

TECHNICAL NOTE

TuBA: Tunable biclustering algorithm reveals clinically relevant tumor transcriptional profiles in breast cancer

Amartya Singh ^{1,2}, Gyan Bhanot ^{1,2,3} and Hossein Khiabani ^{1,2,3,4,*}¹Department of Physics and Astronomy, Rutgers University, 136 Frelinghuysen Rd, Piscataway, NJ 08854;²Center for Systems and Computational Biology, Rutgers Cancer Institute, Rutgers University, 195 LittleAlbany St, New Brunswick, NJ 08903; ³Department of Molecular Biology and Biochemistry, Rutgers University,604 Allison Rd, Piscataway, NJ 08854 and ⁴Department of Pathology and Laboratory Medicine, Rutgers Robert Wood Johnson Medical School, Rutgers University, One Robert Wood Johnson Place, New Brunswick, NJ, 08903

*Correspondence address. Hossein Khiabani, Rutgers Cancer Institute of New Jersey, Rutgers University, 195 Little Albany St, New Brunswick, NJ 08903-2681. Tel: +(732) 235 7554; E-mail: h.khiabani@rutgers.edu  <http://orcid.org/0000-0003-1446-4394>

Abstract

Background: Traditional clustering approaches for gene expression data are not well adapted to address the complexity and heterogeneity of tumors, where small sets of genes may be aberrantly co-expressed in specific subsets of tumors.

Biclustering algorithms that perform local clustering on subsets of genes and conditions help address this problem. We propose a graph-based Tunable Biclustering Algorithm (TuBA) based on a novel pairwise proximity measure, examining the relationship of samples at the extremes of genes' expression profiles to identify similarly altered signatures. **Results:** TuBA's predictions are consistent in 3,940 breast invasive carcinoma samples from 3 independent sources, using different technologies for measuring gene expression (RNA sequencing and Microarray). More than 60% of biclusters identified independently in each dataset had significant agreement in their gene sets, as well as similar clinical implications.

Approximately 50% of biclusters were enriched in the estrogen receptor–negative/HER2-negative (or basal-like) subtype, while >50% were associated with transcriptionally active copy number changes. Biclusters representing gene co-expression patterns in stromal tissue were also identified in tumor specimens. **Conclusions:** TuBA offers a simple biclustering method that can identify biologically relevant gene co-expression signatures not captured by traditional unsupervised clustering approaches. It complements biclustering approaches that are designed to identify constant or coherent submatrices in gene expression datasets, and outperforms them in identifying a multitude of altered transcriptional profiles that are associated with observed genomic heterogeneity of diseased states in breast cancer, both within and across tumor subtypes, a promising step in understanding disease heterogeneity, and a necessary first step in individualized therapy.

Keywords: clustering; gene co-expression; tumor heterogeneity; copy number aberrations; breast invasive carcinoma

Background

The first step in organizing and analyzing high-throughput gene expression datasets is to group together (cluster) genes, or samples based on some mathematical measure of similarity between the respective entities of interest. Because a priori knowledge about both the relevant genes and the unique phenotypic characteristics of samples is usually limited, clustering is often

performed in an unsupervised manner [1–5]. Quite frequently, measures of similarity (e.g., the Pearson correlation coefficient, Spearman correlation coefficient, mutual information) are used to quantify the level of similarity between every pair of genes (or samples) across all samples (or genes). Such an approach is known as global clustering. In case of datasets with a heterogeneous assortment of samples, only a small subset of genes in a fraction of the total set of samples may be co-regulated in spe-

Received: 20 February 2019; Revised: 17 April 2019; Accepted: 6 May 2019

© The Author(s) 2019. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

cific cellular processes. This is especially true for diseases like cancer that manifest a plethora of diseased phenotypes. In gene expression datasets comprising tumor samples, depending on the heterogeneity of the diseased states, there may be multiple distinct transcriptional alterations that are exhibited by multiple (not necessarily exclusive) subsets of the tumor samples. Moreover, it is well known that even in normal cells, the same genes can regulate and participate in multiple distinct pathways, depending on the context. Therefore, global clustering is not an optimal approach to identify co-expressed sets of genes or samples in gene expression datasets associated with heterogeneous diseases.

