



OPEN

Novel cancer subtyping method based on patient-specific gene regulatory network

Mai Adachi Nakazawa¹, Yoshinori Tamada^{1,4}, Yoshihisa Tanaka^{2,3}, Marie Ikeguchi¹, Kako Higashihara¹ & Yasushi Okuno^{1,3}

The identification of cancer subtypes is important for the understanding of tumor heterogeneity. In recent years, numerous computational methods have been proposed for this problem based on the multi-omics data of patients. It is widely accepted that different cancer subtypes are induced by different molecular regulatory networks. However, only a few incorporate the differences between their molecular systems into the identification processes. In this study, we present a novel method to identify cancer subtypes based on patient-specific molecular systems. Our method realizes this by quantifying patient-specific gene networks, which are estimated from their transcriptome data, and by clustering their quantified networks. Comprehensive analyses of The Cancer Genome Atlas (TCGA) datasets applied to our method confirmed that they were able to identify more clinically meaningful cancer subtypes than the existing subtypes and found that the identified subtypes comprised different molecular features. Our findings also show that the proposed method can identify the novel cancer subtypes even with single omics data, which cannot otherwise be captured by existing methods using multi-omics data.

Cancer is a highly heterogeneous disease and is known to differ among patients. This heterogeneity renders one cancer type to be composed of multiple subtypes, which are characterized by different molecular features. Clinical identification of these cancer subtypes is currently one of the major challenges in cancer research. Identifying the subtypes can provide an understanding of the underlying molecular mechanisms and thereby design precise treatment strategies for efficient cancer management. In recent years, advances in high-throughput sequencing technologies have generated large amounts of data on various cancer types. For example, The Cancer Genome Atlas (TCGA) contains multi-omics data, including gene expression, mutation, methylation, and copy number, of over 34 cancer types. These multi-omics data allow improvements in cancer subtyping via computational methods^{1–3}. However, most studies do not identify the cancer subtypes based on differences in the molecular systems, but they are based only on the differences in the numerical patterns of the omics data.

Network representation of molecule-to-molecule relationships is a key to understanding a fundamental molecular system, and it plays an important role in understanding each biological process and the molecular mechanisms of cancer⁴. This idea can also be applied to obtain an understanding of complex human diseases, and is known as *network medicine* where the diseases are rarely caused by single molecular defects but are more likely driven by combinations of various biological processes^{4–11}. This concept has already been employed in recent cutting-edge research for discovering cancer-related genes^{12–14}. Therefore, knowledge of such networks could also be a promising data source for cancer subtyping. Some well-known types of biological networks are gene regulatory networks and protein–protein interaction networks¹⁵. Although the importance of molecular systems has been shown in recent years, only a few studies have incorporated the knowledge of molecular networks into their clustering processes^{16–18}. However, these methods do not sufficiently express the molecular systems for two reasons. First, the networks used do not contain a large number of genes that are supposed to be expressed in cells. In fact, only those genes that are already known to be involved in certain cancer types have been included in the networks^{16–18}. Second, the networks used are constructed from public databases that do not include condition-dependent networks^{19,20}. Recent studies have revealed that biological networks vary between

¹Graduate School of Medicine, Kyoto University, Kyoto 606-8507, Japan. ²Graduate School of Pharmaceutical Sciences, Kyoto University, Kyoto 606-8507, Japan. ³Biomedical Computational Intelligence Unit, HPC- and AI-driven Drug Development Platform Division, RIKEN Center for Computational Science, Kobe 650-0047, Japan. ⁴Present address: Innovation Center for Health Promotion, Hirosaki University, Hirosaki 036-8562, Japan. ✉email: y.tamada@hirosaki-u.ac.jp; okuno.yasushi.4c@kyoto-u.ac.jp

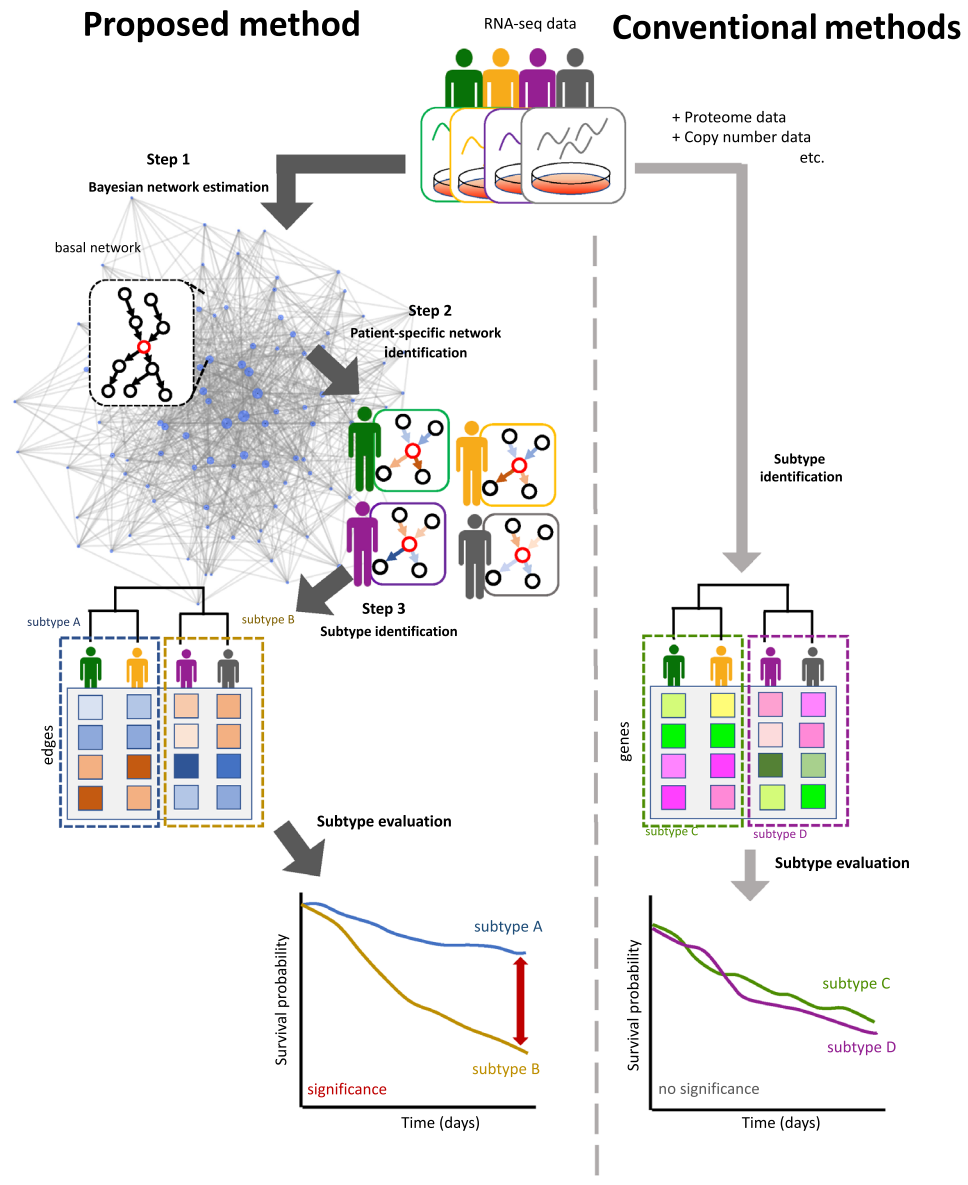


Figure 1. Overview of our method.

normal and the disease states^{21,22}. Because genetic interactions are condition-specific, the networks of particular types of cancers are different from those found in these databases.

