

BRANENet: Embedding Multilayer Networks for Omics Data Integration

Surabhi Jagtap

CentraleSupélec

Aurélie Pirayre

IFP Energies nouvelles (IFPEN)

Frédérique Bidard

IFP Energies nouvelles (IFPEN)

Laurent Duval

IFP Energies nouvelles (IFPEN)

Fragkiskos D. Malliaros (✉ fragkiskos.malliaros@centralesupelec.fr)

CentraleSupélec

Research Article

Keywords: multi-omics data, multilayer network, biological network integration, graph representation learning, regulatory network inference

Posted Date: April 27th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1588328/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

RESEARCH

BRANENet: Embedding Multilayer Networks for Omics Data Integration

Surabhi Jagtap^{1,2}, Aurélie Pirayre², Frédérique Bidard², Laurent Duval² and Fragkiskos D. Malliaros^{1*}

*Correspondence:

fragkiskos.malliaros@centralesupelec.fr

¹CentraleSupélec, Inria, Université Paris-Saclay, Gif-Sur-Yvette, France

²IFP Energies nouvelles, 1 et 4 avenue de Bois-Préau, 92852 Rueil-Malmaison, France

Full list of author information is available at the end of the article

Abstract

Background: Gene expression is regulated at different molecular levels, including chromatin accessibility, transcription, RNA maturation, and transport. These regulatory mechanisms have strong connections with cellular metabolism. In order to study the cellular system and its functioning, omics data at each molecular level can be generated and efficiently integrated. Here, we propose BRANENET, a novel multi-omics integration framework for multilayer heterogeneous networks. BRANENET is an expressive, scalable, and versatile method to learn node embeddings, leveraging random walk information within a matrix factorization framework. Our goal is to efficiently integrate multi-omics data to study different regulatory aspects of multilayered processes that occur in organisms. We evaluate our framework using multi-omics data of *Saccharomyces cerevisiae*, a well-studied yeast model organism.

Results: We test BRANENET on transcriptomics (RNA-seq) and targeted metabolomics (NMR) data for wild-type yeast strain during a heat-shock time course of 0, 20, and 120 minutes. Our framework learns features for differentially expressed bio-molecules showing heat stress response. We demonstrate the applicability of the learned features for targeted omics inference tasks: transcription factor (TF)-target prediction, integrated omics network (ION) inference, and module identification. The performance of BRANENET is compared with existing network integration methods. Our model outperforms baseline methods by achieving high prediction scores for a variety of downstream tasks.

Keywords: multi-omics data; multilayer network; biological network integration; graph representation learning; regulatory network inference

Background

Gene expression in eukaryotes is regulated at different molecular levels, including chromatin accessibility, transcription, RNA maturation, and transport. These regulatory processes are effected by specific bio-molecules potentially in separate cellular compartments with interconnected steps [1]. One of the earlier steps consists of epigenetic modifications that can modify DNA accessibility and chromatin structure, ensuring access to the transcriptional machinery [2, 3]. These modifications occur across the genome, regulating the final synthesis of the mRNAs [4]. These newly synthesized transcripts (mRNAs) are extensively modified and translated into proteins [5]. In this process, mRNA molecules are guided by RNA-binding proteins that control their expression [6]. Finally, the encoded protein products (metabolites) participate in numerous processes that regulate cellular metabolism [7, 8]. The metabolites produced in this step are transformed and/or stored. A number of

these metabolites/compounds (for instance, Acetyl-CoA, glucose or methyl groups) successively participate in chromatin modifications (epigenetics) that regulate gene expression [9, 10]. A comprehensive interpretation of such phenomena would help to understand the structure, function, and dynamics of cellular systems [11]. For this purpose, the collective characterization and quantification of data at different molecular levels are essential. Omics technologies give access to these regulatory mechanisms and are thus able to provide large amounts of data at each molecular level. Despite the abundance and diversity of numerous available omics datasets, there are some noticeable challenges regarding their acquisition, processing, efficient integration, and interpretation [12, 11, 13].

To this direction, networks are widely used to represent individual bio-molecules (nodes) and their biological relationships (edges). These relationships may reflect gene co-expression/regulation, protein-protein interactions (PPIs) or information about production or consumption of metabolites [14]. The major challenge thus pertains to effectively integrate the variety of these relationships (layers) encoded in multiple networks. In this study, we are inspired by graph representation learning (GRL) algorithms that allow to encode the high-dimensional graph structure into compact low dimensional embedding vectors [15, 16]. The aim is to optimize the mapping from the original graph-based node representation onto a lower dimensional space such that their geometric relationship in this latent space reflects the structure of the input graph (network). After optimizing the embedding space, the learned embeddings can be used as input features for downstream omics inference tasks, such as bio-marker identification, drug-target prediction, disease-gene associations, gene regulatory network (GRN) inference [17], and protein function prediction [13].

In the related literature, a variety of methods have been proposed for computing network embeddings either based on random walks [18, 19], matrix factorization [20], or neural networks [15]. However, the majority of them has specifically been designed for single-layer networks. Nevertheless, for multilayer biological networks, only a few existing network integration strategies leverage GRL. As one of the proposed models, OHMNET [21] is a hierarchy-aware task-independent, unsupervised method for protein feature learning. MULTI-NET [22] is an extension of the DEEPWALK [18] and NODE2VEC[19] to multilayer graphs. It allows to perform random walks by defining paths to traverse the nodes. DEEPNF [23], a network fusion method, is based on multimodal deep autoencoder (MDA) to integrate different heterogeneous networks. More recently, BRANE-EXP [24] leverages random walks with the concept of exponential family graph embeddings [25, 26]. It generalizes multilayer random walk-based GRL methods to an instance of exponential family distributions. Nevertheless, most of the existing methods are challenged when applied to real and heterogeneous omics data which demands comprehensive handling towards preserving biological relevance [11]. In such cases, it is necessary to obtain knowledge-based representations of the nodes to assist omics data analysis. Here, we propose BRANENET, a novel multi-omics integration framework for multilayer heterogeneous networks. In particular, we aim to embed graph-based information from multi-omics data in a lower-dimensional space. We leverage a properly chosen random walk-based Positive Pointwise Mutual Information (PPMI) matrix, that

is able to capture relevant context around each node of interest. We introduce network integration with multilayered heterogeneous graph embeddings, that perform matrix factorization by approximating the spectrum of this PPMI matrix. We apply BRANENET over a multi-omics dataset for heat-shock response in the well-known yeast model *Saccharomyces cerevisiae* [27] and learn embeddings for genes, transcription factors (TFs), and metabolites. We focus on TF-target prediction, integrated omics network (ION) inference, and module detection. We validate the embeddings on various downstream tasks and demonstrate BRANENET's effectiveness by comparing its performance to existing network integration approaches.

Methods

Data description

To test our framework, we have used recently published multi-omics yeast datasets acquired in the study of [27]. These datasets present three basic layers of the transcriptional circuit, including, one type of epigenetic modification (H4K12ac mark for identification of active promoters obtained from ChIP-Seq), gene expression (RNA-seq), and targeted metabolomics (NMR). The dataset is comprised of 7,126 genes, 1,970 H4K12ac peaks, and 37 metabolites. To obtain this data, yeast culture flask was grown at 30°C until the exponential phase. This culture was split into three different flasks. One flask was maintained at 30°C and labeled as 0 minute (t0). The other two flasks were incubated at 39°C for 20 minutes (t20) and 120 minutes (t120), respectively. Aliquots from all three flasks (t0, t20, and t120) were collected for ChIP-seq (epigenomics), RNA-seq (transcriptomics), and NMR (metabolomics). This process was repeated four times to generate four biological replicates. The datasets were pre-processed using various bioinformatics tools [27]. These consistent datasets, dedicated to the study of the heat stress response, appear as a good candidate to test and evaluate our proposed omics data integration methodology. We recall the experimental setup and summarize the workflow in Figure 1.

