

1 **Inferring novel lncRNA associated with Ventricular septal defect by**
2 **DNA methylation interaction network**

3 Min Zhang¹, Yue Gu¹, Mu Su², Shumei Zhang¹, Chuangeng Chen¹, Wenhua Lv¹, Yan
4 Zhang^{1,2*}

5 ¹College of Bioinformatics Science and Technology, Harbin Medical University, Harbin, 150081,
6 China.

7 ²School of Life Science and Technology, Harbin Institute of Technology, Harbin 150001, China

8 *Corresponding author

9 E-mail: tyozhang@ems.hrbmu.edu.cn)

10

11

WITHDRAWN
see manuscript DOI for details

12 **Abstract**

13 Ventricular septal defect (VSD) is one of the most common types of congenital heart
14 disease. To find more and more molecular alteration is conducive to explore the
15 mechanism and biomarker in VSD. Herein we devised a predictive strategy to
16 uncover novel lncRNA of VSD integrating DNA methylation, gene expression and
17 lncRNA expression of early embryo and VSD by profiles from GEO database. In
18 totally, 175 lncRNAs, 7290 genes and 3002 DNA methylation genes were obtained by
19 logistic regression analysis associated with embryonic development. Moreover, 7304
20 DMGs were significant differential methylated by Wilcoxon rank test and Student's
21 test in VSD. We constructed the lncRNA-mRNA co-expression network in embryo
22 (LMCNe). Then, a reconstructed co-expression weighted network (RCWN) was built
23 integrated LMCNe and the DNA methylation associated network (DMAN) based
24 on the correlation of the DNA methylation level and protein interaction network
25 between embryonic development and VSD. We extracted top 10 lncRNAs with higher
26 score performing DRaWR from the weight network, which as potential VSD related
27 lncRNAs. Six lncRNAs showed a high level of expression in the heart tissue recorded
28 in the NONOCOND database. Furthermore, associated lncRNA genes DCAF8L1,
29 NIT1, SH2D7 and DOCK9-AS2 in validated samples showed a prominently
30 association with VSD. These outcomes provide a reference for lncRNA involved in
31 VSD initialization and a new insight for studies of VSD-associated lncRNAs.

32 **Key words:** Ventricular septal defect; DNA methylation; lncRNA; network

33

34 **Author Summary**

35 Ventricular septal defect (VSD) is one of the most common types of congenital
36 heart disease and has a high mortality rate in infants. Many molecular markers have
37 proved effective as biomarker in VSD like DNA methylation and lncRNA. lncRNA is
38 a type of non-coding RNA which has important effect in regulation gene expression
39 and disease occurrence. VSD is an embryonic stage developmental disease. Therefore
40 we hypothesized that lncRNA which was associated with DNA methylation and
41 mRNA in early embryonic development may also affect the occurrence of VSD. So in
42 this work, from the perspective of embryonic development, we devised a predictive
43 strategy to uncovering novel lncRNA of VSD. In our result, four lncRNA associated
44 genes were found differential expressed in VSD and normal samples by qPCR
45 validation. The identification of lncRNA associated with ventricular septal defect is
46 beneficial to further study the mechanism of VSD from the molecular level and also
47 provides a good molecular marker for clinical therapeutic and diagnosis. At the same
48 time, it also provides a new insight for the researches of lncRNA associated with
49 VSD.

50 **1. Introduction**

51 The heart is the first developing organ in embryonic development which can
52 provide oxygen and nutrients for the baby and many genes and biological processes
53 are involved in this procedure. A slight disruption in the process of embryonic
54 development of the heart could cause damage to the embryo, leading to congenital
55 heart disease(1, 2). Congenital heart disease is one of the most common defects in
56 infants accompanied with a higher mortality rate(3). Some reports show that 4-14 of
57 every 1,000 newborns have congenital heart disease, and congenital heart disease
58 accounts for 40% of infant defects(4-7). Congenital heart disease could be caused by a
59 variety of factors, including environmental, genetic, epigenetic effects(8). VSD is
60 considered to be one of the most common types of congenital heart disease,
61 accounting for about 30% of congenital heart disease(9). VSD could exist in isolation
62 and could as part of double outlet right ventricle and other cardiac anomalies(10).

63 DNA methylation is an important epigenetic modification that plays an important
64 role in cell differentiation, regulation of gene expression, disease occurrence and
65 development and participates in biological processes of progressive organs. DNA
66 methylation often occurs on the fifth carbon atom of cytosine, formed by DNA
67 methyltransferase. Hypermethylation in gene promoter can contribute tumor
68 suppressor silencing leading to disease. DNA methylation is reprogrammed during the
69 zygote phase of early embryos and rewrite into the genome during cell division and
70 differentiation (11, 12). Epigenetic marks have a sustained effect on individuals and
71 progeny and DNA methylation pattern is stable in heredity(13, 14).Furthermore, most

72 genes in cell were arranged by DNA methylation in the period of early embryo
73 development and abnormal DNA methylation could lead to defect in births(15).
74 Folate is a B-vitamin providing methyl group for methylation reactions in embryonic
75 cell including DNA methylation and could reduce the risks of congenital heart disease
76 (16). Recently, many studies have been investigated in understanding the role of DNA
77 methylation in congenital heart disease. For example, higher abundance of
78 methylation biomark S-adenosylhomocysteine (SAH) in blood increased the risk of
79 CHD (17). The hypermethylation of RAB43 and KIAA0310 leads to ER-te-Golgi
80 dysfunction and eventually contributes to the occurrence of VSD and abnormal
81 embryonic heart development(18). What is more, the hypermethylation of NOX5
82 frequently occurs in the pathogenesis of VSD(19).

83 In recent years, the human genome is divided into coding and non-coding
84 regions(20). With the advanced genome sequencing technology, many mRNA-like
85 functional transcripts have been found in noncoding regions which are non-coding
86 RNAs (21-23). Non-coding RNA have vital influence in regulating gene expression
87 instead of being a protein template(24). LncRNA is a type of non-coding RNA whose
88 length is greater than 200bp having divided into five categories: sense, antisense,
89 bidirectional, intronic, and intergenic(25, 26). Accumulating evidence suggested that
90 lncRNA has a wide range of biological functions involved in the regulation of gene
91 expression networks, such as chromosome printing, cell growth and differentiation
92 and tumorigenesis(27-29). However, lncRNA associated with VSD is rarely reported.

93 Herein, by analyzing the DNA methylation, gene expression and lncRNA profile of

94 early embryo and DNA methylation data of VSD. We constructed an integrated
95 RCWN modified by DNA methylation correlation based on protein interaction
96 network. This network not only illustrated intricate relationship between coding genes
97 and non-coding lncRNA in embryonic development, but also revealed the DNA
98 methylation-mediated correlation between embryonic development and VSD. Then
99 we selected the top 10 lncRNAs with higher score as VSD associated lncRNAs after
100 performing DRaWR algorithm in the RCWN. Six lncRNA were found expressed in
101 heart tissue included in NONCODE database. Furthermore, we investigated the
102 expression of lncRNA gene were apparently associated with VSD by comparing the
103 normal and VSD cardiac tissue. These lncRNAs were mainly involved in the function
104 of regulation of growth heart force and so on, which would be a reference for further
105 analysis.

106 **2. Results**

107 **2.1 Identification of differential methylation genes(DMGs) in VSD**

108 VSD is a disease that caused by abnormal events occurring in embryonic
109 development. DNA methylation plays an important role in embryonic development,
110 and may affect VSD. Therefore, we devised an integrated prediction network to
111 predict lncRNA associated with VSD. We integrated MBD-Seq DNA methylation
112 data containing four normal right atrium and eight VSD samples and DNA
113 methylation data lncRNA expression profile and gene expression profile of early
114 embryos. The LMCNe was built in embryonic firstly. Then, a reconstructed
115 co-expression weighted network (RCWN) was built integrated LMCNe and DMAN

116 based on the correlation of the DNA methylation level and protein-protein interaction
117 between embryonic development and VSD. By performing Discriminative Random
118 Walk with Restarts (DRaWR), selecting the highest score of 10 lncRNA as a
119 candidate lncRNA Figure 1.

120 For DNA methylation data of VSD, 18,322,642 CpG sites were obtained by data
121 filtering and deletion (see methods). 741,027 CpG sites were considered as
122 differential methylated sites by Wilcoxon rank test. Subsequently, these differential
123 DNA methylation CpG sites were matched to the promoter region of 9,274 refseq
124 genes and the 7,292 differential DNA methylation genes (DMG) were screened using
125 the Student test with setting threshold $p < 0.05$ and districted DNA methylation level
126 greater than 0.2 between normal and VSD figure 2 A. The overall trend of the
127 differential DNA methylation genes were that DNA methylation levels of most genes
128 were upregulated, with only a small fraction of the downregulation in figure 2 B.
129 Moreover, the DNA methylation level of DMGs in VSDs were significantly higher
130 than that in the normal samples figure 2 C. The functional enrichment of these DMGs
131 showed that they were associated with signal transduction, organ development, Rap1
132 signaling pathways and other basic functions which has great regulatory for bodies
133 (figure 2 D).

134 **2.2 Identification of embryonic development-related lncRNA, mRNA and DNA** 135 **methylation gene**

136 We treated six stages of early embryonic DNA methylation data including oocytes,
137 Zygotes, 2-cell-stage, 4-cell-stages, 8-cell-stage, morulae stage. In the following

138 analysis, methylated promoter of 14,304 refseq gene were reserved. 8,123 lncRNA
139 expression and 10,176 mRNA expression were available from the document with
140 RPKM (see Method). Subsequently, we fitted an Ordered Rogers regression model to
141 screen for lncRNAs and genes associated with embryonic development. In brief,
142 3,002 associated DNA methylation genes (AMGs) were regarded as associated with
143 embryonic development and most of the genes were in a low DNA methylation state.
144 With further embryo developmental stages, DNA methylation level of AMGs
145 exhibited a tendency of demethylation which was consistent with previous reports
146 figure 3 A(30, 31). These AMGs were mainly involved in cell differentiation, cell
147 adhesion, DNA template regulation, protein transcription, RNA transcription
148 pathways, and other pathways and biological process that had played an important
149 role in embryonic development and organ growth (figure 3 B). However, for the 7304
150 gene expression associated with embryonic development, their expression showed a
151 slight increase in figure 3 C. The function of these genes were mostly enriched in
152 transcription, DNA-templated, cell division and other biological processes and
153 pathways (figure 3D). Similarly, the expression of 175 embryonic associated lncRNAs
154 were identified and we studied protein coding genes within upstream and downstream
155 100kbp of lncRNAs and found that they were primarily focused on transcription,
156 DNA-templated, cellular protein localization and signaling pathways of regulating
157 pluripotency of stem cells (figure 3 E and F). Additionally, as displayed in figure S 1,
158 overlapped embryonic development-related DNA methylation genes and gene
159 expression and found 566 overlapping genes indicating that DNA methylation levels

160 of the promoters of these genes have a negative regulatory effect on gene expression,
161 in other words, higher DNA methylation level are obviously associated with
162 down-regulated gene expression.

163 **2.3 Construction of LMCNe**

164 In order to illustrate the relationship between lncRNA and protein coding genes in
165 embryonic development. We established a LMCNe on the basis of embryo-associated
166 genes and lncRNAs in the figure 3 with Pearson correlation coefficient more than 0.8
167 between mRNA and lncRNA. LMCNe contained 3,344,593 linkages and 7,304
168 mRNAs connected with 175 lncRNAs which indicated that many protein-coding were
169 regulated by the same lncRNA and also governed by differential lncRNA at the same
170 time figure S 2. The distribution of node degree evenly appeared to pow-law
171 distribution suggested that LMCNe was a proper network with biology implication
172 figure S2

173 **2.4 Construction of DMAN**

174 In recent years, protein-protein interaction network (PPI) has dominated the
175 network and revealed more complex biological mechanisms and biological processes.
176 It is typically used in network analysis. Here, we constructed the DMAN of early
177 embryos and VSD by two steps in the background of PPI. First, we acquired
178 significant correlation DNA methylation gene pairs in embryo and VSD methylation
179 genes set separately (see Materials and Methods). In the second step, we used the PPI
180 network as the background network, and the two significant co-methylation gene set
181 were to be incorporated into a DMAN through commons nodes and edges.

