

# Including phenotypic causal networks in genome-wide association studies using mixed effects structural equation models

**Running Head:** Structural equation modeling for association studies

Mehdi Momen<sup>1</sup>, Ahmad Ayatollahi Mehrgardi<sup>1\*</sup>, Mahmoud Amiri Roudbar<sup>1</sup>, Andreas Kranis<sup>2</sup>, Renan Mercuri Pinto<sup>3,4</sup>, Bruno D. Valente<sup>4</sup>, Gota Morota<sup>5</sup>, Guilherme J. M. Rosa<sup>4,6</sup>, Daniel Gianola<sup>4,6,7</sup>

<sup>1</sup> Department of Animal Science, Faculty of Agriculture, Shahid Bahonar University of Kerman (SBUK), Kerman, Iran

<sup>2</sup> Roslin Institute, University of Edinburgh, Midlothian, UK, EH25 9PS

<sup>3</sup> Department of Exact Sciences, University of São Paulo - ESALQ, Piracicaba-SP, Brazil

<sup>4</sup> Department of Animal Sciences, University of Wisconsin, Madison, WI, USA

<sup>5</sup> Department of Animal Science, University of Nebraska-Lincoln, Lincoln, Nebraska, USA

<sup>6</sup> Department of Biostatistics and Medical Informatics, University of Wisconsin, Madison, WI, USA

<sup>7</sup> Department of Dairy Science, University of Wisconsin, Madison, WI, USA

\*Corresponding author: mehrgardi@uk.ac.ir

Email addresses:

MM: momenmehdi@yahoo.com

AAM: mehrgardi@uk.ac.ir

MA: mahmoud.amiri225@gmail.com

GM: morota@unl.edu

AK: andreas.kranis@roslin.ed.ac.uk

RMP: rpinto@wisc.edu

BDV: bvalente@wisc.edu

GJMR: grosa@wisc.edu

DG: gianola@ansci.wisc.edu

## 31 **Abstract**

### 32 **Background**

33 Phenotypic networks describing putative causal relationship among multiple phenotypes can be  
34 used to infer single-nucleotide polymorphism (SNP) effects in genome-wide association studies  
35 (GWAS). In GWAS with multiple phenotypes, reconstructing underlying causal structures  
36 among traits and SNPs using a single statistical framework is essential for understanding the  
37 entirety of genotype-phenotype maps. A structural equation model (SEM) can be used for such  
38 purpose.

### 39 **Methods**

40 We applied SEM to GWAS (SEM-GWAS) in chickens, taking into account putative causal  
41 relationships among body weight (BW), breast meat (BW), hen-house production (HHP), and  
42 SNPs. We assessed the performance of SEM-GWAS by comparing the model results with those  
43 obtained from traditional multi-trait association analyses (MTM-GWAS).

### 44 **Results**

45 Three different putative causal path diagrams were inferred from highest posterior density (HPD)  
46 intervals of 0.75, 0.85, and 0.95 using the IC algorithm. A positive path coefficient was estimated  
47 for  $BM \rightarrow BW$ , and negative values were obtained for  $BM \rightarrow HHP$  and  $BW \rightarrow HHP$  in all  
48 implemented scenarios. Further, the application of SEM-GWAS enabled decomposition of SNP  
49 effects into direct, indirect, and total effects, identifying whether a SNP effect is acting directly or  
50 indirectly on given trait. In contrast, MTM-GWAS only captured overall genetic effects on traits,  
51 which is equivalent to combining the direct and indirect SNP effects from SEM-GWAS.

## 52 **Conclusions**

53 Our results suggested that SEM-GWAS provides insights into mechanisms by which SNPs affect  
54 traits through partitioning effects into direct, indirect, and total components. Thus, we provide  
55 evidence that SEM-GWAS captures complex relationships and delivers a more comprehensive  
56 understanding of SNP effects compared to MTM-GWAS.

57 **Key words:** Causal structure, GWAS, multiple traits, path analysis, SEM, SNP effect

58

## 59 **Background**

60 Genome-wide association studies (GWAS) have become a standard approach for investigating  
61 relationships between common genetic variants in the genome (e.g., single-nucleotide  
62 polymorphisms, SNPs) and phenotypes of interest in human, plant, and animal genetics [1-3]. A  
63 typical GWAS is based on univariate linear or logistic regression of phenotypes on genotypes for  
64 each SNP individually while often adjusting for the presence of nuisance covariates [4]. A  
65 statistically significant association indicates that SNPs may be in strong linkage disequilibrium  
66 (LD) with quantitative trait loci (QTLs) that contribute to the trait etiology. Alternatively, multi-  
67 trait model GWAS (MTM-GWAS) can be used to test for genetic associations among a set of  
68 traits [5-7]. It has been established that MTM-GWAS reduces false positives and increases the  
69 statistical power of association tests, explaining the recent popularity of this method. MTM-  
70 GWAS can be used to study genetic associations of multiple traits; however, it does not identify  
71 factors that mediate relationships between the detected effects and dependencies involving  
72 complex traits.

73 Complex traits are the product of various cryptic biological signals that may affect a trait of  
74 interest either directly or indirectly through other intermediate traits [8]. A standard regression  
75 cannot describe such complex relationships between traits and QTLs properly. For instance, some  
76 traits may simultaneously act as both dependent and independent variables. Structural equation  
77 modeling (SEM) is an extended version of Wright's path analysis [9, 10], that offers a powerful  
78 technique for modeling causal networks. In a complex genotype-phenotype setting involving  
79 many traits, a given trait can be influenced not only by genetic and systematic factors but also by  
80 other traits (as covariates) as well. Here, QTLs may not affect the target trait directly; instead, the  
81 effects may be mediated by upstream traits in a causal network. Indirect effects may therefore  
82 constitute a proportion of perceived pleiotropy, and these concepts apply to sets of heritable  
83 traits, organized as networks, are common in biological systems. An example from dairy cattle  
84 production systems, described by Gianola and Sorensen [10] and Rosa, et al. [11], is that higher  
85 milk yield increases the risk of a particular disease, such as mastitis, while the prevalence of the  
86 disease may negatively affect milk yield. In humans, obesity is a key factor influencing insulin  
87 resistance, which subsequently causes type 2 diabetes. A list of causal networks across human  
88 diseases and candidate genes is described in Kumar and Agrawal [12] and Schadt [13].

89 Although MTM-GWAS is a valuable approach, it only captures correlations or associations  
90 among traits and does not provide information about causal relationships. Knowledge of the  
91 causal structures underlying complex traits is essential, as correlation does not imply causation.  
92 For example, a correlation between two traits, T1 and T2, could be attributed to a direct effect of  
93 T1 on T2, T2 on T1, or to additional variables that jointly influence both traits [11]. Likewise, if  
94 we know a "causal" SNP is linked to a QTL, we can imagine three possible scenarios: 1) causal

95 ( $SNP \rightarrow T1 \rightarrow T2$ ), 2) reactive ( $SNP \rightarrow T2 \rightarrow T1$ ), or 3) independent ( $T1 \leftarrow SNP \rightarrow T2$ ).  
96 Scenarios (1) and (2), do not causes pleiotropy but produce association.  
97 A SEM methodology has the ability to handle complex genotype-phenotype maps in GWAS  
98 placing an emphasis on causal networks [14]. Therefore, SEM-based GWAS (SEM-GWAS) may  
99 provide a better understanding of biological mechanisms and of relationships among a set of  
100 traits than MTM-GWAS. SEM can potentially decompose the total SNP effect on a trait into  
101 direct and indirect (i.e., mediated) contributions. However, SEM-derived GWAS has not been  
102 discussed or applied fully in quantitative genetic studies yet. Our objective was to illustrate the  
103 potential utility of SEM-GWAS by using three production traits in broiler chickens genotyped for  
104 a battery of SNP as a case example.

