

Integrative Biology

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

Insight box statement:

The components of living cells interact with each other forming biological networks. The topology of a biological network determines to a large extent its dynamic properties and modes of operation. The controllability of a transcriptional regulatory network can be interpreted as the ability of the cell to control the expression of genes based on control by some transcription factors responding to environmental cues. Here we show that the controllability is a function of the topology and the complexity of the system. Internal loops in the network increase the controllability, but it may cause instability of the system. Thus, there seems to be a trade-off between controllability and stability of regulatory networks.

1 **Controllability analysis of transcriptional regulatory networks** 2 **reveals circular control patterns among transcription factors**

3 Tobias Österlund¹, Sergio Bordel¹, Jens Nielsen^{1,2}

4
5 ¹ Novo Nordisk Foundation Center for Biosustainability, Department of Chemical and
6 Biological Engineering, Chalmers University of Technology, SE-41296 Göteborg, Sweden

7 ² Novo Nordisk Foundation Center for Biosustainability, Technical University of Denmark,
8 Fremtidsvej 3, DK-2970 Hørsholm

11 **Abstract:**

12 Transcriptional regulation is the most committed type of regulation in living cells where
13 transcription factors (TFs) control the expression of their target genes and TF expression is
14 controlled by other TFs forming complex transcriptional regulatory networks that can be
15 highly interconnected.

16 Here we analyze the topology and organization of nine transcriptional regulatory networks for
17 *E.coli*, yeast, mouse and human, and we evaluate how the structure of these networks
18 influences two of their key properties, namely controllability and stability. We calculate the
19 controllability for each network as a measure of the organization and interconnectivity of the
20 network. We find that the number of driver nodes n_D needed to control the whole network is
21 64% of the TFs in the *E.coli* transcriptional regulatory network in contrast to only 17% for the
22 yeast network, 4% for the mouse network and 8% for the human network. The high
23 controllability (low number of drivers needed to control the system) in yeast, mouse and
24 human is due to the presence of internal loops in their regulatory networks where the TFs
25 regulate each other in a circular fashion. We refer to these internal loops as circular control
26 motifs (CCM). The *E.coli* transcriptional regulatory network, which does not have any CCMs,
27 shows a hierarchical structure of the transcriptional regulatory network in contrast to the
28 eukaryal networks. The presence of CCMs also has influence on the stability of these
29 networks, as the presence of cycles can be associated with potential unstable steady-states

30 where even small changes in binding affinities can cause dramatic rearrangements of the state
31 of the network.

32 **Introduction**

33 The number of biological network reconstructions for model organisms like *S.cerevisiae* and
34 *E.coli* has increased dramatically during recent years and include many types of networks, e.g.
35 signaling networks, protein interaction networks¹ and metabolic networks². For the model
36 organisms *S.cerevisiae* and *E.coli* there is a compendium of genome-scale metabolic network
37 reconstructions (GENREs) and genome-scale metabolic models (GEMs) available, which
38 have been used and applied for several different purposes^{3,4}. Even though these models have
39 shown excellent capabilities in predicting different phenotypes, they do have limitations and
40 false predictions is generally due to missing information about regulation of the metabolism⁵.
41 Attempts to incorporate transcriptional regulation of the metabolic genes into FBA
42 simulations have been done both for *E.coli*^{5,6} and *S.cerevisiae*⁷. Transcriptional regulation is
43 condition dependent in the sense that most transcription factors (TFs) bind and recognize
44 specific sequence motifs and a majority of TFs appear to regulate transcription only at
45 specific growth conditions or under specific environmental perturbations^{8,9}. By integrating
46 information about transcriptional regulation to the genome-scale metabolic models we would
47 improve the ability of the model to predict a phenotype from a given genotype. In order to be
48 able to incorporate regulatory information in the models we need to understand more about
49 the organization and structure of transcriptional regulatory networks, as well as how the
50 network behaves under different conditions. The organization of transcriptional regulatory
51 networks has also an impact on the genotype to phenotype relationship in complex diseases as
52 reviewed by Vidal et al.¹⁰

53 The first step towards understanding the regulation of biological processes on a global (i.e.
54 genomic) scale is to reconstruct the transcriptional regulatory network (TRN). For
55 *S.cerevisiae*, TF-DNA interactions have been characterized by ChIP-chip experiments and
56 then been used to construct the yeast TRN^{9,11,12}. In order to model the transcriptional
57 regulation different approaches have been taken. The TRN can be represented as a Boolean
58 model where the TFs are either on or off (1 or 0) based on the activity of other transcription
59 factors and environmental factors. The Boolean modeling approach is implemented in the
60 rFBA framework⁵ where the states of the metabolic genes depend on the states of the
61 controlling TFs. Another approach is probabilistic regulation of metabolism (PROM)¹³ where

62 the probability of TF regulation can be estimated for each TF-target gene pair by counting if
63 the TF and the target gene are expressed in a large number of transcriptome experiments, e.g.
64 microarrays. The PROM method does not use Boolean logic to describe the probability of TF
65 binding; instead the probability is continuous between 0 and 1. However it requires setting a
66 threshold value for a gene to be expressed or not in order to estimate the probabilities of TF
67 binding.

68 Human and mouse transcriptional regulation is even more complex than for *E.coli* and yeast.
69 The human ENCODE project¹⁴ aims to characterize and map functional elements of the
70 human genome, including cis-regulatory elements and non-coding RNAs etc. Two databases,
71 Cscan¹⁵ and Chip Enrichment Analysis (ChEA)¹⁶ have re-analyzed part of the ChIP-seq data
72 from the ENCODE project and from other publications for human and mouse.

73 Here we constructed nine different transcriptional regulatory networks from different ChIP-
74 chip and ChIP-seq experiments for *E.coli*, *S.cerevisiae*, human and mouse. These networks
75 were analyzed in terms of organization, topology and network structure in order to get
76 increased understanding about the transcriptional regulation in these organisms, and how the
77 TRN architecture differs between different species. In order to do this we analyzed each of the
78 networks in terms of network controllability¹⁷ and stability, and calculated how many driver
79 nodes are needed to control the system. For *S.cerevisiae* we also identified TFs that respond to
80 environmental cues by analyzing microarray data from several chemostat studies where the
81 environment was tightly controlled. By controlling the TFs that respond to the environment
82 and calculate how many other TFs in the network that also can be controlled we obtained an
83 understanding of the condition-specific behavior of the yeast transcriptional regulatory
84 network.