Differential expression analysis

Genes, metabolites, TFs are hereafter referred to as bio-molecules. More generally, genes are regulated by TFs. Therefore, we separate TFs (genes coding for TFs) and non-TF (genes not coding for TFs) from transcriptomics data. Now, to obtain differentially expressed bio-molecules, we first take the average of control samples in the four biological replicates. Then, we compute the log₂ of Fold Change (log₂FC) for each bio-molecule in eight test samples (four in t20 and four in t120) by taking its ratio against the average of four control samples (t0). For each test sample, we select non-TFs if log₂FC is higher than 2 (over-expressed) or lower than -2 (under-expressed). However, it is well known that expression TFs do not vary considerably as compared to non-TFs [28]. Therefore, we lower the threshold of log₂FC for genes encoding for TFs (TFs). TFs were considered as differentially expressed, if log₂FC is higher than 1 (over-expressed) or lower than -1 (under-expressed). For metabolites and H4K12ac peaks, we choose the log₂FC threshold similarly to TFs. If a log₂FC value is meaningful with respect to the above thresholds in at least one biological replicate, we consider the corresponding bio-molecule as differentially expressed. Note that, for the published yeast multi-omics data presented in the data descriptor [27], no statistical tests were performed. Therefore, we only rely on the log₂FC value.

Construction of intra-omics and inter-omics networks

Intra-omics networks are constructed using the same type of bio-molecules, for example, gene-gene co-expression, or metabolite-metabolite correlation networks. These networks are built on data obtained from multi-omics experiments for instance: genomics, epigenetics, transcriptomics, proteomics, and metabolomics. We obtain differentially co-expressed bio-molecules by computing the Pairwise Pearson correlation coefficient (ρ) [29] of log2FC profiles described above, i.e., log2FC for the eight samples (four in t20 and four in t120). Two intra-omic elements were said to be correlated if the absolute value of ρ is higher than 0.8. These intra-omic correlation networks are represented as a set of their adjacency matrices.

Inter-omics networks link bio-molecules of different types. They are constructed using biological *a priori* information showing presence of TF binding sites or H4K12ac epigenetic marks in the promoter of gene, biochemical reactions within genes and metabolites. This information can be acquired from various bioinformatics databases such as SGD [30], YEASTRACT (Yeast Search for Transcriptional Regulators And Consensus Tracking) [31], YeastPathways [30], and BioCyc [32]. This *a priori* knowledge bridges the gap to relate two different omics types, for instance gene-metabolite, TF-target, gene-epigenetic mark. For each differentially expressed bio-molecule of one type (e.g., gene), we obtained its relationship with an bio-molecule of another type (e.g., TF and metabolite).

The BRANENET model

In Figure 2, an illustration of BRANENET is presented. The model takes two sets of networks, intra- and inter-omics X and Y respectively. It first builds a multilayer network G using both sets which are represented by $\bar{\mathbf{A}}_{|N| \times |N|}$, where N is a set of all bio-molecules. Then, we compute a random walk-based PPMI matrix \mathbf{S} for $\bar{\mathbf{A}}$ via a closed-form solution. Matrix $\mathbf{S}_{|N| \times |N|}$ is then factorized using Singular Value Decomposition (SVD), and the d -dimensional embedding vectors $\boldsymbol{\Omega}_d \in \mathbb{R}^{|N| \times d}$ ($d \ll |N|$) are given by $\mathbf{U}_d \sqrt{\boldsymbol{\Sigma}_d}$. Next, we describe the workflow in detail.

Construction of a multilayer network

X is a set of x intra-omics networks represented as $\mathbf{X}_{n_1 \times n_1}^{(1)}, \mathbf{X}_{n_2 \times n_2}^{(2)}, \dots, \mathbf{X}_{n_x \times n_x}^{(x)}$, while Y is a set of $y = \frac{x(x-1)}{2}$ inter-omics networks represented as $\mathbf{Y}_{n_1 \times n_2}^{(1)}, \mathbf{Y}_{n_1 \times n_3}^{(2)}, \dots, \mathbf{Y}_{n_{x-1} \times n_x}^{(y)}$. A multilayer network G of $|N|$ nodes (bio-molecules) and $|E|$ edges (interactions) is built using sets X and Y . The network is represented by its supra-adjacency matrix $\bar{\mathbf{A}}_{|N| \times |N|}$ that is defined as:

$$\bar{\mathbf{A}} = \bigoplus_x \mathbf{A}^{(x)} + \mathbf{C}, \quad (1)$$

where $\bigoplus_x \mathbf{A}^{(x)}$ is the intra-omics adjacency matrices and \mathbf{C} is a block matrix with zero diagonal blocks that stores inter-omics connections obtained from elements in Y . The final output of $\bar{\mathbf{A}}$ has intra-omics networks represented as blocks in the main diagonal and inter-omics networks represented as off-diagonal matrices.

Representation learning

To embed nodes from different omics modalities into a common latent space towards integrating inter- and intra-omics relationships, we construct a random walk matrix \mathbf{S} for the multilayer graph G . \mathbf{S} is defined by the random walk transition probabilities to traverse nodes within and across layers. The flexibility of random walks to traverse within and across layers allows us to capture inter- and intra-layer node neighbourhood information. This is an important and useful property to consider while performing data integration from multilayer networks [24, 22]. For instance, starting from node v in G , a random walk traverses the multilayer graph, moving across neighborhood nodes of v chosen uniformly at random. This process repeats for a predefined number of walks per node. Nevertheless, for large networks, simulating random walks is computationally expensive and therefore it is not a recommended approach. To address this limitation, we leverage the relationship between random walk-based GRL algorithms that rely on the SKIP-GRAM model (for instance, DEEPWALK [18]) and matrix factorization [33]. In particular, a multilayer random walk matrix \mathbf{S} is defined by computing the closed-form of a properly normalized PPMI based random walk transition matrix. This PPMI matrix is the well-studied pointwise mutual information (PMI) matrix that represents node similarities, shifted by a global constant. It has been shown that normalized PPMI is better at optimizing SKIP-GRAM’s objective and show better performance than WORD2VEC derived models [18, 19] in Natural Language Processing (NLP tasks) [34, 33]. For any graph G , \mathbf{S} is given by:

$$\mathbf{S} = \log \left\{ \frac{\text{vol}(G)}{bT} \left(\frac{1}{T} \sum_{r=1}^T \mathbf{P}^r \right) \mathbf{D}^{-1} \right\}, \quad (2)$$

where $\mathbf{P} = \mathbf{D}^{-1}\bar{\mathbf{A}}$. Matrices $\bar{\mathbf{A}}$ and \mathbf{D} are respectively the adjacency and diagonal degree matrices of the graph G and $\text{vol}(G)$ is the sum of the node degrees of G . T corresponds to the window size and b is number of negative samples [33]. In order to obtain node embeddings from matrix \mathbf{S} , we perform spectral decomposition using SVD [35], given by $\mathbf{S} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$. Since \mathbf{S} is a real and symmetric matrix, \mathbf{U} and \mathbf{V} correspond to singular vector matrices and $\mathbf{\Sigma}$ is the singular value matrix. The integrated embedding matrix $\mathbf{\Omega}_d$ of dimension $|N| \times d$ is given by the first d eigenvectors of \mathbf{S} , appropriately weighted by the square root of $\mathbf{\Sigma}_d$ as, $\mathbf{\Omega}_d = \mathbf{U}_d\sqrt{\mathbf{\Sigma}_d}$.

Downstream tasks

The embeddings learned from BRANENET can be used for various downstream tasks, for instance, TF-target prediction, ION inference, identification of biomarkers (e.g., heat stress responsive genes/TFs), identification of biologically related clusters, and visualization. Their details are shown below.

