology;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

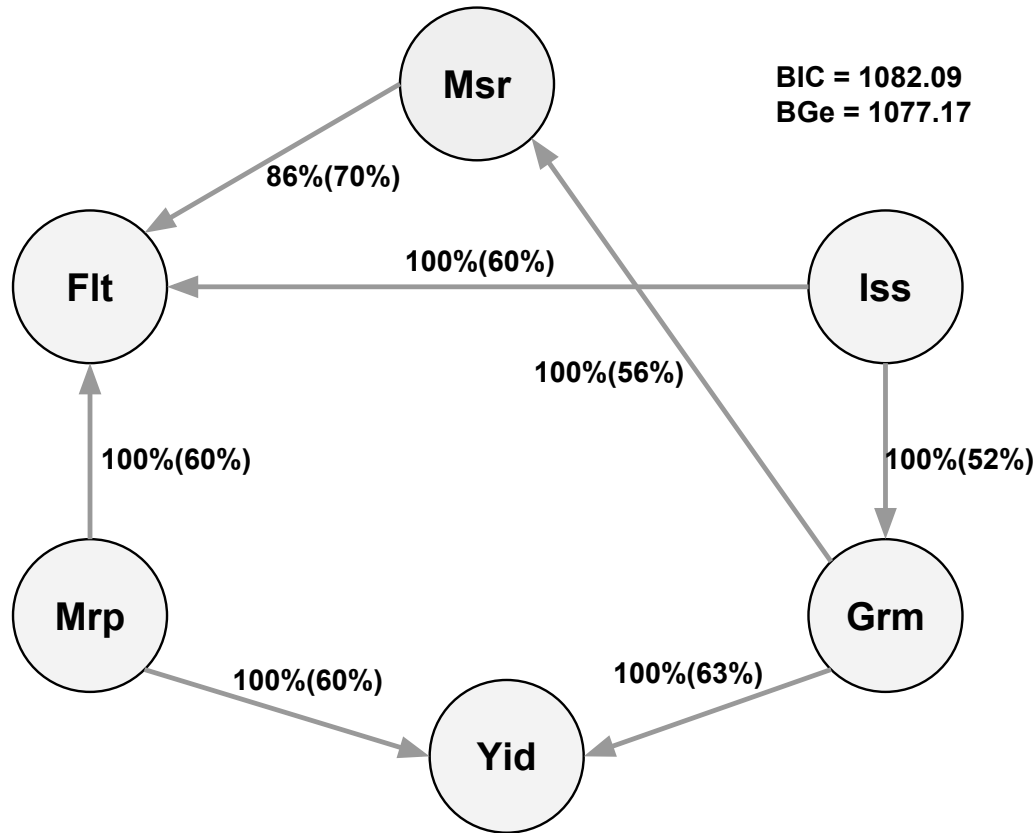


Figure 3: Bayesian network between six latent variables based on the Hill Climbing algorithm. The quality of the structure was evaluated by bootstrap resampling and model averaging across 500 replications. Labels of the edges refer to the strength and direction (parenthesis) which measure the confidence 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.