To address these concerns, a variety of local biclustering algorithms have been proposed that satisfy the following requirements: (i) a cluster of genes is defined with respect to only a subset of conditions (patient samples) and vice versa, and (ii) the clusters are not exclusive and/or exhaustive—i.e., a gene/condition may belong to >1 cluster or to none at all [6–9]. Based on the type of biclusters and the mathematical formulation used to discover them, biclustering techniques are categorized by Oghabian et al. [10] into 4 classes: (i) correlation maximization methods that identify subsets of genes and samples where the expression values of genes (or samples) is highly correlated across samples (or genes) [6], (ii) variance minimization methods that identify biclusters where the expression values have low variance among the selected genes or conditions or both [11], (iii) 2-way clustering methods that iteratively perform 1-way clustering on the genes and samples [12], and (iv) probabilistic and generative methods that employ stochastic approaches to discover genes (or samples) that are similarly expressed in subsets of samples (or genes) [13, 14]. Another classification scheme proposed by Pontes et al. [15, 16] categorizes the generated biclusters based on their gene expression patterns into 4 classes: (i) biclusters with constant values, (ii) biclusters with constant values on rows (genes) or columns (conditions), (iii) biclusters with additive and/or multiplicative relationships between genes and conditions, and (iv) biclusters based on evidence that a subset of genes is up-regulated or down-regulated across a subset of conditions without taking into account actual expression values; data in such biclusters do not follow any mathematical model.

In this paper, we introduce a graph-based method called the Tunable Biclustering Algorithm (TuBA), which discovers biclusters consistent with the latter category. TuBA is based on a novel measure of proximity that identifies aberrantly co-expressed gene sets within subsets of tumor samples that correspond to the expression extremals for the genes. A key feature of the proximity measure used in TuBA is that it does not rely explicitly on the actual gene expression values. We demonstrate the utility of TuBA by applying it to 3 large, independent cohorts of breast invasive carcinoma (BRCA) encompassing 3,940 patients. In addition to detecting known pathways and subtypes associated with breast cancer, TuBA was able to uncover several novel sets of co-expressed genes across subtypes that may be relevant as biomarkers for therapeutic identification and intervention.

Methods

Proximity measure

TuBA's proximity measure addresses the following question: in a given gene expression dataset, which genes exhibit higher (or lower) expression levels in the same subset of samples relative to the rest? In other words, if we only consider the top (or bot-

tom) x percentile samples for every gene, which gene pairs share a significant number of samples between their percentile sets? The number of samples shared between any pair of percentile sets follows the hypergeometric distribution; therefore, we can compute the significance (P -value) of overlaps between pairs of percentile sets based on the numbers of shared samples by using the 1-sided Fisher's exact test. Thus, TuBA's proximity measure between 2 genes is defined by the significance of overlaps between their respective percentile sets (Fig. 1).

In a real biological dataset, we expect the following 2 scenarios to arise: (i) subsets of genes associated with particular biological processes/pathways are co-expressed in all samples. In this case, it is reasonable to expect a significant agreement between the sets of samples that exhibit higher (or lower) expression levels of the involved genes. (ii) Alternatively, subsets of genes may be dysregulated via shared underlying mechanisms, such that their expression levels are higher (or lower) compared with the rest of the samples that are not influenced by that mechanism. The latter case is of particular interest for datasets associated with diseased states, especially cancers, because these gene co-expression signatures and their underlying mechanisms could help us identify potential biomarkers with prognostic and/or predictive value. This is the basic motivation behind standard differential expression analyses as well [17]. However, unlike usual differential co-expression analyses, our proximity measure does not rely on any pre-specification of subtypes.

A salient feature of our proximity measure is that it does not model the distributions of the measured expression levels of genes across samples. Moreover, it does not rely on significant differences between the expression levels of genes in samples comprising the extremal sets vs the rest of the samples. Thus, biologically relevant gene co-expression signatures can be identified without restricting the analysis exclusively to genes that exhibit differential expression across subsets of samples. In the case of tumor datasets, this increases the likelihood for identification of gene co-expression signatures associated with the microenvironment. Another salient feature of our proximity measure is that there is no penalty for relative changes in ranks of samples in the respective percentile sets. This is important because, even if the ranks of matching samples are significantly different in the 2 percentile sets, there is still valuable information to be gleaned by virtue of the fact that these subsets of samples exhibit higher (or lower) expression levels for a given gene pair compared with all the other samples. This feature of our proximity measure makes it less sensitive to noise compared with other proximity measures, such as the Spearman's rank correlation.

