

Supporting Information for DiME: A scalable disease module identification algorithm with application to glioma progression

Yunpeng Liu¹, Daniel A. Tennant², Zexuan Zhu⁴, John K. Heath³, Xin Yao¹, Shan He^{1,3*}

1 School of Computer Science, University of Birmingham, Birmingham, UK

2 School of Cancer Sciences, University of Birmingham, Birmingham, UK

3 Centre for Systems Biology, School of Biological Sciences, University of Birmingham, Birmingham, UK

4 College of Computer Science and Software Engineering, Shenzhen University, China

*** E-mail: s.he@cs.bham.ac.uk**

The B-score Algorithm Pseudo-code

Algorithm 1 The B-score Algorithm

```

1: function BSCORE(module  $C$ )
2:    $t \leftarrow 0, C_0 \leftarrow C, B_0 \leftarrow \emptyset, n_C \leftarrow |C|$ 
3:   while  $t < n_C - 1$  do
4:     for each  $i \in C_t$  do
5:       Calculate  $p_i$  as  $p_i = \sum_{q=k_i^{int}}^{k_i} f(q|C)$ 
6:     end for
7:      $w_{t+1} \leftarrow$  vertex in  $C_t$  with highest  $p_i$ 
8:      $w_{t+2} \leftarrow$  vertex in  $C_t$  with second highest  $p_i$ 
9:      $B_{t+1} \leftarrow B_t \cup \{w_{t+1}\}, C_{t+1} \leftarrow C_t \setminus \{w_{t+1}\}$ 
10:    Recalculate  $p$  values for all nodes currently in  $B_{t+1}$ 
11:     $p_l \leftarrow$  lowest  $p$  value of vertices in  $B_{t+1}$ 
12:    if  $p_{w_{t+2}} > p_l$  then
13:      swap( $p_{w_{t+2}}, p_l$ )
14:    end if
15:    Compute  $Pr(< S_{t+1} | C_{t+1}, B_{t+1}, p_{w_{t+2}})$ , where  $S_{t+1} = \sum_{i \in B_{t+1}} p_i$ 
16:     $t \leftarrow t + 1$ 
17:  end while
18:  Return  $\min_t Pr(< S_t | C_t, B_t, p_{w_{t+1}})$ 
19: end function

```

The B-score measure assumes a null model where edges within the module (community) of interest is held unchanged while the remaining connections in the network are randomly shuffled. A probabilistic measure based on hypergeometric distribution is then calculated for each module member to evaluate the likelihood that the observed number of within-module edges would arise from the null model. Such a probability is then summed over the number of possible within-module connections (from the observed value to the maximum possible value - the total degree of the node) to give a cumulative probability p_i of observing an intra-module degree equal to or larger than the observed value under the null model.

The above p_i is calculated for all nodes in the module and sorted to identify the “worst” node in the module - the node with the highest p_i . The B-score algorithm assumes that, for a truly non-random module, the probability of observing such a worst p_i value as the *minimum* among all nodes currently not belonging to the module is expected to be very low under the null model for its calculation. The original B-score algorithm also incorporated a stochastic element into the calculation of p_i , took into consideration a list of k worst nodes in the module and utilized the probability that the sum of the scores of these worst nodes in a module obtained from a random background model is smaller than the observed value as the final statistical significance measure (the B-score) for the module. The B-score is calculated over multiple runs and the average is used for evaluation of module statistical significance. Such an additional step has been shown to act as a resampling step and guard against possible significant community structure from random graphs.

Calculation of the conservation score

We first define a reference network which can be seen as the ground true network without noise. We extract the modules from the network as reference modules: R_i ($i = 1, \dots, n_r$). We also define a set of m noisy co-expression networks by introducing m different levels of edge noise to the reference network.

We extract the modules N_{i_j} ($i = 1, \dots, n_m$) from the j th noisy network and repeat this for all m noisy networks. Then we perform the following steps to calculate the conservation score.

1. Find the best matching module N_{k_j} in the j th noisy network for the k th reference modules R_k , where k_j is obtained by:

$$k_j = \arg \max_i \frac{|R_k \cap N_{i_j}|}{\min(|R_k|, |N_{i_j}|)}$$

2. Repeat step 1 to find the best matching modules in all m noisy networks: N_{k_j} , $j = 1, \dots, m$.
3. Calculate the conservation score for the reference module R_k using the following formula:

$$\text{ConservationScore}(R_k) = \frac{|R_k \cap_{j=1}^m N_{k_j}|}{|R_k|}$$

4. Repeat steps 1-3 to calculate conservation scores of all reference modules $\text{ConservationScore}(R_i)$, $i = 1, \dots, n_r$.