85 **Materials and methods**

86 **Generating TF-TF regulatory networks**

87 The nine transcriptional regulatory networks used in this study were derived using different
88 ChIP-chip and ChIP-seq datasets as a starting point. The three *S. cerevisiae* networks were
89 derived from the Yeasttract database 2011-10¹⁸, the Harbison TF-DNA interaction data with
90 binding p-value < 0.001 and with binding p-value < 0.005⁹. The TRN from Yeasttract includes
91 TF-gene interactions both with direct evidence (ChIP-chip) and indirect evidence (the gene
92 was transcriptionally changed in a TF knockout). The *E.coli* transcriptional regulatory

93 network was derived from the *E.coli* Regulon DB version 8.2¹⁹. For mouse and human TF-
94 gene interactions were collected from the Chip enrichment analysis database (ChEA)¹⁶.
95 Version 1 of the ChEA database was downloaded from <http://amp.pharm.mssm.edu/chea> in
96 October 2013. Version 2 of the database was downloaded from
97 <http://amp.pharm.mssm.edu/ChEA2/> in November 2013. For human we also downloaded TF-
98 gene interactions from the Cscan database¹⁵. For the controllability analysis we considered
99 only TF-TF interactions, so all non-TF genes were filtered out. For visualization of the
100 networks we used Cytoscape version 3.0.1. The hierarchical structures of the networks were
101 obtained by choosing the Hierarchical layout view in Cytoscape. The simulated scale-free and
102 random networks were constructed using the igraph R-package²⁰ using the Barabasi-Albert²¹
103 and Erdos-Renyi²² models for network growth. The simulated networks all had 100 nodes and
104 a varying number of edges to simulate networks with different average degree.

105 **Controllability analysis**

106 The concept of controllability of complex networks was introduced by Liu et al. 2011¹⁷. The
107 number of driver nodes n_D is defined as the minimum number of nodes that need to be
108 controlled as input to the system to control 100 % of the network. This number is obtained
109 from the maximum number of matched nodes in the network when solving the controllability
110 equation. The maximum matching path is defined in Liu et al. 2011¹⁷ as the maximum set of
111 links in the networks that do not share start or end nodes, i.e. one path that can control all
112 output nodes from all input nodes. Here we used linear programming to retrieve the number
113 of maximum matching nodes. First, the TF-TF interaction network was converted to an $n \times m$
114 matrix, A , where the n rows represents TF nodes and the m columns represent TF-TF
115 interactions (connections). A connection has the value 1 for the TFs that are connected and 0
116 for the TFs that are not involved in this connection. The following linear program was applied
117 to obtain the maximum matching path, i.e. the longest non-overlapping path in the network,
118 connecting the input nodes with the output nodes:

$$119 \text{ Maximize: } \sum_{i=1}^m x_i ,$$
$$120 \text{ Subject to: } Ax = b, x_i \in \{0,1\} \text{ for } i=1,2,\dots,m \quad (1)$$

121 Where x_1, x_2, \dots, x_m are the TF-TF interactions (connections) in the network and b is the vector
122 of input signals which is set to 1 for one of the nodes and 0 for all other nodes in the network.
123 The maximal number of controlled TFs, n_c is calculated as the number of nodes (TFs) that are
124 included in the maximum matching path. The network was constructed and analyzed using

125 the Raven toolbox for Matlab²³. We used the Raven function *getAllSubGraphs()* to calculate
126 the number of internal loops (circular control motifs) and remove the loops that was not
127 covered in the maximum matching paths, i.e. loops without any input or output nodes.

128 **Integrated analysis of transcriptome data**

129 The raw data files (CEL-files) for the selected microarray studies²⁴⁻³⁵ were downloaded from
130 Gene expression omnibus (GEO) and ArrayExpress using the accession numbers given in the
131 papers. All 233 microarrays used the Affymetrix yeast 98 platform which made it possible to
132 normalize all data together. The data were normalized in R using Plier normalization with
133 only perfect match probes. Metadata for each experiment were collected and used to construct
134 the regression model described in Equation 1. The regression model was implemented in R
135 and ANOVA p-values were calculated for each gene and for each of the coefficients β_1, \dots, β_5
136 where the null hypothesis for each gene is that $\beta_i=0$, and the alternative hypothesis is that $\beta_i \neq 0$
137 for $i=1, \dots, 5$. The p-values were corrected for multiple testing using Benjamini and Hochbergs
138 method (FDR).

139 A hypergeometric enrichment test was applied to identify TFs with over-represented
140 significantly changed target genes. The TF-gene interactions were taken from the Yeabstract
141 TRN¹⁸. For the factors oxygen availability, nutrient limitation (N-limited vs. C-limited) and
142 dilution rate (increasing or decreasing) we performed two enrichment tests for each TF, one
143 for up-regulation of the target genes, and one for down-regulation of the target genes using
144 $\text{adj. } p < 0.05$ as cutoff and $\log\text{FC} > 0$ for up-regulation and $\log\text{FC} < 0$ for down-regulation. For
145 the factors with more than two levels (carbon source and extra compound) we only performed
146 one test for each TF.

147

148 **Results**

149 **Controllability of transcriptional regulatory networks**

150 The nine transcriptional regulatory networks (TRNs) that were included in this study are
151 presented in Table 1, including three *S. cerevisiae* TRNs^{9, 18}, one *E. coli* network¹⁹, two mouse
152 networks¹⁶ and three human networks^{15, 16}. Each of these networks were analyzed in terms of
153 network controllability¹⁷ and stability. The networks were constructed using ChIP-chip and
154 ChIP-seq datasets as a base and they contain only TF-TF interactions (i.e. non-TF genes were

155 filtered out). The number of driver nodes n_D is also presented in Table 1, with n_D being
156 defined as the minimum number of input nodes in the network that need to be controlled in
157 order to control 100% of the network¹⁷. As an example, by controlling node A in Figure 1a,
158 we can control all three nodes, while we need to control node A, and C to have control over
159 100% of the TFs in Figure 1b. The example in Figure 1c contains an internal loop which
160 means that 100% of the network can be controlled by controlling any of the nodes A, B or C
161 as input. For a TRN the concept of network controllability corresponds to one or more TFs
162 responding to environmental changes and these TF regulates other TFs by controlling their
163 transcription. The number of driver nodes, n_D , was determined from the “maximal matching”
164 graph (see materials and methods). In terms of stability, the system of control in Figure 1a and
165 Figure 1b will always have asymptotically stable steady states, while the TF system in Figure
166 1c can have either unstable or stable steady states depending on the parameters and the type of
167 regulation (activation or repression). For each network we studied how controllability and
168 stability are related to its topology (random or scale-free) and its average degree, i.e. the
169 average number of connections each TF has to other TFs in the network.

170 **TRNs show circular control motifs**

171 **Error! Reference source not found.** shows the number of driver nodes n_D as a measure of
172 controllability of the nine different TRNs. For the Yeastract network only 17 % of the TFs
173 need to be controlled in order to control all the other TFs. This large controllability is due to a
174 large internal loop in the network containing 78% of the TFs where all the TFs in the loop are
175 controlling each other in a circular fashion. We call this internal loop a circular control motif
176 (CCM) and an example of a CCM can be found in Figure 1c where node A controls node B,
177 node B controls node C and node C control node A in a circular manner. In terms of
178 controllability this means that all the TFs in the loop can be controlled regardless which of the
179 nodes is externally controlled. For the Harbison networks n_D is 37% of the TFs for the
180 $p < 0.005$ network and 48% for the $p < 0.001$ network (these p-values indicate measure the
181 statistical evidence of the interactions included in the network). We see the same trend in
182 these networks; a CCM contains around 36% and 19% of the nodes in these two networks,
183 respectively. For the *E.coli* network derived from Regulon DB, which has a n_D of 64%, we
184 cannot find any large circular control network motif, and it is only possible to control a
185 maximum of 6% of the nodes by controlling a single input node. For the mouse networks n_D
186 is 4% and 5% for the Chea v.1 and Chea v.2 networks, respectively, whereas n_D is 24%, 8%

187 and 54% for the three Human networks. The number of nodes belonging to the CCM for the
188 different networks is presented in Supplementary table 1.