Graph-based algorithm

For each gene, TuBA identifies samples in the uppermost (or lowermost) percentile sets. Pairwise comparison between these percentile sets (using the 1-sided Fisher's exact test) identifies gene pairs that share a statistically significant number of samples. Each significant gene pair is illustrated graphically as a pair of nodes connected by an edge that represents the samples shared between their percentile sets. The complete set of these pairwise graphs generates large graphs from which robust gene co-expression signatures are recovered using the following iterative process (Fig. 2):

1. The graph is pruned such that its elementary units are complete subgraphs (cliques) of size 3 (triangles).

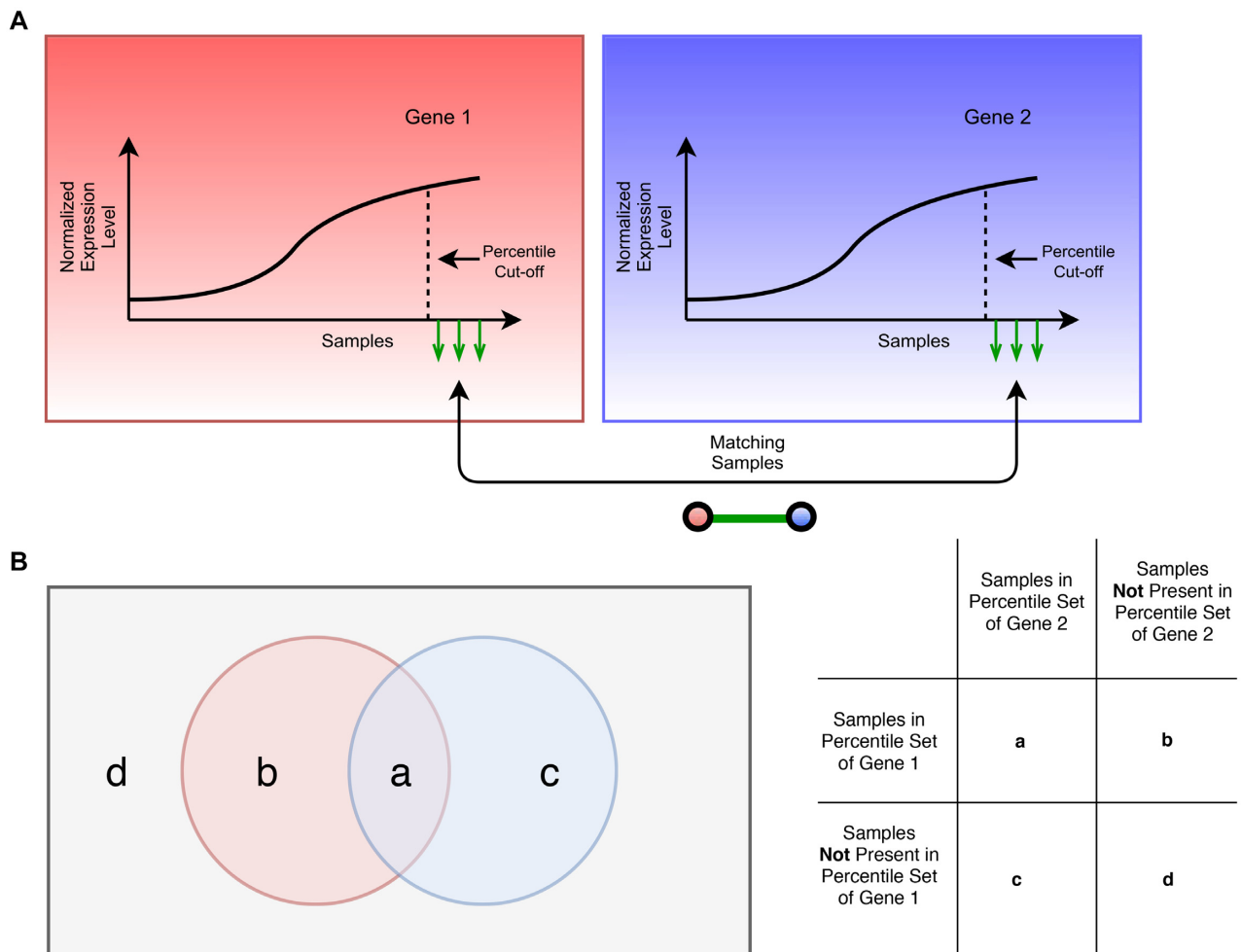


Figure 1: Schematic representation of TuBA's proximity measure. (A) For each gene, samples are arranged in increasing order of expression levels and those corresponding to a fixed percentile set (top or bottom) are compared between each pair of genes as shown. The gene pairs that share a significant number of samples are represented as nodes linked by edges, which represent the samples. (B) The Venn diagram illustrates the set-up of the contingency table for the 1-sided Fisher's exact test. The gray rectangular box represents the set of all samples in the dataset, and the red and blue circles represent the samples in the top (or bottom) percentile sets of gene 1 and gene 2, respectively.

- The largest clique (i.e., the seed) in the pruned graph is identified using the Bron-Kerbosch algorithm [18]. In cases where the largest clique is not unique, the union of all equally large cliques with a non-zero intersection of their nodes is designated as the seed; the remaining largest cliques are identified as new seeds in subsequent iterations.
- The graph is trimmed by removing all the edges that contain any of the nodes in the seed in Step 2. This step significantly reduces the computation time required to identify all the robust cliques in the graph.
