

1 **Application of Modular Response Analysis to Medium- to Large-Size Biological**  
2 **Systems**

3

4 Meriem Mekedem<sup>1,2,3</sup>, Patrice Ravel<sup>1,2,3,4</sup>, Jacques Colinge<sup>1,2,3,5,\*</sup>

5 <sup>1</sup> Université de Montpellier, Montpellier, France

6 <sup>2</sup> Institut de Recherche en Cancérologie de Montpellier, Inserm U1194, Montpellier, France

7 <sup>3</sup> Institut régional du Cancer Montpellier, Montpellier, France

8 <sup>4</sup> Faculté de Pharmacie, Université de Montpellier, Montpellier, France

9 <sup>5</sup> Faculté de Médecine, Université de Montpellier, Montpellier, France

10

11 \*Corresponding author

12 E-mail : [jacques.colinge@inserm.fr](mailto:jacques.colinge@inserm.fr)

13

14 Short title: Application of MRA to large biological systems

15

16

## 17 **Abstract (300 words)**

18 The development of high-throughput genomic technologies associated with recent genetic perturbation  
19 techniques such as short hairpin RNA (shRNA), gene trapping, or gene editing (CRISPR/Cas9) has made it  
20 possible to obtain large perturbation data sets. These data sets are invaluable sources of information  
21 regarding the function of genes, and they offer unique opportunities to reverse engineer gene  
22 regulatory networks in specific cell types. Modular response analysis (MRA) is a well-accepted  
23 mathematical modeling method that is precisely aimed at such network inference tasks, but its use has  
24 been limited to rather small biological systems so far. In this study, we show that MRA can be employed  
25 on large systems with almost 1,000 network components. In particular, we show that MRA performance  
26 surpasses general-purpose mutual information-based algorithms. Part of these competitive results was  
27 obtained by the application of a novel heuristic that pruned MRA-inferred interactions *a posteriori*. We  
28 also exploited a block structure in MRA linear algebra to parallelize large system resolutions.

29

## 30 **Author Summary (150-200 words)**

31 The knowledge of gene and protein regulatory networks in specific cell types, including pathologic cells,  
32 is an important endeavor in the post-genomic era. A particular type of data obtained through the  
33 systematic perturbation of the actors of such networks enables the reconstruction of the latter and is  
34 becoming available at a large scale (networks comprised of almost 1,000 genes). In this work, we  
35 benchmark the performance of a classical methodology for such data called modular response analysis,  
36 which has been so far applied to networks of modest sizes. We also propose improvements to increase  
37 performance and to accelerate computations on large problems.

38

## 39 Introduction

40 The expression and activity of genes and proteins in cells are controlled by highly complex regulatory  
41 networks involving genes and proteins themselves, but also non-coding RNAs, metabolites, etc. Despite  
42 tremendous efforts in research, including all the developments of high-throughput genomic  
43 technologies, a significant portion of this machinery remains uncharted. Moreover, dysregulations in  
44 such networks are related to many diseases, and healthy cells of a same organism feature adjusted  
45 regulatory networks depending on their cell types and states. Techniques, both experimental and  
46 computational methodologies, that enable the inference of regulatory networks for different cells are  
47 obviously of great interest.

48 Reference databases such as Reactome[1], KEGG[2], IntAct[3], or STRING[4] that compile our knowledge  
49 of biological pathways or protein interactions have been established that provide valuable reference  
50 maps. Due to their universal nature, these maps do not reflect natural and pathologic variations of  
51 regulatory networks though some chosen disease pathways might be included. In principle, researchers  
52 should generate data specific to the biological system of interest to assess the actual wiring of its  
53 regulatory network. Specific data can be combined with reference databases in some algorithms, while  
54 others only rely on *de novo* inferences. The field of systems biology has proposed many algorithms for  
55 such a purpose involving different modeling approaches[5–7]. Obviously, algorithms must match the  
56 type of data available to perform the inference such as a transcriptomes or proteomes obtained under  
57 multiple conditions, time series, or perturbation data.

58 In this work, we are interested in the inference of regulatory networks based on systematic perturbation  
59 data. That is, given a biological system of interest, which could be the whole cell, but also a small set of  
60 related genes or proteins such as a pathway or part of a pathway, we have access to information  
61 reporting the activity level of every component (gene/protein). Typical examples are transcript, protein,  
62 or phosphorylated protein abundances. This information is available in basal condition as well as under  
63 the systematic perturbation of each single component. When this type of data are obtained from a  
64 biological system in a steady state, modular response analysis[8] (MRA) has been widely and successfully  
65 applied[9]. The elegance of MRA is that it provides an efficient mathematical framework to estimate a  
66 directed and weighted network representing the system regulatory network. Most applications of MRA  
67 are limited to networks comprised of a modest number of modules (<10). In this study, we want to  
68 explore the application of MRA to medium- (>50) and large-size (>500) systems. It entails a particular

69 implementation of the linear algebra at the heart of MRA to parallelize computations as well as the  
70 introduction of a heuristic to prune the inferred networks *a posteriori* to improve accuracy.

71 As stated above, rewiring of regulatory networks is natural and necessary to yield a multitude of cell  
72 types in higher organisms, and to adapt to distinct environmental conditions. Rewiring is also associated  
73 with several diseases[10,11], an extreme case being cancer[12–14]. For instance, kinase signaling  
74 cascades might be redirected in certain tumors to achieve drug resistance or to foster exaggerated cell  
75 growth. MRA has been applied to a number of such cancer-related investigations[15,16] considering  
76 rather small networks. Here, we take advantage of two published data sets that involve cancer cell lines  
77 and provide systematic perturbation data compatible with MRA requirements. The first – medium-size –  
78 data set[17] reports the transcriptional expression of 55 kinases and 6 non kinases under 11  
79 experimental conditions (unstimulated plus 10 distinct stimulations). Under every condition, the  
80 transcript levels of all the 61 genes were obtained by shallow RNA sequencing, including wild type cells  
81 and cells with individual KOs of each gene. These data hence enable us to infer one network *per*  
82 condition (11 networks) to discover how those 61 genes regulate themselves transcriptionally. The  
83 second – large-size – data set was generated by the next generation of the Connectivity Map (CMap)  
84 using its new L1000 platform[18]. Both shRNA- and CRISPR/Cas9-based systematic perturbations of  
85 roughly 1,000, respectively 350, genes in 9, respectively 5, cell lines were released. These data enable us  
86 to infer  $9+5=14$  networks.

87 We compare the performance of MRA, with and without the proposed pruning heuristic, to mutual  
88 information (MI)-based methods that have found broad acceptance.

89

## 90 Results

### 91 *Network inference algorithms*

92 The availability of large functional genomics data collections (transcriptomes and/or proteomes) has led  
93 to the development of a number of algorithms aimed at inferring interaction networks [7]. An essential  
94 ingredient of most algorithms is the co-expression of genes (or proteins)[19], which can be captured by  
95 simple correlation coefficients[20], mutual information (MI), or diverse statistical models[21]. There are  
96 too many such algorithms to review them all here, but MI-based approaches seem to have provided off-  
97 the-shelf, robust solutions that are widely used. We hence compare MRA to representatives of this  
98 category such as CLR[22], MRNET[23], and ARACNE[24].