In this study, we propose a novel method to identify cancer subtypes by incorporating differences in the molecular systems of the patients. Because a gene network involves gene–gene regulatory relationships and is a fundamental network among the various molecular networks, it can be an adequate representation of molecular systems for our purpose. Therefore, our proposed method is based on the estimated gene network from the gene expression data of patients. The main scheme of the proposed method is illustrated in Fig. 1. Briefly, our method estimates a gene network from a gene expression dataset in the tumors of patients using a Bayesian network. A numerical value is then calculated for every edge of the estimated gene network with respect to each patient's sample. This edge value, also known as the edge contribution value (ECv), is derived by evaluating the contribution of the edge to an expression value with respect to a patient in terms of the estimated molecular system²³. Therefore, differences in ECvs reflect differences in the molecular systems of particular patients. This calculation of ECvs generates a matrix of numerical values consisting of patient-specific networks. Finally, hierarchical clustering was performed using the matrix to categorize patients into subtypes. This simple clustering allows the identification of various subtypes, incorporating complex patient-specific activities of their molecular systems, which cannot be captured by the existing identification methods using the multi-omics data of patients. We used this method to analyze three cancer types from TCGA datasets, namely stomach adenocarcinoma, lung cancer including lung adenocarcinoma and lung squamous cell carcinoma, and breast cancer. Consequently, all three cancer types were grouped into three major novel subtypes, which defined their differential prognoses and distinct molecular properties. Our method identified system-based cancer subtypes using only transcriptome

data, which is more accessible compared to other omics data. Additionally, the proposed method allowed for the extraction of subnetworks to explain the features of the identified subtypes. Collectively, our findings indicate that the proposed method can successfully incorporate cancer-specific gene networks and establish a novel cancer subtype identification that overcomes the limitations of other sophisticated clustering methods, which are based on gene expression data alone.

Methods

In this chapter, we first introduce the gene network estimation method using a Bayesian network. We then explain the ECvs that allow us to quantify the patient-specific characteristics of the gene networks. Finally, the method for subtype identification of cancer patients based on their ECvs is described.

Bayesian network with *B*-spline nonparametric regression. To define the transcriptomic molecular networks, or gene networks, a Bayesian network with *B*-spline nonparametric regression model is used²⁴. A Bayesian network (BN) is a graphical model that represents the conditional independencies among variables as a directed acyclic graph. It can be considered as a hypothetical cause-and-effect regulatory relationship among genes. By representing the gene expressions as the random variables, we can estimate the system-level regulatory relationships from the transcriptome data. Many successful studies have reported on the use of BN for gene network analysis^{25–28}. Assuming p genes, the joint density of the gene expressions in a BN is described as

$$f(x_{i1}, \dots, x_{ip}; \theta_G) = \prod_{j=1}^p f(x_{ij} | pa_{ij}^G; \theta_j), \quad (1)$$

where x_{ij} represents the gene expression of the j -th gene at the i -th sample, θ_G is the parameter vector of the BN represented by G , $pa_{ij}^G = (pa_{i1}^{(j)}, \dots, pa_{i,q_j}^{(j)})$ denotes the gene expression vector of q_j parents of the j -th gene in the i -th sample, and θ_j is the parameter vector for the local density with respect to the j -th gene. The optimal structure of the network is obtained by the maximization of the posterior probability given the observed data as

$$p(G|X) \propto \pi(G) \int \prod_{i=1}^n f(x_{i1}, \dots, x_{ip}; \theta_G) \pi(\theta_G | \lambda) d\theta_G, \quad (2)$$

where X is the observed data matrix, $\pi(G)$ is the prior probability of G , n is the number of samples in X , $\pi(\theta_G | \lambda)$ denotes the prior distribution of θ_G , and λ is the hyperparameter vector. The drawback of the BN is that obtaining the optimal structure for a given dataset is NP-hard. Therefore, the neighbor node sampling and repeat (NNSR) algorithm was used²⁹.

Identification of patients' subtypes based on their molecular networks. Tanaka et al.²³ proposed the edge contribution value (ECv) to extract subnetworks from Bayesian networks, related to specific differences observed for in vitro experiments. Briefly, the *B*-spline nonparametric BN assumes that the gene expression is modeled as

$$x_{ij} = m_1^{(j)}(pa_{i1}^{(j)}) + \dots + m_{q_j}^{(j)}(pa_{i,q_j}^{(j)}) + \epsilon, \quad (3)$$

where $m_k^{(j)}(pa_{ik}^{(j)})$ is a regression function using *B*-spline curves for the k -th parent of the j -th gene, and ϵ is the error term. Because a value of this regression function $m_k^{(j)}(pa_{ik}^{(j)})$ can be considered as a contribution of an edge from the k -th parent to the j -th gene with respect to the i -th sample, Tanaka et al.²³ defined ECv as

$$ECv_{(i)}(jk \rightarrow j) = m_k^{(j)}(pa_{ik}^{(j)}), \quad (4)$$

where j_k represents the index of the k -th parent of the j -th gene. Tanaka et al.²³ considered the differences of ECvs as ΔECv between the control and the TGF β -treated samples, which extract the distinctive edges with a certain threshold for ΔECv . These were defined as the subnetworks characterizing the EMT in lung cancer cell lines. Here, we propose an algorithm that uses ECvs to characterize the patients (samples) and elucidate cancer subtypes. Using ECvs as *quantified* gene networks, patients with similar molecular systems would have similar ECvs, while patients with different molecular systems would have different *quantified* networks. In this context, clustering based on the molecular system differences enables us to identify cancer subtypes. This requires the gene network estimation from the gene expression data of patients and the calculation of ECvs values for the estimated edges.

Identification of cancer subtypes. Assuming E edges in the estimated gene network, the patient's *quantified* network is defined as a vector of E elements (c_{i1}, \dots, c_{iE}) , where $c_{i\nu}$ is an ECv of the ν -th edge for the i -th patient. The *quantified* networks of all the patients are collected and used to construct the columns of a matrix, resulting in an ECv matrix whose (i, ν) element corresponds to an ECv of the ν -th edge for the i -th patient. Using clustering, the patients of this ECv matrix are clustered according to the differences and similarities between their gene networks. As discussed by Tanaka et al.^{23,30}, the use of edges in the estimated network is different in different patients, and such differences of usages can be represented as different edge values using our ECvs. They also discussed that our Bayesian network model can express various patterns of usage or can capture the status of the cellular systems of patients^{23,30}. Except for ECv, there are no topological features for every edge for each patient in the network. Therefore ECvs are suitable for ranking edges and our patient-specific networks based on the ECv matrix are sufficiently different from each other to represent their differences, even though these

networks are derived from the same basal gene network. The ECv matrix consists of all the edges of the network, including approximately 20,000 genes and 150,000 edges. Since the majority of the edges do not represent differences in terms of ECv, parts of the edges are selected prior to clustering. To select the edges for hierarchical clustering based on the ECv matrix, the variance of each edge among patients is used as the ranking edges to represent the differences across the samples. The top N edges showing large variances will be selected. Therefore, hierarchical clustering is performed for the ECv matrix, consisting of the selected N edges, for the categorizing of patients into the different cancer subtypes.