- Steps 2 and 3 are repeated until the graph has no elementary units left.
- The seeds identified in Steps 1–4 are exclusive in their gene sets, i.e., no 2 seeds share a common gene. To create the bicluster, the seeds are reintroduced sequentially into the original pruned graph from Step 2, and nodes that share edges with ≥ 2 nodes in each seed are identified and added to the seed. The resulting graphs are the final biclusters obtained by TuBA.

Note that the requirement of largest cliques as seeds of our biclusters in Step 2 is a key step in our algorithm that enables the

identification of shared altered mechanisms in subsets of samples that exhibit high (or low) expression levels of these genes, while permitting the study of sets of co-expressed genes that are associated with functionally related pathways. Implicit in this requirement is the crucial assumption that the sets of genes making up the largest cliques are co-expressed in a subset of samples that make up the edges. This assumption is not the same as requiring all gene pairs making up the seed to share identical sets of samples, or assuming that all the samples making up the final biclusters co-express all the genes present in the bicluster. Instead, our expectation is that the samples present in the final biclusters are enriched in the top (or bottom) samples for each gene making up the biclusters. We have provided supporting evidence for this expectation in the Results section.

Gene enrichment analysis of the gene sets in the biclusters can be used to identify their functional relevance, and sample enrichment analysis can elucidate potential clinical subtypes, underlying mechanisms of disease, and possible therapeutic approaches. Furthermore, within each bicluster, genes can be assigned degrees, which are the total number of edges that connect them to other genes in the graph. Genes with higher