182 At the first stage, we found 1,464,756 co-methylated DNA methylation gene pairs
183 (DMP) in embryo and 1,515,563 DMPs in VSD set. We observed that some DMPs
184 were attended in both embryo set and VSD set suggested that some DMPs
185 simultaneously affects VSD occurrence and embryo growth. In order to enhance the
186 correlation between two DNA methylation pair sets, we introduced the protein-protein
187 interaction network as the background network, and mapped the DMPs to the PPI
188 network respectively. Only genes that existed in the PPI network were retained and
189 merged into a DMAN by common nodes and co-methylation relationship. DMAN
190 consists of 7,925 nodes and 94,945 edges with three categories of interactions: DMPs
191 occurred in embryo or in VSD's or in both figure 4 A. Moreover, the node-degree
192 distribution of DMAN obeyed the power-law distribution revealed that the DMAN
193 was a scale-free network that only a few nodes were with high connectivity and most
194 nodes are low in connectivity figure 4 B. The interaction in DMAN were established
195 by the absolute value of Pearson correlation coefficient greater than 0.8 and to
196 confirm intensity of DMAN, we randomized the DNA methylation level of the nodes
197 in the network for 100 times and calculated the Pearson correlation coefficient. It was
198 found that the Pearson correlation coefficient of DMAN was much higher than the
199 average stochastic suggested that DMPs were not random events figure 4 C.

200 **2.5 Construction of RCWN**

201 The RCWN was rebuilt by integrating LMCNe and the DMAN which was as the
202 weight of edges attribute and included 3,439,296 links and 14,147 nodes figure 4 D.
203 100 times permutation of Pearson coefficient for each linkage were performed by

204 disorganizing the expression value of lncRNAs and mRNAs. Obviously, the
205 correlation of co-expression networks was greater than the random average
206 correlation coefficient figure 4 F. Obviously, the degree distribution of nodes in
207 this network was also subject to power-law distribution indicating that some genes
208 may regulated by both lncRNA and methylation figure 4 E.

209 **2.6 DRaWR analysis for VSD related lncRNA**

210 With respect to DNA methylation plays an important role in embryonic
211 development and abnormal DNA methylation programming could increase the risk of
212 VSD in neonates (32). Therefore, we hypothesized that VSD was an embryonic stage
213 developmental disease. lncRNAs which were associated with DNA methylation and
214 mRNA in early embryonic development may also affect the occurrence of VSD.
215 Thence we used the DRaWR algorithm to analyze the RCWN and predicted the
216 lncRNA with VSD figure 5 A (see the materials and methods). In our study, we
217 selected 120 genes that simultaneously affect embryonic development and VSD
218 initialization as 'query set' figure 5 B. DNA methylation level of these 120 genes
219 presented a demethylation pattern figure 5 C during developmental stages. We
220 retained the top 10 lncRNAs with higher scores as VSD potential lncRNA
221 lnc-POLE4-8, lnc-DCAF8L1-1, lnc-RAB1A-1, lnc-NIT1-1, lnc-SH2D7-1,
222 lnc-POTEB-15, LINC01467, lnc-ECI2-5, DOCK9-AS2, LINC01622 as showed in
223 figure 6 A . At the same time, we also studied 10 genes, MTO1, DNAJB12,
224 MRFAP1, MAD2L2, MTG1, LPIN2, FBXL22, PEBP1 and POR, which were
225 connected to the 10 lncRNAs in the network with the highest rank figure 6 B. MTO1

226 mutations were associated with hypertrophic cardiomyopathy and PEBP1 could affect
227 heart failure and lactic acidosis and cause respiratory chain deficiency in humans and
228 yeast (38,39). FBXL22 is a cardiac-enriched sarcomere protein and is essential for
229 maintenance of normal contractile function in vivo (40). LPIN2 is likewise a cardiac
230 gene (41). We also analyzed protein-coding genes located 100 kbp near lncRNA, such
231 as APOA2, a gene that encodes high-density lipoprotein binding protein, who was
232 associated with cardiovascular risk(37). The mutation of NDUFS2 and USF1 gene
233 variants was often occurred in cardiomyopathy and has been a biomark in
234 cardiovascular disease(38-40) and its promoter hypermethylation is associated with
235 congenital heart disease which were the neighbors of lnc-NIT1-1(41).

236 **2.7 Validation of lncRNAs in VSD**

237 To explore whether the discovered lncRNAs were indeed associated with VSD,
238 through the NONCODE database, we found that six lncRNAs were highly expressed
239 in the heart tissue and were validated by quantitative PCR in 13 heart samples (see
240 method) figure 7A. Four lncRNAs associated genes exhibited high expression level
241 than control group especially DCAD8L1 and SH2D7 measured by students test.
242 (figure 7) and their function was mainly focused on the regulation of birth, cardiac
243 systolic, cell cycle, and other important pathways related to cardiac development
244 figure 6 C. Some lncRNAs like lnc-POLE4-8 were not only involved in the regulation
245 of growth but also in the contraction of the heart.

246 **3 Discussion**

247 VSD is one of the most common defect in newborns, affected by multiple factors

248 such as genetic factors and epigenetic factors. DNA methylation as the most apparent
249 modification plays an important role in the embryonic development, cell growth,
250 differentiation gradually becoming a diagnostic marker. Hypermethylation of gene
251 promoter region regulated gene expression.

252 In recent years, the crucial role of DNA methylation in congenital heart disease has
253 been gradually found, for example, NOX5 promoter hypermethylation occurred more
254 frequently in VSD, abnormal hypermethylation of gene promoter leads to gene
255 silencing. In this study, we first obtained DNA methylation gene that was
256 significantly different with normal samples by compared with the VSD sample, and
257 also obtained DNA methylation genes associated with embryonic development. By
258 Pearson correlation analysis, the significant DMPs in the VSD were excavated, and
259 the DMPs in embryonic development were also investigated. In order to further
260 increase the relationship of DMPs, we introduced functional correlation
261 protein-protein interaction network. After constructing a DNA methylation associated
262 network, Protein-protein interaction network had increased the robustness of DMAN.
263 Looking for VSD associated related lncRNAs, we used DRaWR algorithm for
264 network analysis and 120 genes were as seed nodes, the methylation of these genes
265 were not only related with embryonic development, but also VSD. We considered top
266 10 lncRNAs as candidate lncRNA. Most of the lncRNAs are lincRAN and adjacent to
267 the heart disease gene. For example, lnc-POLE4-8 is adjacent to HK2 which is
268 Hexokinases phosphorylate glucose to produce glucose-6-phosphat was found to be
269 associated with ventricular function in mice (33).

270 However there also exist some limitations, firstly it did not have enough data of
271 VSD to use to support our study, we downloaded MBD-seq methylation data of VSD
272 containing 8 cases and 4 controls covered only part of CpG sites. Secondly, due to the
273 use of lncRNA and gene expression data from the early embryonic stage, some values
274 were somewhat lower and there may exits false positive rate.

275 **4 Method and Materials**

276 **4.1 Samples and data process**

277 VSD methylation data was download from GEO database
278 (<https://www.ncbi.nlm.nih.gov/geo/>) with the accession number GSE62629 including
279 four normal right atrium and eight VSD samples. Methylation dataset of embryo with
280 six stages was also download from GEO and accession number is GSE49828. As for
281 DNA methylation data only the CpG sites with read coverages more than five times
282 were taken into account. When we quantified the average DAN methylation level of a
283 gene, we calculated the average DNA methylation level of total number of all CpG
284 site located within 1500bp upstream and 500bp downstream of transcription start site
285 as gene promoter DNA methylation level. Gene expression and lncRNA expression
286 profile of embryo with RPKM were reported by Yan (34) . All of genes and lncRNA
287 with an RPKM>0.0 was considered in the following analysis .The hg19 Refseq genes
288 and lncRNA annotation were download from the UCSC Genome Browser
289 (<http://genome.ucsc.edu>) and NONCODE database separately(35). The validate
290 samples were extracted from The first Affiliated Hospital of Harbin Medical
291 University and the study was approved by The Ethics Committee of First Affiliated

292 Hospital of Harbin Medical University.

293 **4.2 Total RNA extraction and purification**

294 The total tissue for RNA extraction varied between 10 and 15 mg from each animal.
295 In order to ensure unbiased analysis of tissue response, total RNA was isolated from
296 randomly sectioned (weighting 10- 15 mg) heart. RNA was isolated from 13 to 8
297 individual samples from each of the treatment and control groups using TRIzol
298 reagent (inveitrogen, Carlsbad, CA, USA) and purified using RNeasy Plus Mini kits
299 (Qiagen, Mississauga, ON, Canada) according to the manufacturer's instruction. Total
300 RNA concentration was measured using a NanoDrop 2000 spectrophotometer
301 (Thermo Fisher Scientific Inc., Wilmington, DE, USA), and RNA quality and
302 integrity were assessed using an Agilent 2100 Bioanalyzer (Agilent Technologies,
303 Inc. Mississauga, ON, Canada) according to the manufacturer's instruction. All
304 samples showed RNA integrity numbers of 7 and above, indicating high quality RNA,
305 and were used to conduct qPCR experiments.

306 **4.3 Real-time quantitative PCR**

307 Total RNA were used to perform reverse transcription using PrimeScript™ RT
308 reagent Kit (TaKaRa, Tokyo, Japan). Real-time polymerase chain reaction
309 amplification of cDNA was performed with SYBR Premix Ex Taq™ II (TaKaRa,
310 Tokyo, Japan). The sequences of the primer sets were listed as follows:

311 DCAF8L1 forward, 5'- TCAAAGTAGTTTCCAGGCAATCTTGTG-3' and reverse,
312 5'- CTGGGATTTGGATGGAGTCGTT-3'; RAB1A forward, 5'-
313 CAGTCACGAGGCTCTCCGAA-3' and reverse, 5'-

314 AGTTCCTCATAGGCTATTGGACTGA-3'; NIT1 forward, 5'-
315 GGAAGACCAGATTATGCGAA-3' and reverse, 5'-
316 CATTATTCCCCAGTTGTCCA-3'; DOCK9-AS2 forward, 5'-
317 AGAGAGCAAAACACCTCCGTT-3' and reverse 5'-
318 AGGAGGAAATGGGTTAGTGTGT-3'; SH2D7 forward, 5'-
319 TCCCCTGGGTTACTCCATTC-3' and reverse 5'-
320 CTCATTACTCTCCCCAAAACCTG-3'; β -actin forward,
321 5'-GCTCCAAGCAGATGCAGCA-3' and reverse
322 5'-CCGGATGTGAGGCAGCAG-3' Quantification of gene expression was
323 determined by comparative quantity, using β -actin gene expression as inner control.

324 4.4 Statistical analysis

325 Due to the multiple sorted stages of embryo, ordered logistic Regression model was
326 used to estimate correlation between development stage of embryo and DNA
327 methylation performing with R packages MASS plor function. Y denoted early
328 stages of embryonic development, X was delegated the DNA methylation level of
329 each gene or lncRNA expression or gene expression level.

$$330 \quad \text{logit}(Y) = \beta_0 + \beta_1 X \quad (1)$$

331 For DNA methylation data in VSD, significant differential methylated CpG sites were
332 identified using Wilcoxon rank test in R software. Bedtools was used to calculate total
333 CpG sites in Refseq gene promoter. To analyze differences in DNA methylation
334 between normal tissue and VSD, p value was calculated by Student's t test.
335 Co-expression relationships among the lncRNA, gene expression and DNA

336 methylation were estimated by Pearson correlation test and correlation coefficient
337 greater than 0.8 were used in following analysis.