Extraction of edges. Although hierarchical clustering categorizes patients into cancer subtypes, the part of the network that is affected by the clustering result is unknown. Therefore, distinctive edges with significant ECv differences need to be extracted, as in Tanaka et al.²³. In their study, they extracted the edges by calculating the Δ ECv between two conditions. However, since their method cannot be applied for more than two groups, we extended their scheme. The following proposed method allows us to extract distinctive edges with significant ECv differences using ECvs in multiple groups. Suppose that there are M groups of patients R_1, \dots, R_M . We define $\tilde{\Delta}$ ECv with respect to group R_r out of these M groups as

$$\tilde{\Delta}ECv^{R_r}(j_k \rightarrow j) = \left| \frac{1}{|R_r|} \sum_{i \in R_r} ECv_{(i)}(j_k \rightarrow j) - \frac{1}{\sum_{s \neq r} |R_s|} \sum_{t \in R_s, s \neq r} ECv_{(t)}(j_k \rightarrow j) \right|. \quad (5)$$

The $\tilde{\Delta}$ ECv of every single edge is then calculated with respect to each subtype, where significant $\tilde{\Delta}$ ECv edges are regarded as distinctive edges of specific subtypes.

Results

Dataset. In this study, our proposed method was applied to TCGA RNA-seq datasets of four types of cancer: stomach adenocarcinoma (STAD)³¹, lung cancer, including lung adenocarcinoma (LUAD)³² and lung squamous cell carcinoma (LUSC)³³; and breast cancer (BRCA)³⁴, which project is supported by the National Cancer Institute (NCI) and National Human Genome Research Institute (NHGRI). LUAD and LUSC were regarded as one dataset and referred to LUNG as they are both lung cancers, and we test whether they are split into different subtypes. These datasets were preprocessed as described in Supplementary information (Supplementary S2.1).

Identification of patients' subtypes based on ECv matrix. We calculated the ECv of every single edge in the estimated network for each patient, respectively. This ECv calculation generated a matrix of numerical values consisting of patient-specific molecular systems. To select the edges for hierarchical clustering based on the ECv matrix, the variances of edges were calculated as described in Method section. For TCGA datasets, we selected the top $N = 250$ edges with the highest variances in ECv among the patients. These 250 edges corresponded to approximately 0.01% of the total number of edges. Categorizing the patients by the small number of edges supposed to be a potentially better identification method. To categorize the patients into network-based subtypes, we performed hierarchical clustering for the ECv matrix consisting of the selected edges in three types of cancer. The ECv heatmap revealed a high variance among the intrinsic ECv matrix of the samples (Fig. 2a, Fig. S1, Fig. S8). The clustering results of the ECv heatmap indicate that patients of each cancer type were categorized into three major subtypes, namely subtype 1, subtype 2, and subtype 3, according to the similarities and differences between the patient networks (Fig. 2a, Fig. S1, Fig. S8). In our dataset, $N = 250$ was approximately the minimum number of edges that produced biologically and clinically meaningful results, as we described in Result section later. For the comparison, we also performed hierarchical clustering for the RNA-seq datasets (Fig. 2b), which is discussed in the later section.

Extraction subtype-specific edges. The Δ ECv value was calculated for every single edge in each subtype across three types of cancer. Edges with a high represent significant differences between subtypes. The distribution of suggested that only limited edges showed significant differences (Fig. 2c–e, Fig. S2). Based on the distributions of Δ ECv, the top 1.0% of the total edges in the estimated network were found to differ significantly. Therefore, we selected the corresponding edges from each subtype and removed any edges that were also selected for other subtypes (Fig. 2f, Fig. S3). We denoted the extracted edges as subtype-specific edges. Networks consisting of these subtype-specific edges were considered as a subnetwork characterizing the identified subtypes.

Applications in TCGA stomach cancer datasets. Stomach cancer is one of the most common leading causes of cancer-related death worldwide³⁵. In the original TCGA paper, the authors demonstrated that stomach cancer is a heterogeneous disease with four molecular subtypes—Epstein-Barr virus (EBV), microsatellite instability (MSI), genomically stable (GS), and chromosomal instability (CIN)³¹. According to the paper, the EBV subtype is enriched for high EBV burden and showed hypermethylation of the DNA, the MSI subtype showed elevated mutation rates and hypermethylation, the CIN subtype showed somatic copy-number aberrations, and the GS subtype did not show any somatic copy-number aberrations³¹. These subtypes were identified using iCluster and are based on the six platforms of the multi-omics molecular signature: somatic mutation, mRNA expression, miRNA expression, promoter methylation, somatic copy number alteration, and protein expression³¹. However, as previously mentioned³¹, no significance was observed between the prognoses of the subtypes (log-rank test p -value = 0.10 > 0.05) (Fig. 3a). This suggests that the multi-omics-based subtypes did not account for the clinical significance, such that subtyping may not provide an opportunity to improve thera-

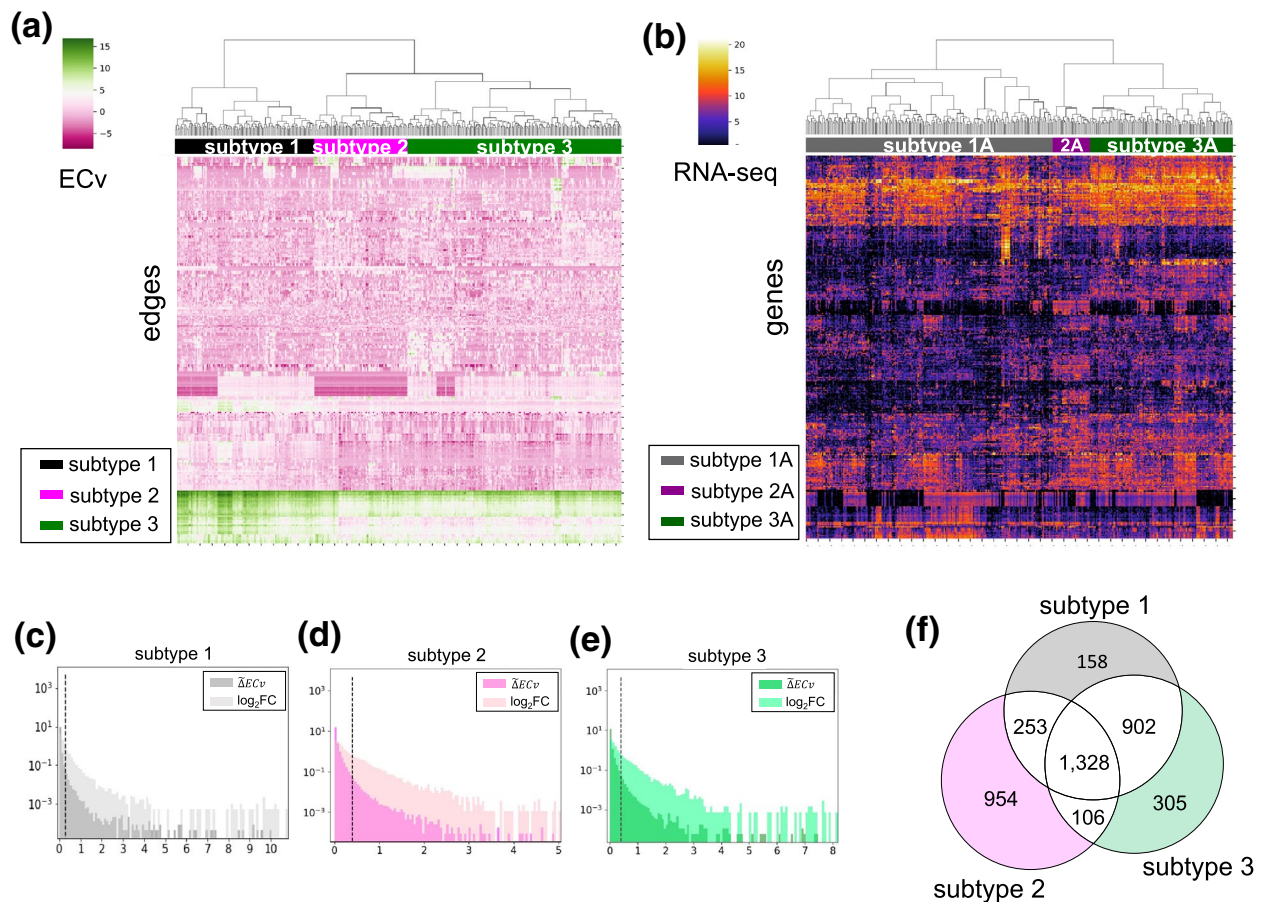


Figure 2. (a) Heatmap showing hierarchical clustering for the ECv matrix in the STAD dataset. (b) Heatmap showing hierarchical clustering for the RNA-seq matrix in the STAD dataset. (c–e) The distribution of ΔECv of edges and absolute \log_2 fold change in genes in the STAD dataset (See supplementary S2.3). Dashed lines represent of the top 1.0% of the total edges in every subtype. (f) The Venn diagram represents the number of edges in the STAD dataset. Colored areas in the Venn diagram represent subtype-specific edges in each subtype.