TF-target prediction

To predict TF-targets, we adapt the traditional link prediction task [36] to TF-target networks. We use the largest connected component of the TF-target network. Then, we split the targets of each TF into two parts to form positive training and test sets by randomly removing 50% of them. The same number of TF-target

pairs that do not exist are sampled to generate negative instances for each training and test sets. The learned embeddings Ω_d are used to compute edge features. In particular, the embeddings of node i and j of size d , given by $\Omega_d[i]$ and $\Omega_d[j]$ respectively, are converted into edge feature vectors using element-wise operations [19, 26] (i) *average*: $\Omega_d[i] + \Omega_d[j]/2$; (ii) *weighted L2*: $|\Omega_d[i] - \Omega_d[j]|^2$. Now, for each positive and negative test and training datasets generated above, edge features are computed. Then, we perform prediction using the logistic regression classifier with L2 regularization [37]. The performance is measured using the area under the precision-recall curve (AUPR) [38]. The performance of BRANENET for TF-target prediction is compared with baseline methods.

ION inference

To infer an ION from the learned embeddings, the pairwise similarity score for nodes i and j is defined as:

$$\delta_{i,j} = \Omega_d[i] \cdot \Omega_d[j] = \sum_{i=1}^d \Omega_d[i] \Omega_d[j]. \quad (3)$$

To validate this network, we compare it with the gold-standard (GS) network of yeast that is built by combining networks from multiple databases, such as BIOGRID [39], STRING [40], and YEASTRACT [41]. The performance of ION inference is measured by computing the Precision@ k and the Matthews Correlation Coefficient (MCC). Precision@ k gives us the fraction of true edges among the inferred top k ones. To compute Precision@ k , we sort based on the pairwise similarity score, then take the top k ones. Then, we compute Precision ($\frac{\text{True Positive (TP)}}{\text{TP} + \text{False Positive (FP)}}$) at each top k edges. We chose $k \in \{1, 10, 50, 100, \dots, 500\}$. MCC is another way to evaluate performance. It measures the differences between the actual values and the predicted ones. The advantages of MCC over Precision-based metrics are shown in recent articles [42]. For the edges with similarity score $\delta_{i,j}$ higher than a threshold ($\theta \in \{0.1, 0.2, \dots, 0.9\}$), we compute MCC ($\frac{\text{True Negative (TN)} \times \text{TP} - \text{False Negative (FN)} \times \text{FP}}{\sqrt{(\text{TP} + \text{FP})(\text{TP} + \text{FN})(\text{TN} + \text{FP})(\text{TN} + \text{FN})}}$). We compare the performance of BRANENET for ION inference with baseline methods.

Module detection

Interestingly in biological networks, the clustering or community structure property is present, under which the graph topology is organized into modules commonly called communities or clusters. To obtain these modules, we first select the top scoring edges ($\theta = 0.7$). Then we find clusters using a greedy modularity maximization algorithm [43]. We select the obtained modules having more than 10 nodes. To validate the obtained clusters, we investigated their biological meaningfulness by performing functional annotation enrichment analysis [44, 45].

Comparison with baseline methods

We compare the performance of link prediction and ION inference using the embeddings learned by BRANENET with existing multilayer network embedding methods. We choose OHMNET [21], MULTINET [22], DEEPPNF [46], and BRANE-EXP

[24] asv our baseline methods. These network integration methods are not specifically developed for multi-omics integration considering biological *a priori* knowledge to learn node embeddings. Therefore, we adapt the existing methods for learning embeddings and performing downstream tasks by keeping the same empirical parameter settings as of BRANENET. Apart from DEEPNF, all baseline models mainly depend on the window size (T) and embedding dimension (d). For DEEPNF, we choose to keep the default model architecture configuration proposed by the authors [46].

Results & Discussion

We present the results of BRANENET applied to the yeast multi-omics dataset [27]. We have identified differentially expressed (DE) bio-molecules as mentioned in the “Methods” section. We have obtained 333 DE genes (non-TF) out of which 310 are upregulated and 23 are downregulated, 55 DE TFs (50: over-expressed; 5: under-expressed), 30 DE metabolites (28: increased concentration; 2: decreased concentration). The list of all DE bio-molecules with their FC and SGD annotation is provided in the [Supplementary data](#). For the epigenetics data, we have observed that no H4K12ac peaks were differentially expressed. Therefore, we discard ChIP-Seq data and use only variable genes, TFs, and metabolites for the study of heat shock response. We then compute intra- and inter-omics networks as described in “Construction of a multilayer network” section. To construct inter-omics relationships, we obtain known TF-target interactions from the YEASTRACT database [31]. Gene-metabolite and TF-metabolite associations were given by the participation of genes or TFs in production and consumption of metabolites in biochemical reactions. This information was acquired from the YeastPathways database [30]. Embeddings are then learned for each node, as discussed in the “Methods” section. We use these embeddings to study different aspects of multi-omics data integration, namely, TF-target Prediction, ION inference, and module detection.

TF-target prediction

We have performed TF-target prediction as mentioned in the “Methods” section. First, we compute node embeddings using BRANENET with parameter $T = 3$, $b = 1$, and $d = 128$. For the same value of d , we learn node embeddings using each baseline method. The edge features are computed using the operators mentioned in the “Methods” section. TF-target prediction is then performed using logistic regression and its performance is measured using the AUPR score. Since, we randomly remove 50% of targets for each TF, we repeat this process 10 times and report the average AUPR scores with standard deviation computed across 10 runs. The results for BRANENET compared to baselines models are summarized in Figure 3. The average AUPR of BRANENET is 10% improved compared to *average* (87%) to *weighted L2* (97.9%) operators. For the empirical parameter settings, the performance of BRANENET is higher in both operators as compared to the baselines methods. *Weighted L2* score of BRANENET’s is 20% higher than BRANE-EXP, the second best performing model. Whereas BRANENET’s *average* score is 10% higher than MULTINET, the second best performing model for *average*. The standard deviation of BRANENET for 10 runs is notably lower (except MULTINET

with *average*) than all the other methods. Overall from Figure 3, we observe that BRANENET outperforms the baseline methods for both operators (i.e., *average* and *weighted L2*). With the exception of DEEPNF, all the rest algorithms constitute mostly extensions of SKIP-GRAM-like GRL models (e.g., DEEPWALK) to multilayer graphs. For single layer graphs, the computational advantages of the spectral SVD algorithm over stochastic gradient training, have been well-studied [47, 34]. First, it is exact, therefore does not require extensive hyper-parameter tuning. Second, unlike SKIP-GRAM’s training procedure, BRANENET does not require layer-wise observations of each node and its respective neighborhood. Hence, it is easily trainable. Previous studies have demonstrated that by explicitly factorizing the PPMI random walk transition matrix using SVD, achieves better performance on downstream tasks compared to SKIP-GRAM derived models [33]. Referring to these studies, we state a hypothesis that our multilayer network embedding approach has similar computational advantages over other SKIP-GRAM-like network embedding models for multilayer graphs [22, 24, 21].

Integrated Omics Network (ION) inference

We infer an ION using the embeddings learned by BRANENET. To validate the performance of ION inference, we reconstruct the gold-standard (GS) using the learned embeddings. The performance is measured by computing Precision@ k and MCC (Mathews Correlation Coefficient). We choose to study the top 500 edges of the inferred ION. We compare the performance of our model to the performance of the baseline methods used. The results are shown in Figure 5. The results of Precision@ k show that BRANENET scores, up to top 320 edges, are higher than DEEPNF. BRANENET outperforms OHMNET, MULTINET, DEEPNF, and BRANE-EXP (Figure 5a). The results of MCC@threshold shows that BRANENET’s performance for different thresholds (θ) is higher than OHMNET, MULTINET, DEEPNF, and BRANE-EXP. As shown in Figure 5b, the MCC of BRANENET was gradually improved with increasing threshold, and began to drop quite sharply at 0.6. Overall, for both metrics, DEEPNF is the second best performing model. Moreover, we have also observed that, to obtain a well-trained model, methods like DEEPNF require large training data involving tuning of numerous parameters. This may cause overfitting/underfitting problems. The best performance of such baseline models can be achieved by extensive hyper-parameter tuning. As an advantage, BRANENET comes with very few parameters.