189 **Controllability of simulated networks**

190 In order to test if the circular control motif (i.e. TFs controlling each other in a circular
191 manner) in yeast, mouse and human, but not in *E.coli*, has evolved as a result of higher
192 connectivity and higher average degree in these networks we simulated random Erdos-Renyi
193 networks²² with different average degree and scale-free Barabasi-Albert networks²¹ with
194 various average degree. The controllability for random Erdos-Renyi networks, when
195 controlling one input node is dependent on the average degree of the network and for the
196 simulated random networks with average degree 10 or higher we only need to control 1 TF as
197 input (approx. 1%) to control 100% of the network (Figure 2). However, for the simulated
198 scale-free networks the controllability is less dependent on the average degree of the network.
199 The n_D for the scale-free network decreases slightly when increasing the average degree but
200 does not exceed 44% for the simulated networks. In Figure 2 we have also plotted the number
201 of driver nodes n_D as a function of the average network degree for the seven real networks
202 included in this study. For the three *S.cerevisiae* networks the trend is that the networks with
203 higher average degree also have lower n_D . The Yeastract network for example has an average
204 degree of 18.28 and a n_D of 17%. The two mouse networks also have high average degrees
205 and low n_D .

206 **Stability analysis of transcriptional regulatory networks**

207 To illustrate the stability of a TRN we will consider a simple network with 3 TFs. For each of
208 these TFs we let n_i denote the number of copies of the i 'th transcription factor in the network
209 and r_i denote the rate of production of the TF. Assuming that the i 'th TF does not regulate its
210 own transcription the time evolution of n_i can be described by the equation $\frac{dn_i}{dt} = r_i - \delta_i n_i$

211 where δ_i is the specific degradation rate of the i 'th TF. The rate r_i can be seen as a function of
212 the abundances of all TFs that are regulating the transcription of the i 'th TF in the regulatory
213 network.

214 The differential matrix of the system can be expressed as $D = \frac{\partial \vec{r}}{\partial \vec{n}} - \Delta$ where $\frac{\partial \vec{r}}{\partial \vec{n}}$ is the
215 Jacobean matrix for the transcription rates and Δ is a diagonal matrix with the specific
216 degradation rates δ . To perform a stability analysis of the system we can calculate the

217 eigenvectors and eigenvalues of the differential matrix D , in the neighborhood of a steady
 218 state point. For the example for the network in Figure 1a the topological differential matrix
 219 will be:

$$220 \quad D = \begin{pmatrix} -\delta_A & \frac{\partial r_B}{\partial n_A} & 0 \\ 0 & -\delta_B & \frac{\partial r_C}{\partial n_B} \\ 0 & 0 & -\delta_C \end{pmatrix}$$

221 Note that we have inserted zeros in the matrix for the instances where we do not have an edge
 222 in the regulatory network, e.g. TF C is not directly controlled by TF A so $\frac{\partial r_C}{\partial n_A} = 0$, etc. The
 223 Eigenvalues of this matrix is obtained by solving the equation $\det(D-\lambda)=0$ and thus the
 224 eigenvalues for the network in Figure 1a becomes $\lambda_1 = -\delta_A$, $\lambda_2 = -\delta_B$, $\lambda_3 = -\delta_C$. If we
 225 assume that the degradation rates δ_A , δ_B and δ_C will be positive, all eigenvalues for the
 226 system will be negative, independent on the parameters and type of regulation, but clearly
 227 dependent only on the topological structure of the network. The fact that we will always have
 228 negative eigenvalues implicates that the system in example 1A will always be stable. Also for
 229 the example in figure 1B the eigenvalues will be $\lambda_1 = -\delta_A$, $\lambda_2 = -\delta_B$, $\lambda_3 = -\delta_C$, and this
 230 network is therefore also stable. For the example in Figure 1c the eigenvalues of the system
 231 will be a function of $\frac{\partial r_B}{\partial n_A}$, $\frac{\partial r_C}{\partial n_B}$ and $\frac{\partial r_A}{\partial n_C}$ and they can be either positive or negative depending
 232 on the parameters and the type of regulation. If all (real parts of the) eigenvalues are negative,
 233 the system will be stable, but if one or more of the real parts of the eigenvalues are positive,
 234 the system will be unstable. A more general derivation of this concept called topological
 235 stability analysis and another example with three transcription factors is presented in
 236 Electronic supplementary information (ESI) and Supplementary figure S1. This approach can
 237 be applied for a general network to calculate the supports of the eigenvectors (components
 238 with non-zero values) of the Jacobean matrix of the system, as a function only of the topology
 239 and regardless of the system parameters.

240 **Hierarchical structure of the TRNs**

241 The hierarchical structure of both the *E.coli* and *S.cerevisiae* transcriptional regulatory
 242 network has previously been reported^{36,37}. However, Figure 3a and b shows that the TRN for
 243 *E.coli* is more hierarchical than yeast when it comes to regulation. The *E.coli* regulation is
 244 controlled by *cRP* at the top of the hierarchical tree (Figure 3b), whereas for yeast there are

245 many TFs on the top hierarchical level (*SWI6*, *MBP1*, *FKH2*, *FKH1*, *ABF1*, *NRG1*, *INO4* and
246 *SKN7*) (Figure 3a). Another difference between the *E.coli* and yeast network structure is the
247 presence of the large circular control motif (CCM) in the yeast network where some of the
248 TFs in this internal loop are even present in the top hierarchical level (*ABF1* and *MBP1*).

249 Figure 3c shows the human network (ChEA v.2) ordered in a hierarchical way, where four
250 TFs are on top of the hierarchical structure (*TP63*, *SOX2*, *AR*, *GABP*). However, since 87% of
251 the TFs in this network are part of the CCM this means that by controlling one of the top
252 nodes, 87% of the network can be controlled. The TFs included in the CCM are marked in
253 green in Figure 3c. Figure 3d shows the same hierarchical structure for the mouse (ChEA v.2)
254 network and here is 92% of the TFs are part of the CCM.

255

256 **Integrated analysis reveals condition-specific regulation in yeast**

257 In order to identify TFs that respond to specific environmental cues that can represent input
258 nodes in the network controllability analysis we analyzed transcriptome data from 233 yeast
259 microarrays collected from 11 separate studies²⁴⁻³⁵. The experiments were all carried out in
260 chemostat cultures where the environment could be tightly controlled, i.e. the specific growth
261 rate, oxygen availability, carbon source etc. The data were analyzed using a regression model
262 in order to describe the expression of each gene in terms of the effect of the environment.