338 **4.5 Go and pathways analysis**

339 Enrichment analysis was analyzed in DAVID database. For associated-lncRNA in
340 embryo, protein-coding genes between 100kbp upstream and 100kbp downstream for
341 each lncRNA were used for enrichment analysis. Go terms and pathways with the
342 threshold $p\text{-value} < 0.05$ were selected as significant enrichment results and visualized
343 top 10 terms of biological process, cellular component, molecular function and
344 pathways of KEGG in Cytoscape3.5 enrichment plug-in.

345 **4.6 DNA Methylation associated network (DMAN)**

346 In this study, we introduced a human protein-protein interaction network as
347 background network for the further analysis. The network integrates the protein
348 interaction relationship for six database, including HPRD ,DIP ,MINT,BIND, IntAct
349 (36) and STRING database. A DNA methylation associated network was extracted by
350 mapping DMG pairs into the background network. DMGs were used as the node of
351 the network and interaction between DMGs serves as network edges formed by
352 Pearson correlation analysis with threshold $p\text{ value} < 0.05$ and coefficient > 0.8 . Then
353 1000 random permutations were performed for DMG pairs calculating Pearson
354 correlation coefficient.

355 **4.7 RCWN**

356 In order to identify lncRNA associated with VSD, we merged LMCNe and DMAN as
357 RCWN. The degree of node refers to the number of neighbor nodes and the weight of

358 an edge was defined as following principles:

359 M_E refers to co-methylation gene pair dataset in embryo, G_E is defined as

360 co-expression gene pair in embryo. M_V donated significant co-methylation gene pairs

361 in VSD. GL_E means significant gene-lncRNA co-expression pair in embryo.

$$362 \text{ weight(edge)} = \begin{cases} 1 & \text{if gene pair occur in } M_E \text{ or } G_E \text{ or } M_V \text{ or } GL_E \\ 2 & \text{if gene pair occur in } M_E \text{ and } G_E \\ 2 & \text{if gene pair occur in } M_E \text{ and } M_V \\ 2 & \text{if gene pair occur in } G_E \text{ and } M_V \\ 3 & \text{if gene pair occur in } M_E, G_E \text{ and } M_V \end{cases}$$

363 4.8 Discriminative random walk with Restarts (DRaWR)

364 Discriminative random walk with Restarts is an algorithm to rank genes for their

365 relevance to a given ‘query’ gene set based on “guilt by association”. The method

366 performs two times of Random walking with restart (RWR). In the first stage of

367 DRaWR, a proper subnetwork was extracted from a large original network including

368 only seed-relevant node. Then ranking all genes according to the relevance to query

369 gene set in the second stage of DraWR in the light of subnetwork(37).

370 In this study, gene set $Q = \{q_1, q_2, \dots, q_n\}$ refers to the seed gene set and $G = \{g_1, g_2, \dots,$

371 $g_m\}$ is the total node of the RCWN. D denotes an adjacency matrix of network and

372 D_{ij} is the link between g_i and g_j .

$$373 D = \begin{bmatrix} D_{11} & D_{12} & \dots & D_{1m} \\ D_{21} & D_{22} & \dots & D_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ D_{n1} & D_{n2} & \dots & D_{nm} \end{bmatrix} \quad (2)$$

374 Then normalized the adjacency matrix N :

$$375 N_{ij} = \frac{D_{ij}}{\sum_j D_{ij}} \quad (3)$$

376 Next normalized each columns of then N to create a transition probability matrix T

377 Using following equation:

378
$$T_{ij} = \frac{N_{ij}}{\sum_{ij} N_{ij}}. \quad (4)$$

379 T_{ij} is the probability the walker according to an edge transition from node i to node
380 j . Moreover follows the equation to perform the RWR :

381
$$W_{t+1} = (1 - r)TW_t + r\alpha \quad (5)$$

382 where r is the restart parameter, the walker reset the walk by moving directly to
383 other genes in query gene set. W_t refers to the probability distribution of all node in
384 the network since t steps of RWR algorithm. Given a cutoff $|W_{t+1} - W_t| < 0.05$
385 ranking all node in W by probability distribution. Notably, in the first stage, W_Q
386 and W_G are calculated respectively, and the nodes with the largest difference
387 between W_Q and W_G are reserved to form a subnet. In the second stage of DRaRW ,
388 all nodes were ranked by repeating RWR on the subnet until $|W'_{t+1} - W'_t| < 0.05$.
389 DRaWR is a R packages and performed in R3.3.1 software.

390 **Acknowledgments**

391 This work was supported by National Natural Science Foundation of China [grant
392 number 81573021].

393 **References**

- 394 1. Satou Y, Satoh N. Gene regulatory networks for the development and evolution of the
395 chordate heart. *Genes & development*. 2006;20(19):2634-8.
- 396 2. McQuillen PS, Miller SP. Congenital heart disease and brain development. *Annals of the*
397 *New York Academy of Sciences*. 2010;1184:68-86.
- 398 3. Hoffman JI, Kaplan S. The incidence of congenital heart disease. *Journal of the*
399 *American College of Cardiology*. 2002;39(12):1890-900.

- 400 4. Mendieta-Alcantara GG, Santiago-Alcantara E, Mendieta-Zeron H, Dorantes-Pina R,
401 Ortiz de Zarate-Alarcon G, Otero-Ojeda GA. [Incidence of congenital heart disease and
402 factors associated with mortality in children born in two Hospitals in the State of Mexico].
403 *Gaceta medica de Mexico*. 2013;149(6):617-23.
- 404 5. Acharya G, Sitras V, Maltau JM, Dahl LB, Kaaresen PI, Hanssen TA, et al. Major
405 congenital heart disease in Northern Norway: shortcomings of pre- and postnatal diagnosis.
406 *Acta obstetricia et gynecologica Scandinavica*. 2004;83(12):1124-9.
- 407 6. Ransom J, Srivastava D. The genetics of cardiac birth defects. *Seminars in cell &*
408 *developmental biology*. 2007;18(1):132-9.
- 409 7. Jenkins KJ, Correa A, Feinstein JA, Botto L, Britt AE, Daniels SR, et al. Noninherited risk
410 factors and congenital cardiovascular defects: current knowledge: a scientific statement from
411 the American Heart Association Council on Cardiovascular Disease in the Young: endorsed
412 by the American Academy of Pediatrics. *Circulation*. 2007;115(23):2995-3014.
- 413 8. Eskedal LT, Hagemo PS, Eskild A, Frosli KF, Seiler S, Thaulow E. A population-based
414 study relevant to seasonal variations in causes of death in children undergoing surgery for
415 congenital cardiac malformations. *Cardiology in the young*. 2007;17(4):423-31.
- 416 9. Ko JM. Genetic Syndromes associated with Congenital Heart Disease. *Korean*
417 *circulation journal*. 2015;45(5):357-61.
- 418 10. de Hemptinne Q, De Becker B, Unger P. Ventricular septal defect: an unusual cause of
419 paradoxical low-gradient aortic stenosis. *European heart journal cardiovascular Imaging*.
420 2017;18(5):609.
- 421 11. Altorok N, Tsou PS, Coit P, Khanna D, Sawalha AH. Genome-wide DNA methylation

- 422 analysis in dermal fibroblasts from patients with diffuse and limited systemic sclerosis reveals
423 common and subset-specific DNA methylation aberrancies. *Annals of the rheumatic diseases*.
424 2015;74(8):1612-20.
- 425 12. Gilsbach R, Preissl S, Gruning BA, Schnick T, Burger L, Benes V, et al. Dynamic DNA
426 methylation orchestrates cardiomyocyte development, maturation and disease. *Nature*
427 *communications*. 2014;5:5288.
- 428 13. Hanson M, Godfrey KM, Lillycrop KA, Burdge GC, Gluckman PD. Developmental
429 plasticity and developmental origins of non-communicable disease: theoretical considerations
430 and epigenetic mechanisms. *Progress in biophysics and molecular biology*.
431 2011;106(1):272-80.
- 432 14. Reik W. Stability and flexibility of epigenetic gene regulation in mammalian development.
433 *Nature*. 2007;447(7143):425-32.
- 434 15. Kim KC, Friso S, Choi SW. DNA methylation, an epigenetic mechanism connecting folate
435 to healthy embryonic development and aging. *The Journal of nutritional biochemistry*.
436 2009;20(12):917-26.
- 437 16. Hernandez-Diaz S, Werler MM, Walker AM, Mitchell AA. Folic acid antagonists during
438 pregnancy and the risk of birth defects. *The New England journal of medicine*.
439 2000;343(22):1608-14.
- 440 17. Obermann-Borst SA, van Driel LM, Helbing WA, de Jonge R, Wildhagen MF, Steegers
441 EA, et al. Congenital heart defects and biomarkers of methylation in children: a case-control
442 study. *European journal of clinical investigation*. 2011;41(2):143-50.
- 443 18. Zhu C, Yu ZB, Chen XH, Pan Y, Dong XY, Qian LM, et al. Screening for differential

- 444 methylation status in fetal myocardial tissue samples with ventricular septal defects by
445 promoter methylation microarrays. *Molecular medicine reports*. 2011;4(1):137-43.
- 446 19. Zhu C, Yu ZB, Chen XH, Ji CB, Qian LM, Han SP. DNA hypermethylation of the NOX5
447 gene in fetal ventricular septal defect. *Experimental and therapeutic medicine*.
448 2011;2(5):1011-5.
- 449 20. Morris KV, Mattick JS. The rise of regulatory RNA. *Nature reviews Genetics*.
450 2014;15(6):423-37.
- 451 21. Guttman M, Garber M, Levin JZ, Donaghey J, Robinson J, Adiconis X, et al. Ab initio
452 reconstruction of cell type-specific transcriptomes in mouse reveals the conserved
453 multi-exonic structure of lincRNAs. *Nature biotechnology*. 2010;28(5):503-10.
- 454 22. Cabili MN, Trapnell C, Goff L, Koziol M, Tazon-Vega B, Regev A, et al. Integrative
455 annotation of human large intergenic noncoding RNAs reveals global properties and specific
456 subclasses. *Genes & development*. 2011;25(18):1915-27.
- 457 23. Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, et al. The GENCODE
458 v7 catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and
459 expression. *Genome research*. 2012;22(9):1775-89.
- 460 24. Qiu MT, Hu JW, Yin R, Xu L. Long noncoding RNA: an emerging paradigm of cancer
461 research. *Tumour biology : the journal of the International Society for Oncodevelopmental*
462 *Biology and Medicine*. 2013;34(2):613-20.
- 463 25. Mattick JS. The genetic signatures of noncoding RNAs. *PLoS genetics*.
464 2009;5(4):e1000459.
- 465 26. Ponting CP, Oliver PL, Reik W. Evolution and functions of long noncoding RNAs. *Cell*.