99 MI is often preferred over correlation for its ability to detect nonlinear relationships. With a network  
100 involving  $n$  genes whose expression levels are measured in  $m$  transcriptomes, we write  $X_i$  the discrete  
101 distribution representing gene  $i$  expression. The MI between genes  $i$  and  $j$  is given by

$$102 \quad MI_{i,j} = H(X_i) + H(X_j) - H(X_i, X_j),$$

103 where  $H(X) = -\sum_{k \in X} p(x_k) \ln(p(x_k))$  is the entropy of a discrete random variable  $X$ . There exist  
104 different estimators for  $H(X)$  that use the  $m$  available transcriptomes[25]. Networks of interactions  
105 identified through MI, imposing a minimal threshold on MI values, are commonly called relevance  
106 networks[26,27]. The CLR algorithm improves over relevance networks by introducing a row- and  
107 column-wise z-score-like transformation of  $MI_{i,j}$  to normalize the MI matrix into a  $Z = (z_{i,j})$  matrix  
108 before thresholding. Namely, for each gene  $i$  CLR computes

$$z_i = \max \left\{ 0; \frac{MI_{i,j} - \text{mean}(MI_{i,\cdot})}{\text{sd}(MI_{i,\cdot})} \right\}$$

109 and then

$$110 \quad z_{i,j} = \sqrt{z_i^2 + z_j^2}.$$

111

112 MRNET applies a greedy maximum relevance strategy to link each gene  $i$  to the gene  $j$  that has  
113 maximum MI with it ( $j = \arg \max MI_{i,j}$ ). Additional links are added recursively maximizing MI with both

114 the gene  $i$  and the already linked genes until a stop criterion based on redundancy is met. A further  
115 approach by pruning was proposed by ARACNE authors, where as in relevance networks a common  
116 threshold is applied to all the  $M_{i,j}$  followed by the application of a pruning rule. This rule states that, if  
117 gene  $i$  interacts with gene  $j$  through gene  $k$ , then  $M_{i,j} \leq \min\{M_{i,k}; M_{k,j}\}$ . Consequently, among each  
118 triplet of nonzero MI after initial thresholding, the weakest interaction is removed.

119

## 120 ***The MRA and MRA+CLR algorithms***

121 Due to its ability to model biological systems at various resolutions, the MRA terminology for a system  
122 component is a module. We follow this terminology and consider that the  $n$  modules composing the  
123 system have their activity levels denoted by  $x \in \mathbb{R}^n$ . Here, modules are genes and  $x_i$  stands for gene  $i$   
124 transcript abundance. If we make the rather nonrestrictive assumption that relationships between  
125 modules are modeled by a dynamical system

$$\dot{x} = f(x)$$

126 ( $f(\cdot)$  must exist but it does need to be known), and the system is in a steady state at the time of  
127 experimental measurements ( $\dot{x} = 0$ ), MRA theory lets us compute an  $n \times n$  matrix of local interaction  
128 strengths  $r = (r_{i,j})$  from a gene  $j$  to a gene  $i$  ( $r_{i,j} = \frac{\partial x_i}{\partial x_j}$ ). The matrix  $r$  is obtained from linear  
129 algebraic computations based on the observed activity of each module in an unperturbed state, and  
130 under the individual, successive perturbations of each module. Details are provided in MRA original  
131 publication[8], reviews of MRA developments[9], or in our recent publication[15]. We use the notations  
132 of this recent paper. In Materials and Methods, we provide a brief overview of MRA along with a  
133 description of the particular way we implemented the linear algebra to take advantage of parallel  
134 computing.

135 Returning to the regulatory network inference problem, the MRA local interaction matrix  $r$  provides us  
136 with a direct estimate of this network. Interactions are signed with positive coefficients representing  
137 activation and negative coefficients representing inhibition. Given the fact that we want to apply MRA to  
138 large systems, where every module does not necessarily have a direct influence on all the others, we  
139 also face the problem of thresholding or pruning. Within the context of this study, we call MRA the  
140 direct use of MRA computations followed by a threshold on the absolute values of  $r$  coefficients (values  
141 below a given threshold in absolute values are set to 0). We also adapted CLR heuristic (z-score-like

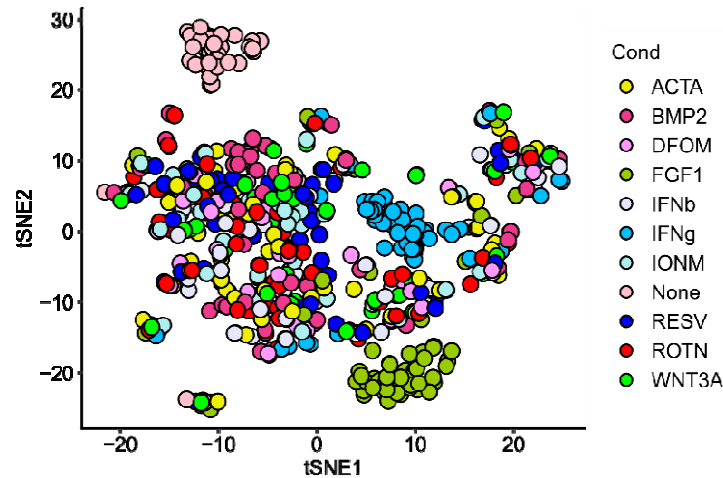
142 computation) to bring  $r$  coefficients to a more uniform scale before thresholding. We call this algorithm  
143 MRA+CLR, see Materials and Methods for details.

144

#### 145 ***Application to a medium-size data set***

146 Gapp *et al.*[17] published a data set, where they studied the transcriptional impact of the full knockouts  
147 (KOs) of 55 tyrosine kinases and 6 non-kinases. We call this data set K61. The systematic perturbations  
148 (KOs) of each gene as well as the unperturbed transcriptomes obviously constitute a *bona fide* MRA data  
149 set. The transcriptomes were acquired under 11 conditions: no stimulation (None), FGF1, ACTA, BMP2,  
150 IFN $\beta$ , IFN $\gamma$ , WNT3A, ionomycin (IONM), resveratrol (RESV), rotenone (ROTN), and deferoxamine (DFOM)  
151 stimulation. Stimulations were applied for 6 hours allowing the cells to adapt and reach a steady state or  
152 near steady state. To facilitate the generation of full-KOs, human HAP1 haploid cells[28] were utilized.  
153 The published transcriptomes were not limited to the expression of the 61 perturbed genes, but here,  
154 due to the specifics of MRA, we limited the data to those 61 genes. Replicates were essentially averaged  
155 (see Materials and Methods), resulting in a 61 $\times$ 61 matrix for each of the 11 conditions. Interestingly,  
156 considering the complete transcriptomes, K61 authors showed in their publication that those clustered  
157 primarily after the stimulatory condition. That is transcriptomes of different KOs obtained under the  
158 same stimulation were closer to each other than transcriptomes of the same KO but under different  
159 conditions. When reduced to the 61 genes of the network, this picture was less pronounced. In Fig. 1,  
160 we see that None-, WNT3A-, and to a certain extent IFN $\gamma$ -stimulated transcriptomes clustered  
161 separately thus potentially indicating rather different network wiring. The other conditions were not  
162 really separated suggesting that more similar networks could take place.