263 The environmental factors included in the regression model are presented in **Error!**

264 **Reference source not found.** We were specifically interested in determination of the TFs
265 that respond to these environmental factors. To identify condition-specific transcriptional
266 regulation we considered three different features: i) TFs whose target genes are up- or down-
267 regulated under different conditions, ii) TFs that do not change in expression between
268 conditions and iii) TFs that are reported in literature to respond to a specific environmental
269 cue.

270 For each environmental factor in **Error! Reference source not found.** we determined the
271 genes where the environmental factor had a significant effect on the gene expression (adjusted
272 p-value <0.05). We then performed a hypergeometric enrichment test for each TF in order to
273 identify the TFs where the expression of the target genes was influenced by the environment.
274 The target genes for each TF were defined by the TRN from Yeasttract. The results for the

275 hypergeometric enrichment test for the factors oxygen availability, nutrient limitation and
276 dilution rate and carbon source are presented in Supplementary Figure S2 and Table S2.

277 Figure 4 shows a comparison of feature i) and ii) above, i.e. the ability of the TF to change the
278 expression of the target genes and the transcriptional change of the TF itself. The genes
279 marked with red in the four plots are the TFs where the adj. p-value of the TF gene itself > 0.8
280 and the p-value of the hypergeometric test (transcriptional change of the target genes) < 0.05 ,
281 i.e. TFs that don't change in expression themselves but their target genes are significantly
282 regulated as a function of the environmental factor. These TFs must be regulated in another
283 way as a response to the environment, either by activation through signaling or through
284 interaction with other TFs. In response to oxygen availability these transcription factors are
285 *HAA1* and *FKH2*, to nutrient limitation *MOT3*, and to altered specific growth rate (dilution
286 rate) the TFs are *PHO4*, *FKH2* and *MGAI*, and in response to changes in carbon source they
287 are *OAF1* and *CST6*.

288 We calculated the controllability of the network when controlling the TFs that were identified
289 to respond to the environment in the TRN from Yeasttract. When we control the nodes marked
290 with red in Figure 4 for each environmental factor it is possible to control between 78-82% of
291 the network, e.g. when controlling both *HAA1* and *FKH2* responding to oxygen limitation the
292 controllability is 81%, 78% for nutrient limitation, 82% for dilution rate and 82% when
293 controlling *OAF1* and *CST6* for the carbon source environmental factor.

294 Discussion

295 Regulation of cellular processes is complex and may occur on different levels in the cell. Here
296 we consider transcriptional regulation, which can be thought of as consisting of three different
297 layers. The first layer is the environment which can be thought of as different environmental
298 cues, e.g. high or low oxygen levels, carbon or nitrogen limitation etc. The second layer
299 consists of the transcription factors where some TFs respond to the environment and the TFs
300 regulate the transcription of other TFs as well as other genes. The third layer consists of genes
301 which encodes for proteins that carry out different functions in the cell. Here we analyzed the
302 structure and topology of transcriptional regulatory networks (TRNs) by calculating the
303 controllability, i.e. how many driver nodes are needed to control the networks.

304 Understanding the regulation of a cellular system on a global level can have large
305 implications in for example metabolic engineering where an organism is engineered to

306 produce high amounts of a chemical of interest^{38, 39}. If we have information about the
307 regulation of metabolism this can help to identify bottlenecks in the route to over-production
308 of the product of interest. Also, knowledge about the organization of transcriptional regulation
309 will help in understanding of development of complex diseases and other complex traits.

310 The controllability analysis of the TRNs for *E.coli*, yeast, human and mouse resulted in
311 several novel findings. The controllability of the *E. coli* TRN is low, i.e. we need to control as
312 much as 64% of the input nodes to be able to control all output nodes, whereas the
313 controllability for yeast, mouse and human is much higher. The organization of the *E.coli*
314 TRN is hierarchical in the sense that most of the TFs do not control any other TFs but are
315 controlled by one or more TFs and the number of TFs that are involved in one route of
316 regulation from *cRP* on the top hierarchical level to the leaves of the tree is maximum 6 as
317 can be seen in the *E.coli* hierarchical tree in Figure 3b. For the yeast, mouse and human
318 networks we identify so called circular control motifs (CCMs) where we can control a large
319 part of the network just by controlling one input node. The presence of these internal loops in
320 the networks also means that the hierarchical organization of the network is less prominent
321 since all nodes in the CCMs can be controlled just by controlling one of the nodes in the loop
322 as input. For the human and mouse network around 90% of the network are interconnected
323 and belongs to CCMs.

324 From analysis of the simulated networks presented in Figure 2 we can see that for a perfect
325 scale-free network there are no CCMs, even for networks with a high average degree, since
326 the corresponding network will follow a hierarchical tree where most of the nodes have a
327 degree of 1 (i.e. controlled by one TF but not controlling any TFs) and there is one or few
328 nodes that have a high degree (network hubs). The degree distribution of the scale free
329 network follows a power law distribution, e.g. $P(k) \sim k^{-\alpha}$, where the parameter α typically has
330 a value between 2 and 3. For the random network the degree distribution is uniformly
331 distributed, i.e. the chance for having a node with a low degree is approximately equal to the
332 chance of having a node with a high degree. For simulated random networks a high
333 connectivity (i.e. average degree > 10) means that one can control 100% of the network
334 through a single input node.

335 The behavior of the *E.coli* network seems to follow the behavior of the simulated scale-free
336 network based on the controllability analysis and the results presented in Figure 2. We also
337 investigated the degree distribution of the *E.coli* TRN and the network is scale-free in the

338 sense that it is possible to fit a power law distribution to the degree distribution
339 (Supplementary figure S1). The mouse and human networks seem to be more random in the
340 way the TFs are interconnected. Also, it is not possible to fit a power law distribution to these
341 networks, which points to the fact that these networks are not scale-free. In the yeast networks
342 the number of nodes with one or two neighbors is much less than expected for a scale-free
343 network, especially for the TRN from Yeastract. The three different yeast networks included
344 in this study differ in the confidence used to consider if a TF-gene interaction should be
345 included in the network or not. The Harbison network uses either $p < 0.001$ or $p < 0.005$ for the
346 binding probability cutoff from the ChIP-chip experiment. The Harbison $p < 0.001$ network is
347 the most conservative network and also has the lowest average degree. The Yeastract network
348 is the least conservative yeast network and has the highest number of TF-gene interactions. It
349 is not possible to compare TRNs from the different organisms (and even the different yeast
350 TRNs), since the TRNs were constructed in different ways, with different confidence scores
351 etc. However, still for all the different eukaryal TRNs there is a clear difference in structure
352 and by using the different TRNs we could investigate the controllability of the network as a
353 function of the average degree of the network, and hence also as a function of the confidence
354 of TF binding scores, or how conservative the network construction process has been. The
355 yeast network seems to behave in between the simulated scale-free and random networks, but
356 as we increase the average degree (and become less conservative in what is considered as a
357 TF binding event) the network behaves more random, and the maximum controllability when
358 controlling one input node increases.