## 105 **Methods**

### 106 **Data set**

107 The analysis included records for 1,351 broiler chickens provided by Aviagen Ltd. (Newbridge,  
108 Scotland) for three phenotypic traits: body weight (BW), ultrasound of breast muscle (BM) at 35  
109 days of age, and hen-house egg production (HHP), defined as the total number of eggs laid  
110 between weeks 28 and 54 per bird. The sample consisted of 274 full-sib families, 326 sires, and  
111 592 dams. More details regarding population and family structure were provided by Momen et al.  
112 [15]. A pre-correction procedure was performed on the phenotypes to account for systematic  
113 effects such as sex, hatch week, pen, and contemporary group for BW, BM, and HHP.

114 Each bird was genotyped for 580,954 SNP markers with a 600k Affymetrix SNP [16] chip  
115 (Affymetrix, Inc., Santa Clara, CA, USA). The Beagle software [17] was used to impute missing  
116 SNP genotypes, and quality control was performed using PLINK version 1.9 [18]. After

117 removing markers that did not fulfill the criteria of minor allele frequencies  $< 1\%$ , call  
118 rate  $> 95\%$ , and Hardy–Weinberg equilibrium (Chi-square test p-value threshold was  $10^{-6}$ ),  
119 354,364 autosomal SNP markers were included in the analysis.

## 120 **Multiple-trait model for GWAS**

121 MTM-GWAS is a single-trait GWAS model extended to multi-dimensional responses. When  
122 only considering additive effects of SNPs, the phenotype of a quantitative trait using the single-  
123 trait model can be described as:

$$124 \quad y_i = \sum_{q=1}^k x_{iq}\beta_q + w_{ij}s_j + e_i \quad (1)$$

125 where  $y_i$  is the phenotypic trait of individual  $i$ ,  $x_{iq}$  is the incidence value for the  $i$ th phenotype in  
126 the  $q$ th level of systematic environmental effect,  $\beta_q$  is fixed effect of the  $q$ th systemic  
127 environmental effect on the trait,  $w_j = (w_1, \dots, w_p)$  is the number of A alleles (i.e.,  $w_j \in \{0,1,2\}$ )  
128 in the genotype of SNP marker  $j$ , and  $s_j$  is the allele substitution effect for SNP marker  $j$ . Strong  
129 LD between markers and QTLs coupled with an adequate marker density increases the chance of  
130 detecting marker and phenotype associations. Hypothesis testing is typically used to evaluate the  
131 strength of the evidence of a putative association. Typically, a  $t$ -test is applied to obtain p-values,  
132 and the statistic is  $T_{ij} = \hat{s}_j/se(\hat{s}_j)$ , where  $\hat{s}$  is the point estimate of the  $j$ th SNP effect and  $se(\hat{s}_j)$   
133 is its standard error.

134 The single locus model described above is naïve for a complex trait because the data typically  
135 contain hidden population structure and individuals have varying degrees of genetic similarity  
136 [19, 20]. Therefore, accounting for covariance structure induced by genetic similarity is expected

137 to produce better inferences [21]. Ignoring effects that reveal genetic relatedness inflates the  
138 residual terms, compromises the ability to detect association. A random effect  $g_i$  including a  
139 covariance matrix reflecting pairwise similarities between additive genetic effects of individuals  
140 can be included to control population stratification. The similarity metrics can be derived from  
141 pedigree information or from whole-genome marker genotypes. This model extended for analysis  
142 of  $t$  traits is given by:

143

$$144 \quad y_{il} = \sum_{q=1}^k x_{iq}\beta_{ql} + w_{ij}s_{jl} + g_{il} + e_{il} \quad (2)$$

145 for  $i = 1, 2, \dots, n$ ,  $l = 1, 2, \dots, t$ . In this extension,  $y_{il}$  is the phenotypic value of the  $l$ th trait for the  
146  $i$ th subject,  $\beta_{qj}$  is the systematic effect of the  $q$ th environmental factor  $x_{iq}$  on the  $l$ th trait,  $s_{jl}$  is  
147 the additive effect of the  $j$ th marker on the  $l$ th trait,  $w_{ij}$  is as previously defined, and  $g_{il}$  and  $e_{il}$   
148 are the random polygenic effect and model residual assigned to individual  $i$  for trait  $l$ ,  
149 respectively. Random effects within a trait follow the multivariate normal distribution,

$$150 \quad \begin{bmatrix} g_l \\ e_l \end{bmatrix} \sim N \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K}\sigma_{g_l}^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{I}\sigma_{e_l}^2 \end{bmatrix} \right), \text{ where } \mathbf{K} \text{ is genetic relationship matrix, } \sigma_{g_l}^2 \text{ is the additive genetic}$$

151 variance of trait  $l$ ,  $\mathbf{I}$  is an identity matrix, and  $\sigma_{e_l}^2$  is the residual variance for trait  $l$ . The multiple-  
152 trait model accounts for the additive genetic ( $\rho_{ll'}$ ) and residual correlation ( $\lambda_{ll'}$ ) between a pair of  
153 traits  $l$  and  $l'$ .

154 The positive definite matrix  $\mathbf{K}$  may be a genomic relationship matrix ( $\mathbf{G}$ ) computed from marker  
155 data, or a pedigree-based matrix ( $\mathbf{A}$ ) computed from genealogical information. The  $\mathbf{A}$  matrix  
156 describes the expected additive similarity among individuals, while the  $\mathbf{G}$  measures the realized

157 fraction of alleles shared. Genomic relationship matrices can be derived in several ways [22-24].

158 Here, we used the form proposed by VANRADEN 2008:

$$159 \quad \mathbf{G} = \frac{\mathbf{M}\mathbf{M}'}{2 \sum p_j q_j} \quad (3)$$

160 where  $\mathbf{M}$  is an  $n \times p$  matrix of centered SNP genotypes and  $p_j$  and  $q_j = 1 - p_j$  are the allele  
161 frequencies at marker locus  $j$ . We evaluated both  $\mathbf{A}$  and  $\mathbf{G}$  in the present study.

162

### 163 **Structural equation model association analysis**

164 A SEM consists of two essential parts: a measurement model and a structural model. The  
165 measurement model depicts the connections between observable variables and their  
166 corresponding latent variables. The measurement model is also known as confirmatory factor  
167 analysis. The critical part of a SEM is the structural model, which can have three forms. The first  
168 consists of observable exogenous and endogenous variables. This model is a restricted version of  
169 a SEM known as path analysis [9]. The second form explains the relationship between exogenous  
170 and endogenous variables that are only latent. The third type is a model consisting of both  
171 manifest and latent variables.