- 466 2009;136(4):629-41.
- 467 27. Vance KW, Sansom SN, Lee S, Chalei V, Kong L, Cooper SE, et al. The long non-coding
468 RNA Paupar regulates the expression of both local and distal genes. *The EMBO journal*.
469 2014;33(4):296-311.
- 470 28. Clark MB, Mattick JS. Long noncoding RNAs in cell biology. *Seminars in cell &*
471 *developmental biology*. 2011;22(4):366-76.
- 472 29. Li J, Xuan Z, Liu C. Long non-coding RNAs and complex human diseases. *International*
473 *journal of molecular sciences*. 2013;14(9):18790-808.
- 474 30. Smith ZD, Meissner A. DNA methylation: roles in mammalian development. *Nature*
475 *reviews Genetics*. 2013;14(3):204-20.
- 476 31. Lee MT, Bonneau AR, Giraldez AJ. Zygotic genome activation during the
477 maternal-to-zygotic transition. *Annual review of cell and developmental biology*.
478 2014;30:581-613.
- 479 32. Wijnands KP, Chen J, Liang L, Verbiest MM, Lin X, Helbing WA, et al. Genome-wide
480 methylation analysis identifies novel CpG loci for perimembranous ventricular septal defects
481 in human. *Epigenomics*. 2017;9(3):241-51.
- 482 33. Del Duca D, Tadevosyan A, Karbassi F, Akhavein F, Vaniotis G, Rodaros D, et al.
483 Hypoxia in early life is associated with lasting changes in left ventricular structure and function
484 at maturity in the rat. *International journal of cardiology*. 2012;156(2):165-73.
- 485 34. Yan L, Yang M, Guo H, Yang L, Wu J, Li R, et al. Single-cell RNA-Seq profiling of human
486 preimplantation embryos and embryonic stem cells. *Nature structural & molecular biology*.
487 2013;20(9):1131-9.

488 35. Bu D, Yu K, Sun S, Xie C, Skogerbo G, Miao R, et al. NONCODE v3.0: integrative
489 annotation of long noncoding RNAs. *Nucleic acids research*. 2012;40(Database
490 issue):D210-5.

491 36. Liu H, Su J, Li J, Liu H, Lv J, Li B, et al. Prioritizing cancer-related genes with aberrant
492 methylation based on a weighted protein-protein interaction network. *BMC systems biology*.
493 2011;5:158.

494 37. Blatti C, Sinha S. Characterizing gene sets using discriminative random walks with
495 restart on heterogeneous biological networks. *Bioinformatics*. 2016;32(14):2167-75.

496
497 Figure legends

498 Figure 1. An overview of predicting strategy of VSD-associated lncRNA.

499 Figure 2 Differential DAN methylation gene in VSD. (A) Volcano plot of P value as a
500 function of the differential DNA methylation level in normal tissue and VSD sample.
501 Red dots represent that DNA methylation level of genes are not significantly
502 different. Green dots represent that DNA methylation level of genes are significantly
503 different. (B) Heatmap of DNA methylation level for 7292 DMGs. (C) Distinction
504 DNA methylation level between normal tissue and VSD sample. Significance was
505 tested using a Wilcoxon test. (D) DAVID enrichment analysis for DMGs in VSD

506 Figure 3. An overview of DNA methylation lncRNA and gene expression to
507 embryonic development. (A) Heatmap and DAVID enrichment analysis for
508 embryonic methylation genes.(B) Heatmap and DAVID enrichment analysis for
509 embryonic gene expression. (C) Heatmap and DAVID enrichment analysis for

510 embryonic lncRNA.

511 Figure 4. RCWN. (A) DMAN. (B) The degree distribution of DMAN follows a
512 power-law distribution. (C) Comparison of Pearson correlation coefficient between
513 real network and random network. (D) DNA methylation-lncRNA-mRNA
514 co-expression weighted network. (E-F) The degree distribution of co-expression
515 network and comparison of Pearson correlation coefficient with permutation.

516 Figure 5 Prediction of VSD-associated lncRNA by DRaWR. (A) An overview of
517 DRaWR. (B) Overlap of gene expression and DNA methylation of embryo and VSD
518 as the 'query set'. (C) The DNA methylation level of 'query set' in six stages of
519 embryo VSD and normal tissue.

520 Figure 6 VSD-associated lncRNAs. (A) The genomic position of the top 10 lncRNAs
521 and genes, the green line means the interaction between gene and gene. The blue line
522 represent the interaction between gene and lncRNA and red ones mean
523 lncRNA-lncRNA co-expression. (B) The subnetwork of top 10 lncRNAs with highest
524 score genes extracted from global co-expression network. (C) The
525 function of lncRNA. Black blocks represent that the lncRNA has the function in
526 NONCODE database.

527 Figure 7 Validation of lncRNA related genes in VSD cardiac tissue. (A)
528 Electropherogram of validated samples. (B) Bar plot of lncRNA-related gene
529 expression in normal and VSD.

530 Figure S1 : Negative regulation of DNA methylation to gene expression.

531 (A) Venn diagram of detected embryonic DNA methylation gene and embryonic gene

532 expression. (B)The scatterplot model shows that the methylation negatively regulates

533 the gene.

534 Figure S2: lncRNA-mRNA co-expression network in embryo(LMCNe) (A) An

535 overview of (LMCNe) (B) The degree distribution of (LMCNe) follows a power-low

536 distribution

537

538

WITHDRAWN
see manuscript DOI for details

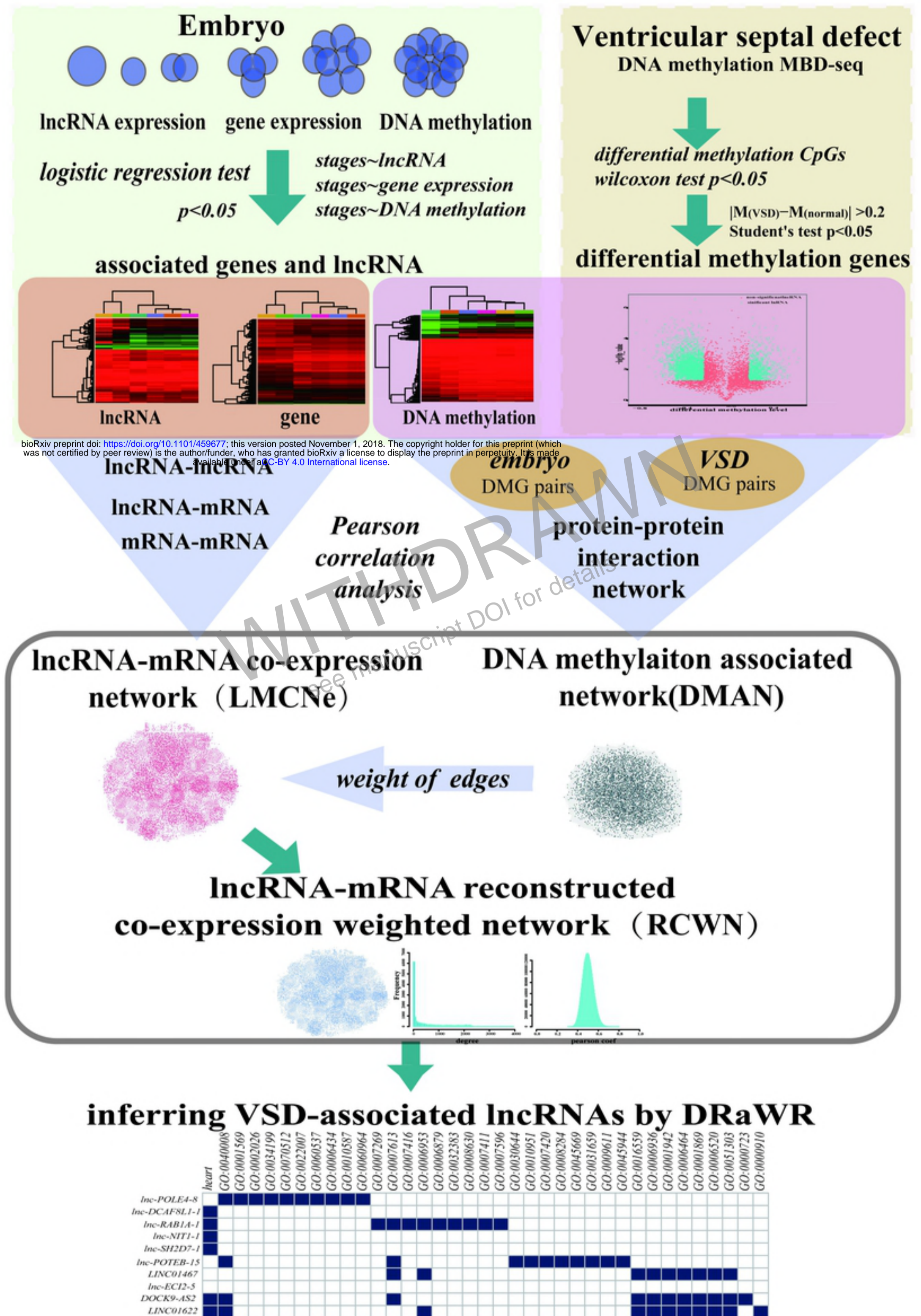
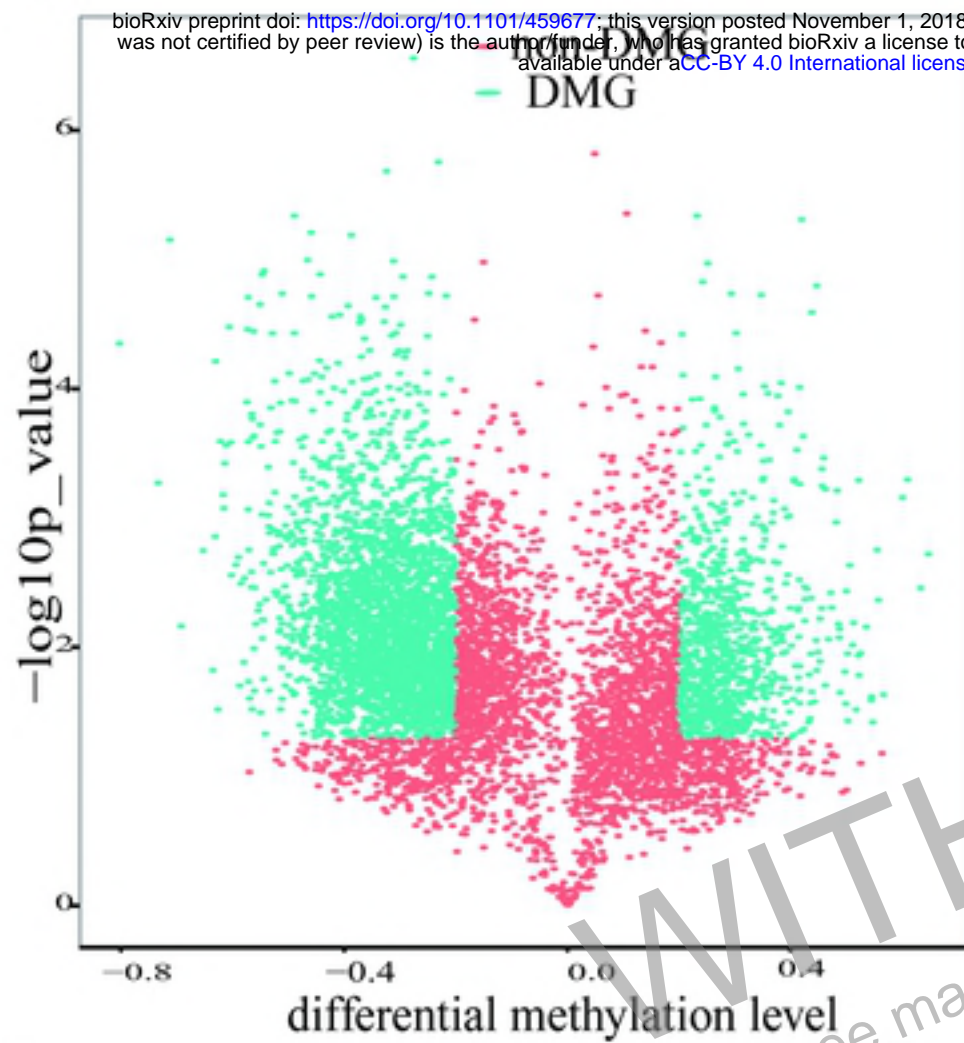
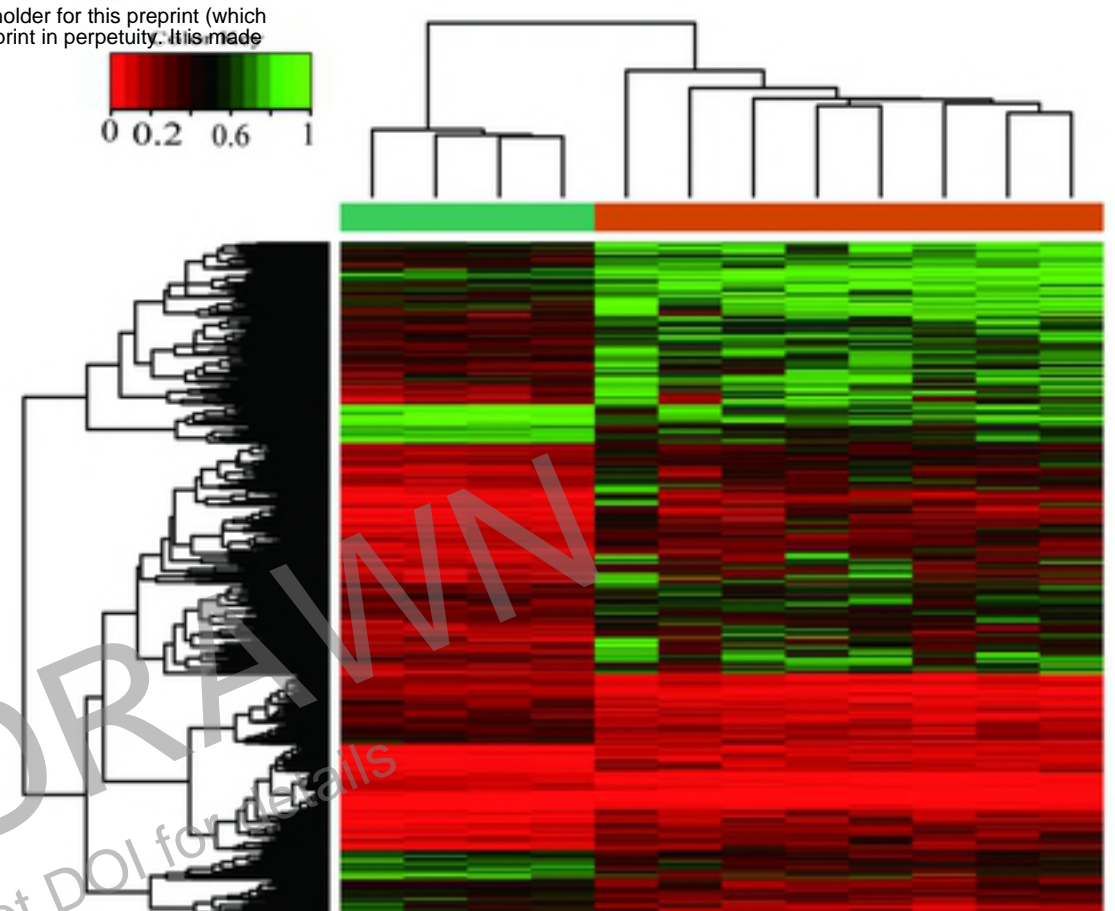