The network inferred by BRANENET with $\theta = 0.7$ is shown in Figure 4. Node colour represents over- (pink) and under- (blue) expressed bio-molecules. Node shape and label color represent gene (circle, black), TF (triangle, purple), and metabolites (square, orange). Edges existing in GS are given in red whereas newly inferred edges are given in green. Edge width is represented by the similarity scores, while the node label size is proportional to its degree. Using the inferred ION, we narrow down the search space from all differentially expressed bio-molecules and identify potential biomarkers in heat stress response. We rank nodes based on their degree. Table 1 shows the obtained 21 bio-molecules that could be potential biomarkers in heat stress response. We have investigated the participation of these genes during heat-shock response in published literature. Using the BRANENET

integrated tool, we are able to recover information from 11 different heat shock response studies. The references of these articles are given in Table 1. We have also validated our results by comparing to another study of heat shock response [48]. We could find the potential bio-markers (1) in the heat stress responsive gene clusters that were identified in this study.

Biologically meaningful modules

To identify modules from the inferred ION, we perform community detection using the Clauset-Newman-Moore greedy modularity maximization algorithm [49]. We select modules with size more than 10 nodes. We have obtained 6 modules. To know if the obtained modules are biologically meaningful, we perform functional enrichment analysis on the two largest modules. We select the terms with p -value lower than 0.05. Their enrichment results are shown in Figure 6. We can clearly see that module A is enriched with catabolic processes including HSP90 and chaperone binding activity related terms while module B is enriched with transport and sporulation. The terms enriched in both these clusters have been discussed over years in yeast heat-shock response studies [48, 50, 51].

Parameter sensitivity analysis

To investigate the robustness in the performance of BRANENET, we performed parameter sensitivity analysis. We used grid-search to assess the uncertainty in the model outputs that is attributed to different values of the window size T and dimension d . We choose $T \in \{1, 2, 3, 4, 5\}$, $d \in \{32, 64, 128, 256\}$, and perform TF-target prediction and ION inference. The results for TF-target prediction are shown in Table 2. The mean AUPR in the given table for *average* and *weighted L2* is 83.8% and 96.1%, respectively, with standard deviation of 2 percent. On the other hand, the results for ION inference are shown in Figure 7. The performance is measured by MCC at different thresholds. From Figure 7, the optimal threshold for θ is between 0.6 and 0.8. We also see that the performance of MCC is increased with respect to d and slightly decreased with T . From the parameter sensitivity analysis for both tasks, we see that our model has lower variance in the results with respect to different parameter settings. Therefore for the new datasets, we recommend users to consider the default parameter settings ($T = 3$ and $d = 128$).

Conclusions

The recent wide application of high-throughput experimental techniques has provided complex high-dimensional heterogeneous data. In turn, the wide availability of these omics data has driven the need to develop methods that can efficiently analyse and integrate these data. We have presented an integrative analysis of multilayer heterogeneous networks for learning low-dimensional features for bio-molecules. BRANENET relies on a GRL technique to learn embeddings that can capture relevant features from complex networks built from single omics data. Considering the biological context and need for omics integration, methods like BRANENET enable us to study the cross-talk between genes, transcription factors, and metabolites at transcriptional and metabolic levels of regulation.

Besides, we have compared the performance of our approach for various downstream tasks with existing network integration methods. To summarize the empirical analysis, the performance of BRANENET is especially appealing mainly because of four reasons. First, our approach integrates experimental data with biological *a priori* knowledge which facilitates the inference of inter-omics relationships. Second, it can generate meaningful embeddings by preserving the inter- and intra-omics interactions. Third, its objective function is independent of the downstream tasks, thereby it is adaptable to various omics data inference tasks. And fourth, it has a low number of parameters to tune. For the predicted bio-markers, experimental as well as *in silico* investigation can help in identifying their functions and biological significance [52]. As a future work, we intend to test BRANENET on multiple omics datasets including epigenetics, proteomics, and metagenomics. We would also like to explore promising directions of multilayer network embedding that can support fast approximations of higher-order random walks transition matrices. BRANENET is a versatile tool that combines data-driven information and biological *a priori* knowledge and provides an effective and scalable network integration framework with diverse downstream integrative omics applications. We consider BRANENET to be a valuable method of biological network integration for experimental follow-up, as well as a baseline candidate for further development of the graph-based multi-omics data integration field.

Funding

This work has been supported in part by ANR (French National Research Agency) under the JCJC project GraphIA (ANR-20-CE23-0009-01).

Availability of data and materials

All data generated or analyzed during this study are included in this published article, its supplementary file and the GitHub repository. BRANET is freely available as a Jupyter notebook present at <https://github.com/Surabhivj/BRANenet>.

Authors' contributions

SJ, FM, AP, and LD designed the study, participating in the mathematical modeling, experimental setup, and downstream bioinformatics tasks. SJ performed the experiments and drafted the manuscript. SJ, FM, AP, LD, and FB collaborated on result interpretation. FM, AP and LD reviewed and edited the manuscript. All authors read and approved the final manuscript.

Declarations

Competing interests

The authors declare that they have no competing interests.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Author details

¹CentraleSupélec, Inria, Université Paris-Saclay, Gif-Sur-Yvette, France. ²IFP Energies nouvelles, 1 et 4 avenue de Bois-Préau, 92852 Rueil-Malmaison, France.

References

- Rodríguez-Navarro, S., Hurt, E.: Linking gene regulation to mRNA production and export. *Current opinion in cell biology* **23**(3), 302–309 (2011)
- Jaenisch, R., Bird, A.: Epigenetic regulation of gene expression: how the genome integrates intrinsic and environmental signals. *Nature genetics* **33**(3), 245–254 (2003)
- Li, D., Yang, Y., Li, Y., Zhu, X., Li, Z.: Epigenetic regulation of gene expression in response to environmental exposures: from bench to model. *Science of The Total Environment*, 145998 (2021)
- Woo, H., Ha, S.D., Lee, S.B., Buratowski, S., Kim, T.: Modulation of gene expression dynamics by co-transcriptional histone methylations. *Experimental & molecular medicine* **49**(4), 326–326 (2017)
- Zhao, B.S., Roundtree, I.A., He, C.: Post-transcriptional gene regulation by mRNA modifications. *Nature reviews Molecular cell biology* **18**(1), 31–42 (2017)