172 SEM can be applied to GWAS as an alternative to MTM-GWAS to study how different causal  
173 paths mediate SNP effects on each trait. The following SEM model was considered:

$$174 \quad y_{il} = \mu_l + \sum_{m \in C} y_m \lambda_{lm} + w_{j(l)} S_{j(l)} + g_{il} + \varepsilon_{il} \quad (4)$$

175 where  $C$  is the set of phenotypic traits that directly affect the trait  $l$ ,  $\lambda_{lm}$  is a structural coefficient  
176 representing the effect of trait  $m$  on trait  $l$ , and  $g_l \sim N(0, \mathbf{K}\sigma_l^2)$  is the polygenic effect of the  $l$ th  
177 trait. The remaining terms are as presented earlier with one important difference: the SNP effects



178 are not interpreted as overall effects on trait  $l$  but instead represent direct effects on trait  $l$ .  
179 Additional indirect effects from the same SNP may be mediated by phenotypic traits in  $C$ . Each  
180 marker is entered into equation (4) one at a time, and its significance is tested. For a discussion of  
181 how SEM represents genetic signals on each trait through multiple causal paths, see Valente, et  
182 al. [25]. Despite the difference in interpretation, the distribution of the vector of polygenic effects  
183 is assumed to be the same as in the MTM-GWAS model. The same applies to residual terms  
184 within a trait. We also consider trait-specific residuals to be independent within an individual.  
185 This restriction is required to render structural coefficients likelihood-identifiable. In addition, the  
186 interpretation of inferences as having a causal meaning requires imposing the restriction that the  
187 residuals' joint distribution be interpreted as the causal sufficiency assumption [26]. In the  
188 present study, all exogenous and endogenous variables were observable, and there was no latent  
189 variable. In hence, causal structure was assumed between the endogenous variables BM, BW, and  
190 HHP.

191 We considered the following GWAS models, which their causal structures were recovered by the  
192 inductive causation (IC) algorithm [26]: (1) MTM-GWAS with pedigree-based kinship  $\mathbf{A}$  (MTM-  
193 A) or marker-based kinship  $\mathbf{G}$  (MTM-G), and (2) SEM-GWAS with  $\mathbf{A}$  (SEM-A) or  $\mathbf{G}$  (SEM-G).  
194 Although nuisance covariates such as environmental factors can be omitted in the graph, they  
195 may be incorporated into the models as exogenous variables. The SEM representation allowed us  
196 to decompose SNP effects into direct, indirect, and total effects.

197 A direct SNP effect is the path coefficient between a SNP as an exogenous variable and a  
198 dependent variable without any causal mediation by any other variable. The indirect effects of a  
199 SNP are those mediated by at least one other intervening endogenous variable. Indirect effects are  
200 calculated by multiplying path coefficients for each path linking the SNP to associated variable,

201 and then summing over all such paths [9]. The overall effect is the sum of all direct and indirect  
202 effects. By explicitly accounting for complex relationship structure among traits in such a way,  
203 SEM provides a better understanding of a genome-wide SNP analysis by allowing us to  
204 decompose effects into direct, indirect, and overall effects within a predefined casual framework.  
205 MTM-GWAS and SEM-GWAS were compared with the logarithm of the likelihood function  
206 ( $\log L$ ), Akaike's Information Criterion (AIC), and the Bayesian Information Criterion (BIC).  
207 The model providing the lowest values for these information criteria is considered to fit the data  
208 better [27]. MTM-GWAS and SEM-GWAS were fitted using the SNP Snappy strategy, which is  
209 implemented in the Wombat software program [28].

## 210 **Searching for a phenotypic causal network in a mixed model**

211 In the SEM-GWAS formulation described earlier, the structure of the underlying causal  
212 phenotypic network needs to be known. Because this is not so in practice, we used a causal  
213 inference algorithm to infer the structure. Residuals are assumed to be independent in all SEM  
214 analyses, so associations between observed traits are viewed as due to causal links between traits  
215 and by correlations among genetic values (i.e.,  $g_1$ ,  $g_2$ , and  $g_3$ ). Thus, to eliminate confounding  
216 problem when inferring the underlying network among traits, we used the approach of Valente, et  
217 al. [29] to search for acyclic causal structures through conditional independencies on the  
218 distribution of the phenotypes, given the genetic effects. A causal phenotypic network was  
219 inferred in two stages: 1) a MTM model [30] was employed to estimate covariance matrices of  
220 additive genetic effects and of residuals, and 2) the causal structure among phenotypes from the  
221 covariance matrix between traits, conditionally on additive genetic effects inferred by the IC  
222 algorithm. The residual (co)variance matrix was inferred using Bayesian MCMC [29, 31], with  
223 samples drawn from the posterior distribution. For each query testing statistical independence

224 between traits  $y_l$  and  $y_{l'}$ , the posterior distribution of the residual partial correlation  $\rho_{y_l, y_{l'}} | S$  was  
225 obtained, where  $S$  is a set of variable (traits) that are independent. Three highest posterior density  
226 (HPD) intervals of 0.75, 0.85, and 0.95 were used to make statistical decisions for SEM-GWAS.  
227 We thus considered SEM-A75 (HPD > 0.75), SEM-A85 (HPD > 0.85), SEM-A95 (HPD > 0.95),  
228 and SEM-G75 (HPD > 0.75). An HPD interval that does not contain zero declares  $y_l$  and  $y_{l'}$  to  
229 be conditionally dependent.

## 230 Results

231 Figure 1 shows phenotypic relationship structures recovered by the IC algorithm for the three  
232 different HPD intervals. Edges connecting two traits represent non-null partial correlations as  
233 indicated by HPD intervals. We compared the two MTM-GWAS and four SEM-GWAS by using  
234 the three chicken traits (BW, BM, and HHP). Only causal structures among the three traits are  
235 shown in Figure 1, because other parts were the same across the different SEM models. Fully  
236 recursive SEM-A75 and SEM-G75 revealed direct effects of BM on BW and HHP, and those of  
237 BW on HHP, as well as an indirect effect of BM on HHP. In addition, SEM-A85 detected a  
238 direct effect of BM on BW, the direct effect of BW on HHP, and the indirect effect of BM on  
239 HHP mediated by BW. Finally, SEM-A95 only identified a direct effect of BM on BW because  
240 of a statistically stringent HPD cutoff imposed.

241 Given the causal structures inferred from the IC algorithm, the following SEM was fitted:

$$242 \begin{cases} \mathbf{y}_1 = \mu + \mathbf{Z}_i \mathbf{g}_1 + W_{ij} S_j + \boldsymbol{\varepsilon}_i \\ \mathbf{y}_2 = \mu + \lambda_{21} \mathbf{y}_1 + \mathbf{Z}_i \mathbf{g}_2 + W_{ij} S_j + \boldsymbol{\varepsilon}_i \\ \mathbf{y}_3 = \mu + \lambda_{31} \mathbf{y}_1 + \lambda_{32} \mathbf{y}_2 + \mathbf{Z}_i \mathbf{g}_3 + W_{ij} S_j + \boldsymbol{\varepsilon}_i \end{cases} \quad (5)$$

243 Note that only a small number of the entries in the structural coefficient matrix ( $\lambda$  in equation 5)  
244 are nonzero due to sparsity. These nonzero entries specify the effect of one phenotype on other

245 phenotypes. The corresponding directed acyclic graph is shown in Figure 2 assuming the causal  
246 relationships among the three traits, where  $y_1$ ,  $y_2$ , and  $y_3$  represent BM, BW, and HHP,  
247 respectively;  $SNP_j$  is the genotype of the  $j$ th SNP;  $S_{jl}$  is the direct SNP effect on trait  $l$ ; and the  
248 remaining variables are as presented earlier. This diagram depicts a fully recursive structure in  
249 which all recursive relationships among the three phenotypic traits are shown. Arrows represent  
250 causal connections, whereas double-headed arrows between polygenic effects are correlations.