Subtype name	CIN	EBV	GS	MSI	Unknown	All
Subtype 1	33	10	37	9	24	113
Subtype 2	25	3	5	21	22	76
Subtype 3	58	11	6	15	83	173

Table 1. The relationship between the existing four molecular subtypes and our identified subtypes.

peutic treatments. We hypothesized that multi-omics data provide limited information on tumor subtyping. Rather, the differences in molecular systems may explain the differences in patients' prognoses.

To address this issue, we applied the proposed method to the preprocessed RNA-seq datasets of STAD. As described in Result section above, the clustering results of the ECv heatmap indicate that stomach cancer is grouped into three major subtypes: subtype 1 (113 samples), subtype 2 (76 samples), and subtype 3 (173 samples) (Fig. 2a). To determine the relationship between the existing multi-omics-based subtypes and our identified subtypes, we summarized the number of patients across them (Table 1) and found that our subtyping was different from the multi-omics-based subtypes. These findings suggest that our proposed method, based on the patient-specific molecular systems, can identify novel cancer subtypes that cannot be captured by existing methods using multi-omics data. To investigate the extent to which our proposed method identifies cancer subtypes, we conducted a survival analysis of the three identified subtypes. A better method is key for the identification of cancer subtypes and different prognoses, since patients with different molecular systems require different drug treatments. The Kaplan-Meier survival probability curves in the identified subtypes indicated that each subtype had a significantly different prognosis pattern (log-rank test p -value = 0.00011 < 0.05) (Fig. 3b).

Furthermore, to confirm whether the gene network information improves the identification of the cancer subtypes, hierarchical clustering was performed using RNA-seq data alone, without network information. The top 322 genes showing the highest variances of the RNA-seq data in STAD were selected, as the 250 edges with the

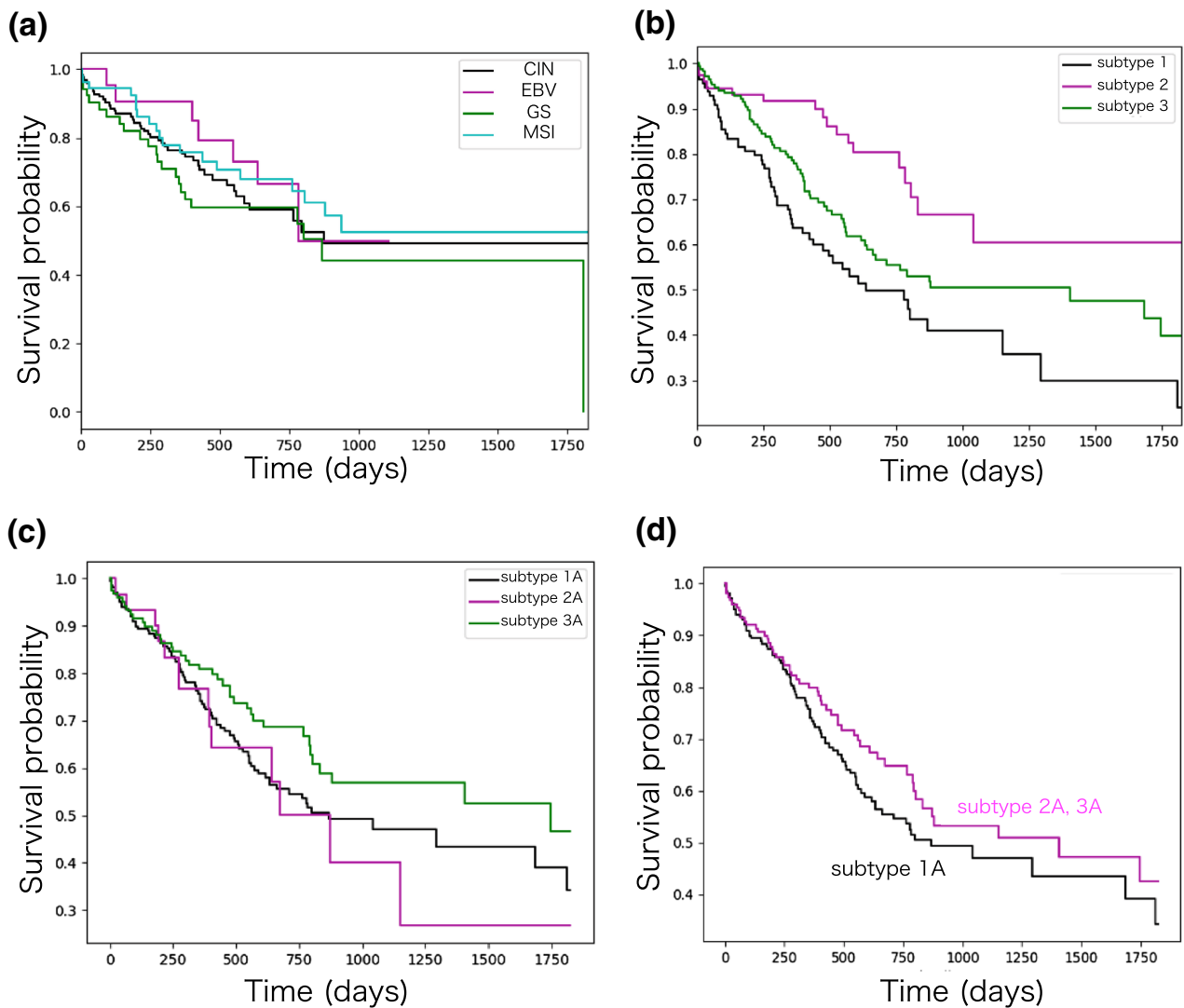


Figure 3. (a) Kaplan–Meier survival probability curves of patients for the multi-omics-based subtypes. The log-rank test p -value = 0.10 for the four subtypes. The log-rank test between two subtypes; 0.45 (CIN vs EBV) > 0.05, 0.30 (CIN vs GS) > 0.05, and 0.50 (CIN vs MSI) > 0.05, 0.31 (EBV vs GS) > 0.05, 0.95 (EBV vs MSI) > 0.05, 0.16 (GS vs MSI) > 0.05. (b) Kaplan–Meier survival probability curves of patients for the identified network-based subtypes. The log-rank test p -value = 0.00011 for the identified three subtypes. The log-rank test between two subtypes; 0.00016 (subtype 1 vs 2) < 0.05, 0.042 (subtype 1 vs 3) < 0.05, and 0.013 (subtype 2 vs 3) < 0.05. (c) Kaplan–Meier survival probability curves of patients for the identified RNA-seq based three subtypes. The log-rank test p -value = 0.06 for the identified three subtypes. The log-rank test between two subtypes; 0.70 (subtype 1 vs 2) > 0.05, 0.091 (subtype 1 vs 3) > 0.05, and 0.14 (subtype 2 vs 3) > 0.05. (d) Kaplan–Meier survival probability curves of patients for the identified RNA-seq based two major subtypes. The log-rank test p -value = 0.19 > 0.05.