163



164

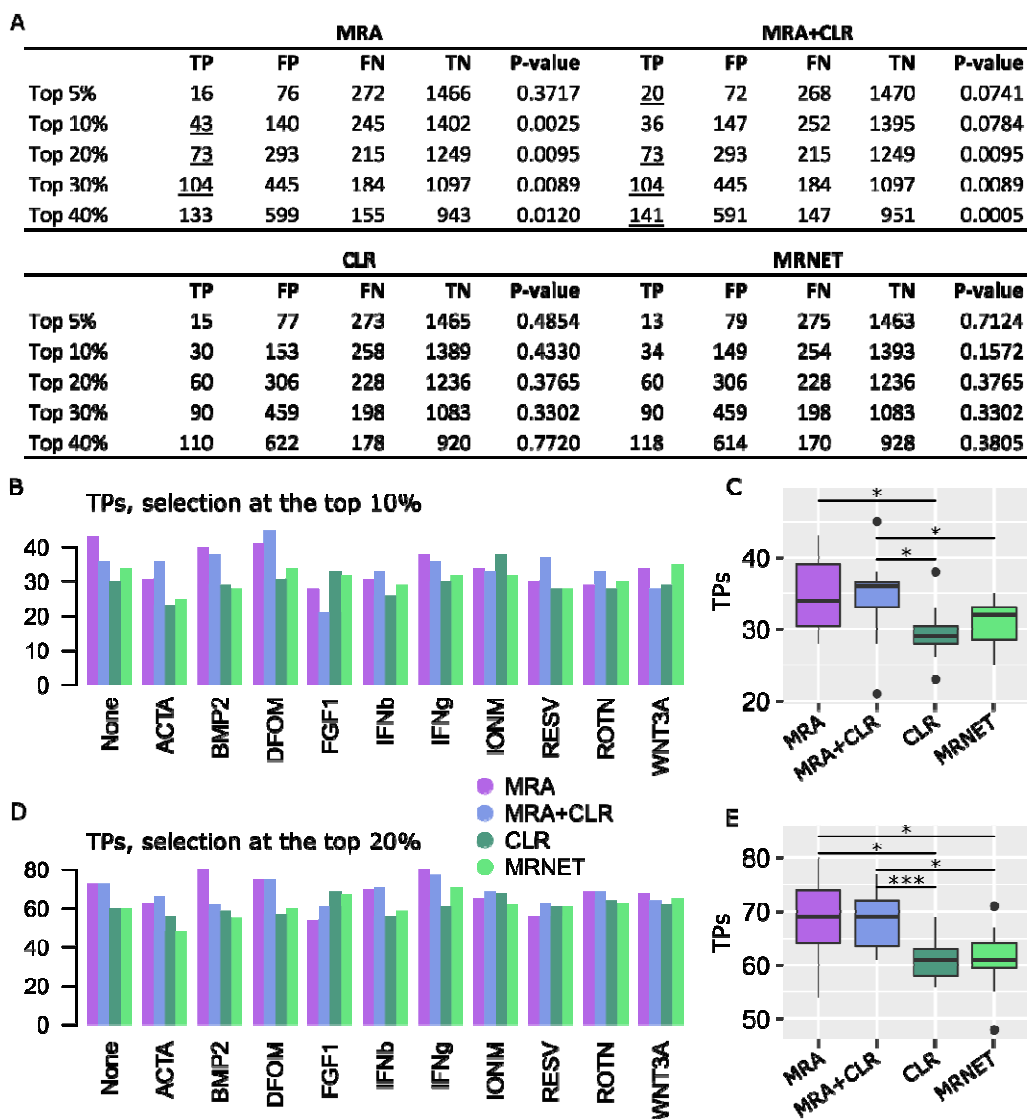
165 **Figure 1.** t-distributed stochastic neighbor embedding (t-SNE) 2D projection of the 61  
166 transcriptomes of the K61 data set.

167

168 We applied MRA, MRA+CLR, CLR, MRNET, and ARACNE to the K61 data set, the later 3 algorithms  
169 implementations were provided by the minet BioConductor package[25]. To estimate performance, we  
170 compared our results with the STRING database[4] due to its broad content. In fact, working with  
171 transcriptomic data, the inferred networks might overlap protein complexes as well as certain parts of  
172 known pathways, but they might also unravel different types of relationships such as genetic  
173 interactions, strong co-regulation, etc. Physical interaction of well-described pathway databases[1,3]  
174 might thus be too restrictive. To apply a uniform selection mechanism to all of the algorithms, we simply  
175 took the top 5%, 10%, 20%, 30% and 40% scores of the returned interaction matrices and determined  
176 the intersection with STRING. This resulted in confusion matrices reporting true/false positives (TPs/FPs)  
177 and true/false negatives (TNs/FNs) along with a P-value for the significance of the STRING intersection  
178 (hypergeometric test). A representative example (None condition) is featured in Fig. 2A, while the  
179 complete results are in Suppl. Table 1. Given the limited overlap between STRING and our data, and the  
180 rather large numbers involved in the confusion matrices, we found the P-values rather unstable (small  
181 differences in confusion matrices might cause important changes in terms of P-values). They should  
182 hence be regarded as indicative only. Because we used a constant reference (STRING), and all the  
183 algorithm scores were selected in identical numbers, reporting the number of TPs gives a clear  
184 indication of the relative algorithm performances. In Fig. 2B-E we provide these numbers at the top 10%  
185 and the top 20% selection levels. ARACNE implementation in minet did not perform well, typically  
186 reaching half of CLR or MRNET TPs. Accordingly, ARACNE performance is not reported in Fig. 2, but in



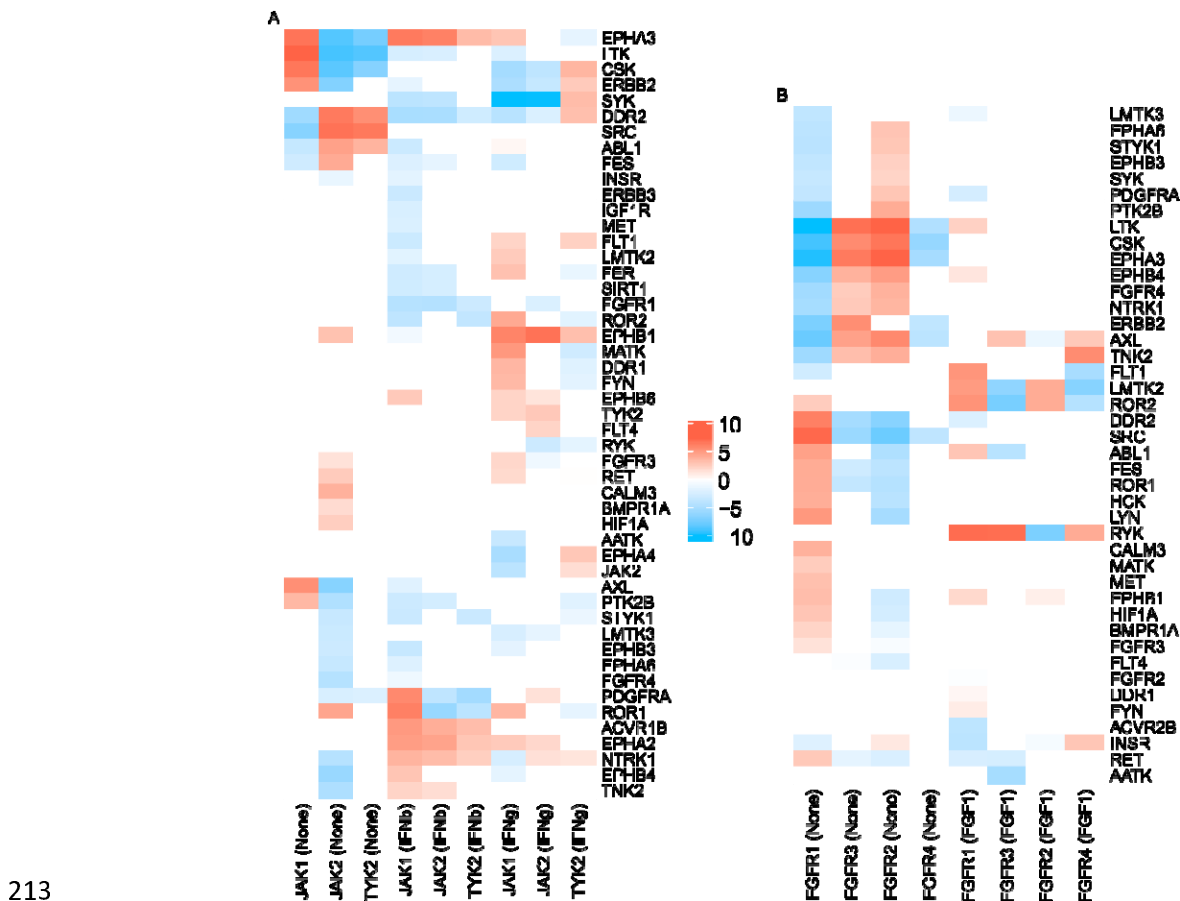
187 Suppl. Table 1 only. The CLR heuristic applied on top of MRA did not provide much performance  
 188 increase, but it resulted in more stable performances thus making it nonetheless an attractive option.



189  
 190 **Figure 2.** Performance on K61 data. **A.** Representative confusion matrices for the None condition. **B.** TP numbers at  
 191 the top 10% selection level. **C.** Comparison between the algorithm TP numbers (Wilcoxon test, 2-sided, \*P < 0.05).  
 192 **D.** TP numbers at the top 20% selection level. **E.** Comparison between the algorithm TP numbers (Wilcoxon test, 2-  
 193 sided, \*P < 0.05, \*\*\*P < 0.005).