251 << **Figure 1 about here**>>

252 << **Figure 2 about here**>>

253 We examined the fit of each model implemented to assess how well it describes the data (Table  
254 1). Valente, et al. [25] showed that re-parametrization and reduction of a SEM mixed model yield  
255 the same joint probability distribution of observation as in MTM suggesting that expected  
256 likelihood of SEM and MTM should be the same. As expected, SEM-GWAS and MTM-GWAS  
257 showed very similar results (e.g., SEM-A75 vs. MTM-A and SEM-G75 vs. MTM-G). Among the  
258 models considered, the ones involving **G** exhibited a better fit. SEM-A85 and SEM-A95, sharing  
259 a subset of the SEM-A75 structure, presented almost identical AIC and BIC values.

260 << **Table 1 about here**>>

## 261 **Structural coefficients**

262 Table 2 presents the causal structural path coefficients for endogenous variables (BM, BW, and  
263 HHP). All models have positive effects for BM→BW, whereas the BM→HHP and BW→HHP  
264 relationships have negative path coefficients. The latter confirmed the fact that chicken breeding  
265 is divided into broiler and layer sections due to the negative genetic correlation between BW and  
266 HHP.

267

<<Table 2 about here>>

268 Also shown in Table 2 are the magnitudes of the SEM structural coefficient reflecting the  
 269 intensity of the causality. The positive coefficient  $\lambda_{21}$  quantifies the (direct) causal effect of BM  
 270 on BW. This suggests that a 1-unit increase in BM results in a  $\lambda_{21}$  -unit increases in BW.  
 271 Likewise, the negative causal effects  $\lambda_{31}$  and  $\lambda_{32}$  offer the same interpretation.

## 272 **Decomposition of SNP effect paths using a fully recursive model**

273 We can decompose the SNP effects into direct and indirect effects using Figure 2. The direct  
 274 effect of the SNP  $j$  on  $y_3$  (HHP) is given by  $d_{SNP_j \rightarrow y_3} : \hat{S}_{j(y_3)}$ , where  $d$  denotes the direct effect.  
 275 Note there are only one direct and many indirect paths. We find three indirect paths from  $SNP_j$  to  
 276  $y_3$  mediated by  $y_1$  and  $y_2$  (i.e., the nodes formed by other traits). The first indirect effect is  
 277  $ind_{(1)SNP_j \rightarrow y_3} : \lambda_{32}(\lambda_{21}\hat{S}_{j(y_1)})$  in the path mediated by  $y_1$  and  $y_2$ , where  $ind$  denotes the indirect  
 278 effect. The second indirect effect  $ind_{(2)SNP_j \rightarrow y_3} : \lambda_{32}\hat{S}_{j(y_2)}$ , is mediated by  $y_2$ . The last indirect  
 279 effect, is  $ind_{(3)SNP_j \rightarrow y_3} : \lambda_{31}\hat{S}_{j(y_1)}$ , mediated by  $y_1$ . Therefore, the overall effect is given by  
 280 summing all four paths,  $T_{SNP_j \rightarrow y_3} : \lambda_{32}(\lambda_{21}\hat{S}_{j(y_1)}) + \lambda_{32}\hat{S}_{j(y_2)} + \lambda_{31}\hat{S}_{j(y_1)} + \hat{S}_{j(y_3)}$ . The fully  
 281 recursive model of the overall SNP effect is then:

$$282 \quad \begin{cases} T_{\hat{S}_{j \rightarrow y_1} : \hat{S}_{j(y_1)}} \\ T_{\hat{S}_{j \rightarrow y_2} : \lambda_{21}(\hat{S}_{j(y_1)}) + \hat{S}_{j(y_2)}} \\ T_{\hat{S}_{j \rightarrow y_3} : \lambda_{32}[\lambda_{21}(\hat{S}_{j(y_1)}) + \hat{S}_{j(y_2)}] + \lambda_{31}(\hat{S}_{j(y_1)}) + \hat{S}_{j(y_3)}} \end{cases} \quad (6)$$

283 For  $y_1$  (BM), there is only one effect, so the overall effect is equal to the direct effect. For  $y_2$   
 284 (BW) and  $y_3$  (HHP), direct and indirect SNP effects are involved. There are two paths for  $y_2$ : one  
 285 indirect,  $ind_{S_j \rightarrow y_2} : \hat{S}_{j(y_1)} \rightarrow y_1 \rightarrow y_2$ , and one direct,  $d_{S_j \rightarrow y_2} : \hat{S}_{j(y_2)} \rightarrow y_2$ . Here, SNP effect is

286 direct and mediated thorough other phenotypes according to causal networks in SEM-GWAS  
287 (Figures 1 and 2). For instance, the overall SNP effect for  $y_3$  into four direct and indirect paths is  
288  $T_{\hat{S}_{j \rightarrow y_3}} : \lambda_{32}\lambda_{21}\hat{S}_{j(y_1)} + \lambda_{32}\hat{S}_{j(y_1)} + \lambda_{31}\hat{S}_{j(y_1)} + \hat{S}_{j(y_3)}$ .

289 The scatter plots in Figure 3 compare the estimated total effects for HHP ( $T_{\hat{S}_{j \rightarrow y_3}}$ ) obtained from  
290 SEM-GWAS and those from MTM-GWAS. We observed good agreement between SEM-GWAS  
291 and MTM-GWAS. The total SNP signals derived from SEM and MTM are the same but SEM  
292 provides biologically relevant additional information.

293 <<Figure 3 about here>>

294 Supplementary Figures S1-S4 present scatter plots of MTM-GWAS and SEM-GWAS signals  
295 (SEM-A75, SEM-G75, SEM-A85, and SEM-A95) for the  $BM \rightarrow BW$  path, which was a common  
296 path across all SEM-GWAS considered. These two traits have a genetic correlation of 0.5 (results  
297 not shown). We break the SEM causal link into direct, indirect, and overall effects based on the  
298 IC algorithm with HPD > 0.85, whereas MTM-GWAS capture an overall SNP effect on BW.  
299 Scatter plots of the overall effects from SEM-GWAS and the total effects from MTM-GWAS  
300 indicated almost perfect agreement (top left plots, Supplementary Figures S1–S4). We observed  
301 concomitance between estimated overall and direct effects (top right plots, Supplementary  
302 Figures S1–S4). In contrast, there was less agreement in the magnitude of the SNP effects when  
303 comparing overall vs. indirect effects (bottom left plots, Supplementary Figures S1–S4). There  
304 was no linear relationship between the indirect and direct SNP effects (bottom right plots,  
305 Supplementary Figures S1–S4). In short, genetic signals detected in SEM-GWAS were close to  
306 those of MTM-GWAS for overall effects because both models are based on a multivariate

307 approach with the same covariance matrix. In all SEM-GWAS, results showed that direct effect  
308 contributed to overall effects than the indirect effects.