ECv matrix were composed of 322 genes. Consequently, we identified three RNA-seq-based subtypes (Fig. 2b). However, these subtypes did not show any significant differences in terms of their prognoses (log-rank test p -value = 0.06 > 0.05) (Fig. 3c). Despite employing two major subtypes in the clustering result, the differences were not significant (Fig. 3d). To determine the relationship between the network-based and the RNA-seq-based subtypes, we summarized the number of patients across them (Table S2) and found that network-based subtypes were different from the RNA-seq-based subtypes. These results further suggest that our network-based method might generate a better cancer subtyping profile. Moreover to confirm whether the gene network information improves the identification of the cancer subtypes, we also applied the iNMF method. We compared our method with the iCluster method through four molecular subtypes identified by the original TCGA paper and applied another clustering method, the iNMF, as it is a successful method for cancer subtyping that can handle multi-omics data^{36,37}. We set three clusters when performing the iNMF as we identified three subtypes in our method. However, although we performed using gene expression data alone and using multi-omics data consisting of gene expression, miRNA expression, copy number, and DNA methylation, in both cases, these subtypes did not show any significant differences in their prognosis (Fig. S4).

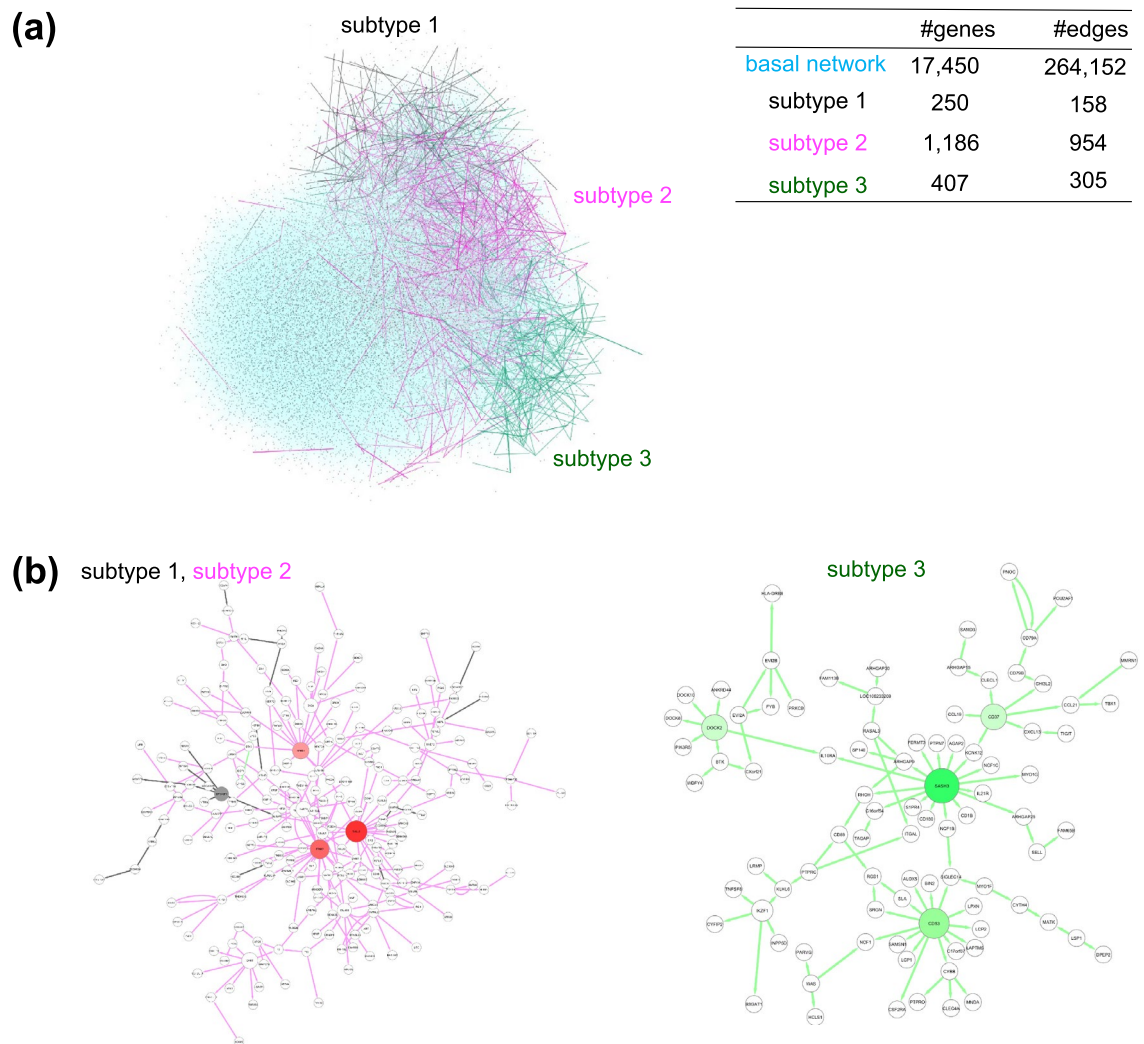


Figure 4. Visualization of subtype-specific subnetworks in the STAD dataset. **(a)** Subnetworks of subtype-specific edges were highlighted with the basal network (blue). **(b,c)** The biggest connected component in the subnetwork of subtype-specific edges in each subtype. Edges and nodes were colored by each subtype: subtype 1 (gray), subtype 2 (magenta), and subtype 3 (green). Colored nodes were hub nodes in each subtype and the color gradient represents the outdegree of hubs.

To characterize the network-based subtypes obtained from the ECv matrix, we highlighted the subnetworks composed of subtype-specific edges in the estimated basal network (Fig. 4a). The subtype-specific edges constituted a module in the basal network, especially in terms of subtype 3 (Figs. 2f, a). In Fig. 4a, the node layout of the basal network was arranged only using its topological structure without ECv information. These networks indicate that the extracted subnetworks form modules without being relaid. This suggests that the differences in the partial modules of the network might affect the identification of the cancer subtypes. Furthermore, to account for the properties of the identified subtypes, we verified their molecular features using gene ontology analysis. The subtype-specific networks were composed of 250, 1186, and 407 genes in subtype 1, subtype 2, and subtype 3, respectively. The ontology analysis results indicated that, according to the top five biological function terms, each subtype had a characteristic molecular feature (Fig. 5). Although most of the biological functions in the subtypes were related to development, the developmental stages or tissues varied between the subtypes. For example, “cardiovascular system development and function” was found in subtype 1, while “embryonic development” was found in subtype 2 and “cellular development” was found in subtype 3. In particular, the top five of biological functions in subtype 3, which were associated with a moderate prognosis, were completely different from those in the other subtypes. Most of the biological functions observed in subtype 1 and subtype 2 were related to development, while “cellular growth and proliferation” and “cell-to-cell signaling and interaction” were observed exclusively for subtype 3. Moreover, to demonstrate the clinical relevance of the identified subtypes, we summarized the number of patients across several clinical indices, such as sex, age at diagnosis, and tumor stage, in TCGA (Table S3-5). These results suggest that subtype 2 is characteristic of females, although the age at diagnosis and the tumor stage are not relevant to the identified subtypes (Table S3-5). Furthermore, we visualized a network composed of subtype-specific edges (Fig. 2f) and extracted the largest connected component