figure1

A



B

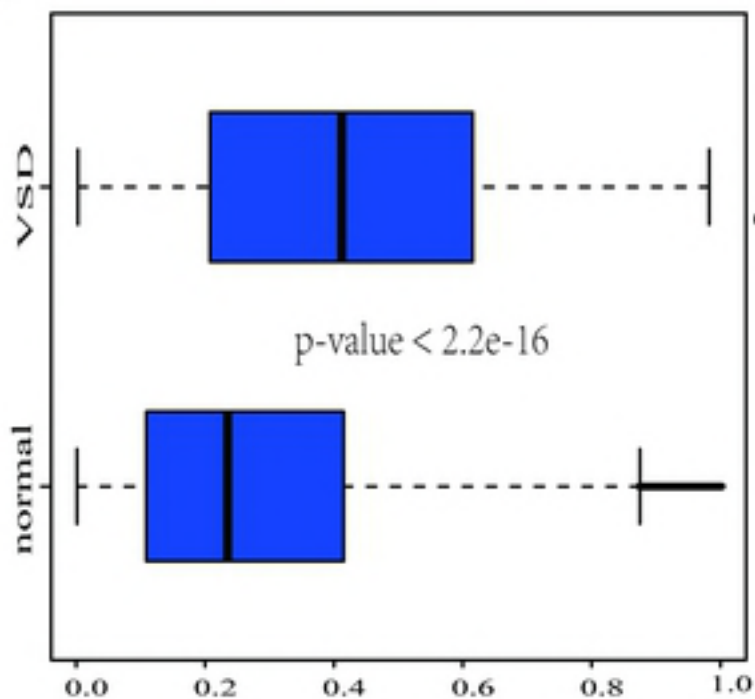


GO:0007165~signal transduction
 GO:0043547~positive regulation of GTPase activity
 GO:0045893~positive regulation of transcription, DNA-templated
 GO:0007275~multicellular organism development
 GO:0007155~cell adhesion
 GO:0008284~positive regulation of cell proliferation
 GO:0043066~negative regulation of apoptotic process
 GO:0035556~intracellular signal transduction
 GO:0006954~inflammatory response
 GO:0042493~response to drug
 hsa04080:Neuroactive ligand-receptor interaction
 hsa04060:Cytokine-cytokine receptor interaction
 hsa04015:Rap1 signaling pathway
 hsa04810:Regulation of actin cytoskeleton
 hsa04510:Focal adhesion
 hsa05164:Influenza A
 hsa04020:Calcium signaling pathway
 hsa04261:Adrenergic signaling in cardiomyocytes
 hsa04380:Osteoclast differentiation
 hsa05142:Chagas disease (American trypanosomiasis)

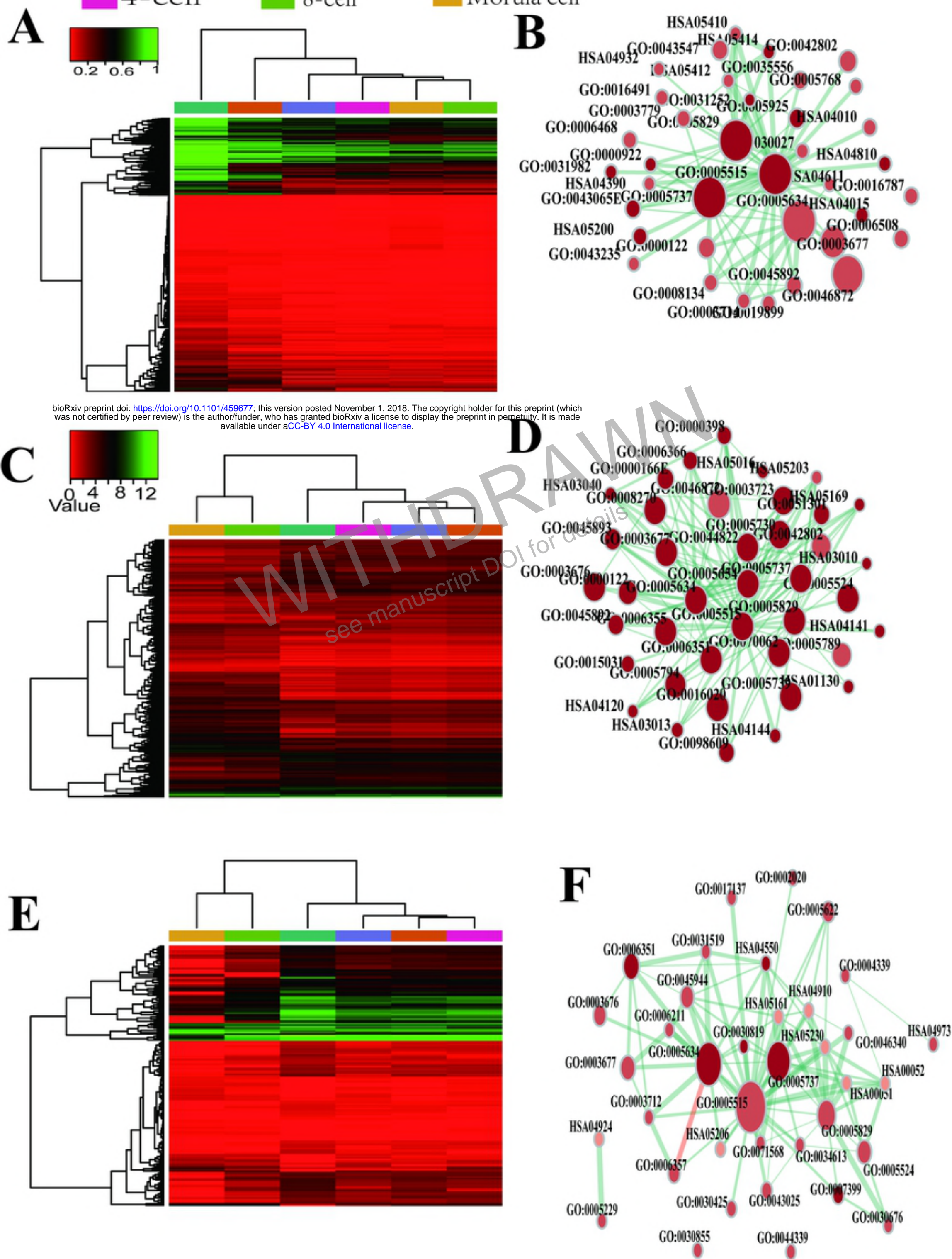
D



C



oocyte cell Zygote cell 2-cell
 4-cell 8-cell Morula cell



bioRxiv preprint doi: <https://doi.org/10.1101/459677>; this version posted November 1, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