Derivation of $\Delta \tilde{W}$

Let x_i be the boolean variable indicating whether the i th node is selected as a community member. Denote the entire set of nodes as V , and N as the total number of nodes in V . Let \mathbf{A} be the adjacency matrix of the entire network. Following the denotations used by Zhao et al. (2011), we have

$$\tilde{W}_S = O_S \cdot \frac{|S_c|}{|S|} + O_S - (O_S + B_S) = \sum_{i,j \in S} A_{ij} x_i x_j \left(\frac{|S_c|}{|S|} + 1 \right) - \sum_{i \in S, j \in V} A_{ij} x_i \quad (1)$$

If node k is in S , the only move that will change \tilde{W} is to move it from S to S_c . The new \tilde{W} after moving will be:

$$\tilde{W}_{S'} = \left(\sum_{i,j \in S} A_{ij} x_i x_j - 2 \sum_{j \in S} A_{kj} x_j \right) \left(\frac{|S_c| + 1}{|S| - 1} + 1 \right) - \left(\sum_{i \in S, j \in V} A_{ij} x_i - \sum_{j \in V} A_{kj} \right); \quad (2)$$

and if node k is in S_c , the only move that will change \tilde{W} is to move it from S_c to S . The new \tilde{W} after moving will be:

$$\tilde{W}_{S'} = \left(\sum_{i,j \in S} A_{ij} x_i x_j + 2 \sum_{j \in S} A_{kj} x_j \right) \left(\frac{|S_c| - 1}{|S| + 1} + 1 \right) - \left(\sum_{i \in S, j \in V} A_{ij} x_i + \sum_{j \in V} A_{kj} \right) \quad (3)$$

Therefore, we can calculate the change in the value of \tilde{W} when node k is in S :

$$\begin{aligned} \Delta \tilde{W}_k &= \tilde{W}_{S'} - \tilde{W} \\ &= \left(O_S - 2 \sum_{j \in S} A_{kj} x_j \right) \left(\frac{|S_c| + 1}{|S| - 1} \right) - O_S \cdot \frac{|S_c|}{|S|} - \sum_{j \in S} A_{kj} x_j + \sum_{j \in S} A_{kj} \end{aligned} \quad (4)$$

$$= O_S \cdot \frac{N}{|S|(|S| - 1)} - 2 \frac{N}{|S| - 1} \sum_{j \in S} A_{kj} x_j + \sum_{j \in S} A_{kj}, \quad (5)$$

Similarly, we can obtain $\Delta\tilde{W}$ when node k is in S_c

$$\begin{aligned}\Delta\tilde{W}_k &= \tilde{W}_{S'} - \tilde{W} \\ &= \left(O_S + 2 \sum_{j \in S} A_{k_j} x_j \right) \left(\frac{|S_c| - 1}{|S| + 1} \right) - O_S \cdot \frac{|S_c|}{|S|} + \sum_{j \in S} A_{k_j} x_j - \sum_{j \in S} A_{k_j}\end{aligned}\quad (6)$$

$$= -O_S \cdot \frac{N}{|S|(|S| + 1)} + 2 \frac{N}{|S| + 1} \sum_{j \in S} A_{k_j} x_j - \sum_{j \in S} A_{k_j}, \quad (7)$$

Combine the above two equations, we finally derive the equation for $\Delta\tilde{W}$:

$$\Delta\tilde{W}_k = \begin{cases} O_S \cdot \frac{N}{|S|(|S|-1)} - 2 \frac{N}{|S|-1} \sum_{j \in S} A_{k_j} x_j + \sum_{j \in S} A_{k_j} & \text{if } k \in S \\ -O_S \cdot \frac{N}{|S|(|S|+1)} + 2 \frac{N}{|S|+1} \sum_{j \in S} A_{k_j} x_j - \sum_{j \in S} A_{k_j} & \text{if } k \in S_c \end{cases}$$

where $O_S = \sum_{i,j \in S} A_{ij} x_i x_j$.

It is easy to observe from (9) and (10) that the change in \tilde{W} for any flipping of node membership can be calculated in linear time (O_S itself can be updated in linear time for each flip too).

Two Unique DiME Modules in the Grade II and IV Glioma Co-expression Networks

Figure Legends

Tables

Table S1. Relative loss of genes under different B-score cutoffs

B-score Cutoff	Technique								
	DiME			MCODE			Modularity		
	0.05	0.001	1×10^{-5}	0.05	0.001	1×10^{-5}	0.05	0.001	1×10^{-5}
Rembrandt Data (GBM)	32.97% (574 / 1741)	50.09% (872 / 1741)	54.68% (952 / 1741)	58.09% (452 / 778)	81.36% (633 / 778)	83.03% (646 / 778)	41.16% (1073 / 2607)	49.33% (1286 / 2607)	99.04% (2582 / 2607)
TCGA Data (GBM)	30.19% (358 / 1186)	42.50% (504 / 1186)	51.85% (615 / 1186)	33.75% (188 / 557)	39.14% (218 / 557)	45.60% (254 / 557)	2.14% (36 / 1681)	45.63% (767 / 1681)	90.78% (1526 / 1681)
Rembrandt Data (grade II Glioma)	47.27% (1230 / 2602)	62.95% (1638 / 2602)	68.14% (1773 / 2602)	61.96% (728 / 1175)	65.79% (773 / 1175)	69.96% (822 / 1175)	94.97% (3546 / 3734)	98.23% (3668 / 3734)	98.93% (3694 / 3734)
GEO Data (grade II Glioma)	42.46% (1106 / 2605)	66.64% (1736 / 2605)	71.48% (1862 / 2605)	56.59% (466 / 822)	67.76% (557 / 822)	74.70% (614 / 822)	94.76% (3255 / 3435)	99.33% (3412 / 3435)	99.71% (3425 / 3435)