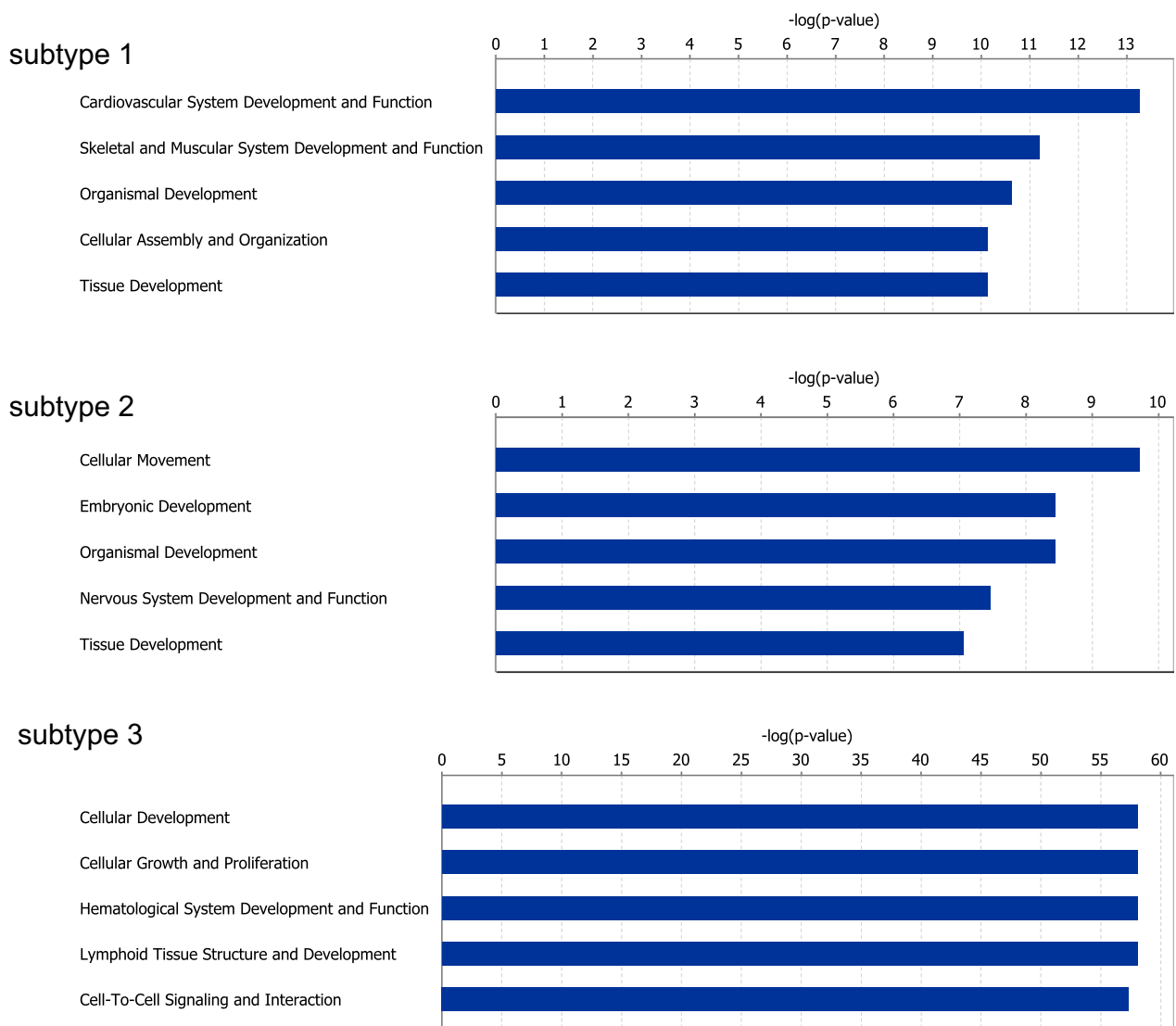


Figure 5. The top five terms of biological functions in the STAD dataset.

from the visualized network in each subtype (Fig. 4b,c). The results suggested that the subtype-specific edges in subtypes 2 and 3 were composed of a large connected subnetwork. Moreover, we found *SALL2*, *ETNK2*, and *APBB1* were located as the top hub genes in subtype 2. These genes are implied as cancer-related genes^{38–40}, and thus may play an important role in characterizing these subtypes.

Identification of cancer subtypes in lung cancer. To test the effectiveness of our method in other datasets, we applied it to the lung cancer dataset (LUNG). As shown in Fig. S1, the clustering results of the ECv heatmap indicate that lung cancer is grouped into three subtypes: subtype 1 (227 samples), subtype 2 (121 samples), and subtype 3 (343 samples). Then, survival analysis was conducted, similar to STAD analysis. The Kaplan-Meier survival probability curves in the identified subtypes indicate that each subtype has a significantly different pattern of prognosis (log-rank test p -value = $1.3e-14 < 0.05$) (Fig. S5a). Next, the differences in the molecular features between the identified subtypes were determined. The subtype-specific networks were found to be composed of 158, 1049, and 582 genes in subtype 1, subtype 2, and subtype 3, respectively. Gene ontology analysis of the subtypes indicate that all of the top five biological functions varied between them (Fig. S6). These findings suggest that our proposed method might also work for different cancer types. Moreover, hierarchical clustering followed by survival analysis, was performed for RNA-seq data as shown in the STAD section. Consequently, three RNA-seq-based subtypes were identified and the prognoses of these subtypes were significantly different (log-rank test p -value = $1.5e-15 < 0.05$) (Fig. S5b and S5c). While the patients in network-based subtype 1 were almost identical with those in RNA-seq-based subtype 1A, those in network-based subtype 2 and subtype 3 were different from RNA-seq-based subtypes (Table S6). Furthermore, the network-based and RNA-seq-based clustering could almost completely categorize LUAD and LUSC (Table S6). This may suggest that patients with various molecular features can be categorized even without network information, as LUAD and LUSC have characteristic molecular features that are significant for categorizing them using transcriptome data^{41,42}. Gene

ontology analysis indicates that network-based subtypes could reveal completely different characteristic molecular features among the subtypes. Thus, this suggests that our method can identify novel subtypes that cannot be detected using RNA-seq clustering. We also visualized a network composed of subtype-specific edges (Fig. S7a) and extracted the largest connected component from the visualized network in each subtype (Fig. S7b-d).

Identification of cancer subtypes in breast cancer. To illustrate the generalization of our method, we applied it to breast cancer datasets (BRCA). Breast cancer is well known as a heterogeneous disease with different molecular subtypes⁴³. As shown in Fig. S8, the clustering results of the ECV heatmap indicate that breast cancer is grouped into three subtypes: subtype 1 (192 samples), subtype 2 (464 samples), and subtype 3 (439 samples). Similar to the STAD and LUNG datasets, we then conducted a survival analysis of BRCA. The Kaplan-Meier survival probability curves for the identified subtypes indicated that subtype 3 had a different prognosis pattern and is almost statistically significant (log-rank test p -value = 0.0504) (Fig. S9). The differences in the molecular features between the identified subtypes were then determined. The subtype-specific networks were found to be composed of 1052, 463, and 953 genes in subtypes 1, 2, and 3, respectively. Gene ontology analysis of the subtypes indicated that all of the top five biological functions varied between them (Fig. S10). The results of the analysis indicated that each subtype, especially subtype 1, had a characteristic molecular feature. Next, hierarchical clustering followed by survival analysis, similar to the STAD and LUNG datasets, was performed for RNA-seq data. Consequently, three RNA-seq-based subtypes were identified and the prognoses of subtype 3A were significantly different (log-rank test p -value = 0.005 < 0.05) (Fig. S9b and S9c). To determine the relationship between the network-based and the RNA-seq-based subtypes, we summarized the number of patients across them (Table S7) and found that the patients in network-based subtype 1 were almost identical to those in RNA-seq-based subtype 1A. Consistent with the analysis of the LUNG datasets, these results suggest that patients with different molecular features can be grouped regardless of with network information or without it, and these subtypes show significantly different prognoses. According to the analysis of the LUNG datasets, network-based subtypes are similar to the RNA-seq-based subtypes in the patients with various molecular features; therefore, the results suggest that our method can also be applied to BRCA datasets. Furthermore, we also visualized a network composed of subtype-specific edges (Fig. S11a) and extracted the largest connected component from the visualized network for each subtype (Fig. S11b-d).