### 309 **Manhattan plot of direct, indirect, and overall SNP effects**

310 Figure 4 depicts a Manhattan plot summarizing the magnitude of direct, indirect, and overall SNP  
311 effects. We plotted the decomposed SNP effects on BW along chromosomes to visualize  
312 estimated marker effects from SEM-GWAS. The indirect and direct effects provide a view of  
313 SNP effects from a perspective that is not available for the total effect of MTM-GWAS. For  
314 instance, many pleiotropic QTLs have positive direct effects on BW but negative effects on BM.  
315 There were two estimated SNP effects on chromosomes 1 and 2 that deserve particular attention.  
316 These two SNPs are highlighted with black circles and red ovals. The overall effect of the first  
317 SNP consisted of large indirect and small direct effects on BM, whereas the opposite pattern was  
318 observed for the second SNP, which showed large direct and small indirect effects. Although the  
319 overall effects of these SNPs were similar (top Manhattan plot, Figure 4), use of decomposition  
320 allowed us to find out that the trait of interest is affected in different manners: the second SNP  
321 effect acted directly on BW without any mediation by BM, whereas the first SNP reflected a  
322 large effect mediated by BM on BW. Collectively, new insight regarding the direction of SNP  
323 effects can be obtained using the SEM-GWAS methodology.

324 It should also be noted that the estimated additive SNP effects obtained from the four SEM-  
325 GWAS can be used for inferring pleiotropy. For instance, a pleiotropic QTL may have a large  
326 positive direct effect on BW but may exhibit a negative indirect effect coming from BM, which  
327 in turn reduces the total QTL effect on BW. Arguably, the methodology employed here would be  
328 most effective when the direct and indirect effects of a QTL are in opposite directions. If the

329 direct and indirect QTL effects are in the same direction, the power of SEM-GWAS may be the  
330 same as the overall power of MTM-GWAS.

331 <<Figure 4 about here>>

## 332 Discussion

333 It is becoming increasingly common to analyze a set of traits simultaneously in GWAS by  
334 leveraging genetic correlations between traits [32, 33]. In the present study, we illustrated the  
335 potential utility of a SEM-based GWAS approach, which has the potential advantage of  
336 embedding a pre-inferred causal structure across phenotypic traits [29]. SEM-GWAS accounts  
337 for the relationship of mediating variables that could be either dependent or independent with  
338 restriction on a residual covariance. This is a useful approach when multiple mediators interplay  
339 influencing the final outcomes [34, 35]. SEM-GWAS is achieved by first inferring the structure  
340 of network between phenotypic traits. For this purpose, we used a modified version of the IC  
341 algorithm described by Valente, et al. [29]. The IC algorithm was used to explore putative causal  
342 links among phenotypes obtained from a residual covariance matrix, in a model that accounted  
343 for systematic and genetic confounding factors such as polygenic additive effects. It then  
344 produced a posterior distribution of partial residual correlations between any possible pairs of  
345 variables. Three different causal path diagrams were inferred from HPD intervals of 0.75, 0.85,  
346 and 0.95. We observed that the number of identified paths decreased with an increase in the HPD  
347 interval value. Only a path connecting BM and BW was present in all HPD intervals considered.  
348 Moreover, we found that the partial residual correlation between BM and HHP was weaker than  
349 that between BM and BW. This may explain why the path between BM and HHP was not  
350 detected with HPD intervals larger than 0.75.



351 Estimated path coefficients reflect the strength of each causal link. For instance, a positive path  
352 coefficient from BM to BW suggests that a unit increase in BM directly results in an increase in  
353 BW. Our results showed that MTM-GWAS and SEM-GWAS were similar in terms of the  
354 goodness of fit as per the AIC and BIC criteria. This finding is in agreement with theoretical  
355 work of Valente, et al. [25] showing the equivalence between models. Thus, MTM-GWAS and  
356 SEM-GWAS produced the same marginal phenotypic distributions and goodness of fit values. A  
357 similar approach has been proposed by Li, et al. [14], Mi, et al. [36], and Wang and van Eeuwijk  
358 [37]. The main difference between our approach and theirs is that they used SEM in the context  
359 of standard QTL mapping, whereas our SEM-GWAS is developed for GWAS based on a linear  
360 mixed model.

361 The advantage of SEM-GWAS over MTM-GWAS is that the former decomposes SNP effects by  
362 tracing inferred causal networks. Our results showed that by partitioning SNP effect into direct,  
363 indirect, and total components, an alternative perspective of SNP effects can be obtained. As  
364 shown in Figure 4, direct and indirect effects may differ in magnitude and sign, acting in the  
365 same direction or even antagonistic manners. Note that the total SNP effects inferred from SEM-  
366 GWAS were the same as the estimated SNP effects from MT-GWAS (Figure 3). However,  
367 knowledge derived from the decomposition of SNP effects may be critical for animal and plant  
368 breeders to breaking unfavorable indirect QTL effects, or to obtain better SNP effects estimates  
369 than those from MTM-GWAS [e.g., 36].

## 370 **Conclusion**

371 SEM offers insights into how phenotypic traits relate to each other. We illustrated potential  
372 advantages of SEM-GWAS relative to the commonly used standard MTM-GWAS by using three

373 chicken traits as an example. SNP effects pertaining to SEM-GWAS have a different meaning  
374 than those in MTM-GWAS. Our results showed that SEM-GWAS enabled the identification of  
375 whether a SNP effect is acting directly or indirectly, i.e. mediated, on given trait. In contrast,  
376 MTM-GWAS only captures overall genetic effects on traits, which is equivalent to combining  
377 direct and indirect SNP effects from SEM-GWAS together. Thus, SEM-GWAS offers more  
378 information and provides an alternative view of putative causal network, enabling a better  
379 understanding the genetic quiddity of traits at the genomic level.

380

## 381 **Conflict of Interest**

382 The authors do not have any conflict of interest.

## 383 **Author's contributions**

384 MM carried out the study and wrote the first draft of the manuscript. GJMR and DG designed the  
385 experiment, supervised the study and critically contributed to the final version of manuscript. GM  
386 contributed to the interpretation of results, provided critical insights, and revised the manuscript.  
387 BDV and AAM participated in discussion and reviewed the manuscript. MA, AK and RMP  
388 contributed materials and revised the manuscript. All authors read and approved the final  
389 manuscript.

## 390 **Acknowledgment**

391 The first author wishes to acknowledge the Ministry of Science, Research and Technology of Iran  
392 for financially supporting his visit to the University of Wisconsin-Madison. Work was partially

393 supported by the Wisconsin Agriculture Experiment Station under hatch grant 142-PRJ63CV to  
394 DG.