figure 3

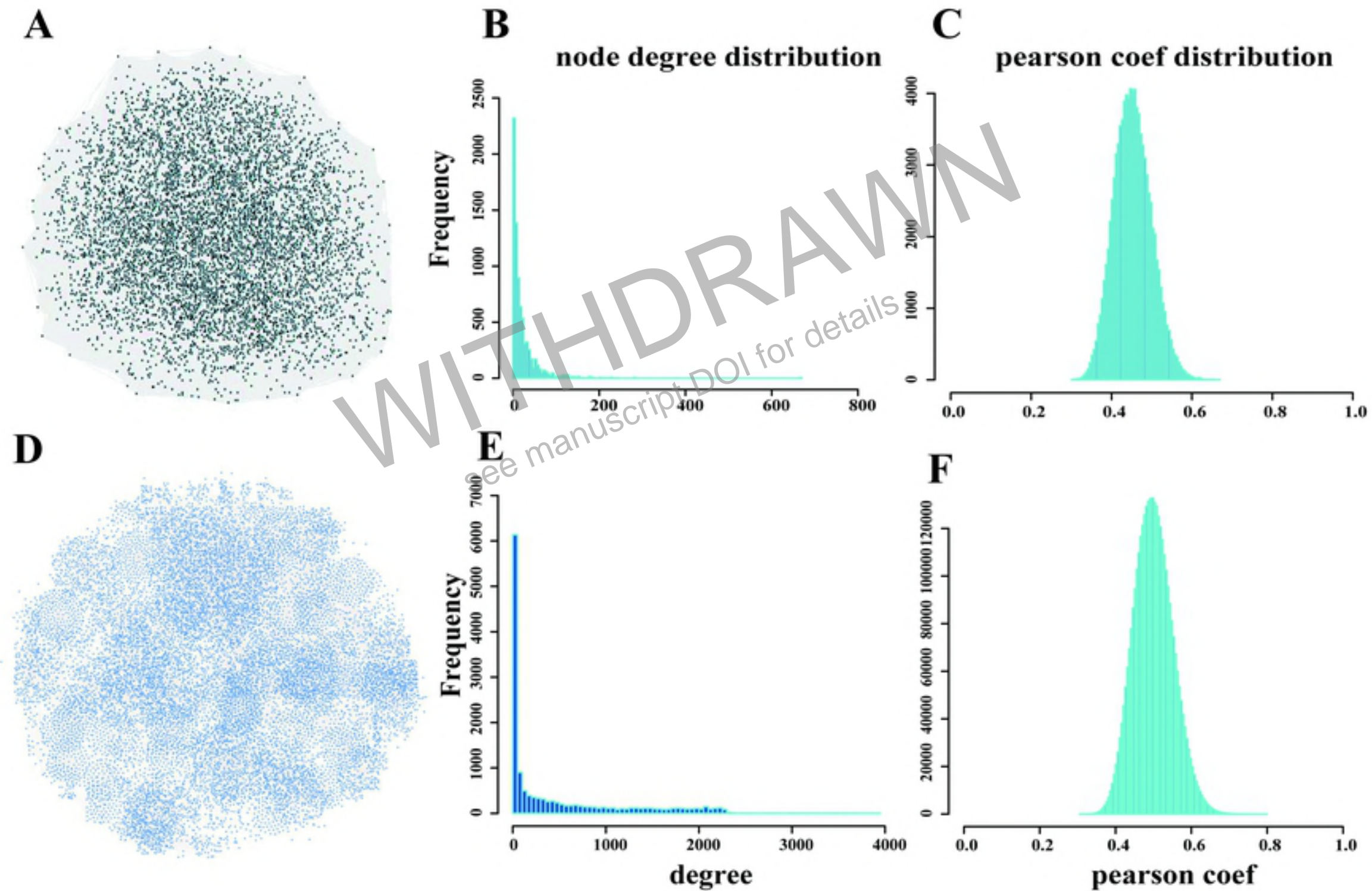


figure 4

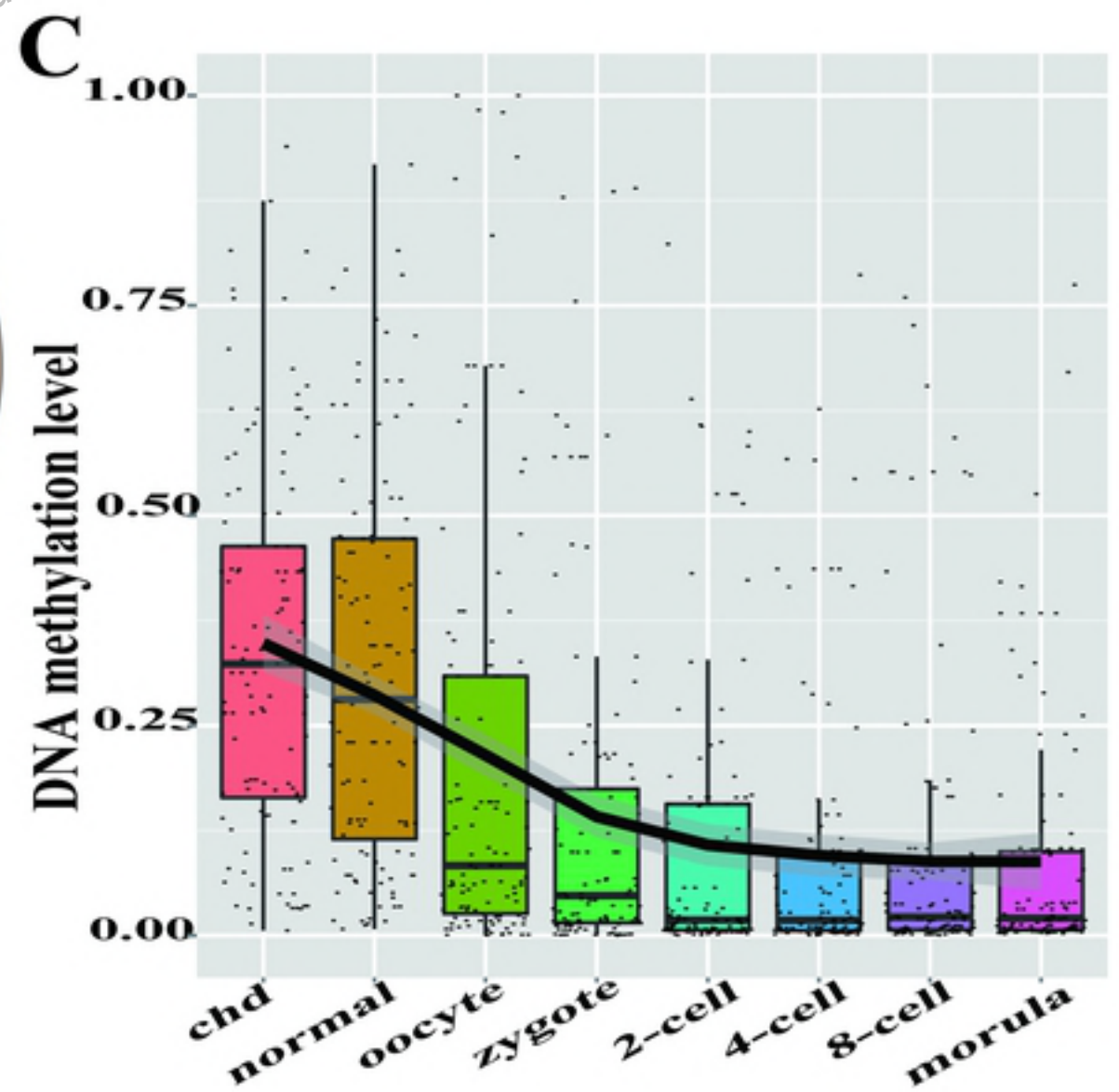
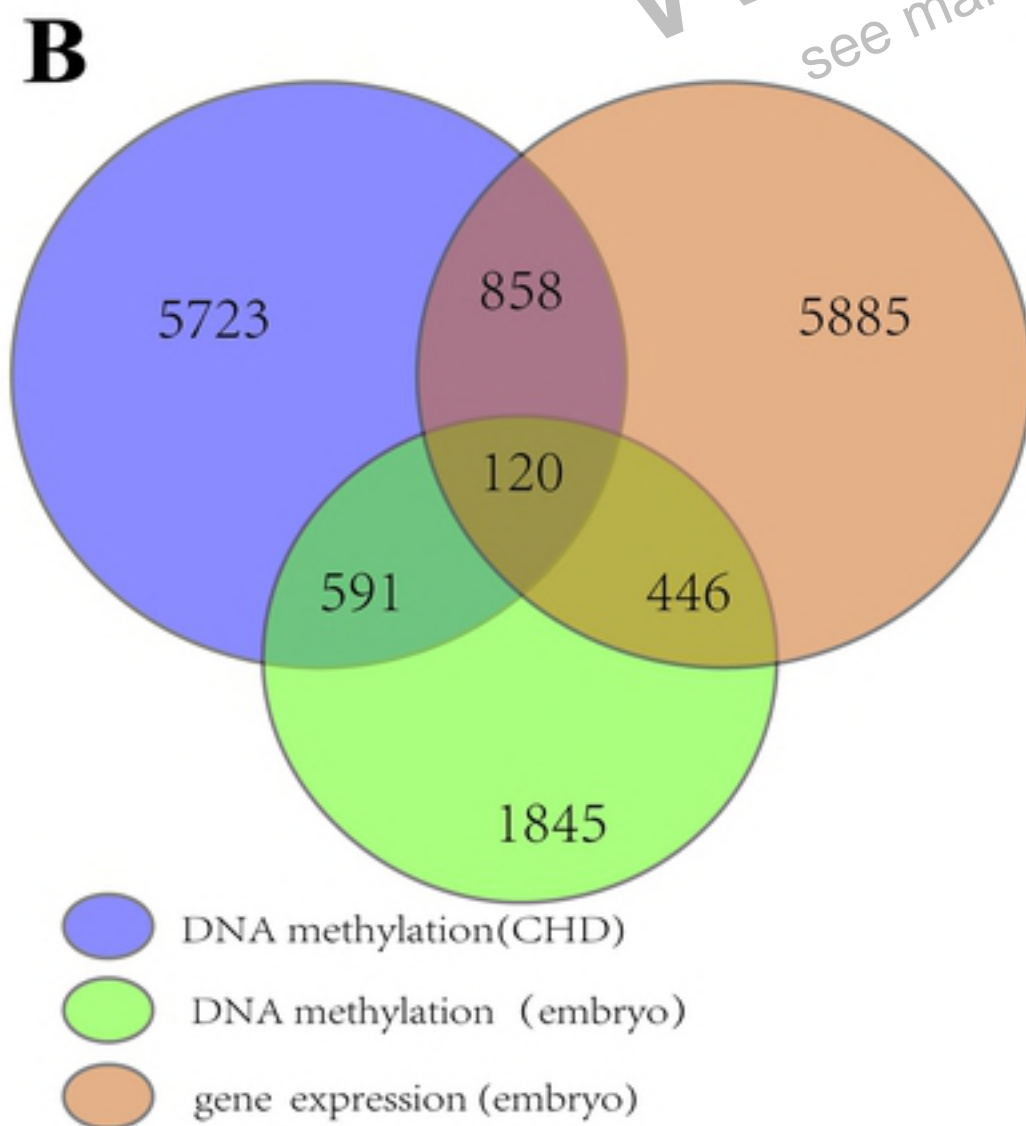
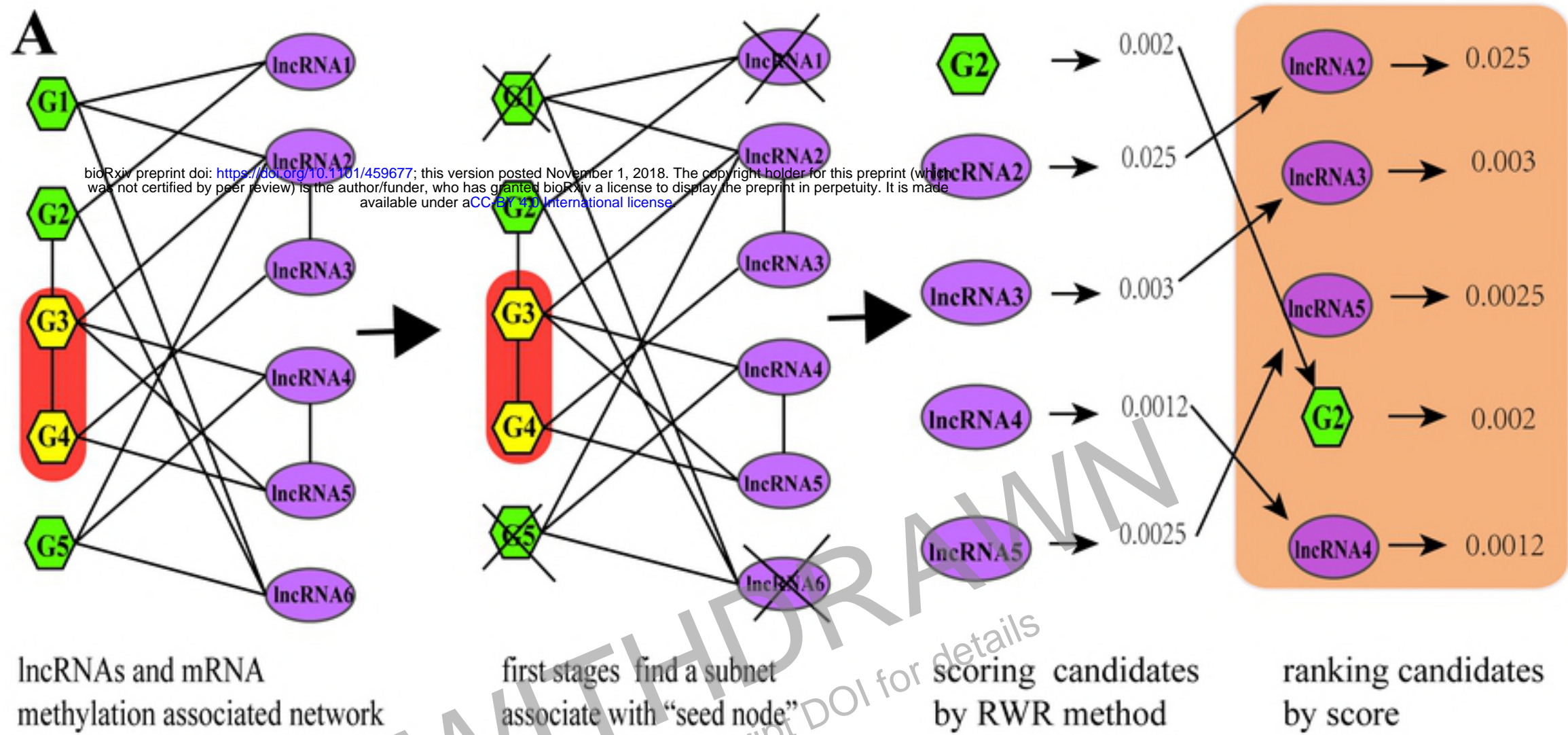