Discussion

In this study, we proposed a novel method for the identification of cancer subtypes based on patient-specific molecular systems. The proposed method is able to identify novel subtypes with different prognoses, as well as the differences in molecular properties between stomach cancer, lung cancer, and breast cancer. Differences in molecular systems are not necessarily associated with the prognoses of patients. However, it is likely to affect the effectiveness and/or medical treatment options available for these patients. For this reason, our novel subtypes may be related to prognosis of patient.

Although many types of omics data are currently available, it remains difficult to integrate multi-omics data in research. Each type of omics data can be used to categorize cancers into various subtypes in terms of prognosis, pathological findings, and others. However, since our proposed method uses only transcriptome data, even though our gene network-based method was successful, it may not be sufficient to obtain an in-depth understanding of the molecular systems. Despite this, changes in the different layers of omics networks influence the transcriptome profile at some level. This could explain why our proposed method, based on the gene network, was able to identify novel cancer subtypes using only the transcriptome data. There, however, remains room for improvement in the method reported in this study, wherein identification using multi-omics data based on estimated systems represents an informative strategy for the identification of cancer subtypes.

Data availability

All the patient lists generated in this study are provided in the supplementary data. All the networks are available at NDEX (The basal network in STAD; <https://www.ndexbio.org/viewer/networks/1dabd135-8bab-11eb-9e72-0ac135e8bacf>, The subtype-specific network in STAD; <https://www.ndexbio.org/viewer/networks/4e61c7cf-8889-11eb-9e72-0ac135e8bacf>, The basal network in LUNG; <https://www.ndexbio.org/viewer/networks/0e943431-8e00-11eb-9e72-0ac135e8bacf>, The subtype-specific network in LUNG; <https://www.ndexbio.org/viewer/networks/be019ba4-8e01-11eb-9e72-0ac135e8bacf>, The basal network in BRCA; <https://www.ndexbio.org/viewer/networks/6c7d71dd-1ebf-11ec-9fe4-0ac135e8bacf>, The subtype-specific network in BRCA; <https://www.ndexbio.org/viewer/networks/e9435f90-1ec5-11ec-9fe4-0ac135e8bacf>).

Received: 14 April 2021; Accepted: 12 November 2021

Published online: 08 December 2021

References

- Shen, R., Olshen, A. B. & Ladanyi, M. Integrative clustering of multiple genomic data types using a joint latent variable model with application to breast and lung cancer subtype analysis. *Bioinformatics* **25**, 2906–2912. <https://doi.org/10.1093/bioinformatics/btp543> (2009).
- Mo, Q. *et al.* Pattern discovery and cancer gene identification in integrated cancer genomic data. *Proc. Natl. Acad. Sci. USA* **110**, 4245–4250. <https://doi.org/10.1073/pnas.1208949110> (2013).
- Gao, Y. & Church, G. Improving molecular cancer class discovery through sparse non-negative matrix factorization. *Bioinformatics* **21**, 3970–3975. <https://doi.org/10.1093/bioinformatics/bti653> (2005).
- Barabási, A. L., Gulbahce, N. & Loscalzo, J. Network medicine: A network-based approach to human disease. *Nat. Rev. Genet.* **12**, 56–68. <https://doi.org/10.1038/nrg2918> (2011).