Table S2. Module Members in A Unique DiME Module (Grade II Glioma) Larger than 10 Genes

Module Name	Gene Symbol	Gene Product
Mesenchyme morphogenesis and cell division / differentiation (grade II glioma)	<i>MEX3A</i>	Mex-3 RNA Binding Family Member A
	<i>EBF4</i>	Early B-Cell Factor 4
	<i>HEY1</i>	Hairy/Enhancer-Of-Split Related With YRPW Motif 1
	<i>KCNQ2</i>	Potassium Voltage-Gated Channel, KQT-Like Subfamily, Member 2
	<i>SLC13A3</i>	Solute Carrier Family 13 Member 3
	<i>SLC22A15</i>	Solute Carrier Family 22, Member 15
	<i>ABCA5</i>	ATP-Binding Cassette, Sub-Family A (ABC1), Member 5
	<i>CAMKMT</i>	Calmodulin-Lysine N-Methyltransferase
	<i>PNPLA4</i>	Patatin-Like Phospholipase Domain Containing 4
	<i>RGN</i>	Regucalcin
	<i>ECHDC2</i>	Enoyl CoA Hydratase Domain Containing 2
	<i>PRMT5</i>	Protein Arginine Methyltransferase 5
	<i>MAML2</i>	Mastermind-Like 2
	<i>ZNF22</i>	Zinc Finger Protein 22
	<i>C16orf89</i>	(Chromosome 16 Open Reading Frame 89)
	<i>KLHL26</i>	Kelch-Like Family Member 26
	<i>CLQL1</i>	Complement Component 1, Q
	<i>KLRC3</i>	Killer Cell Lectin-Like Receptor Subfamily C, Member 3
	<i>S100A13</i>	S100 Calcium Binding Protein A13
	<i>SELENBP1</i>	Selenium Binding Protein 1
	<i>TPPP3</i>	Tubulin Polymerization-Promoting Protein Family Member 3
	<i>MYC</i>	Proto-Oncogene C-Myc
	<i>FAM50B</i>	Family With Sequence Similarity 50, Member B
	<i>CHST6</i>	Carbohydrate (N-Acetylglucosamine 6-O) Sulfotransferase 6
	<i>RYR3</i>	Ryanodine Receptor 3
	<i>RCAN3</i>	Calcipressin-3
	<i>PHLDA1</i>	Pleckstrin Homology-Like Domain, Family A, Member 1
	<i>SULF2</i>	Sulfatase 2
<i>RPLP0</i>	Ribosomal Protein, Large, P0	
<i>CDKL2</i>	Cyclin-Dependent Kinase-Like 2 (CDC2-Related Kinase)	
<i>C11orf96</i>	Chromosome 11 Open Reading Frame 96	
<i>STPG1</i>	Sperm-Tail PG-Rich Repeat Containing 1	

Table S3. Module Members in A Unique DiME Module (Grade IV Glioma) Larger than 10 Genes

Module Name	Gene Symbol	Gene Product
Regulation of vesicle-related processes (grade IV glioma)	<i>MGAT1</i>	Mannosyl (Alpha-1,3-)-Glycoprotein Beta-1,2-N-Acetylglucosaminyltransferase
	<i>RAB32</i>	Ras-Related Protein Rab-32
	<i>ANKRD13B</i>	Ankyrin Repeat Domain-Containing Protein 13B
	<i>SBK1</i>	SH3 Domain Binding Kinase 1
	<i>SGMS2</i>	Sphingomyelin Synthase 2
	<i>PLBD1</i>	Phospholipase B Domain Containing 1
	<i>LRRTM2</i>	Leucine Rich Repeat Transmembrane Neuronal 2
	<i>FSD1</i>	Fibronectin Type III And SPRY Domain Containing 1
	<i>WIPI1</i>	WD Repeat Domain, Phosphoinositide Interacting 1
	<i>BICC1</i>	Bicaudal C Homolog 1
	<i>SOX8</i>	SRY (Sex Determining Region Y)-Box 8
	<i>RRAS</i>	Related RAS Viral (R-Ras) Oncogene Homolog
	<i>LMF1</i>	Lipase Maturation Factor 1
<i>GATAD2B</i>	GATA Zinc Finger Domain Containing 2B	