6. Dreyfuss, G., Kim, V.N., Kataoka, N.: Messenger-RNA-binding proteins and the messages they carry. *Nature reviews Molecular cell biology* **3**, 195–205 (2002)
7. Metallo, C.M., Vander Heiden, M.G.: Understanding metabolic regulation and its influence on cell physiology. *Molecular cell* **49**(3), 388–398 (2013)
8. Carthew, R.W.: Gene regulation and cellular metabolism: An essential partnership. *Trends in Genetics* **37**(4), 389–400 (2021)
9. Du, J., Johnson, L.M., Jacobsen, S.E., Patel, D.J.: DNA methylation pathways and their crosstalk with histone methylation. *Nature reviews Molecular cell biology* **16**(9), 519–532 (2015)
10. Sabari, B.R., Zhang, D., Allis, C.D., Zhao, Y.: Metabolic regulation of gene expression through histone acylations. *Nature reviews Molecular cell biology* **18**(2), 90–101 (2017)
11. Subramanian, I., Verma, S., Kumar, S., Jere, A., Anamika, K.: Multi-omics data integration, interpretation, and its application. *Bioinformatics and biology insights* **14**, 1–24 (2020)
12. Ma, T., Zhang, A.: Integrate multi-omics data with biological interaction networks using multi-view factorization autoencoder (mae). *BMC genomics* **20**(11), 1–11 (2019)
13. Yue, X., Wang, Z., Huang, J., Parthasarathy, S., Moosavinasab, S., Huang, Y., Lin, S.M., Zhang, W., Zhang, P., Sun, H.: Graph embedding on biomedical networks: methods, applications and evaluations. *Bioinformatics* **36**(4), 1241–1251 (2020)
14. Huber, W., Carey, V.J., Long, L., Falcon, S., Gentleman, R.: Graphs in molecular biology. *BMC Bioinformatics* **8**(S8) (2007)
15. Hamilton, W.L., Ying, R., Leskovec, J.: Representation learning on graphs: Methods and applications. *IEEE Data Engineering Bulletin* **40**(3), 52–74 (2017)
16. Ietswaart, R., Gyori, B.M., Bachman, J.A., Sorger, P.K., Churchman, L.S.: Genewalk identifies relevant gene functions for a biological context using network representation learning. *Genome biology* **22**(1), 1–35 (2021)
17. Kc, K., Li, R., Cui, F., Yu, Q., Haake, A.R.: Gne: a deep learning framework for gene network inference by aggregating biological information. *BMC systems biology* **13**(2), 1–14 (2019)
18. Perozzi, B., Al-Rfou, R., Skiena, S.: Deepwalk: Online learning of social representations. In: *Proc. 20th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, pp. 701–710 (2014)
19. Grover, A., Leskovec, J.: node2vec: Scalable feature learning for networks. In: *Proc. 22nd ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, pp. 855–864 (2016)
20. Kuang, D., Ding, C., Park, H.: Symmetric nonnegative matrix factorization for graph clustering. In: *Proc. 2012 SIAM International Conference on Data Mining*, pp. 106–117 (2012). SIAM
21. Zitnik, M., Leskovec, J.: Predicting multicellular function through multi-layer tissue networks. *Bioinformatics* **33**(14), 190–198 (2017)
22. Bagavathi, A., Krishnan, S.: Multi-Net: a scalable multiplex network embedding framework. In: *Proc. Int. Conf. on Complex Networks and Their Applications*, pp. 119–131 (2018)
23. Gligorijević, V., Barot, M., Bonneau, R.: deepNF: deep network fusion for protein function prediction. *Bioinformatics* **34**(22), 3873–3881 (2018)
24. Jagtap, S., Çelikkanat, A., Pirayre, A., Bidard, F., Duval, L., Malliaros, F.D.: Multiomics data integration for gene regulatory network inference with exponential family embeddings. In: *29th European Signal Processing Conference (EUSIPCO)*, pp. 1221–1225 (2021)
25. Rudolph, M., Ruiz, F., Mandt, S., Blei, D.: Exponential family embeddings. In: *Proc. 30th Conf. on Neural Information Processing Systems*, pp. 478–486 (2016)
26. Çelikkanat, A., Malliaros, F.D.: Exponential family graph embeddings. In: *Proc. AAAI Conference on Artificial Intelligence*, vol. 34, pp. 3357–3364 (2020)
27. Nuño-Cabanes, C., Ugidos, M., Tarazona, S., Martín-Expósito, M., Ferrer, A., Rodríguez-Navarro, S., Conesa, A.: A multi-omics dataset of heat-shock response in the yeast RNA binding protein Mip6. *Scientific data* **7**(69) (2020)
28. Dalman, M.R., Deeter, A., Nimishakavi, G., Duan, Z.-H.: Fold change and p-value cutoffs significantly alter microarray interpretations. In: *BMC Bioinformatics*, vol. 13, pp. 1–4 (2012)
29. Zhang, B., Horvath, S.: A general framework for weighted gene co-expression network analysis. *Statistical applications in genetics and molecular biology* **4**(1) (2005)
30. Cherry, J.M., Adler, C., Ball, C., Chervitz, S.A., Dwight, S.S., Hester, E.T., Jia, Y., Juvik, G., Roe, T., Schroeder, M., *et al.*: SGD: *Saccharomyces Genome Database*. *Nucleic acids research* **26**(1), 73–79 (1998)
31. Teixeira, M.C., Monteiro, P.T., Palma, M., Costa, C., Godinho, C.P., Pais, P., Cavalheiro, M., Antunes, M., Lemos, A., Pedreira, T., *et al.*: YEASTRACT: an upgraded database for the analysis of transcription regulatory networks in *Saccharomyces cerevisiae*. *Nucleic acids research* **46**(D1), 348–353 (2018)
32. Karp, P.D., Billington, R., Caspi, R., Fulcher, C.A., Latendresse, M., Kothari, A., Keseler, I.M., Krummenacker, M., Midford, P.E., Ong, Q., *et al.*: The BioCyc collection of microbial genomes and metabolic pathways. *Briefings in bioinformatics* **20**(4), 1085–1093 (2019)
33. Qiu, J., Dong, Y., Ma, H., Li, J., Wang, K., Tang, J.: Network embedding as matrix factorization: Unifying DeepWalk, LINE, PTE, and node2vec. In: *Proc. 11th ACM International Conference on Web Search and Data Mining*, pp. 459–467 (2018)
34. Levy, O., Goldberg, Y.: Neural word embedding as implicit matrix factorization. *Advances in neural information processing systems* **27**, 2177–2185 (2014)
35. Bisgard, J.: *Analysis and Linear Algebra: The Singular Value Decomposition and Applications*, 1st edn. Student Mathematical Library, p. 217. American Mathematical Society, Providence, RI, USA (2020)
36. Liben-Nowell, D., Kleinberg, J.: The link-prediction problem for social networks. *Journal of the American society for information science and technology* **58**(7), 1019–1031 (2007)
37. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Vanderplas, J.: scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12** (2011)
38. Flach, P., Kull, M.: Precision-recall-gain curves: PR analysis done right. *Advances in neural information processing systems* **28** (2015)