194  
 195 In their article, K61 authors discussed interesting differences in JAK1 *versus* JAK2 and TYK2 signaling,  
 196 three members of the JAK family. In particular, they found that JAK1 KO cells were insensitive to IFNb  
 197 and IFNg stimulation, while JAK2 and TYR2 KO cells responded normally although, in general, all these

198 proteins are known to contribute to transcriptional response upon type I and II interferon stimuli[29]. To  
 199 illustrate how network inference might provide some clue on such differences, we report in Fig. 3A the  
 200 MRA+CLR-inferred transcriptional interaction strengths between those three genes and their targets  
 201 under the unstimulated (None), IFN $\beta$ , and IFN $\gamma$  conditions. In the absence of stimulation, we clearly  
 202 notice opposed influences of JAK1 on its targets compared to JAK2 and TYR2 (first three columns), which  
 203 already indicate different signal transduction capabilities. Upon IFN $\beta$  stimulation, the interactions are  
 204 closer with opposed action on ROR1 and PDGFRA. JAK2 and TYR2 remained highly similar in this  
 205 condition. IFN $\gamma$  stimulation induced three different patterns with ROR1 transcriptional inhibition  
 206 remaining a specific mark of JAK1. Gapp *et al.* also found differences in FGF receptors. FGF-induced  
 207 response was attenuated in FGFR1 and FGFR3 KO cells, but preserved in FGFR2 and FGFR4 KO cells. In  
 208 Figure 3B, we notice an almost perfect inversion of the activation/inhibition pattern between FGFR1  
 209 *versus* FGFR2 and FGFR3. FGFR4 adopted a very different configuration with limited interactions. This  
 210 observation already indicates a distinct role for FGFR1. Upon FGF stimulation, the interactions are more  
 211 patchy, but certain oppositions can be found such as a strong inhibitory action of FGFR1 and FGFR3 on  
 212 RYK transcription.



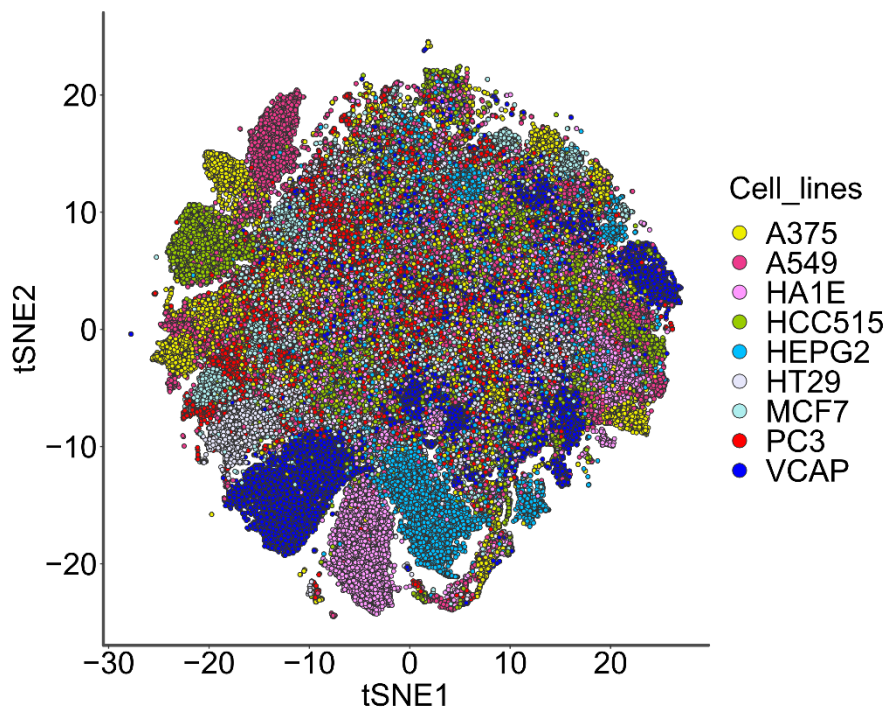
214 **Figure 3.** MRA+CLR-inferred interactions (top 20% selected). **A.** Interaction strengths (in  $\log_2$  with sign preserved)  
215 between JAK1, JAK2, and TYR2 and their targets. Stimulatory conditions are in brackets (None, IFN $\beta$ , IFN $\gamma$ ) **B.**  
216 Interaction strengths between FGFR1, FGFR2, FGFR3, and FGFR4 and their targets.

217

### 218 ***Application to a large-size data set***

219 CMap next generation platform L1000[18] has recently released (December 2020) a new batch of data.  
220 These data are in majority comprised of transcriptomes obtained in reference cancer cell lines under a  
221 large number of perturbations with chemical agents, but most importantly shRNA-induced knockdowns  
222 and CRISPR/Cas9 KOs. L1000 cost effective design entailed the identification of roughly 1,000 *hallmark*  
223 genes from which a large proportion of the whole transcriptome can be inferred. The L1000 platform  
224 only measures the expression of the hallmark genes experimentally. Two subsets of these data interest  
225 us.

226 A first data set is composed of the almost systematic shRNA perturbation of all the hallmark genes, thus  
227 providing an expression matrix close to 1,000 $\times$ 1,000 in size for 9 human cell lines: A375 (metastatic  
228 melanoma), A549 (lung adenocarcinoma), HCC515 (non-small cell lung cancer, adenocarcinoma), HT29  
229 (colorectal adenocarcinoma), HEPG2 (hepatocellular carcinoma), MCF7 (breast adenocarcinoma), PC3  
230 (metastatic prostate adenocarcinoma), VCAP (metastatic prostate cancer), and HA1E (normal kidney  
231 cells). To alleviate shRNA off-target effects, L1000 employed multiple hairpins, which were integrated  
232 into a consensus gene signature (CSG) that the authors showed to be essentially devoid of off-target  
233 consequences[18]. Cells were harvested 96 hours after shRNA perturbation leaving time to reach a  
234 steady state that is compatible to shRNA common use. Due to variation in data production, the actual  
235 matrix sizes ranged from 815 $\times$ 815 (MCF7) to 938 $\times$ 938 (A375). Interestingly, the t-SNE 2D projection of  
236 all the L1000 shRNA transcriptomes used here clearly indicate cell line specific subnetworks as well as  
237 shared, core parts (Fig. 4).