## 395 **References**

- 396 1. Visscher Peter m, Brown Matthew a, Mccarthy Mark i, Yang J, Five Years of GWAS  
397 Discovery, *The American Journal of Human Genetics*. 2012; 90: 7-24.
- 398 2. Hayes B, Goddard M, Genome-wide association and genomic selection in animal breeding  
399 This article is one of a selection of papers from the conference “Exploiting Genome-wide  
400 Association in Oilseed Brassicas: a model for genetic improvement of major OECD crops for  
401 sustainable farming”, *Genome*. 2010; 53: 876-883.
- 402 3. Brachi B, Morris G P, Borevitz J O, Genome-wide association studies in plants: the missing  
403 heritability is in the field, *Genome Biology*. 2011; 12: 232.
- 404 4. Sikorska K, Lesaffre E, Groenen P F, Eilers P H, GWAS on your notebook: fast semi-parallel  
405 linear and logistic regression for genome-wide association studies, *BMC Bioinformatics*.  
406 2013; 14: 1-11.
- 407 5. Zhou X, Stephens M, Genome-wide efficient mixed-model analysis for association studies,  
408 *Nat Genet*. 2012; 44: 821-824.
- 409 6. Korte A, Vilhjálmsson B J, Segura V, Platt A, Long Q, Nordborg M, A mixed-model approach  
410 for genome-wide association studies of correlated traits in structured populations, *Nature*  
411 *genetics*. 2012; 44: 1066-1071.
- 412 7. O’reilly P F, Hoggart C J, Pomyen Y, Calboli F C, Elliott P, Jarvelin M-R, Coin L J,  
413 MultiPhen: joint model of multiple phenotypes can increase discovery in GWAS, *PloS one*.  
414 2012; 7: e34861.
- 415 8. Falconer D S, Mackay T F, Introduction to quantitative genetics (4th edn), *Trends in Genetics*.  
416 1996; 12: 280.
- 417 9. Wright S, Correlation and causation, *Journal of agricultural research*. 1921; 20: 557-585.
- 418 10. Gianola D, Sorensen D, Quantitative genetic models for describing simultaneous and  
419 recursive relationships between phenotypes, *Genetics*. 2004; 167: 1407-1424.

- 420 11. Rosa G J, Valente B D, De Los Campos G, Wu X-L, Gianola D, Silva M A, Inferring causal  
421 phenotype networks using structural equation models, *Genetics Selection Evolution*. 2011; 43:  
422 6.
- 423 12. Kumar S, Agrawal S, Disease-oriented Causal Networks, *Encyclopedia of Systems Biology*,  
424 Springer, 2013, pp. 593-594.
- 425 13. Schadt E E, Chapter 10 - Reconstructing Causal Network Models of Human Disease A2 -  
426 Lehner, Thomas, in: Miller B L, State M W (Eds.), *Genomics, Circuits, and Pathways in*  
427 *Clinical Neuropsychiatry*, Academic Press, San Diego, 2016, pp. 141-160.
- 428 14. Li R, Tsaih S-W, Shockley K, Stylianou I M, Wergedal J, Paigen B, Churchill G A,  
429 Structural model analysis of multiple quantitative traits, *PLoS genetics*. 2006; 2: e114.
- 430 15. Momen M, Mehrgardi A A, Sheikhy A, Esmailizadeh A, Fozi M A, Kranis A, Valente B D,  
431 Rosa G J M, Gianola D, A predictive assessment of genetic correlations between traits in  
432 chickens using markers, *Genetics Selection Evolution*. 2017; 49: 16.
- 433 16. Kranis A, Gheyas A A, Boschiero C, Turner F, Yu L, Smith S, Talbot R, Pirani A, Brew F,  
434 Kaiser P, Development of a high density 600K SNP genotyping array for chicken, *BMC*  
435 *genomics*. 2013; 14: 59.
- 436 17. Browning S R, Browning B L, Rapid and accurate haplotype phasing and missing-data  
437 inference for whole-genome association studies by use of localized haplotype clustering, *Am J*  
438 *Hum Genet*. 2007; 81.
- 439 18. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira M A, Bender D, Maller J, Sklar P, De  
440 Bakker P I, Daly M J, PLINK: a tool set for whole-genome association and population-based  
441 linkage analyses, *The American Journal of Human Genetics*. 2007; 81: 559-575.
- 442 19. Listgarten J, Lippert C, Kadie C M, Davidson R I, Eskin E, Heckerman D, Improved linear  
443 mixed models for genome-wide association studies, *Nat Meth*. 2012; 9: 525-526.
- 444 20. Gianola D, Fariello M I, Naya H, Schön C-C, Genome-Wide Association Studies with a  
445 Genomic Relationship Matrix: A Case Study with Wheat and Arabidopsis, *G3: Genes,*  
446 *Genomes, Genetics*. 2016; 6: 3241-3256.
- 447 21. Galesloot T E, Van Steen K, Kiemeny L A, Janss L L, Vermeulen S H, A comparison of  
448 multivariate genome-wide association methods, *PloS one*. 2014; 9: e95923.
- 449 22. Vanraden P, Efficient methods to compute genomic predictions, *J Dairy Sci*. 2008; 91: 4414 -  
450 4423.

- 451 23. Yang J, Benyamin B, Mcevoy B, Gordon S, Henders A, Nyholt D, Common SNPs explain a  
452 large proportion of the heritability for human height, *Nat Genet.* 2010; 42: 565 - 569.
- 453 24. Forni S, Aguilar I, Misztal I, Different genomic relationship matrices for single-step analysis  
454 using phenotypic, pedigree and genomic information, *Genet Sel Evol.* 2011; 43: 1.
- 455 25. Valente B D, Rosa G J, Gianola D, Wu X-L, Weigel K, Is structural equation modeling  
456 advantageous for the genetic improvement of multiple traits?, *Genetics.* 2013; 194: 561-572.
- 457 26. Pearl J, Causal inference in statistics: An overview, *Statistics surveys.* 2009; 3: 96-146.
- 458 27. Akaike H, Information theory and an extension of the maximum likelihood principle,  
459 *Selected Papers of Hirotugu Akaike*, Springer, 1998, pp. 199-213.
- 460 28. Meyer K, Tier B, Churchill G A, “SNP Snappy”: A Strategy for Fast Genome-Wide  
461 Association Studies Fitting a Full Mixed Model, *Genetics.* 2012; 190: 275-277.
- 462 29. Valente B D, Rosa G J, De Los Campos G, Gianola D, Silva M A, Searching for recursive  
463 causal structures in multivariate quantitative genetics mixed models, *Genetics.* 2010; 185:  
464 633-644.
- 465 30. Henderson C, Quaas R, Multiple trait evaluation using relatives' records, *Journal of Animal*  
466 *Science.* 1976; 43: 1188-1197.
- 467 31. Wu X L, Heringstad B, Gianola D, Bayesian structural equation models for inferring  
468 relationships between phenotypes: a review of methodology, identifiability, and applications,  
469 *Journal of Animal Breeding and Genetics.* 2010; 127: 3-15.
- 470 32. Gao H, Zhang T, Wu Y, Wu Y, Jiang L, Zhan J, Li J, Yang R, Multiple-trait genome-wide  
471 association study based on principal component analysis for residual covariance matrix,  
472 *Heredity.* 2014; 113: 526-532.
- 473 33. Wu B, Pankow J S, Genome-wide association test of multiple continuous traits using imputed  
474 SNPs, *Statistics and its interface.* 2017; 10: 379.
- 475 34. Bellavia A, Valeri L, Decomposition of the total effect in the presence of multiple mediators  
476 and interactions, *American journal of epidemiology.* 2017.
- 477 35. Barfield R, Shen J, Just A C, Vokonas P S, Schwartz J, Baccarelli A A, Vanderweele T J, Lin  
478 X, Testing for the indirect effect under the null for genome-wide mediation analyses, *Genetic*  
479 *epidemiology.* 2017; 41: 824-833.