359 Network controllability has earlier been introduced by Liu et al.¹⁷ where they also analyzed
360 many different real networks in terms of controllability, including the yeast and *E.coli* TRNs.
361 However the networks that were included in this study were different in the way they were
362 constructed and the number of TF-gene interactions that were included. Here we use network
363 controllability as a tool to study the topology and organization of the network and we identify
364 CCMs in the yeast, mouse and human networks, but not in the *E.coli* network. In terms of
365 evolution and the ability of the organism to adapt to changes in the environment it is
366 interesting to see that the eukaryal TRNs included in this study seem to contain CCMs, but the
367 bacterial TRN does not. Furthermore, based on our analysis it seems like more complex
368 organisms have more random organization of their transcriptional regulatory network.

369 We find the presence of circular motifs in yeast, mouse and human interesting, as this may
370 improve robustness in the transcriptional regulation when it comes to e.g. adapting to
371 environmental changes.

372 In terms of stability a linear network motif will always be stable. A circular network motif,
373 like the example presented in Figure 1c, on the other hand can be unstable, but also be more
374 stable depending on the parameters and types of interaction. If two of the interactions in a
375 triangular motif are positive (i.e. activation) and the third interaction is negative (i.e.
376 repression) the steady state of the system will be less sensitive to changes in environmental
377 parameters compared to a linear pathway. However, gene deletions or perturbations might
378 disturb the system and make the steady state unstable, causing a rearrangement of the system
379 similar to a non-equilibrium phase change, if an environmental or genetic parameter passes a
380 certain threshold. A detailed discussion of these phenomena is contained in the supplementary
381 material. According to our observations, more complex organisms contain more potentially
382 unstable steady states of their TRN. This can contribute to a higher capability to maintain
383 homeostasis under environmental changes, but also to sudden rearrangements of the state of
384 the system as a result of mutations. Using the approach of topological stability analysis
385 introduced in this paper we can identify the non-zero elements of the Jacobean matrix and
386 find potentially unstable motifs in the TRN, without knowing the parameters of the system,
387 but only as a function of the topology.

388 For the yeast TRN we identified TFs that respond to environmental cues by analyzing a large
389 amount of microarray data from different controlled environments. The transcriptional
390 regulation was found to be highly condition-specific and if we could identify condition-
391 specific responses that allowed us to see how the yeast TRN looks under specific conditions
392 or when changing environment from one state to another. This analysis shows that key TFs
393 seem to exert a high degree of controllability in response to different environmental cues, i.e.
394 a few TFs can control a large number of other TFs in the regulatory network.

395 In conclusion we perform analysis of the topology of TRNs for different species and find that
396 there is an increasing complexity in terms of connectivity and controllability when moving
397 from bacteria to yeast and further to mouse and human. Whereas the TRN for *E.coli* is scale-
398 free the TRNs for eukaryotes seems to be more random, mainly due to the presence of circular
399 control motifs (CCMs) involving a large number of TFs. These large CCMs enable control of
400 a large fraction of the TRN through control of many single TFs, which may have been

401 important for establishing increased robustness towards different environmental changes. On
402 the other hand the presence of CCMs can cause instability due to changes in binding affinity
403 for the TFs in the motif, and this may result in large reprogramming of cellular function, e.g.
404 resulting in disease development.

405 Acknowledgements

406 The authors thank Intawat Nookaew for valuable input regarding construction of the
407 regulatory networks. Simulations were performed at Chalmers Centre for Computational
408 Science and Engineering (C3SE). We acknowledge funding from the European Research
409 Council project INSYSBIO (grant no. 247013), the Novo Nordisk Foundation and the
410 Department of Energy (grant no. 5710003389).

411 Author's contributions

412 TÖ analyzed the data, constructed the networks, performed the controllability analysis and
413 simulations and wrote the paper. SB derived the topological stability analysis. SB and JN
414 supervised the work and edited the paper. All authors read and approved the manuscript.

415 References

- 416 1. P. Bork, L. J. Jensen, C. von Mering, A. K. Ramani, I. Lee and E. M. Marcotte, *Current opinion*
417 *in structural biology*, 2004, 14, 292-299.
- 418 2. I. Thiele and B. Ø. Palsson, *Nature protocols*, 2010, 5, 93-121.
- 419 3. T. Österlund, I. Nookaew and J. Nielsen, *Biotechnol Adv*, 2012, 30, 979-988.
- 420 4. D. McCloskey, B. Ø. Palsson and A. M. Feist, *Molecular systems biology*, 2013, 9.
- 421 5. M. W. Covert and B. Ø. Palsson, *Journal of Biological Chemistry*, 2002, 277, 28058-28064.
- 422 6. M. W. Covert, E. M. Knight, J. L. Reed, M. J. Herrgard and B. O. Palsson, *Nature*, 2004, 429,
423 92-96.
- 424 7. M. Herrgård, B. Lee, V. Portnoy and B. Palsson, *Genome research*, 2006, 16, 627.
- 425 8. G. Chua, M. D. Robinson, Q. Morris and T. R. Hughes, *Current opinion in microbiology*, 2004,
426 7, 638-646.
- 427 9. C. T. Harbison, D. B. Gordon, T. I. Lee, N. J. Rinaldi, K. D. Macisaac, T. W. Danford, N. M.
428 Hannett, J.-B. Tagne, D. B. Reynolds and J. Yoo, *Nature*, 2004, 431, 99-104.
- 429 10. M. Vidal, M. E. Cusick and A.-L. Barabasi, *Cell*, 2011, 144, 986-998.
- 430 11. T. I. Lee, N. J. Rinaldi, F. Robert, D. T. Odom, Z. Bar-Joseph, G. K. Gerber, N. M. Hannett, C. T.
431 Harbison, C. M. Thompson and I. Simon, *Science*, 2002, 298, 799-804.
- 432 12. T. R. Hughes and C. G. de Boer, *Genetics*, 2013, 195, 9-36.
- 433 13. S. Chandrasekaran and N. D. Price, *Proceedings of the National Academy of Sciences*, 2010,
434 107, 17845-17850.
- 435 14. I. Dunham, E. Birney, B. R. Lajoie, A. Sanyal, X. Dong, M. Greven, X. Lin, J. Wang, T. W.
436 Whitfield and J. Zhuang, 2012.
- 437 15. F. Zambelli, G. M. Prazzoli, G. Pesole and G. Pavesi, *Nucleic acids research*, 2012, 40, W510-
438 W515.