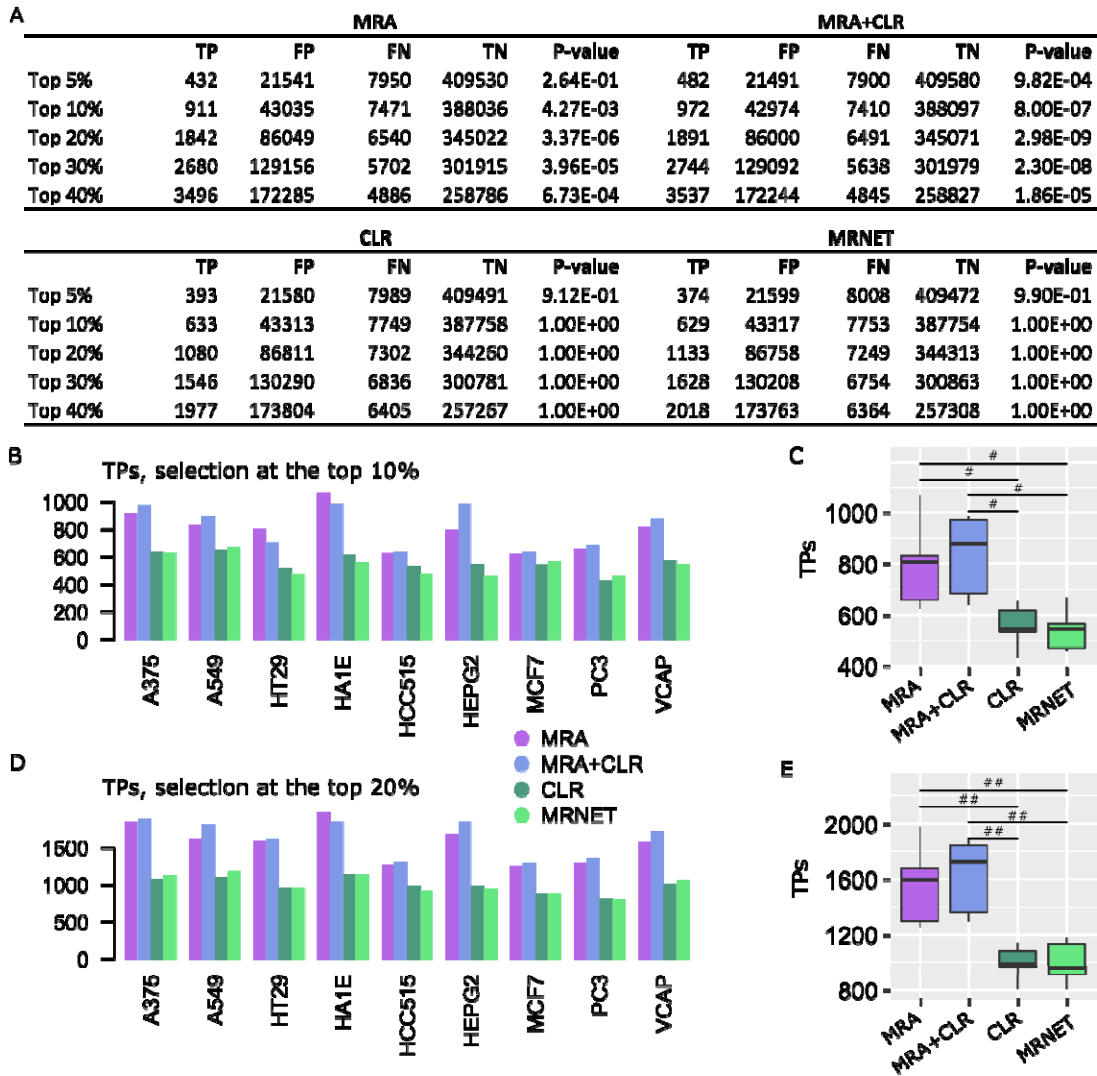


238

239 **Figure 4.** t-SNE projection of L1000 shRNA data. We note well-separated clusters that are specific to certain cell  
240 lines, *e.g.*, HA1E, VCAP, HCC515, HEPG2, A549, A375, as well as shared undistinguishable profile. This indicates  
241 potential common and specific subnetworks across the cell lines.

242 We followed the same performance evaluation procedure as above for K61. A representative (A375  
243 cells) confusion matrix is reported in Fig. 5A (full results in Suppl. Table 2), followed by TP numbers at  
244 the top 10% and top 20% selection levels in Fig. 5B-E. With these larger matrices, but also knockdown  
245 perturbations instead of KOs, MRA and MRA+CLR advantage was much augmented. Moreover, the CLR  
246 heuristic not only attenuated performance variability, but it almost systematically outperformed MRA  
247 alone.

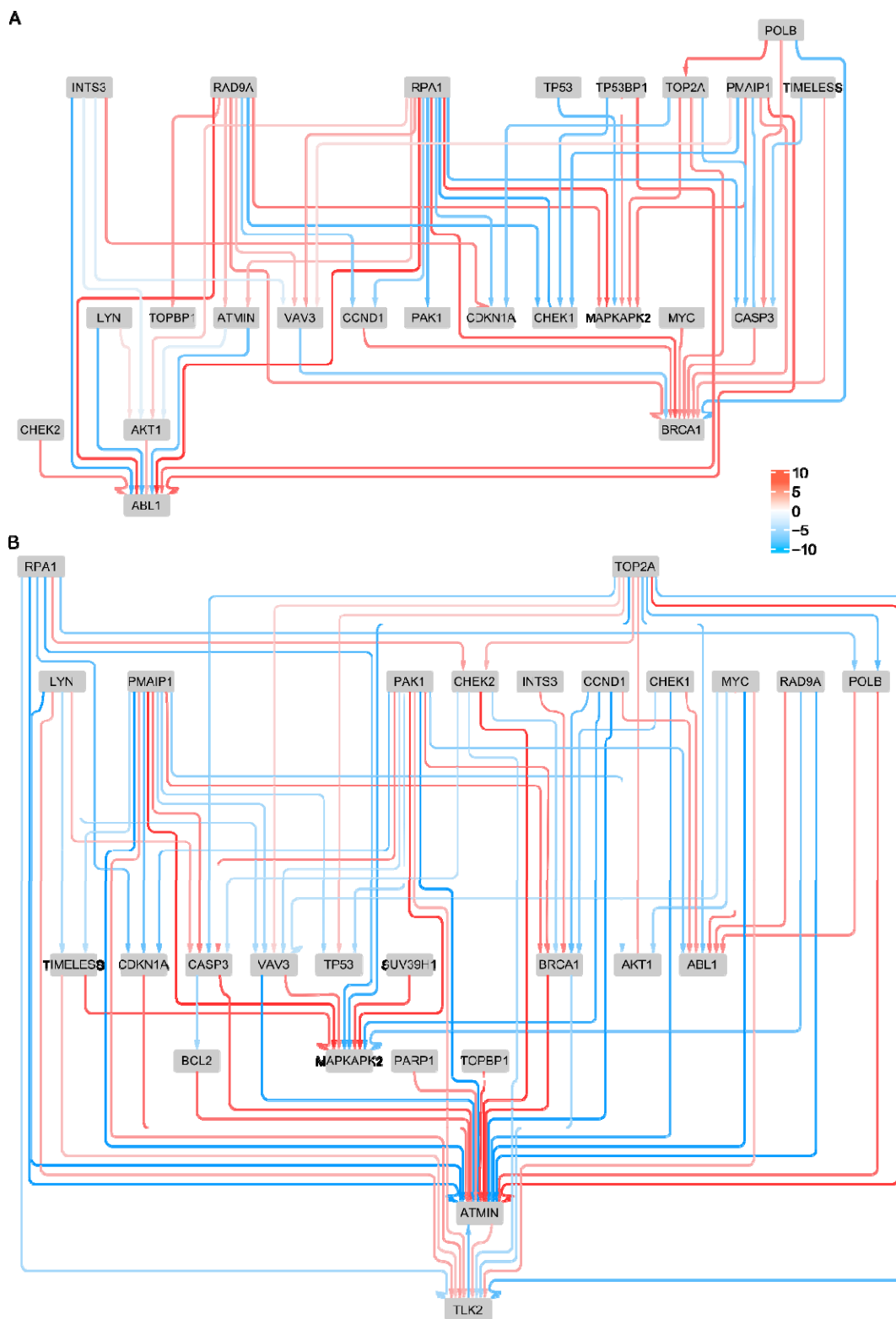
248 To illustrate the interest of network inference at this scale, we intersected MRA+CLR inferences in  
249 normal kidney HA1E and melanoma A375 cells with a Gene Ontology term, *i.e.*, GO:0006974 cellular  
250 response to DNA damage stimulus. In Fig. 6, we can notice the difference in connectivity between  
251 normal cells and cells where this process is obviously exacerbated, in particular the regulation of ATMIN  
252 a key molecule in DNA repair. This result is in agreement with the known rewiring of genetic networks in  
253 response to DNA damage[30].



254

255 **Figure 5.** Performance on L1000 shRNA data. **A.** Representative confusion matrices for A375 cells. **B.** TP numbers at  
 256 the top 10% selection level. **C.** Comparison between the algorithm TP numbers (Wilcoxon test, 2-sided, #P <  
 257 0.001). **D.** TP numbers at the top 20% selection level. **E.** Comparison between the algorithm TP numbers (Wilcoxon  
 258 test, 2-sided, #P < 0.001, ##P < 0.00005).

259



260

261 **Figure 6.** Networks inferred with MRA+CLR (top 10% selection) in normal kidney cells (A) and melanoma cells (B)  
 262 for genes involved in cellular response to DNA damage stimulus (GO:0006974).

263

264 The second L1000 data set of interest is the CRISPR/Cas9 collection of KOs. These data were only  
 265 available for five cell lines: A375, A549, HT29, MCF7, and PC3. The matrix sizes ranged from 343 343  
 266 (MCF7) to 359 359 (A375). Performance results are featured in Fig. 7 and Suppl. Table 3. Although MRA  
 267 and MRA+CLR again dominated the other algorithms, their advantage was less pronounced on these  
 268 large, full KO data.