39. Oughtred, R., Rust, J., Chang, C., Breitkreutz, B.-J., Stark, C., Willems, A., Boucher, L., Leung, G., Kolas, N., Zhang, F., et al.: The biogrid database: A comprehensive biomedical resource of curated protein, genetic, and chemical interactions. *Protein Science* **30**(1), 187–200 (2021)
40. Szklarczyk, D., Gable, A.L., Nastou, K.C., Lyon, D., Kirsch, R., Pyysalo, S., Doncheva, N.T., Legeay, M., Fang, T., Bork, P., Jensen, L.J., von Mering, C.: The STRING database in 2021: customizable protein–protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic Acids Research* **49**(D1), 605–612 (2020)
41. Monteiro, P.T., Oliveira, J., Pais, P., Antunes, M., Palma, M., Cavalheiro, M., Galocha, M., Godinho, C.P., Martins, L.C., Bourbon, N., et al.: YEASTRACT+: a portal for cross-species comparative genomics of transcription regulation in yeasts. *Nucleic acids research* **48**(D1), 642–649 (2020)
42. Chicco, D., Jurman, G.: The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics* **21**(6) (2020)
43. Clauset, A., Newman, M.E., Moore, C.: Finding community structure in very large networks. *Physical review E* **70**(6), 066111 (2004)
44. Dennis, G., Sherman, B.T., Hosack, D.A., Yang, J., Gao, W., Lane, H.C., Lempicki, R.A.: David: database for annotation, visualization, and integrated discovery. *Genome biology* **4**(9), 1–11 (2003)
45. Bindea, G., Mlecnik, B., Hackl, H., Charoentong, P., Tosolini, M., Kirilovsky, A., Fridman, W.-H., Pagès, F., Trajanoski, Z., Galon, J.: Clueo: a cytoscape plug-in to decipher functionally grouped gene ontology and pathway annotation networks. *Bioinformatics* **25**(8), 1091–1093 (2009)
46. Gligorijevic, V., Barot, M., Bonneau, R.: deepNF: deep network fusion for protein function prediction. *Bioinformatics* **34**(22), 3873–3881 (2018)
47. Qiu, J., Dong, Y., Ma, H., Li, J., Wang, K., Tang, J.: Network embedding as matrix factorization: Unifying deepwalk, line, pte, and node2vec. In: *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, pp. 459–467 (2018)
48. Castells-Roca, L., García-Martínez, J., Moreno, J., Herrero, E., Bellí, G., Pérez-Ortín, J.E.: Heat shock response in yeast involves changes in both transcription rates and mRNA stabilities. *PLOS One* **6**(2), 17272 (2011)
49. Clauset, A., Newman, M.E.J., Moore, C.: Finding community structure in very large networks. *Physical Review E* **70**(6) (2004)
50. Morano, K.A., Grant, C.M., Moye-Rowley, W.S.: The response to heat shock and oxidative stress in *Saccharomyces cerevisiae*. *Genetics* **190**(4), 1157–1195 (2012)
51. Verghese, J., Abrams, J., Wang, Y., Morano, K.A.: Biology of the heat shock response and protein chaperones: budding yeast (*Saccharomyces cerevisiae*) as a model system. *Microbiology and Molecular Biology Reviews* **76**(2), 115–158 (2012)
52. Lee, D., Redfern, O., Orengo, C.: Predicting protein function from sequence and structure. *Nature reviews molecular cell biology* **8**(12), 995–1005 (2007)
53. Pastor-Flores, D., Ferrer-Dalmau, J., Bahí, A., Boleda, M., Biondi, R.M., Casamayor, A.: Depletion of yeast PDK1 orthologs triggers a stress-like transcriptional response. *BMC genomics* **16**(1), 1–21 (2015)
54. Oromendia, A.B., Dodgson, S.E., Amon, A.: Aneuploidy causes proteotoxic stress in yeast. *Genes & development* **26**(24), 2696–2708 (2012)
55. Segal, E., Shapira, M., Regev, A., Pe'er, D., Botstein, D., Koller, D., Friedman, N.: Module networks: identifying regulatory modules and their condition-specific regulators from gene expression data. *Nature genetics* **34**(2), 166–176 (2003)
56. Yamamoto, A., Mizukami, Y., Sakurai, H.: Identification of a novel class of target genes and a novel type of binding sequence of heat shock transcription factor in *Saccharomyces cerevisiae*. *Journal of Biological Chemistry* **280**(12), 11911–11919 (2005)
57. Matsumoto, R., Akama, K., Rakwal, R., Iwahashi, H.: The stress response against denatured proteins in the deletion of cytosolic chaperones SSA1/2 is different from heat-shock response in *Saccharomyces cerevisiae*. *BMC genomics* **6**(1), 1–15 (2005)
58. Düvel, K., Santhanam, A., Garrett, S., Schnepfer, L., Broach, J.R.: Multiple roles of Tap42 in mediating rapamycin-induced transcriptional changes in yeast. *Molecular cell* **11**(6), 1467–1478 (2003)
59. Berry, D.B., Gasch, A.P.: Stress-activated genomic expression changes serve a preparative role for impending stress in yeast. *Molecular biology of the cell* **19**(11), 4580–4587 (2008)
60. O'Duibhir, E., Lijnzaad, P., Benschop, J.J., Lenstra, T.L., van Leenen, D., Groot Koerkamp, M.J., Margaritis, T., Brok, M.O., Kemmeren, P., Holstege, F.C.: Cell cycle population effects in perturbation studies. *Molecular systems biology* **10**(6), 732 (2014)
61. Shivaswamy, S., Iyer, V.R.: Stress-dependent dynamics of global chromatin remodeling in yeast: dual role for SWI/SNF in the heat shock stress response. *Molecular and cellular biology* **28**(7), 2221–2234 (2008)
62. Spedale, G., Meddens, C.A., Koster, M.J., Ko, C.W., van Hooff, S.R., Holstege, F.C., Timmers, H.T.M., Pijnappel, W.P.: Tight cooperation between Mot1p and NC2 β in regulating genome-wide transcription, repression of transcription following heat shock induction and genetic interaction with SAGA. *Nucleic acids research* **40**(3), 996–1008 (2012)
63. Gasch, A.P., Spellman, P.T., Kao, C.M., Carmel-Harel, O., Eisen, M.B., Storz, G., Botstein, D., Brown, P.O.: Genomic expression programs in the response of yeast cells to environmental changes. *Molecular biology of the cell* **11**(12), 4241–4257 (2000)

Figures

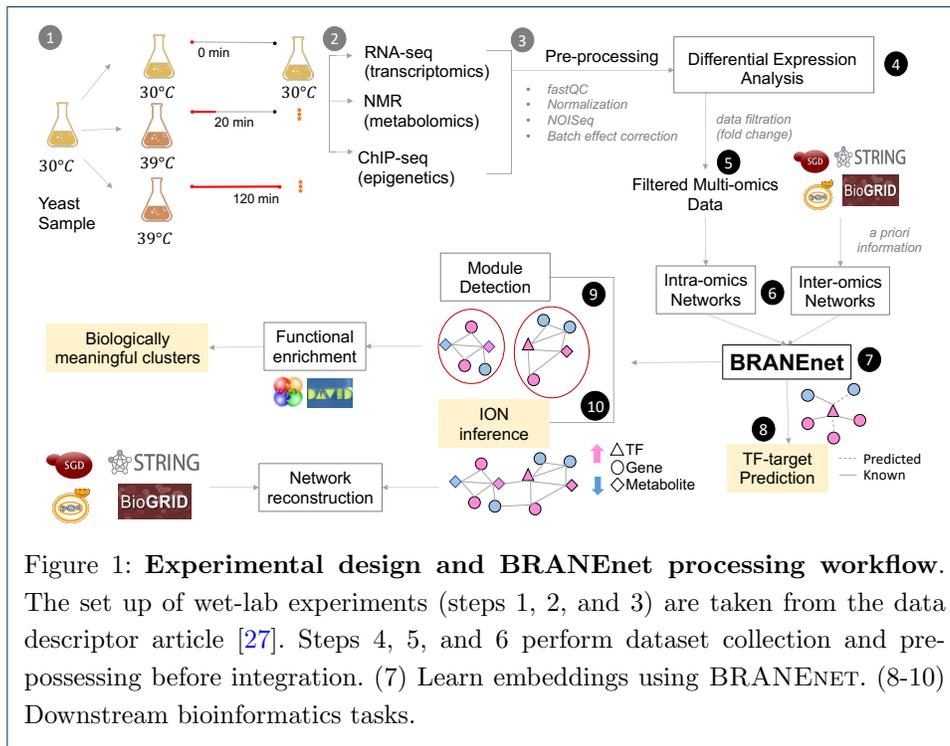
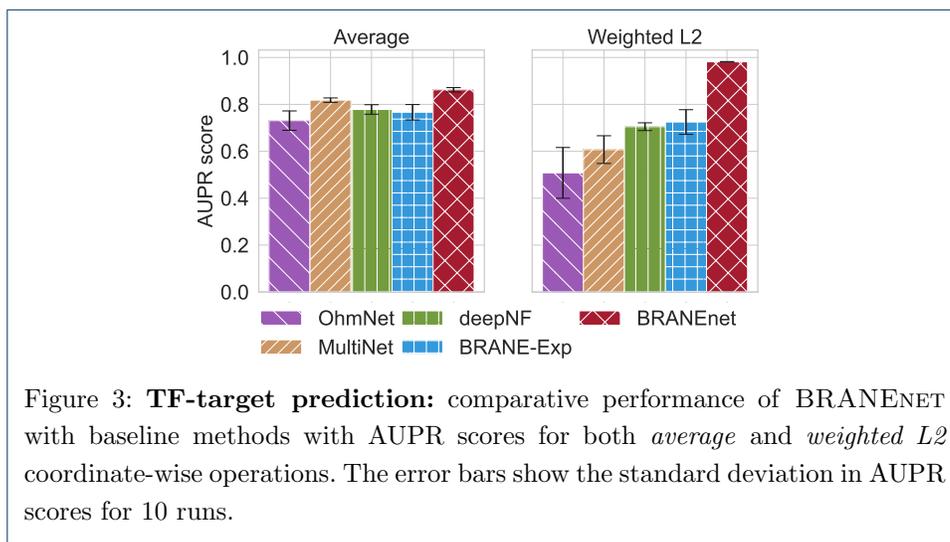
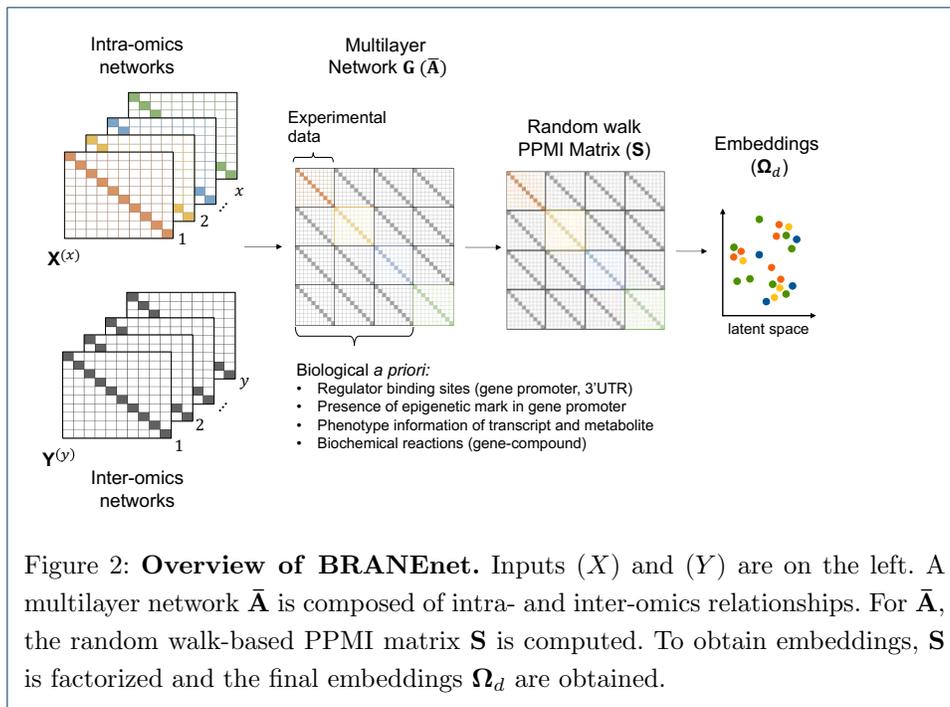