- 439 16. A. Lachmann, H. Xu, J. Krishnan, S. I. Berger, A. R. Mazloom and A. Ma'ayan, *Bioinformatics*,
440 2010, 26, 2438-2444.
- 441 17. Y.-Y. Liu, J.-J. Slotine and A.-L. Barabási, *Nature*, 2011, 473, 167-173.
- 442 18. M. C. Teixeira, P. Monteiro, P. Jain, S. Tenreiro, A. R. Fernandes, N. P. Mira, M. Alenquer, A. T.
443 Freitas, A. L. Oliveira and I. Sá-Correia, *Nucleic acids research*, 2006, 34, D446-D451.
- 444 19. S. Gama-Castro, H. Salgado, M. Peralta-Gil, A. Santos-Zavaleta, L. Muñiz-Rascado, H. Solano-
445 Lira, V. Jimenez-Jacinto, V. Weiss, J. S. García-Sotelo and A. López-Fuentes, *Nucleic acids*
446 *research*, 2011, 39, D98-D105.
- 447 20. G. Csardi and T. Nepusz, *InterJournal, Complex Systems*, 2006, 1695.
- 448 21. A.-L. Barabási and R. Albert, *science*, 1999, 286, 509-512.
- 449 22. P. Erdős and A. Renyi, *Acta Mathematica Hungarica*, 1961, 12, 261-267.
- 450 23. R. Agren, L. Liu, S. Shoaie, W. Vongsangnak, I. Nookaew and J. Nielsen, *PLoS computational*
451 *biology*, 2013, DOI: 10.1371/journal.pcbi.1002980.
- 452 24. T. Knijnenburg, J.-M. Daran, M. van den Broek, P. Daran-Lapujade, J. de Winde, J. Pronk, M.
453 Reinders and L. Wessels, *BMC genomics*, 2009, 10, 53.
- 454 25. C. Cipollina, J. van den Brink, P. Daran-Lapujade, J. T. Pronk, M. Vai and J. H. de Winde,
455 *Microbiology*, 2008, 154, 337-346.
- 456 26. D. A. Abbott, T. A. Knijnenburg, L. M. de Poorter, M. J. Reinders, J. T. Pronk and A. J. van
457 Maris, *FEMS yeast research*, 2007, 7, 819-833.
- 458 27. J. Aguilera, T. Petit, J. H. de Winde and J. T. Pronk, *FEMS yeast research*, 2005, 5, 579-593.
- 459 28. R. De Nicola, L. A. Hazelwood, E. A. De Hulster, M. C. Walsh, T. A. Knijnenburg, M. J. Reinders,
460 G. M. Walker, J. T. Pronk, J. M. Daran and P. Daran-Lapujade, *Appl Environ Microbiol*, 2007,
461 73, 7680-7692.
- 462 29. V. M. Boer, S. L. Tai, Z. Vuralhan, Y. Arifin, M. C. Walsh, M. D. Piper, J. H. de Winde, J. T. Pronk
463 and J. M. Daran, *FEMS yeast research*, 2007, 7, 604-620.
- 464 30. T. A. Knijnenburg, J. H. de Winde, J. M. Daran, P. Daran-Lapujade, J. T. Pronk, M. J. Reinders
465 and L. F. Wessels, *BMC Genomics*, 2007, 8, 25.
- 466 31. S. L. Tai, V. M. Boer, P. Daran-Lapujade, M. C. Walsh, J. H. de Winde, J. M. Daran and J. T.
467 Pronk, *Journal of Biological Chemistry*, 2005, 280, 437-447.
- 468 32. P. Daran-Lapujade, M. L. Jansen, J. M. Daran, W. van Gulik, J. H. de Winde and J. T. Pronk, *The*
469 *Journal of biological chemistry*, 2004, 279, 9125-9138.
- 470 33. P. Daran-Lapujade, S. Rossell, W. M. van Gulik, M. A. Luttkik, M. J. de Groot, M. Slijper, A. J.
471 Heck, J. M. Daran, J. H. de Winde, H. V. Westerhoff, J. T. Pronk and B. M. Bakker, *Proc Natl*
472 *Acad Sci U S A*, 2007, 104, 15753-15758.
- 473 34. A. Fazio, M. Jewett, P. Daran-Lapujade, R. Mustacchi, R. Usaite, J. Pronk, C. Workman and J.
474 Nielsen, *BMC genomics*, 2008, 9, 341.
- 475 35. B. Regenber, T. Grotkjær, O. Winther, A. Fausbøll, M. Åkesson, C. Bro, L. K. Hansen, S.
476 Brunak and J. Nielsen, *Genome biology*, 2006, 7, R107.
- 477 36. H. Yu and M. Gerstein, *Proceedings of the National Academy of Sciences*, 2006, 103, 14724-
478 14731.
- 479 37. H.-W. Ma, B. Kumar, U. Ditges, F. Gunzer, J. Buer and A.-P. Zeng, *Nucleic acids research*, 2004,
480 32, 6643-6649.
- 481 38. E. J. Kerkhoven, P. J. Lahtvee and J. Nielsen, *FEMS yeast research*, 2014.
- 482 39. M. Li and I. Borodina, *FEMS yeast research*, 2014.

483

484

485 **Figure legends**

486 **Figure 1 – Network controllability and stability.** Concept of network controllability. (A) To
487 control all three nodes, it is enough to control node A in the left network ($n_D=1$). (B) two
488 driver nodes are needed to control the middle network ($n_D=2$). (C) Since node A, B and C are
489 connected in a circular control motif (CCM) we can choose any of the three nodes as driver
490 node and still have 100% control over the network. In terms of stability network (A) and (B)
491 will always be stable but network (C) can possibly be unstable.

492 **Figure 2 – Controllability when controlling one input node for simulated and real**
493 **networks.** The x-axis shows the average degree of the network and the y-axis shows the
494 controllability (relative amount of driver nodes needed to control 100% of the network).

495 **Figure 3 – Hierarchical structure of the transcriptional regulatory networks.** (A)
496 *S.cerevisiae* Harbison network. (B) *E.coli* RegulonDB network (C) Human ChEA v.2 network
497 (D) Mouse ChEA v.2 network

498 **Figure 4 – TFs responding to environment.** The x-axis in each plot shows the hyper-
499 geometric p-value for each TF based on the expression of the target genes. A low hyper-
500 geometric p-value indicates that the environmental factor has an effect on the expression of
501 the target genes. The y-axis in each plot shows the adjusted p-value of the TF gene, a high p-
502 value indicates that the gene is not changed in expression as a function of the environmental
503 change. The TFs marked with red have a hyper-geometric p-value for the target genes less
504 than 0.05 and a adjusted p-value for the TF gene greater than 0.8 (A) Oxygen availability. (B)
505 Nutrient limitation. (C) Dilution rate. (D) Carbon source.