- 480 36. Mi X, Eskridge K, Wang D, Baenziger P S, Campbell B T, Gill K S, Dweikat I, Bovaird J,  
481 Regression-based multi-trait QTL mapping using a structural equation model, *Statistical*  
482 *applications in genetics and molecular biology*. 2010; 9.
- 483 37. Wang H, Van Eeuwijk F A, A New Method to Infer Causal Phenotype Networks Using QTL  
484 and Phenotypic Information, *PLOS ONE*. 2014; 9: e103997.

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

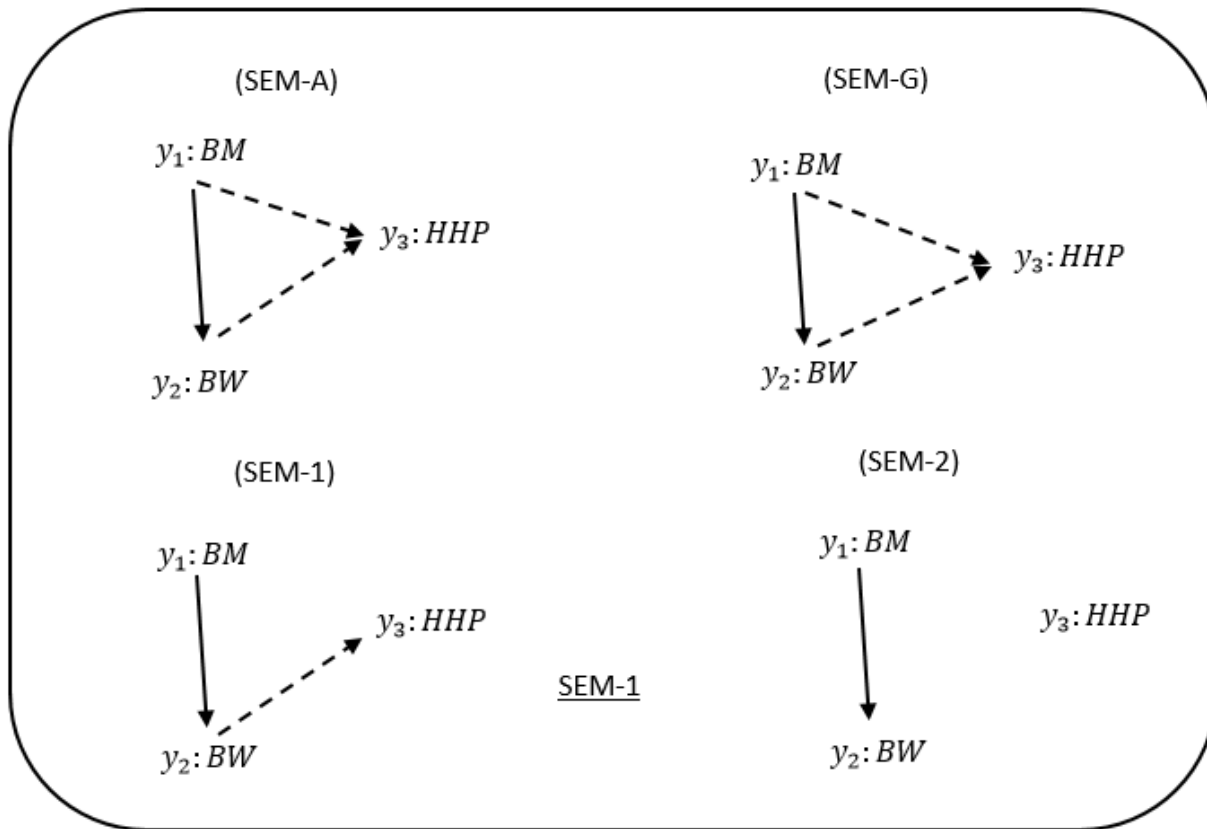
500

501

502

503

504 **Figures**



505

506 **Figure 1 Causal graphs inferred using the IC algorithm among three traits: breast meat**  
507 **(BM), body weight (BW) and hen-house production (HHP) in the chicken data.** SEM-A75  
508 and SEM-G75 were the inferred fully recursive causal structures with HPD > 0.75 and corrected  
509 for genetic confounder using **A** (pedigree-based) and **G** (marker-based) matrices. SEM-A85 and  
510 SEM-A95 were obtained with HPD > 0.85 and HPD > 0.95, respectively, corrected with **A**.  
511 Arrows indicate direction of causal relationships. Dashed lines indicate negative coefficients, and  
512 the continuous arrows indicate positive coefficients.

513

514

515

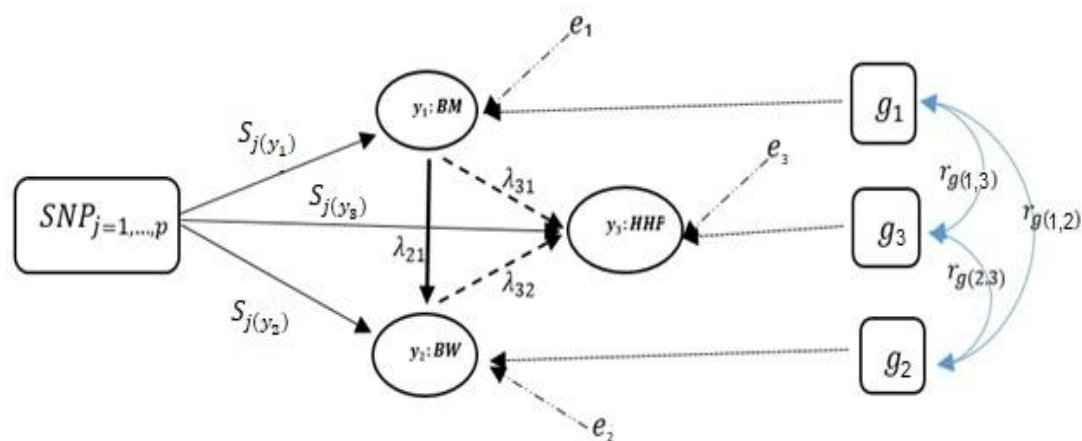
516

517

518

519

520



521

522 **Figure 2 A diagram for causal path analysis of SNP effects in a fully recursive structural**  
 523 **equation model for three traits,  $p$  exogenous independent SNP variables, and three**  
 524 **correlated polygenic effects.** Arrows indicate the direction of causal effects and dashed lines  
 525 represent associations among the three phenotypes. Genetic correlation between traits ( $r_g$ ),  
 526 polygenic effects ( $g_l$ ), environmental effect on trait  $l$  ( $e_l$ ), effects of  $j$  th SNP on  $l$  th trait ( $S_{j(y_l)}$ ),  
 527 and recursive effect of phenotype  $l'$  on phenotype  $l$  ( $\lambda_{l,l'}$ ).