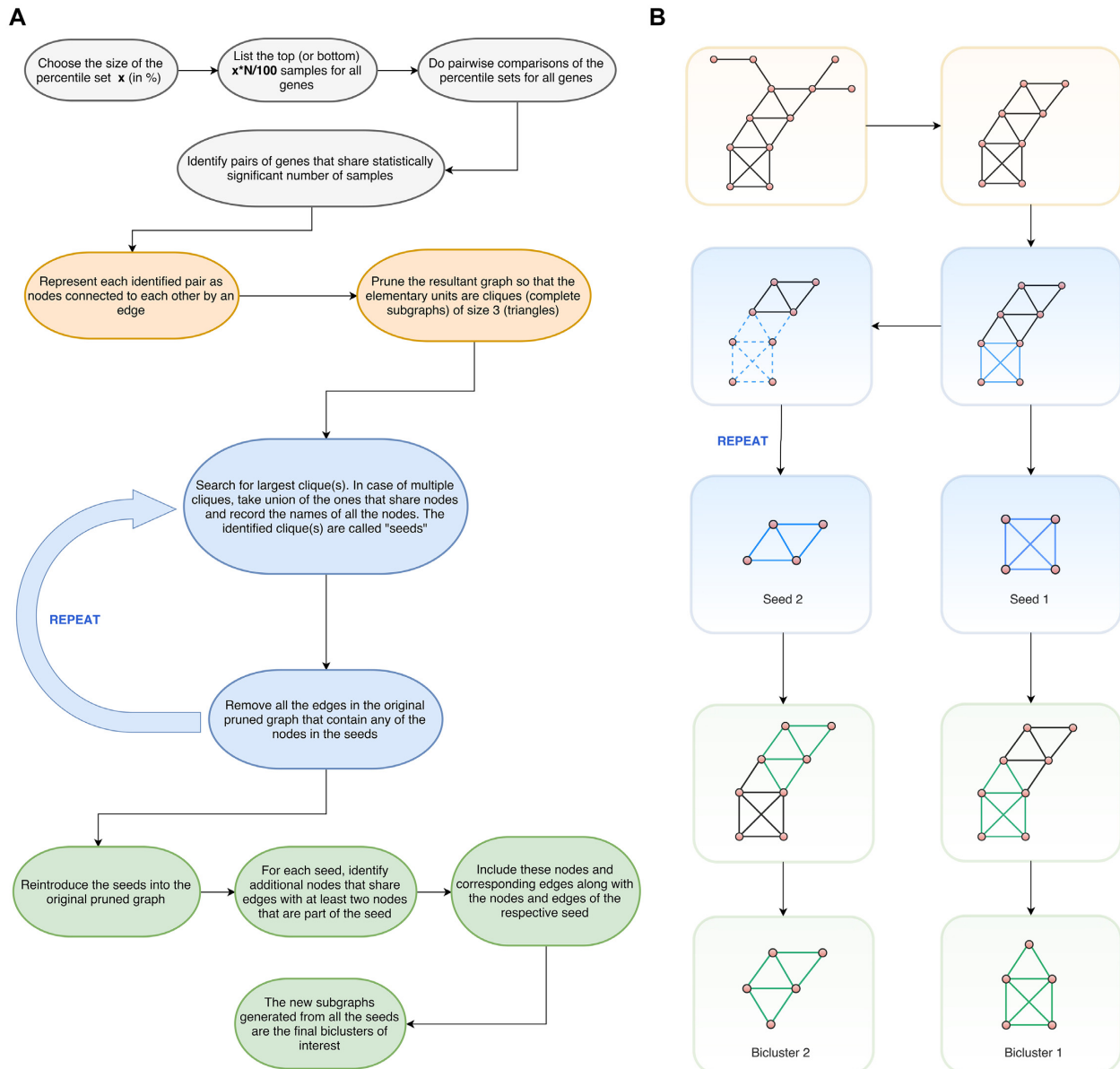


Figure 2: TuBA's schematics. (A) Flow chart of the pipeline for TuBA. (B) Schematic representation of the graph-based approach to discover biclusters.

degrees exhibit co-expression with other genes in a greater proportion of samples in the bicluster. These could be candidate driver genes.

Tuning TuBA

TuBA has 2 adjustable parameters:

1. The percentile cut-off: This parameter controls the number of top or bottom samples (based on expression levels) considered for comparison between genes.
2. The overlap significance cut-off: The P -value threshold used to assess significance of the overlap of samples between percentile sets for each gene pair. This parameter controls the minimum number of samples that must be shared between percentile sets for an association to be considered significant, and to be represented in the graph.

The parameters can be seen as “knobs” that can be tuned to probe different levels of heterogeneity in the population. For a given dataset, the choice of these 2 parameters determines the number, as well as the composition of the final biclusters. The choice of the first parameter determines the level of heterogeneity and/or the extent of prevalence of genomic alterations in tumors that may be of interest to the investigator. To illustrate how the choice of percentile cut-off could affect the identification of co-expressed gene pairs, we consider a hypothetical dataset consisting of 200 samples. Assume that there is a gene pair in this dataset that is up-regulated in 5% of the samples such that the top 5% percentile sets (i.e., top 10 samples) are identical for the 2 genes. Fig. 3A shows the significance values for overlaps (P -values calculated using 1-sided Fisher's exact test) as a function of the fraction of samples that overlap. At 5% percentile cut-off, the significance value for an overlap fraction of 1 (complete match/overlap between percentile sets) is $P = 4.45e-$