506

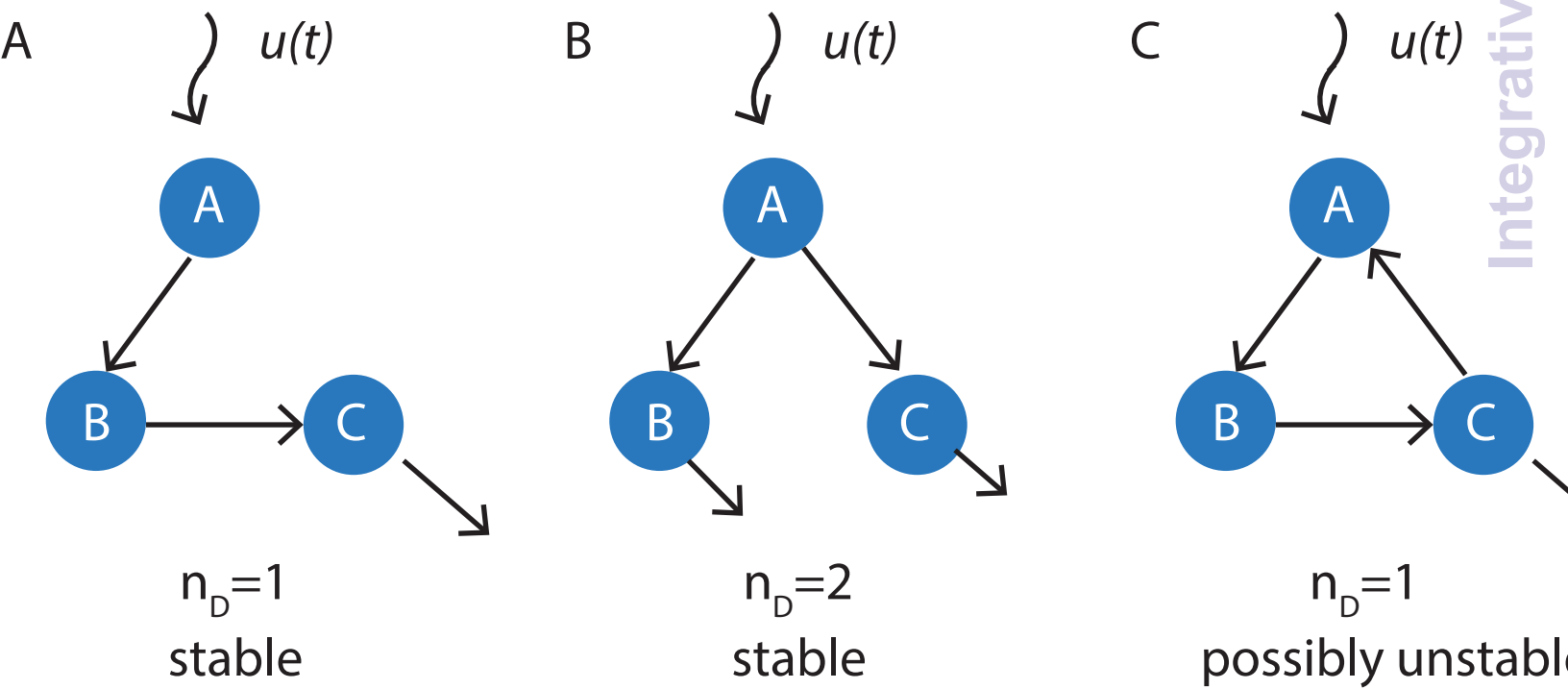
507

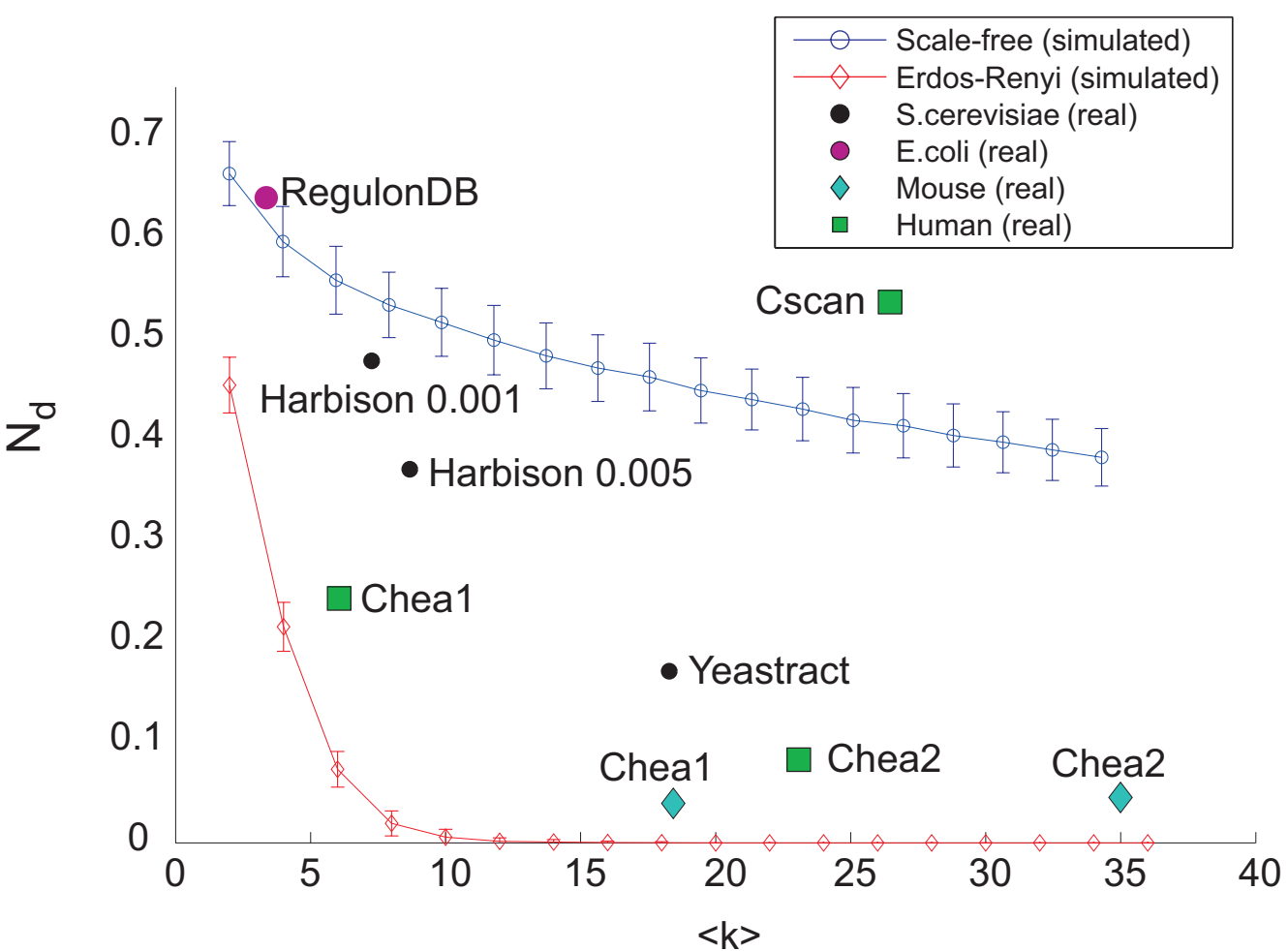
508

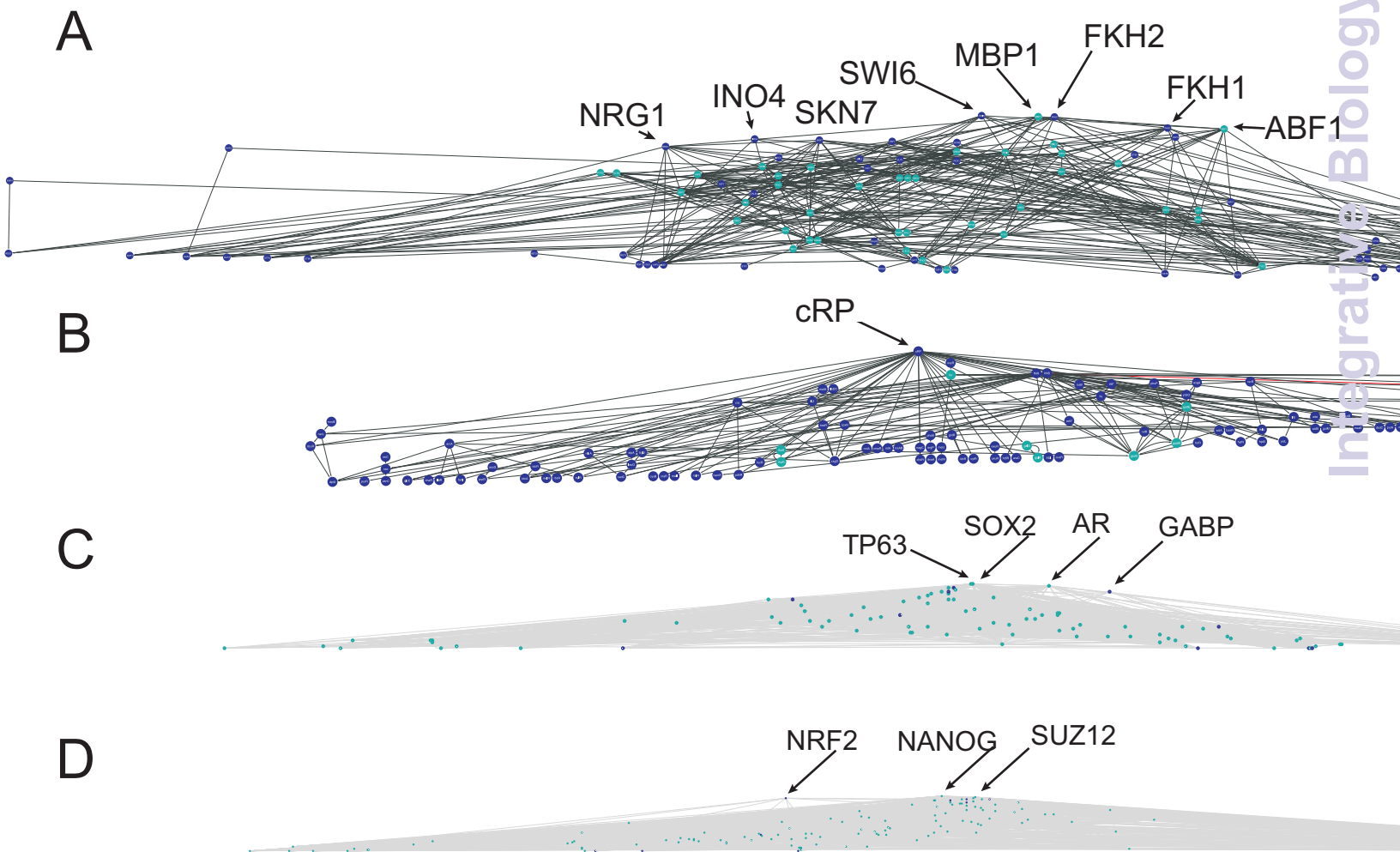
509

510