528

529

530

531

532

533

534

535

536

537

538

539

540

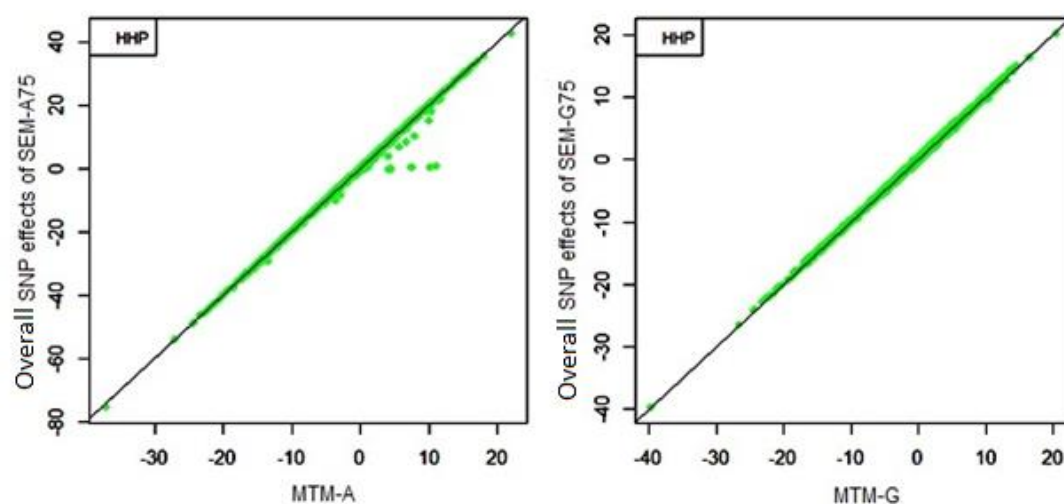


541

542

543

544



545

546 **Figure 3 Comparison of multiple trait (MTM) and fully recursive overall SNP effects**  
547 **obtained with A (pedigree-based) and G (marker-based) from structural equation modeling**  
548 **(SEM)-based GWAS. Overall effects in SEM are the sum of all direct and indirect effects. HHP:**  
549 **hen-house egg production.**

550

551

552

553

554

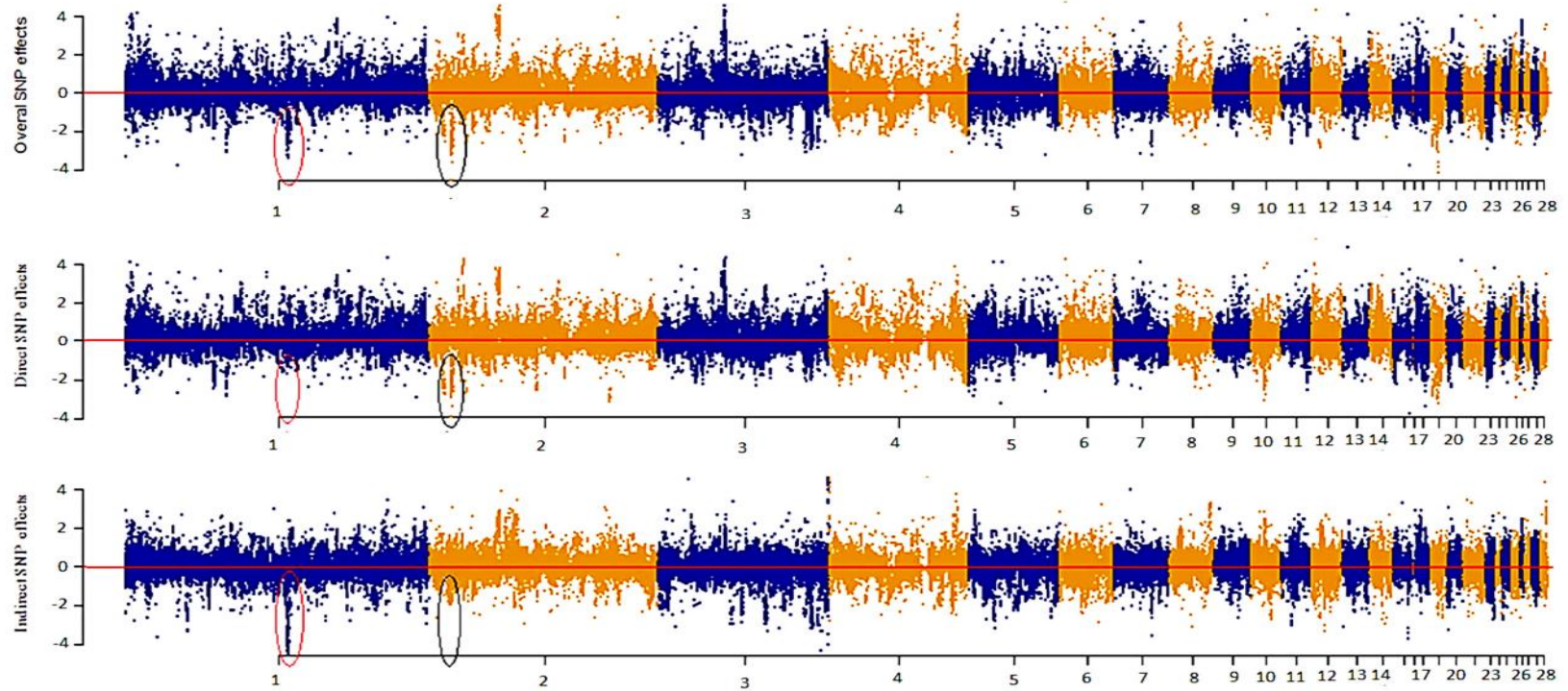
555

556

557

558

559



560

561 **Figure 4** Manhattan plot showing overall, direct, and indirect SNP effects using a full recursive model based on G matrix for  
 562 **body weight (BW).**

563

564 **Tables**

565

**Table 1 Model comparison criteria: logarithm of the restricted maximum likelihood function (log L), Akaike's information criteria (AIC), Schwarz Bayesian information criteria (BIC) to evaluate model fit for two MTM and four SEM models.**

Model	Maximum log L	-1/2 AIC	-1/2 BIC
MTM-A	-7093.480	-7105.48	-7142.436
SEM-A75	-7098.370	-7110.415	-7147.321
SEM-A85	-7095.188	-7107.188	-7144.143
SEM-A95	-7097.517	-7109.517	-7146.470
MTM-G	-6529.270	-6541.276	-6578.232
SEM-G75	-6537.391	-6549.391	-6586.34

A: pedigree-based relationship matrix, G: VanRaden's marker-based relationship matrix

566

567

568

569

570

571

572

573

574

575

**Table 2 Estimates of three causal structural coefficients ( $\lambda$ ) derived from four different structural models. BM: breast meat. BW: body weight. HHP: hen-house production. SEM-75: HPD > 0.75. SEM-G75: HPD > 0.75. SEM-A85: HPD > 0.85. SEM-A95: HPD > 0.95.**

Path	Structural Models			
	SEM-75	SEM-G75	SEM-A85	SEM-A95
$\lambda_{BM \rightarrow BW}(\lambda_{21})$	2.13	2.19	2.14	2.14
$\lambda_{BM \rightarrow HHP}(\lambda_{31})$	-0.17	-0.28	***	***
$\lambda_{BW \rightarrow HHP}(\lambda_{32})$	-0.27	-0.096	-0.31	***

576

577

578

579

580

581

582

583