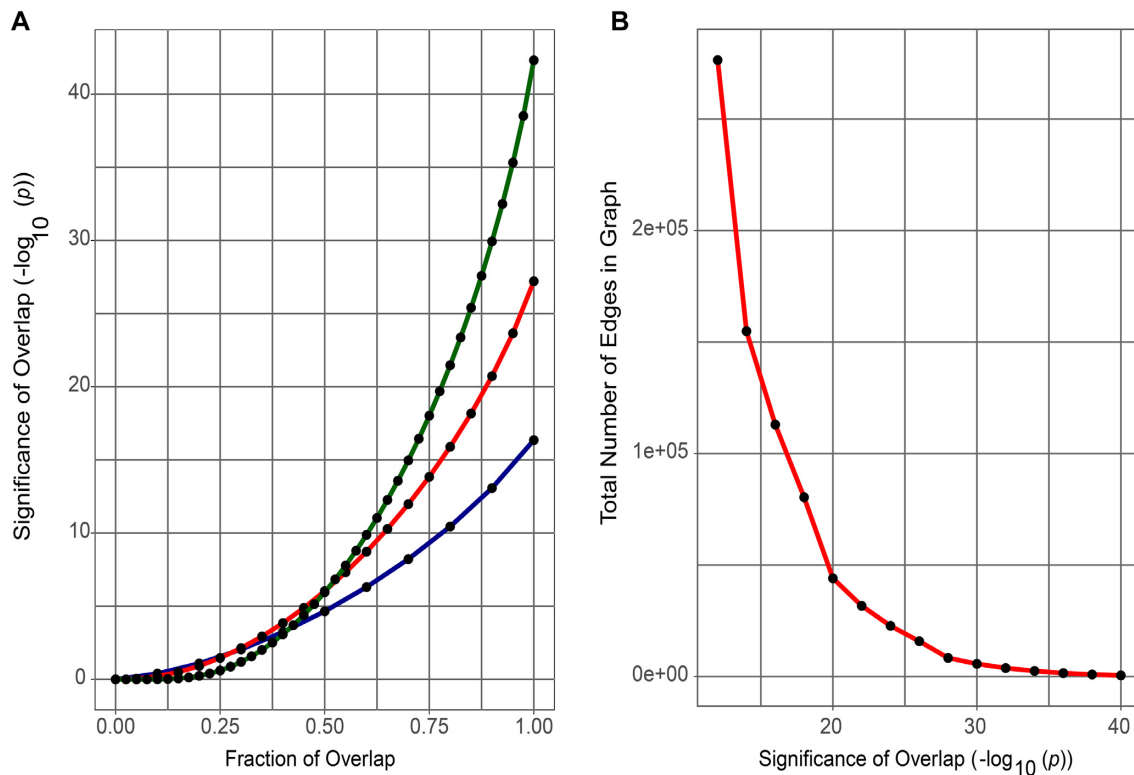


Figure 3: Tuning TuBA's parameters. (A) Significance of overlap corresponding to fraction of overlap between 0 (no matches/overlap) and 1 (all samples match/overlap) for percentile set size of: (i) top 20% (green), (ii) top 10% (red), and (iii) top 5% (dark blue) for a hypothetical dataset consisting of 200 samples. (B) Divergence of the total number of edges in the graph for the TCGA RFS dataset as we lower the cut-off for the significance of overlap.

17 (dark blue curve). If instead we had chosen an upper 10% percentile cut-off, we would have an overlap fraction of 0.5 (overlap of 10 out of 20 samples) corresponding to a significance value for overlap between $1e-10$ and $1e-5$. Thus, an increase in the size of the percentile set results in loss of significance for aberrant co-expression signatures found in smaller subsets of samples. This does not imply that it is generally better to choose smaller percentile sets. In fact, a reduction in the size of the percentile set increases the likelihood that a number of samples match purely by chance. (Note the P -values corresponding to the overlap fraction of 1 for the 3 cases in Fig. 3A, further demonstrated by permutation tests in the Results.) Thus, as we vary the size of the percentile set, there exists a trade-off between the sensitivity (identification of altered transcriptional profiles in small subsets of population) on one hand and the overlap significance on the other.

The choice of the second parameter—the extent of patient/sample overlap between percentile sets—determines the gene pairs that will be represented in the graph that is explored iteratively to identify sets of co-expressed genes. As we lower the significance of overlap (increasing P -values), new genes and samples get added to the graph, resulting in an increase in the number of edges (Fig. 3B). Further lowering of the overlap significance results in the addition of many more edges to the graph; however, this addition is not accompanied by a proportional increase in the samples or genes added to the graph. As more edges get added to the graph, the computational effort required for finding maximal cliques increases. Because the maximal clique problem is NP-hard [19], it can take exponential time to find all maximal cliques. Thus, the cut-off for the significance of overlap is informed by the trade-off between the gain of new information in the biclusters in terms of new samples and genes,

and the number of edges added to the graph that leads to a disproportionate increase in the computational effort. We propose the following heuristic for choosing the cut-off value: the cut-off for the significance level of overlap should be such that a decrease in the significance level by an order of magnitude leads to a 40–60% increase in the number of edges that get added to the graph (note the number of edges that get added to the graph at overlap significance P -values $>1e-20$ in Fig. 3B).