Figure 1: **Experimental design and BRANENet processing workflow.** The set up of wet-lab experiments (steps 1, 2, and 3) are taken from the data descriptor article [27]. Steps 4, 5, and 6 perform dataset collection and pre-processing before integration. (7) Learn embeddings using BRANENET. (8-10) Downstream bioinformatics tasks.



Tables

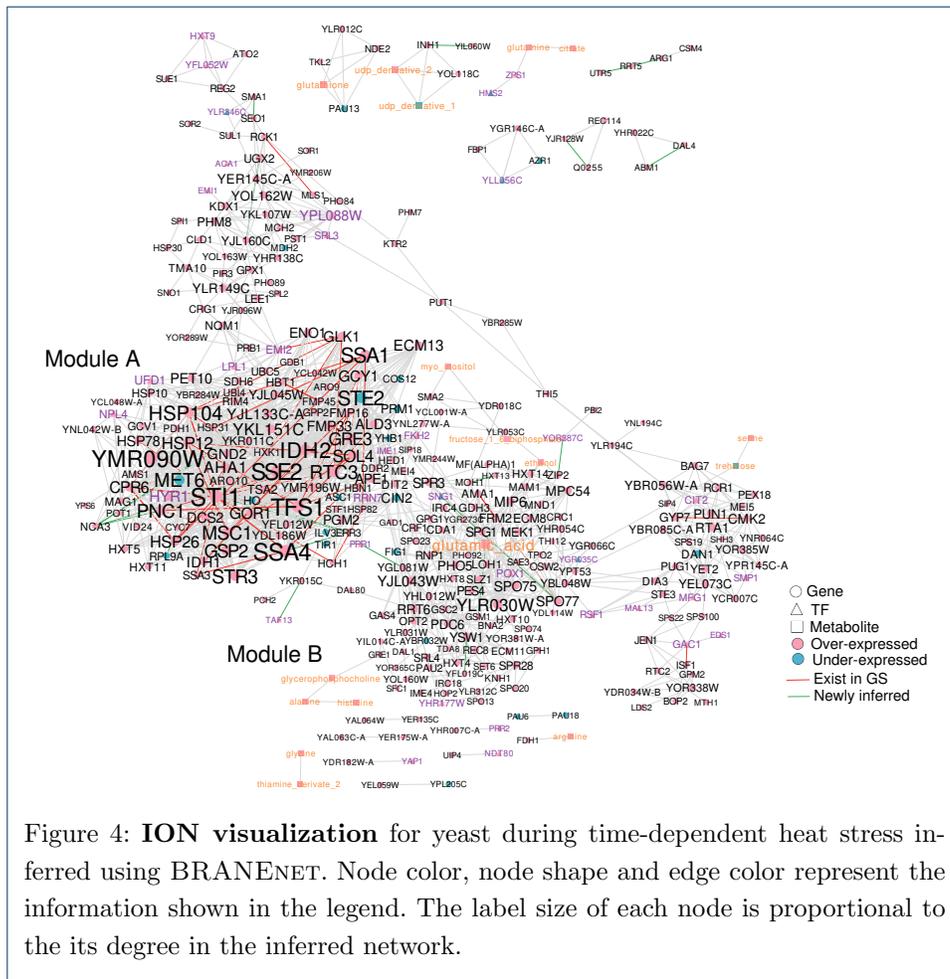


Figure 4: **ION visualization** for yeast during time-dependent heat stress inferred using BRANENET. Node color, node shape and edge color represent the information shown in the legend. The label size of each node is proportional to the its degree in the inferred network.

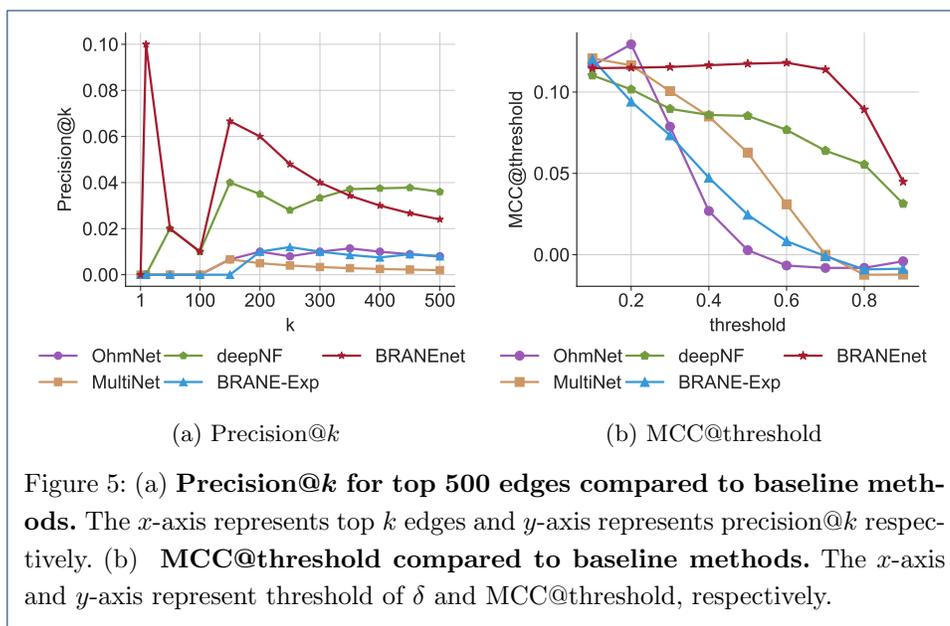


Figure 5: (a) **Precision@k** for top 500 edges compared to baseline methods. The x -axis represents top k edges and y -axis represents precision@ k respectively. (b) **MCC@threshold** compared to baseline methods. The x -axis and y -axis represent threshold of δ and MCC@threshold, respectively.