figure 5

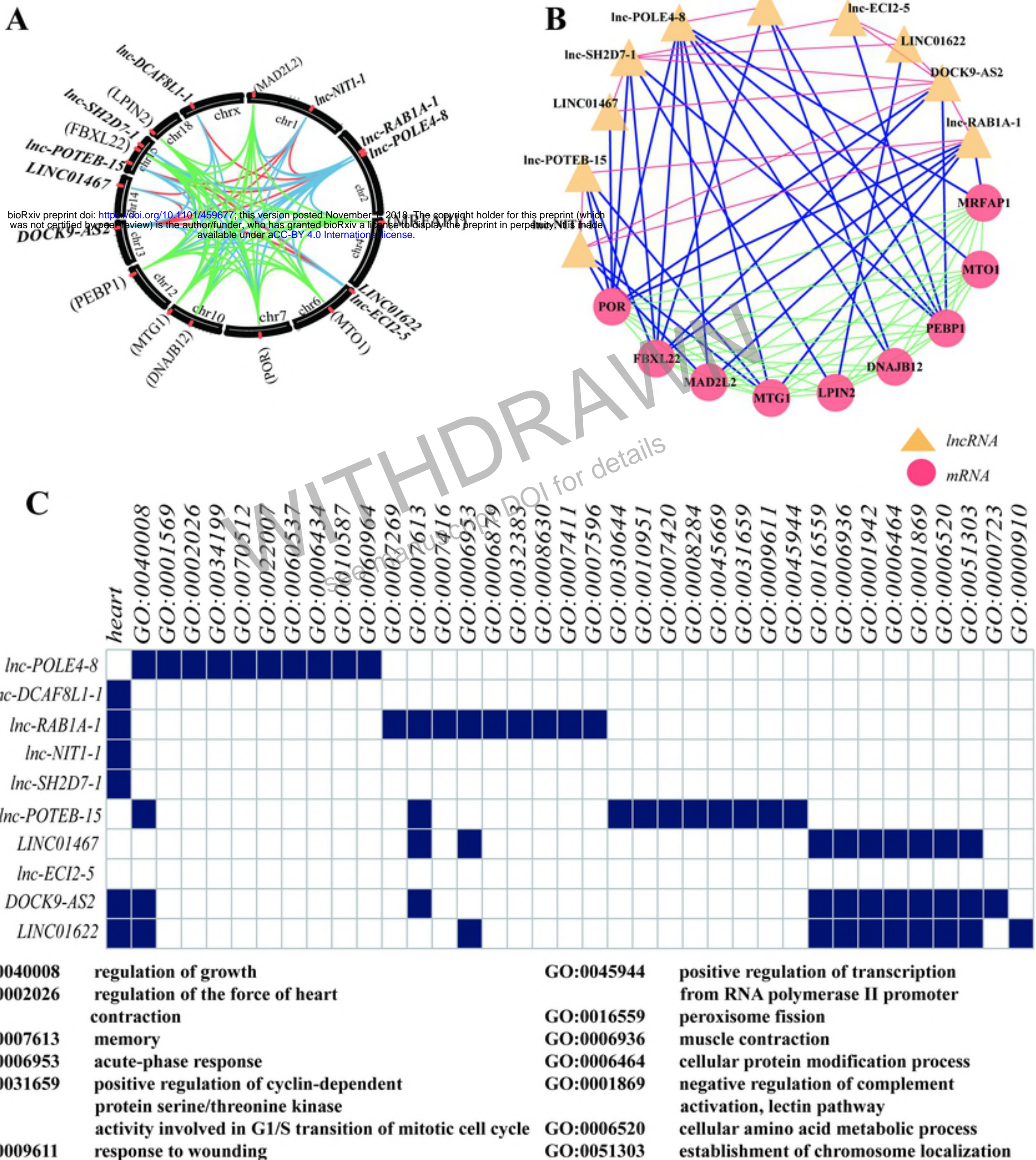


figure 6

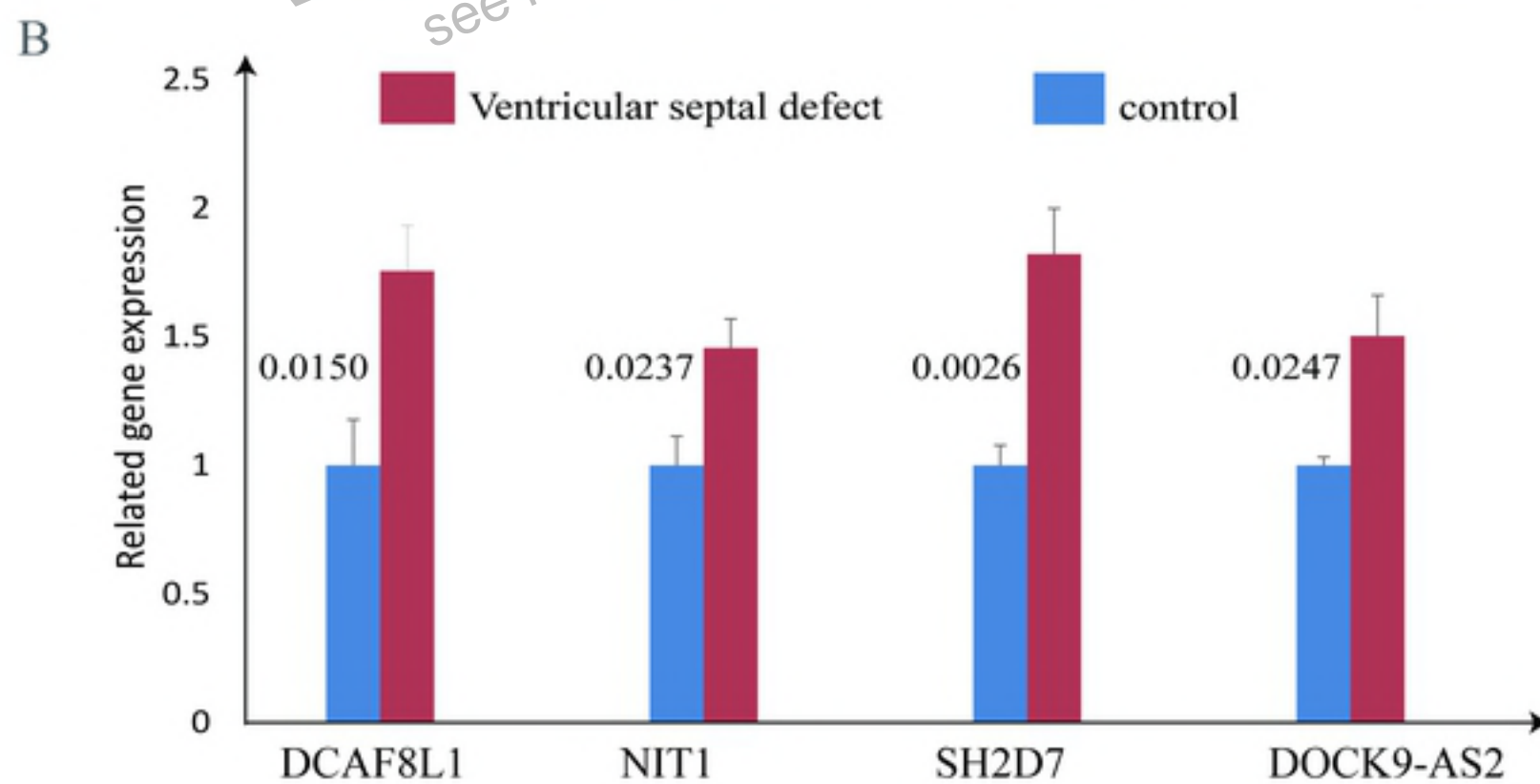
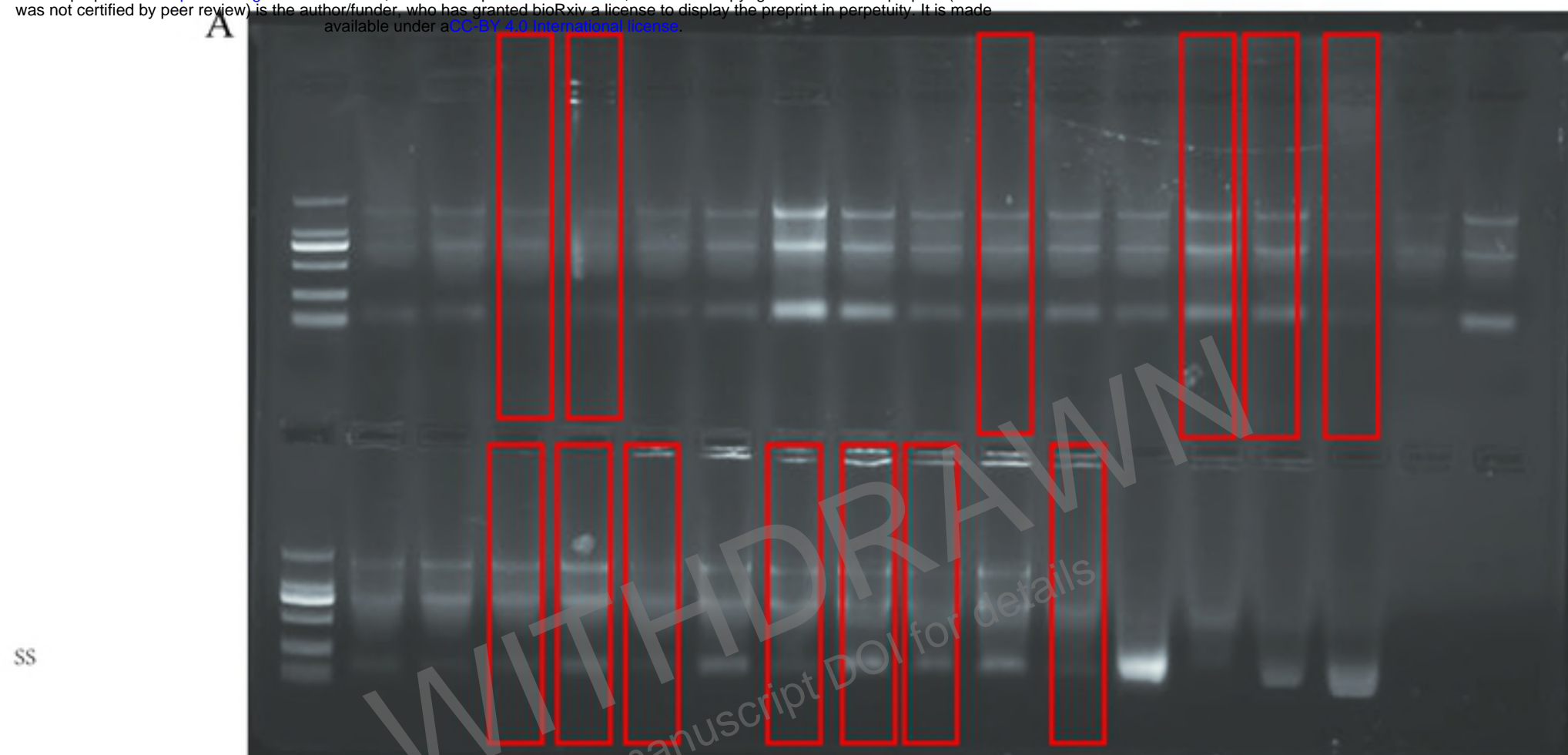