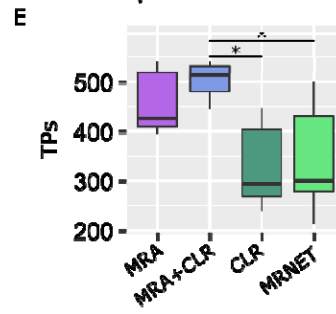
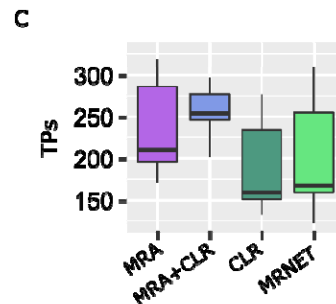
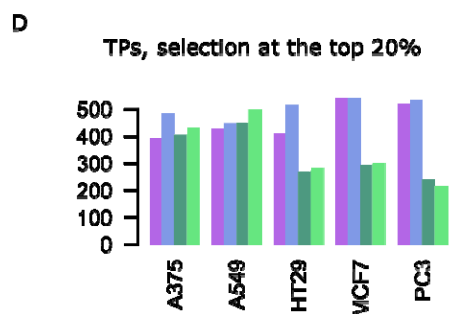
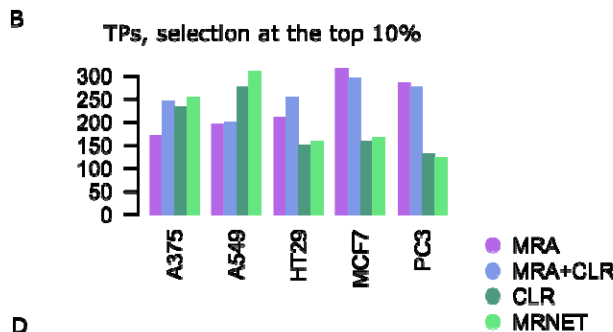
269

**A**

	MRA					MRA+CLR				
	TP	FP	FN	TN	P-value	TP	FP	FN	TN	P-value
Top 5%	76	3138	2463	58584	1.00E+00	114	3100	2425	58622	8.96E-01
Top 10%	171	6256	2368	55466	1.00E+00	246	6181	2293	55541	7.14E-01
Top 20%	394	12459	2145	49263	1.00E+00	<u>482</u>	12371	2057	49351	9.09E-01
Top 30%	632	18647	1907	43075	1.00E+00	<u>705</u>	18574	1834	43148	9.95E-01
Top 40%	868	24837	1671	36885	1.00E+00	<u>943</u>	24762	1596	36960	9.99E-01

	CLR					MRNET				
	TP	FP	FN	TN	P-value	TP	FP	FN	TN	P-value
Top 5%	131	3083	2408	58639	3.68E-01	<u>151</u>	3063	2388	58659	1.62E-02
Top 10%	234	6193	2305	55529	9.17E-01	<u>255</u>	6172	2284	55550	4.82E-01
Top 20%	404	12449	2135	49273	1.00E+00	431	12422	2108	49300	1.00E+00
Top 30%	574	18705	1965	43017	1.00E+00	581	18698	1958	43024	1.00E+00
Top 40%	727	24978	1812	36744	1.00E+00	723	24982	1816	36740	1.00E+00



270

271 Figure 7. Performance on L1000 CRISPR/Cas9 data. **A**. Representative confusion matrices for A375 cells. **B**. TP  
 272 numbers at the top 10% selection level. **C**. Comparison between the algorithm TP numbers. **D**. TP numbers at the  
 273 top 20% selection level. **E**. Comparison between the algorithm TP numbers (Wilcoxon test, 2-sided, \*P < 0.05).

274

275

## 276 **Discussion**

277 We presented a particular application of MRA to large biological systems and showed its competitive  
278 performance compared to first-in-class MI-based inference methods. Obviously, MI-based methods  
279 have a much broader spectrum of application, as they do not need specific and systematic perturbations  
280 on the components of the biological system whose network is inferred. Nevertheless, when  
281 perturbation data are available, our results suggest that a dedicated method, relying on a modeling  
282 approach might deliver good performance in a robust fashion. The simple heuristic we proposed to  
283 prune MRA inferences, which was adapted from the CLR algorithm, provided more stability in MRA  
284 performance. In many cases, especially with very large systems ( $n \approx 1,000$ ), this heuristic boosted  
285 performance.

286 Although the number of data sets was limited, we could notice much superior improvement over MI-  
287 based methods with L1000 shRNA knockdown perturbation data compared to the two full KO data sets.  
288 This might relate to the linearization at the heart of MRA modeling, where the error depends on the  
289 magnitude of perturbations (see our derivation of MRA through Taylor series expansion[15]). Very  
290 strong perturbation such as full KOs might bring the data away from MRA area of safe application.

291



## 292 Materials and Methods

293

### 294 *Modular response analysis*

295 We briefly recall the main MRA equations to facilitate the reading of this text, and to explain the  
 296 particular way we implemented the linear algebra. We assume that the biological system is comprised of  
 297  $n$  modules whose activity levels are denoted by  $x \in \mathbb{R}^n$ . We further admit the existence of  $n$  intrinsic  
 298 parameters,  $p \in \mathbb{R}^n$ , one per module, and each of them can be perturbed by an elementary  
 299 perturbation. One can imagine  $x$  reporting mRNA abundances and perturbations induced by shRNAs for  
 300 instance. Lastly, we assume that there exist  $S \subset \mathbb{R}^n \times \mathbb{R}^n$ , an open subset, and  $f: S \rightarrow \mathbb{R}^n$  of class  $\mathcal{C}^1$ ,  
 301 *i.e.*, continuously differentiable, such that

$$302 \quad \dot{x} = f(x, p). \quad (1)$$

303 We do not need to know  $f(x, p) = (f_1(x, p), \dots, f_n(x, p))^t$  explicitly, but we need the existence of a  
 304 time  $T > 0$  such that all the solutions, for any  $p$  and initial conditions of  $x$ , have reached a steady state,  
 305 *i.e.*,

$$\dot{x} = 0, \forall t > T.$$

306 The unperturbed, basal state of the modules is denoted  $x(p^0) \in \mathbb{R}^n$  and it has corresponding  
 307 parameters  $p^0 \in \mathbb{R}^n$ . By the application of the implicit function theorem and Taylor expansion at the  
 308 first order [8,15], MRA relates the experimental observations of the global effect of perturbations to  
 309 local interaction strengths, *i.e.*, the matrix  $r = (r_{i,j}) = \left( \frac{\partial x_i}{\partial x_j} \frac{x_j}{x_i} \right)$  that we mentioned in Results. Such local  
 310 interactions are obviously signed and non-symmetric. To compute  $r$ , we need to compute the relative  
 311 global change induced by each elementary perturbation in each module. These values are compiled in a  
 312  $n \times n$  matrix denoted  $R = (R_{i,k})$  with

$$313 \quad R_{i,k} = \left( \frac{\Delta x_i}{x_i} \right)_{q_k},$$

314 the relative difference in activity of module  $i$  upon  $\Delta p_k$  change induced by an elementary perturbation  
 315  $q_k$  that touches module  $k$  only. The relationship between observational data in  $R$  and the local  
 316 interactions we want to estimate in  $r$  are provided by the following equations

$$317 \quad \left( \frac{\Delta x_i}{x_i} \right)_{q_k} = \sum_{j \neq i} r_{i,j} \left( \frac{\Delta x_j}{x_j} \right)_{q_k}, \quad k \neq i, \quad (2)$$

$$318 \quad \left( \frac{\Delta x_i}{x_i} \right)_{q_i} = \sum_{j \neq i} r_{i,j} \left( \frac{\Delta x_j}{x_j} \right)_{q_i} + \frac{\partial x_i}{\partial p_i}(p^0) \left( \frac{\Delta p_i}{x_i} \right). \quad (3)$$

319 By setting  $r_{i,i} = -1$ , Eqs (2) and (3) can be put together in matrix form and we obtain

$$320 \quad rR = -P, \quad (4)$$

321 where  $P$  is a diagonal  $n \times n$  matrix with

$$322 \quad P_{i,i} = \frac{\partial x_i}{\partial p_i}(p^0) \left( \frac{\Delta p_i}{x_i} \right), \quad i \in \{1, \dots, n\}. \quad (5)$$

323 Eq. (3) can be solved in two steps:  $r = -PR^{-1}$  and  $r_{i,i} = -1$  imply  $P_{i,i}(R^{-1})_{i,i} = 1$ , thus

$$324 \quad P_{i,i} = \frac{1}{(R^{-1})_{i,i}}.$$

325 Therefore,

$$326 \quad r = -[\text{diag}(R^{-1})]^{-1}R^{-1}. \quad (6)$$

327 In practice, relative differences in  $R$  are often estimated with the more stable formula

$$328 \quad R_{i,k} = 2 \left( \frac{x_i(p^0 + \Delta p_k) - x_i(p^0)}{x_i(p^0 + \Delta p_k) + x_i(p^0)} \right), \quad (7)$$

329 where we denote  $x(p^0 + \Delta p)$  the steady-state corresponding to the changed parameters  $p^0 + \Delta p$ , i.e.,  
 330 the solution of  $\dot{x}(p^0 + \Delta p) = f(x(p^0 + \Delta p), p^0 + \Delta p)$ .