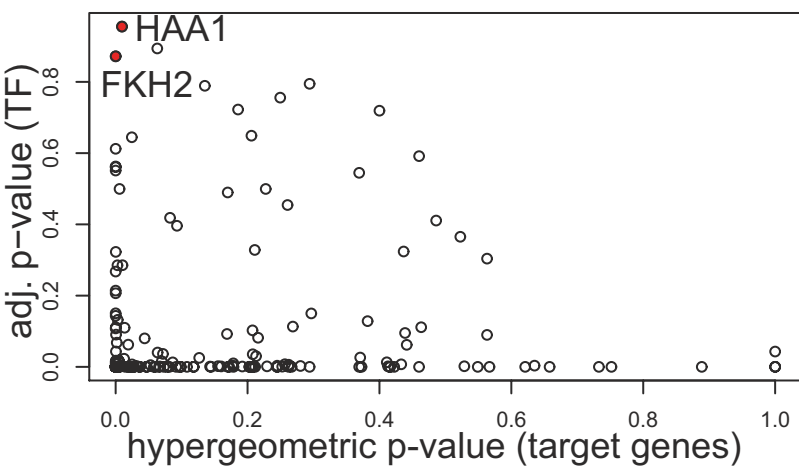
511



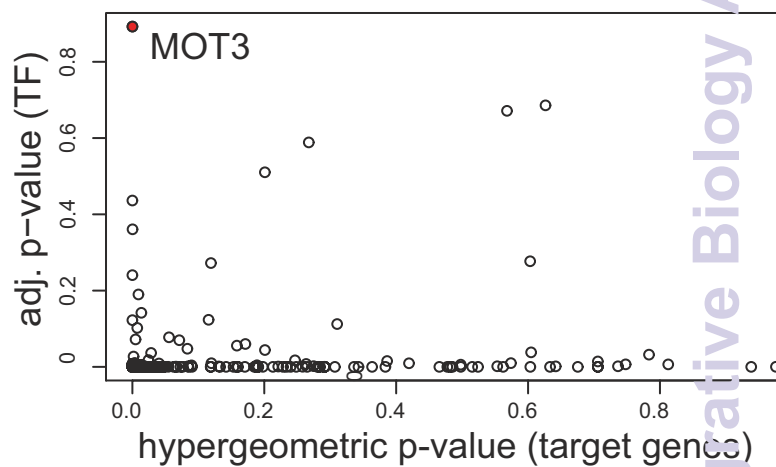




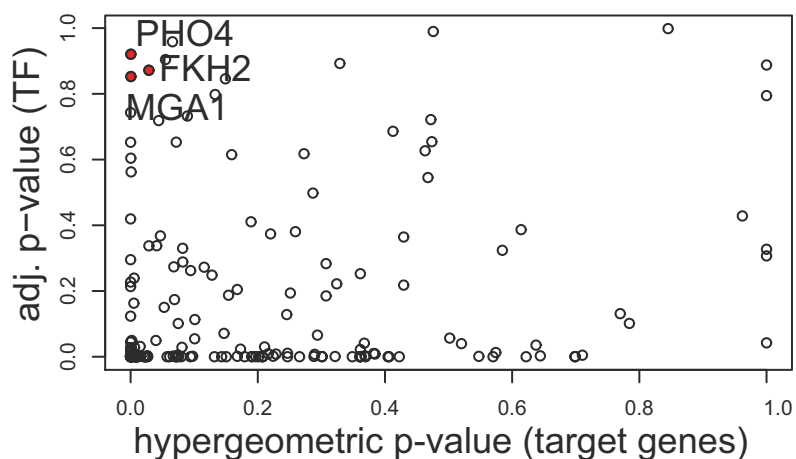
A



B



C



D