5. Conte, F. *et al.* A paradigm shift in medicine: A comprehensive review of network-based approaches. *Biochimica et Biophysica Acta - Gene Regulatory Mechanisms* **1863**, 194416. <https://doi.org/10.1016/j.bbagrm.2019.194416> (2020).
6. Paci, P. *et al.* Gene co-expression in the interactome: moving from correlation toward causation via an integrated approach to disease module discovery. *npj Syst. Biol. Appl.* **7**, 1–11. <https://doi.org/10.1038/s41540-020-00168-0> (2021).
7. Fiscon, G. & Paci, P. SAveRUNNER: An R-based tool for drug repurposing. *BMC Bioinf.* **22**, 150. <https://doi.org/10.1186/s12859-021-04076-w> (2021).
8. Fiscon, G., Conte, F., Farina, L. & Paci, P. SAveRUNNER: A network-based algorithm for drug repurposing and its application to COVID-19. *PLoS Comput. Biol.* **17**, e1008686. <https://doi.org/10.1371/JOURNAL.PCBI.1008686> (2021). [arXiv:2006.03110](https://arxiv.org/abs/2006.03110).
9. Yu, X., Zeng, T., Wang, X., Li, G. & Chen, L. Unravelling personalized dysfunctional gene network of complex diseases based on differential network model. *J. Transl. Med.* **13**, 189. <https://doi.org/10.1186/s12967-015-0546-5> (2015).
10. Zhang, W., Zeng, T., Liu, X. & Chen, L. Diagnosing phenotypes of single-sample individuals by edge biomarkers. *J. Mol. Cell Biol.* **7**, 231–241. <https://doi.org/10.1093/jmcb/mjv025> (2015).
11. Kuijjer, M. L., Tung, M. G., Yuan, G. C., Quackenbush, J. & Glass, K. Estimating sample-specific regulatory networks. *iScience* **14**, 226–240. <https://doi.org/10.1016/j.isci.2019.03.021> (2019). [arXiv:1505.06440](https://arxiv.org/abs/1505.06440).
12. Fiscon, G., Conte, F., Licursi, V., Nasi, S. & Paci, P. Computational identification of specific genes for glioblastoma stem-like cells identity. *Sci. Rep.* **8**, 7769. <https://doi.org/10.1038/s41598-018-26081-5> (2018).
13. Panebianco, V. *et al.* Prostate cancer screening research can benefit from network medicine: An emerging awareness. *npj Syst. Biol. Appl.* **6**, 13. <https://doi.org/10.1038/s41540-020-0133-0> (2020).
14. Falcone, R. *et al.* BRAF V600E -mutant cancers display a variety of networks by SWIM analysis: Prediction of vemurafenib clinical response. *Endocrine* **64**, 406–413. <https://doi.org/10.1007/s12020-019-01890-4> (2019).
15. Yu, D., Kim, M., Xiao, G. & Hwang, T. H. Review of biological network data and its applications. *Genom. Inf.* **11**, 200–210. <https://doi.org/10.5808/gi.2013.11.4.200> (2013).
16. Xu, T. *et al.* Identifying cancer subtypes from miRNA-TFmRNA regulatory networks and expression data. *PLoS ONE* **11**, e0152792. <https://doi.org/10.1371/journal.pone.0152792> (2016).
17. Guo, Y., Qi, Y., Li, Z. & Shang, X. Improvement of cancer subtype prediction by incorporating transcriptome expression data and heterogeneous biological networks. *BMC Med. Genom.* **11**, 119. <https://doi.org/10.1186/s12920-018-0435-x> (2018).
18. Liu, Y., Gu, Q., Hou, J. P., Han, J. & Ma, J. A network-assisted co-clustering algorithm to discover cancer subtypes based on gene expression. *BMC Bioinf.* **15**, 37. <https://doi.org/10.1186/1471-2105-15-37> (2014).
19. Kanehisa, M., Goto, S., Sato, Y., Furumichi, M. & Tanabe, M. KEGG for integration and interpretation of large-scale molecular data sets. *Nucleic Acids Res.* **40**, D109–D114. <https://doi.org/10.1093/nar/gkr988> (2012).
20. Schaefer, C. F. *et al.* PID: The pathway interaction database. *Nucleic Acids Res.* **37**, D674–D679. <https://doi.org/10.1093/nar/gkn653> (2009).
21. Singh, A. J., Ramsey, S. A., Filtz, T. M. & Kiuoussi, C. Differential gene regulatory networks in development and disease. *Cell. Mol. Life Sci.* **75**, 1013–1025. <https://doi.org/10.1007/s00018-017-2679-6> (2018).
22. Ideker, T. & Krogan, N. J. Differential network biology. *Mol. Syst. Biol.* **8**, 565. <https://doi.org/10.1038/msb.2011.99> (2012).
23. Tanaka, Y., Tamada, Y., Ikeguchi, M., Yamashita, F. & Okuno, Y. System-based differential gene network analysis for characterizing a sample-specific subnetwork. *Biomolecules* **10**, 306. <https://doi.org/10.3390/biom10020306> (2020).
24. Imoto, S., Goto, T. & Miyano, S. Estimation of genetic networks and functional structures between genes by using Bayesian networks and nonparametric regression. *Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing* 175–186, https://doi.org/10.1142/9789812799623_0017 (2002).
25. Wang, L. *et al.* Cell cycle gene networks are associated with melanoma prognosis. *PLoS ONE* **7**, e34247. <https://doi.org/10.1371/journal.pone.0034247> (2012).
26. Arima, C. *et al.* Lung adenocarcinoma subtypes definable by lung development-related miRNA expression profiles in association with clinicopathologic features. *Carcinogenesis* **35**, 2224–2231. <https://doi.org/10.1093/carcin/bgu127> (2014).
27. Gendelman, R. *et al.* Bayesian network inference modeling identifies TRIB1 as a novel regulator of cell-cycle progression and survival in cancer cells. *Can. Res.* **77**, 1575–1585. <https://doi.org/10.1158/0008-5472.CAN-16-0512> (2017).
28. Creixell, P. *et al.* Pathway and network analysis of cancer genomes. *Nat. Methods* **12**, 615–621. <https://doi.org/10.1038/nmeth.3440> (2015).
29. Tamada, Y. *et al.* Estimating genome-wide gene networks using nonparametric bayesian network models on massively parallel computers. *IEEE/ACM Trans. Comput. Biol. Bioinf.* **8**, 683–697. <https://doi.org/10.1109/TCBB.2010.68> (2011).
30. Tanaka, Y. *et al.* Dynamic changes in gene-to-gene regulatory networks in response to SARS-CoV-2 infection. *Sci. Rep.* **11**, 11241. <https://doi.org/10.1038/s41598-021-90556-1> (2021).
31. The Cancer Genome Atlas Research Network. Comprehensive molecular characterization of gastric adenocarcinoma. *Nature* **513**, 202–209. <https://doi.org/10.1038/nature13480> (2014).
32. The Cancer Genome Atlas Research Network. Comprehensive molecular profiling of lung adenocarcinoma. *Nature* **511**, 543–550. <https://doi.org/10.1038/nature13385> (2014).
33. The Cancer Genome Atlas Research Network. Comprehensive genomic characterization of squamous cell lung cancers. *Nature* **489**(7417) 519–525 <https://doi.org/10.1038/nature11404> (2012).
34. The Cancer Genome Atlas Research Network. Comprehensive molecular portraits of human breast tumours. *Nature* **490**, 61–70. <https://doi.org/10.1038/nature11412> (2012).
35. Rawla, P. & Barsouk, A. Epidemiology of gastric cancer: Global trends, risk factors and prevention. *Przegląd Gastroenterologiczny* **14**, 26–38. <https://doi.org/10.5114/pg.2018.80001> (2019).
36. Yang, Z. & Michailidis, G. A non-negative matrix factorization method for detecting modules in heterogeneous omics multi-modal data. *Bioinformatics* **32**, 1–8. <https://doi.org/10.1093/bioinformatics/btv544> (2016).
37. Schlicker, A. *et al.* Subtypes of primary colorectal tumors correlate with response to targeted treatment in colorectal cell lines. *BMC Med. Genom.* **5**, 66. <https://doi.org/10.1186/1755-8794-5-66> (2012).
38. Hermosilla, V. E. *et al.* Developmental SALL2 transcription factor: A new player in cancer. *Carcinogenesis* **38**, 680–690. <https://doi.org/10.1093/carcin/bgx036> (2017).
39. Lee, J. H. *et al.* APBB1 reinforces cancer stem cell and epithelial-to-mesenchymal transition by regulating the IGF1R signaling pathway in non-small-cell lung cancer cells. *Biochem. Biophys. Res. Commun.* **482**, 35–42. <https://doi.org/10.1016/j.bbrc.2016.11.030> (2017).
40. Li, L., Mou, Y. P., Wang, Y. Y., Wang, H. J. & Mou, X. Z. miR-199a-3p targets ETNK1 to promote invasion and migration in gastric cancer cells and is associated with poor prognosis. *Pathol. Res. Pract.* **215**, 152511. <https://doi.org/10.1016/j.prp.2019.152511> (2019).
41. Chen, M., Liu, X., Du, J., Wang, X. J. & Xia, L. Differentiated regulation of immune-response related genes between LUAD and LUSC subtypes of lung cancers. *Oncotarget* **8**, 133–144. <https://doi.org/10.18632/oncotarget.13346> (2017).
42. Wang, C. *et al.* RNA-Seq profiling of circular RNA in human lung adenocarcinoma and squamous cell carcinoma. *Mol. Cancer* **18**, 134. <https://doi.org/10.1186/s12943-019-1061-8> (2019).
43. Grimaldi, A. M. *et al.* The new paradigm of network medicine to analyze breast cancer phenotypes. *Int. J. Mol. Sci.* **21**, 6690. <https://doi.org/10.3390/ijms21186690> (2020).

Acknowledgements

We thank Dr. K. Fukuyama for advice on the clinical interpretation and members of the Okuno Lab. and Project 20 in Life Intelligence Consortium for their helpful discussions. Computational resources were provided by the Super Computer System, Human Genome Center, Institute of Medical Science, University of Tokyo. This work was supported by Cabinet Office, Government of Japan, Public/Private R&D Investment Strategic Expansion Program (PRISM). This work was supported by RIKEN Junior Research Associate Program. This work used computational resources of the supercomputer Fugaku provided by RIKEN through the HPCI System Research Project (Project ID: hp210145). The results here are in whole or part based upon data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>.

Author contributions

M.A.N, Y.T. (Yoshinori Tamada), Y.O. conceived the experiments, analyzed the result, and wrote the manuscript; Y.T. (Yoshihisa Tanaka), M.I., and K.H. conducted experiments; Y.T. (Yoshihisa Tanaka) and Y.O. reviewed and edited the manuscript.

Competing interests

Y. Tamada and Y. Okuno have a patent application on the method for identification of patient-specific network used in this study through the technology licensing organization in Kyoto University. Other authors declare no conflict of interest.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-02394-w>.

Correspondence and requests for materials should be addressed to Y.T. or Y.O.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021