Because the principal goal of TuBA is to identify subsets of genes that are co-expressed at high (or low) levels within subsets of samples, the exact number of biclusters is not biologically relevant despite possible small variations in their total number as the algorithm is tuned. We investigated the consistency of TuBA's biclusters across different choices of the parameters. We used the hypergeometric test to identify biclusters that share significant fractions of their genes and observed that despite a 5-fold difference in the significance level of overlap, there is $>80\%$ agreement between the sets of biclusters obtained for different choices of the overlap significance cut-off. The results of the analysis are presented and discussed under Robustness of TuBA's Biclusters in Supplementary Methods and Supplementary Table 4.

BRCA datasets

We primarily applied TuBA to 3 independent BRCA datasets that used distinct methods for measuring transcript levels: (i) The Cancer Genome Atlas (TCGA) RNA sequencing (RNA-Seq) gene expression dataset using the Illumina HiSeq 2000 RNA-Seq platform, (ii) METABRIC gene expression dataset using the Illumina HT-12 v3 microarray platform, and (iii) 6 cohorts with gene

expression data from Gene Expression Omnibus (GEO) using the Affymetrix HGU133A microarray platform. To compare results among datasets, we applied TuBA to only their common gene sets. For clinical association analysis, we prepared 2 separate datasets for patients with known recurrence-free survival (RFS) status (908 patients) and patients with known Prosigna Breast Cancer Prognostic Gene Signature Assay (PAM50) subtype annotation (522 patients), respectively. Henceforth, we refer to TCGA RFS, METABRIC RFS, and GEO RFS datasets simply as TCGA, METABRIC, and GEO, respectively, and PAM50 datasets are indicated specifically.

1. TCGA—BRCA: The $\log_2(x + 1)$ transformed RSEM normalized counts of Level 3 data (16 August 2016 version), the clinical data (including relapse status and PAM50 subtype annotation from the 2012 *Nature* study [20]) (27 April 2016 version), and gene-level copy number variation (CNV) data, as estimated by GISTIC2 [21] (16 August 2016 version), were downloaded from the University of California Santa Cruz (UCSC) Xena Portal [22]. Genes with zero expression in all samples, as well as the samples with NAs for any gene, were removed from the analysis.
2. METABRIC: Normalized gene expression data, the clinical file, and the copy number file for the METABRIC study were downloaded from the cBioPortal [23] on 14 May 2017 [24, 25]. A gene expression dataset of 1,970 samples that had both relapse status and PAM50 subtype annotation was used in this study.
3. GEO: MAS5 normalized gene expression data and the clinical data with relapse status were downloaded from Gyorffy & Schafer [26] on 10 May 2017. The dataset comprises samples from 6 independent cohorts. After processing, our gene expression dataset consisted of 1,062 patients with relapse status.

Additional validation datasets

To compare TuBA's performance with other biclustering algorithms, we applied it to the following datasets:

1. breastCancerNKI: Gene expression data from a breast cancer study published by van't Veer et al. [27], and van de Vijver et al. [28] were downloaded in the form of an eSet using the breastCancerNKI package [29] in R. The dataset was further processed by removing probes with >10% missing values, and imputing the missing values for the included probes [30].
2. ESTIMATE: Scores for the level of stromal cells present and the infiltration level of immune cells in tumor tissues for 906 of 908 samples for the TCGA—BRCA RNA-SeqV2 dataset using the ESTIMATE algorithm were downloaded from ESTIMATE [31] on 2017-10-12.
3. GTEx: RNA-Seq raw counts data from the Genotype-Tissue Expression (GTEx) portal [32] were downloaded on 15 June 2017. The dataset comprises all tissue samples currently available in the GTEx database. The 214 breast tissue samples were identified and normalized using the DESeq package in R [33].