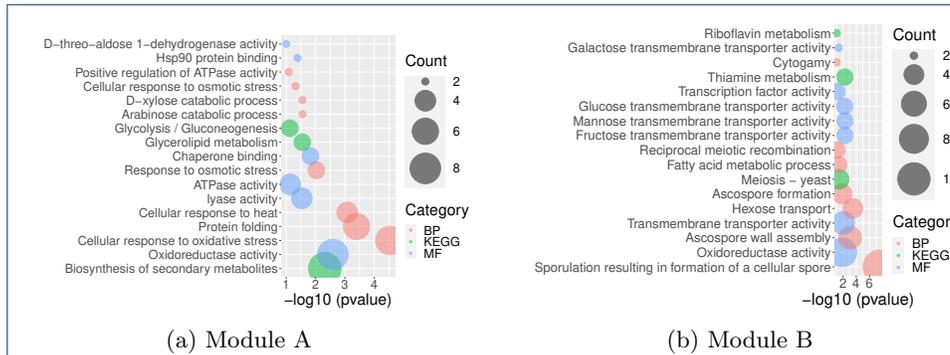


Figure 6: **Functional enrichment of modules A and B.** The y -axis represents the list of significantly enriched terms, while the x -axis shows their significance value ($-\log_{10}(p\text{-value})$). Different colors of circles indicate types of functional annotations: biological process (BP) in pink, molecular function (MF) in blue, and KEGG pathway in green. The size of circle represents the number of differentially expressed genes/TFs.

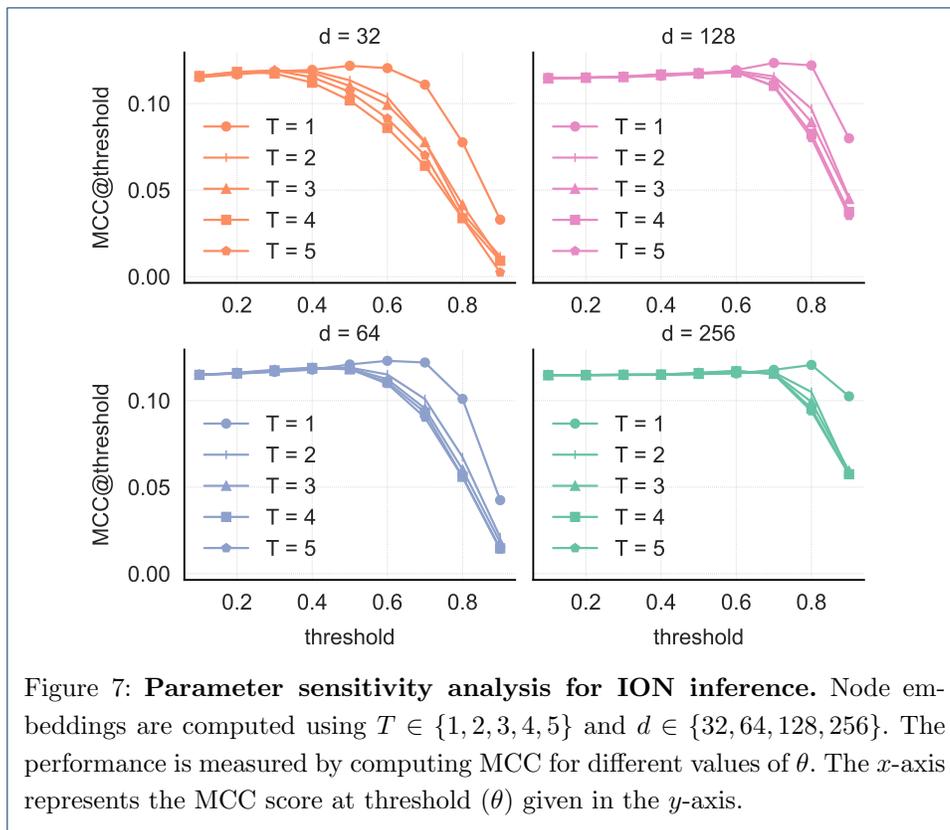


Table 1: **ION based identification of potential bio-markers.** The table provides the names, over- (\uparrow) or under- (\downarrow) expressed, node degree in ION (\mathcal{D}), function, cross references, and BRANENET module information comparison with external studies of potential bio-markers during heat stress response in yeast.

Name	\uparrow/\downarrow	\mathcal{D}	Function	Reference	BRANENET modules
STI1	\uparrow	40	Hsp90 cochaperone	[53, 54, 55, 56, 57, 58, 59]	A
SSA4	\uparrow	39	Heat shock protein	[59, 58, 57, 60, 55, 61, 62]	A
TFS1	\uparrow	46	Inhibitor of carboxy-peptidase Y , Ras GAP	[53, 54, 55, 57, 58, 59, 61]	A
YMR090W	\uparrow	50	Unknown function	[53, 58, 55]	A
SSE2	\uparrow	36	Hsp110 family member	[59, 58, 57, 54, 53, 55, 62, 61]	A
IDH2	\uparrow	37	Oxidative decarboxylation of isocitrate	[59, 58, 54, 55]	A
SSA1	\uparrow	40	ATPase	[59, 58, 54, 55, 61, 56]	A
STE2	\downarrow	36	Receptor for α -factor pheromone	[61]	A
HSP104	\uparrow	33	Disaggregase	[59, 58, 57, 54, 53, 55, 62, 61]	A
STI1	\uparrow	41	Hsp90 cochaperone	[53, 54, 55, 56, 57, 58, 59]	A
MET6	\downarrow	31	Cobalamin-independent methionine synthase	[58, 54]	A
STR3	\uparrow	30	Peroxisomal cystathionine beta-lyase	[58, 54, 61]	A
RTC3	\uparrow	28	Unknown function	[58, 54, 53, 55, 61]	A
MSC1	\uparrow	27	Unknown function	[59, 58, 57, 54, 53, 55, 61, 62]	A
PNC1	\uparrow	27	Nicotinamidase acid	[59, 58, 54, 55, 61, 62]	A
GSP2	\uparrow	30	GTP binding protein	[54, 55]	A
GRE3	\uparrow	31	Aldose reductase	[59, 58, 63, 57, 54, 61]	A
YLR030W	\uparrow	31	Unknown function	-	B
SOL4	\uparrow	32	6-phosphogluconolactonase	[58, 54, 53, 55, 61]	A
HSP12	\uparrow	28	Heat shock protein	[59, 58, 57, 54, 53, 55, 61, 62]	A
IDH1	\uparrow	25	Oxidative decarboxylation of isocitrate	[54, 55]	A

Table 2: **Parameter sensitivity analysis for TF-target prediction.** Node embeddings are computed using $T \in \{1, 2, 3, 4, 5\}$ and $d \in \{32, 64, 128, 256\}$. The performance is measured by computing AUPR score for *average* and *weighted L2* coordinate-wise operations.

T	Average				Weighted L2			
	d				d			
	32	64	128	256	32	64	128	256
1	0.700	0.880	0.880	0.870	0.980	0.982	0.983	0.983
2	0.790	0.820	0.850	0.860	0.916	0.945	0.966	0.968
3	0.830	0.850	0.870	0.870	0.956	0.967	0.979	0.979
4	0.780	0.810	0.850	0.860	0.852	0.938	0.966	0.968
5	0.820	0.840	0.860	0.870	0.952	0.968	0.981	0.981

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [BRANEnetsupplementaryfile.xlsx](#)
- [BRANetSurabhiJBMCBioinformatics.zip](#)