331

### 332 **Parallelized and stable linear algebra**

333 Eq. (6) requires the computation of the inverse of the matrix  $R$ , which is less efficient and less stable  
 334 than LU decomposition with pivot search[31]. These technical issues are usually irrelevant with small  
 335 systems, but in applications of MRA to larger biological systems they should be addressed.

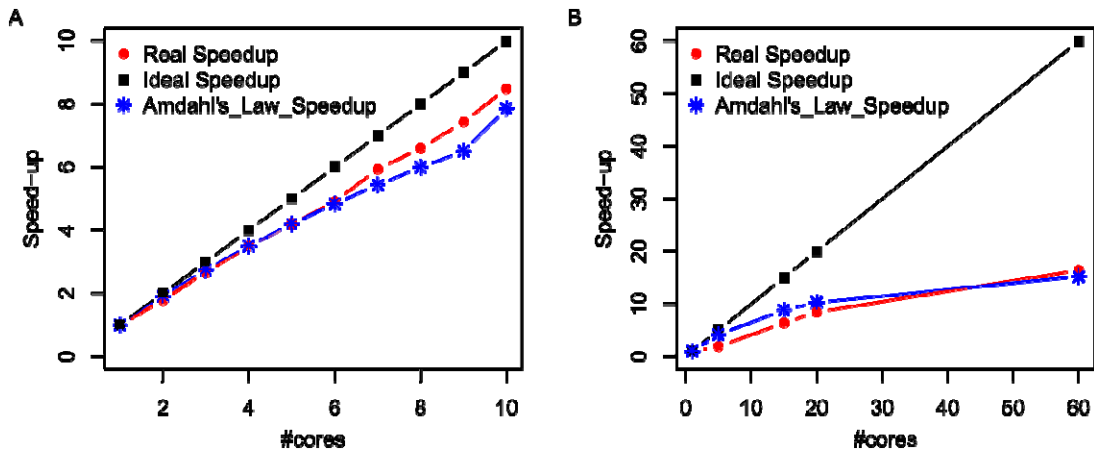
336 As several authors noticed, including in MRA original publication[8], the homogeneous Eq. (2) is  
 337 sufficient to compute  $r$ . Moreover, letting  $i$  take the values  $1, \dots, n$ , we remark that Eq. (2) defines  $n$   
 338 systems of linear equations of dimension  $n - 1$ , which can be solved independently. In particular, those  
 339 systems can be solved on independent processors by performing the LU decomposition with pivot  
 340 search. Illustrative speedup curves are featured in Fig. 8. Depending on the size of  $n$ , each such  
 341 subsystem could itself benefit from a parallel solver if enough processors were available.

342 When Eq. (2) is solved for each value of  $i$ , it is straightforward to solve Eq. (3) to find  $P_{i,i}$  values in case  
 343 those are required:

$$344 \quad \left( \frac{\Delta x_i}{x_i} \right)_{q_i} = \sum_{j \neq i} r_{i,j} \left( \frac{\Delta x_j}{x_j} \right)_{q_i} + P_{i,i} \Leftrightarrow P_{i,i} = \sum_{j \neq i} r_{i,j} \left( \frac{\Delta x_j}{x_j} \right)_{q_i} - \left( \frac{\Delta x_i}{x_i} \right)_{q_i},$$

345 where Eq. (4) was used for the definition of  $S_{ideal}$ .

346



347

348 **Figure 8.** Speedup curves. **A.** K61 data (None condition, 61 × 61 matrix). **B.** L1000 shRNA data (A375 cells, 938 × 938  
 349 matrix).

350

### 351 ***CLR, MRNET, and ARACNE computations***

352 We used the implementation of these algorithms provided by the BioConductor R package minet[25].

353 The performance reported here reflects the performance of this specific implementation.

354

### 355 ***CLR heuristic adapted to MRA***

356 We adapted the CLR normalization scheme by means of z-score computation to MRA matrix content.

357 From  $\mu_{i,j}$  we thus derive a  $\mu_{i,j}^z$  defined as follow:

358 
$$\mu_{i,j}^z = \frac{\mu_{i,j} - \mu_{i,\cdot}}{\sigma_{i,\cdot}},$$
 with  $\sigma_{i,\cdot}$  the standard deviation of  $\mu_{i,\cdot}$ 's  $i$ -th row,

359 
$$\mu_{i,j}^z = \frac{\mu_{i,j} - \mu_{\cdot,j}}{\sigma_{\cdot,j}},$$
 with  $\sigma_{\cdot,j}$  the standard deviation of  $\mu_{\cdot,j}$ 's  $j$ -th column,

360 
$$\mu_{i,j}^z = \frac{\mu_{i,j} - \mu_{\cdot,\cdot}}{\sigma_{\cdot,\cdot}},$$
 and

361

362

### 363 ***Data sets preparation***

364 TK61 data were obtained on multiple 96-well plates. Accordingly, we tried to stick to this format  
365 preparing data for MRA computations. We computed an  $R$  matrix for each plate and then simply  
366 averaged the relevant  $R$ 's for each experimental condition to obtain the averaged  $R$  used in MRA. For  
367 MI-based inferences, we averaged all the relevant values.

368 L1000 shRNA data were extracted at level 5 (L1000 terminology) where CGSs (integration of multiple  
369 shRNA hairpins to alleviate off-target effects) were transformed into z-scores for normalization purposes  
370 by the authors of the data. Consequently, values representing the abundance of a gene were no longer  
371 positive numbers but just real numbers. Eq. (7) above was adapted to compute the relative changes in  
372 MRA  $R$  matrices according to

$$R_{i,k} = 2 \left( \frac{\text{CGS}_i(p^0 + \Delta p_k) - \text{CGS}_i(p^0)}{|\text{CGS}_i(p^0 + \Delta p_k)| + |\text{CGS}_i(p^0)|} \right)$$

373 avoiding potential divisions by 0 in case of small values with opposed signs.

374 L1000 CRISPR/Cas9 data were averaged over replicates (also level 5).

375

### 376 ***Performance evaluation***

377 STRING as well as MI-based inference are devoid of direction of interaction and a sign. Therefore, the  
378 intersection of inferences with STRING content only used the upper triangular part of matrices  
379 representing the inferences (such matrices are symmetric anyway). To provide a fair comparison with  
380 MRA and MRA+CLR, we filled the upper triangular part of  $r$  according to  $r_{i,j} = \max\{|r_{i,j}|; |r_{j,i}|\}$ ,  $i < j$ .