figure 7

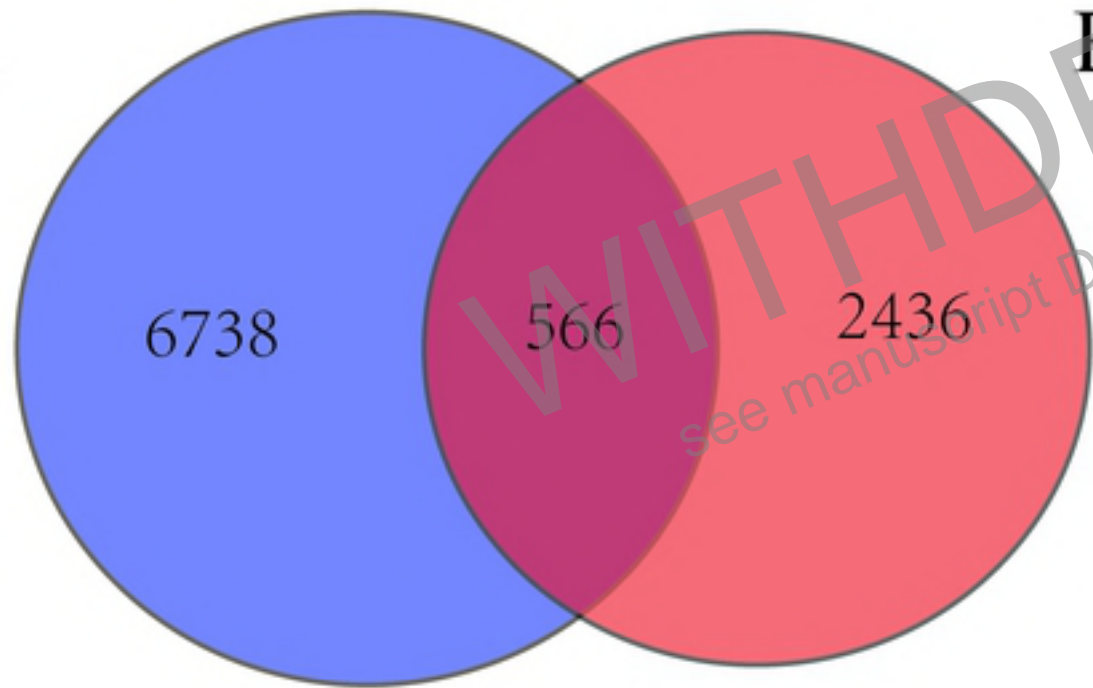
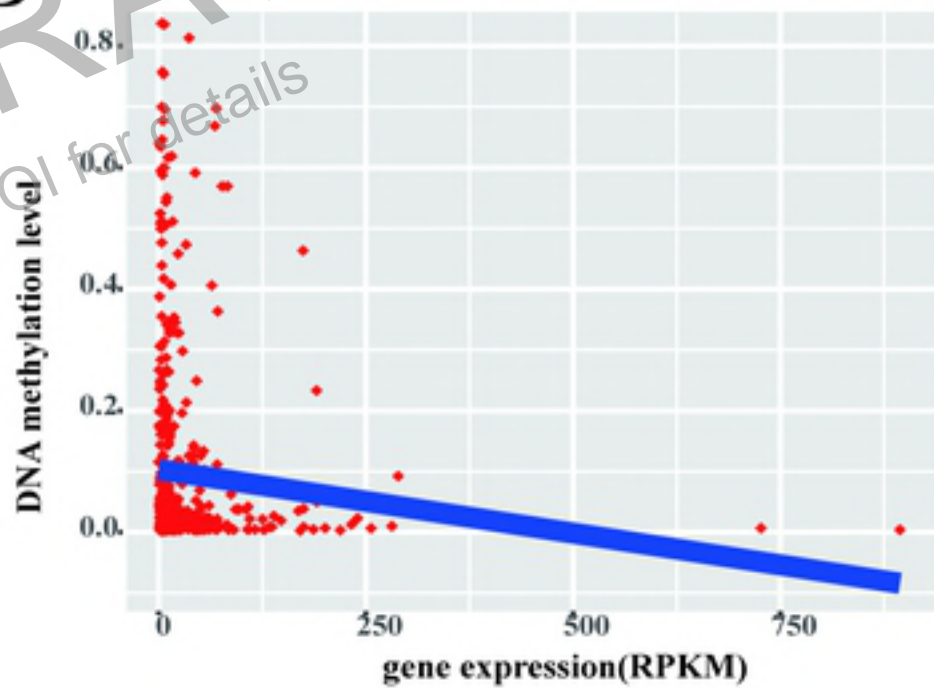
A**B**

figure S1

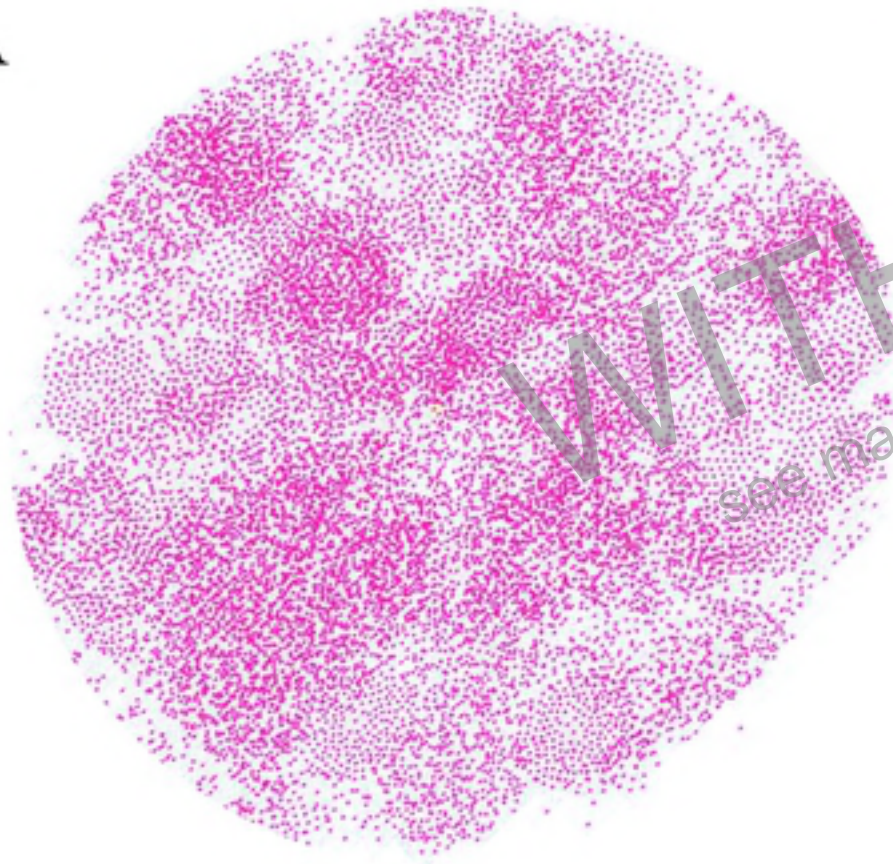
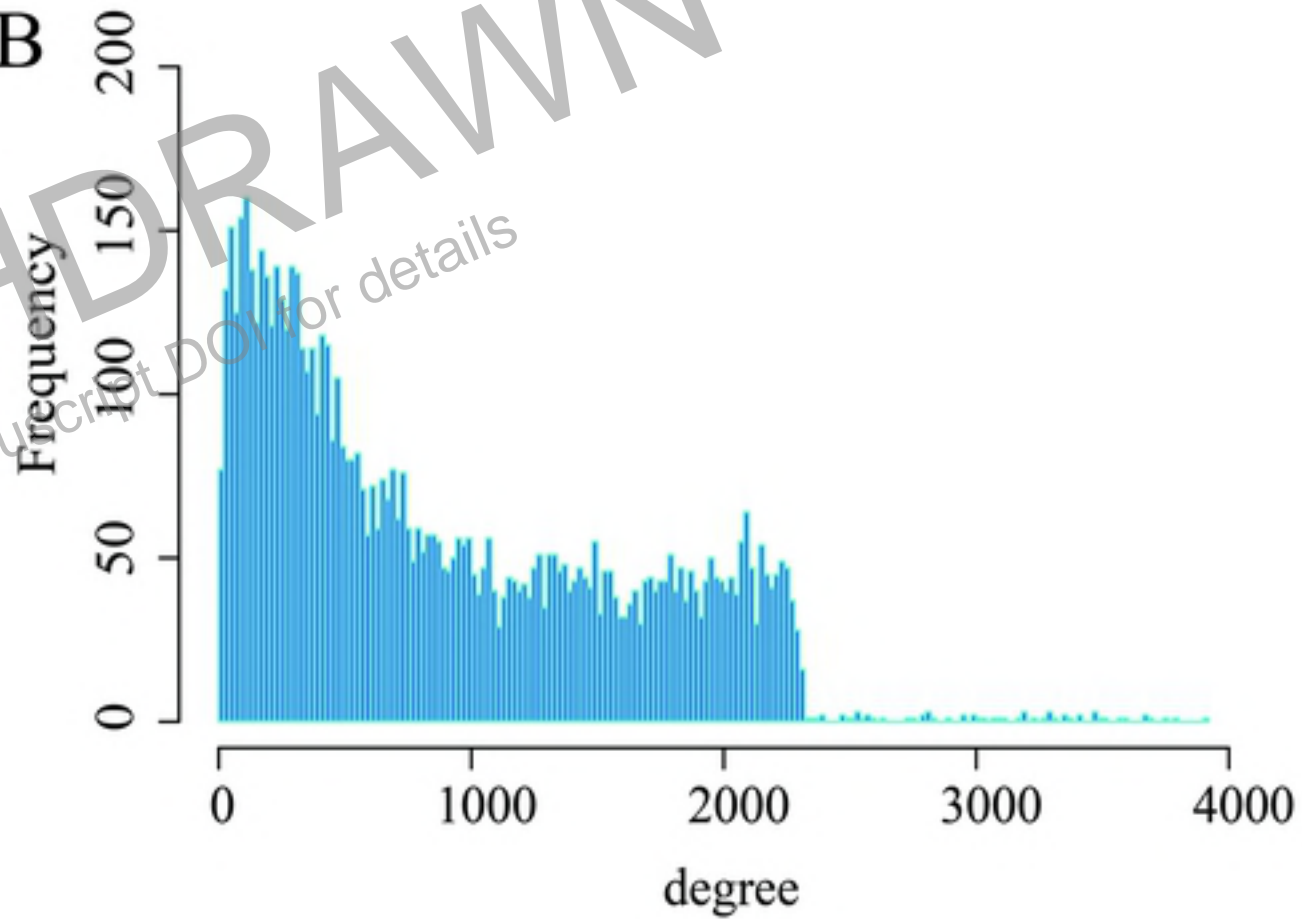
A**B**

figure S2