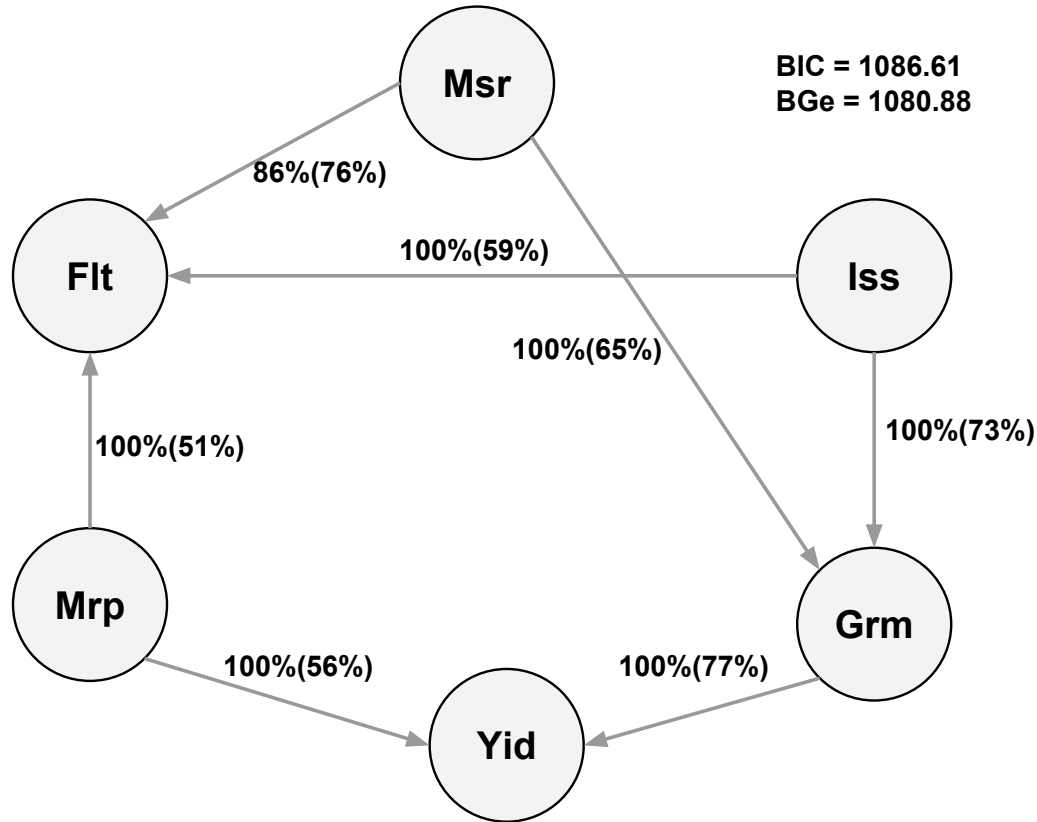


Figure 4: Bayesian network between six latent variables based on the Tabu algorithm. The quality of the structure was evaluated by bootstrap resampling and model averaging across 500 replications. Labels of the edges refer to the strength and direction (parenthesis) which measure the confidence of the 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.