381

### 382 **Acknowledgements**

383 MM was supported by a PhD fellowship of the Algerian government.

384

385

## 386 References

- 387 1. Matthews L, Gopinath G, Gillespie M, Caudy M, Croft D, de Bono B, et al. Reactome knowledgebase  
388 of human biological pathways and processes. *Nucleic Acids Res.* 2009;37: D619-22.  
389 doi:10.1093/nar/gkn863
- 390 2. Kanehisa M, Araki M, Goto S, Hattori M, Hirakawa M, Itoh M, et al. KEGG for linking genomes to life  
391 and the environment. *Nucleic Acids Res.* 2008;36: D480-4.
- 392 3. Kerrien S, Alam-Faruque Y, Aranda B, Bancarz I, Bridge A, Derow C, et al. IntAct--open source  
393 resource for molecular interaction data. *Nucleic Acids Res.* 2007;35: D561-5.  
394 doi:10.1093/nar/gkl958
- 395 4. Szklarczyk D, Franceschini A, Kuhn M, Simonovic M, Roth A, Minguéz P, et al. The STRING database  
396 in 2011: functional interaction networks of proteins, globally integrated and scored. *Nucleic Acids*  
397 *Res.* 2011;39: D561-8. doi:10.1093/nar/gkq973
- 398 5. Bansal M, Belcastro V, Ambesi-Impiombato A, di Bernardo D. How to infer gene networks from  
399 expression profiles. *Molecular Systems Biology.* 2007;3: 78. doi:10.1038/msb4100120
- 400 6. Babbie AC, Stumpf MPH, Thorne T. Gene Regulatory Network Inference. In: Wolkenhauer O, editor.  
401 *Systems Medicine.* Oxford: Academic Press; 2021. pp. 86–95. doi:10.1016/B978-0-12-801238-  
402 3.11346-7
- 403 7. Emmert-Streib F, Dehmer M, Haibe-Kains B. Gene regulatory networks and their applications:  
404 understanding biological and medical problems in terms of networks. *Front Cell Dev Biol.* 2014;2.  
405 doi:10.3389/fcell.2014.00038
- 406 8. Kholodenko BN, Kiyatkin A, Bruggeman FJ, Sontag E, Westerhoff HV, Hoek JB. Untangling the wires:  
407 a strategy to trace functional interactions in signaling and gene networks. *Proc Natl Acad Sci U S A.*  
408 2002;99: 12841–6. doi:10.1073/pnas.192442699
- 409 9. Santra T, Rukhlenko O, Zhernovkov V, Kholodenko BN. Reconstructing static and dynamic models of  
410 signaling pathways using Modular Response Analysis. *Current Opinion in Systems Biology.* 2018;9:  
411 11–21. doi:10.1016/j.coisb.2018.02.003
- 412 10. Hu JX, Thomas CE, Brunak S. Network biology concepts in complex disease comorbidities. *Nat Rev*  
413 *Genet.* 2016;17: 615–629. doi:10.1038/nrg.2016.87
- 414 11. Huttlin EL, Bruckner RJ, Paulo JA, Cannon JR, Ting L, Baltier K, et al. Architecture of the human  
415 interactome defines protein communities and disease networks. *Nature.* 2017;545: 505–509.  
416 doi:10.1038/nature22366
- 417 12. Assi SA, Imperato MR, Coleman DJL, Pickin A, Potluri S, Ptasinska A, et al. Subtype-specific regulatory  
418 network rewiring in acute myeloid leukemia. *Nat Genet.* 2019;51: 151–162. doi:10.1038/s41588-  
419 018-0270-1
- 420 13. Pawson T, Warner N. Oncogenic re-wiring of cellular signaling pathways. *Oncogene.* 2007;26: 1268–  
421 1275. doi:10.1038/sj.onc.1210255

- 422 14. Weinstein IB, Joe A. Oncogene addiction. *Cancer Res.* 2008;68: 3077–3080; discussion 3080.  
423 doi:10.1158/0008-5472.CAN-07-3293
- 424 15. Jimenez-Dominguez G, Ravel P, Jalaguier S, Cavallès V, Colinge J. An R package for generic modular  
425 response analysis and its application to estrogen and retinoic acid receptor crosstalk. *Sci Rep.*  
426 2021;11: 7272. doi:10.1038/s41598-021-86544-0
- 427 16. Klinger B, Sieber A, Fritsche-Guenther R, Witzel F, Berry L, Schumacher D, et al. Network  
428 quantification of EGFR signaling unveils potential for targeted combination therapy. *Mol Syst Biol.*  
429 2013;9: 673. doi:10.1038/msb.2013.29
- 430 17. Gapp BV, Konopka T, Penz T, Dalal V, Bürckstümmer T, Bock C, et al. Parallel reverse genetic  
431 screening in mutant human cells using transcriptomics. *Molecular Systems Biology.* 2016;12: 879.  
432 doi:10.15252/msb.20166890
- 433 18. Subramanian A, Narayan R, Corsello SM, Peck DD, Natoli TE, Lu X, et al. A Next Generation  
434 Connectivity Map: L1000 Platform and the First 1,000,000 Profiles. *Cell.* 2017;171: 1437-1452.e17.  
435 doi:10.1016/j.cell.2017.10.049
- 436 19. Horvath S, Dong J. Geometric interpretation of gene coexpression network analysis. *PLoS Comput*  
437 *Biol.* 2008;4: e1000117. doi:10.1371/journal.pcbi.1000117
- 438 20. Obayashi T, Hayashi S, Shibaoka M, Saeki M, Ohta H, Kinoshita K. COXPRESdb: a database of  
439 coexpressed gene networks in mammals. *Nucleic Acids Res.* 2008;36: D77-82.  
440 doi:10.1093/nar/gkm840
- 441 21. Wang YXR, Huang H. Review on statistical methods for gene network reconstruction using  
442 expression data. *J Theor Biol.* 2014;362: 53–61. doi:10.1016/j.jtbi.2014.03.040
- 443 22. Faith JJ, Hayete B, Thaden JT, Mogno I, Wierzbowski J, Cottarel G, et al. Large-scale mapping and  
444 validation of *Escherichia coli* transcriptional regulation from a compendium of expression profiles.  
445 *PLoS Biol.* 2007;5: e8. doi:10.1371/journal.pbio.0050008
- 446 23. Meyer PE, Kontos K, Lafitte F, Bontempi G. Information-theoretic inference of large transcriptional  
447 regulatory networks. *EURASIP J Bioinform Syst Biol.* 2007; 79879. doi:10.1155/2007/79879
- 448 24. Margolin AA, Nemenman I, Basso K, Wiggins C, Stolovitzky G, Favera RD, et al. ARACNE: An  
449 Algorithm for the Reconstruction of Gene Regulatory Networks in a Mammalian Cellular Context.  
450 *BMC Bioinformatics.* 2006;7: S7. doi:10.1186/1471-2105-7-S1-S7
- 451 25. Meyer PE, Lafitte F, Bontempi G. minet: A R/Bioconductor Package for Inferring Large  
452 Transcriptional Networks Using Mutual Information. *BMC Bioinformatics.* 2008;9: 461.  
453 doi:10.1186/1471-2105-9-461
- 454 26. Butte AJ, Kohane IS. Mutual information relevance networks: functional genomic clustering using  
455 pairwise entropy measurements. *Pac Symp Biocomput.* 2000; 418–429.  
456 doi:10.1142/9789814447331\_0040

- 457 27. Butte AJ, Tamayo P, Slonim D, Golub TR, Kohane IS. Discovering functional relationships between  
458 RNA expression and chemotherapeutic susceptibility using relevance networks. PNAS. 2000;97:  
459 12182–12186.
- 460 28. Carette JE, Guimaraes CP, Varadarajan M, Park AS, Wuethrich I, Godarova A, et al. Haploid genetic  
461 screens in human cells identify host factors used by pathogens. Science. 2009;326: 1231–1235.  
462 doi:10.1126/science.1178955
- 463 29. Rane SG, Reddy EP. Janus kinases: components of multiple signaling pathways. Oncogene. 2000;19:  
464 5662–5679. doi:10.1038/sj.onc.1203925
- 465 30. Bandyopadhyay S, Mehta M, Kuo D, Sung M-K, Chuang R, Jaehnig EJ, et al. Rewiring of genetic  
466 networks in response to DNA damage. Science. 2010;330: 1385–1389. doi:10.1126/science.1195618
- 467 31. Golub GH, Loan CFV. Matrix Computations. JHU Press; 2013.

468

469

## 470 **Data availability**

471 Data used in this work were made publicly available by their respective authors.

472

## 473 **Supporting information caption**

474 **Supplementary Table 1.** Confusion matrices on the K61 data set.

475 **Supplementary Table 2.** Confusion matrices on the L1000 shRNA data set.

476 **Supplementary Table 3.** Confusion matrices on the L1000 CRISPR/Cas9 data set.

477