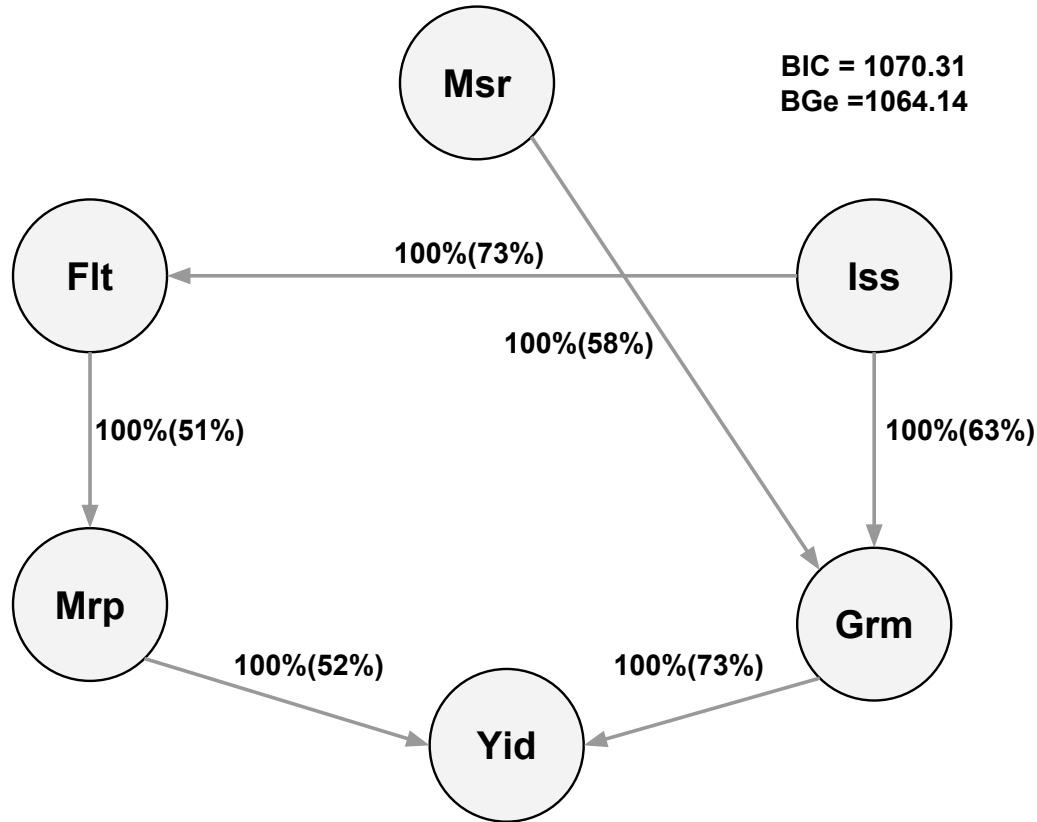


Figure 5: Bayesian network between six latent variables based on the Max-Min Hill Climbing algorithm. The quality of the structure was evaluated by bootstrap resampling and model averaging across 500 replications. Labels of the edges refer to the strength and direction (parenthesis) which measure the 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.



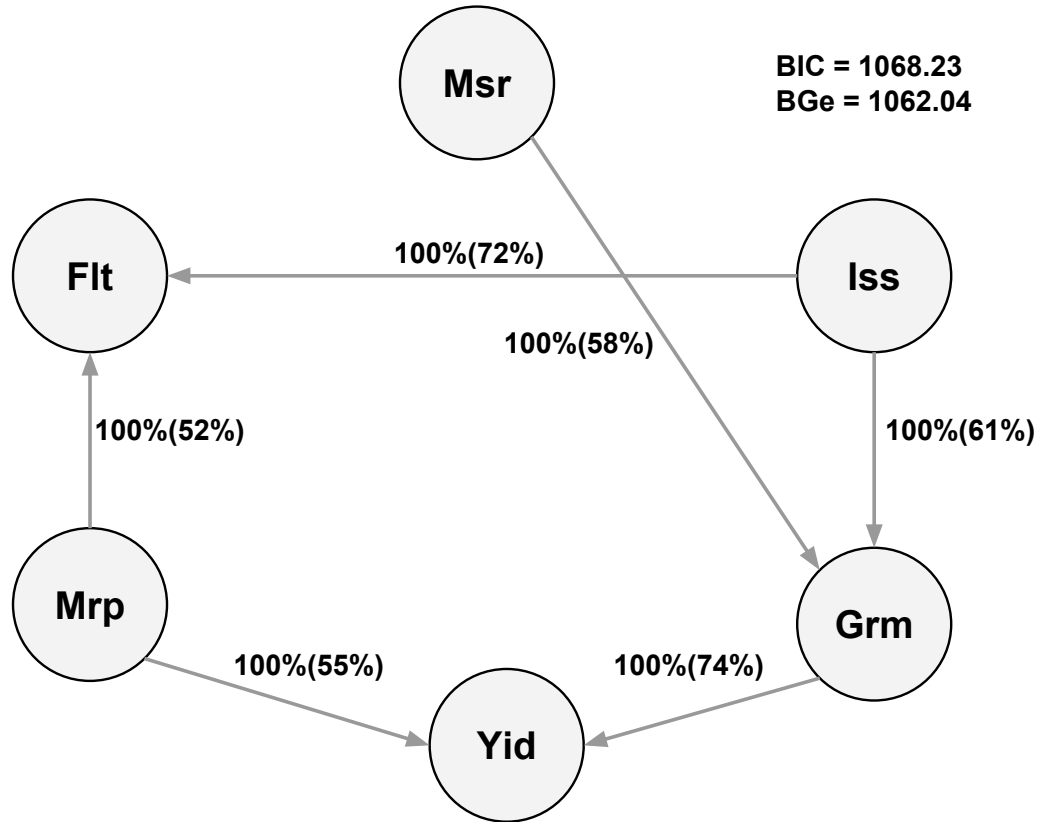


Figure 6: Bayesian network between six latent variables based on the General 2-Phase Restricted Maximization algorithm. The quality of the structure was evaluated by bootstrap resampling and model averaging across 500 replications. Labels of the edges refer to the strength and direction (parenthesis) 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.