Statistical analysis

All computations were performed with R 3.3.0 [34]. The igraph package [35] was used to perform network/graph computations with some data summary functions performed using the plyr package [36]. The figures with graphs showing the genes for

some of the biclusters were generated using Cytoscape v3.4.0 [37]. Permutation test was performed on the METABRIC dataset (1,970 samples) with upper percentile set size cut-off of 5%. For each gene, we permuted the labels of the samples prior to ascertainment of the samples that corresponded to the top 5%, respectively. The significance values for overlaps between every pair of genes were computed using the Fisher's exact test. We performed 100 iterations of these permutation tests in total. The data.table package was used to handle data files, and the ggplot2 package was used to make plots [38]. GeneSCF [39] was used to perform gene set clustering based on functional annotation and to associate biclusters with specific biological processes. A binary matrix with biclusters along the rows and samples along the columns was generated to perform hierarchical clustering. Samples belonging to respective biclusters were assigned a value of 1. The Hamming distance was used to measure dissimilarity between the biclusters as well as the samples. All tests for enrichment were performed using the Fisher's exact test; where necessary the *P*-values were corrected for multiple hypothesis testing using the Benjamini–Hochberg false discovery rate (FDR) method [40]. The details of the contingency tables for the tests are provided in the Supplementary Methods. All enrichments are reported at $FDR < 0.05$, unless specified otherwise.

Results

Benchmarking TuBA's proximity measure

TuBA's pairwise proximity measure for genes is based on overlaps between their extremal subsets of samples. In comparison, global pairwise linear correlation coefficients such as the Pearson's correlation coefficient and Spearman's rank correlation coefficient are both susceptible to the influence of outliers. Even in the absence of outliers, genes that are co-expressed across all samples are expected to have significant overlaps between the subsets of samples that correspond to their top (or bottom) percentile sets. Given these observations, we tested the hypothesis that gene sets identified by global proximity measures have significant overlap with those identified by TuBA within its biclusters.

We computed the Pearson's correlation coefficients between all possible pairs of genes in the TCGA and METABRIC datasets, respectively. We shortlisted all gene pairs with the correlation coefficient ≥ 0.6 . We then used our graph-based algorithm to identify gene co-expression modules within the graphs. For TCGA and METABRIC, we obtained 569 and 298 gene co-expression modules, respectively (Supplementary Table 1). We investigated the association between the gene sets in biclusters discovered by TuBA's proximity measure vs gene co-expression modules identified by global correlation metrics by performing a hypergeometric test for gene overlaps. The null hypothesis was the absence of significant overlap between their respective sets of genes. Ninety percent (316 of 353) of the biclusters discovered by TuBA in the TCGA dataset comprised gene sets that were enriched in ≥ 1 gene co-expression module ($FDR < 0.001$), while 86% (293 of 340) of the biclusters discovered by TuBA in the METABRIC dataset were enriched in ≥ 1 module.

We performed a similar analysis using Spearman's rank correlation, with a cut-off of 0.6 for the correlation coefficient. We obtained 524 and 232 gene co-expression modules for TCGA and METABRIC, respectively. A total of 81% (285 of 353) of TuBA's biclusters in the TCGA dataset comprised gene sets that were enriched in ≥ 1 module. Overall, we see a significant enrichment of gene sets in TuBA's biclusters with the co-expression modules

obtained by using the 2 global proximity measures. This is in concordance with our hypothesis stated earlier.

As noted earlier, due to samples that exhibit aberrant/outlier expression of some genes, the linear correlation coefficients can often get skewed to reflect greater pairwise correlations between such genes that our graph-based algorithm can identify. However, due to the global nature of these proximity measures, the resulting graphs lack any information on the samples that might be associated with aberrant expression of these genes. In other words, unlike the case for TuBA, the edges in these graphs do not represent any subset of samples; they simply reflect an association between the genes by virtue of their pairwise correlation coefficient being greater than the chosen cut-off. The novel design of our proximity measure enables precise identification of co-expressed gene sets, while discerning the subsets of samples that exhibit higher (or lower) expression levels of these genes relative to the rest of the samples.

Choice of parameters for TuBA

For a given choice of the size of percentile set, TuBA generates plots that illustrate the number of added genes, added edges, and added samples as the overlap significance cut-off is varied. These are used to inform the choice of the overlap significance cut-off based on our proposed heuristic. Because the experimental platform and the total number of samples were different among the analyzed datasets, the choice of the overlap significance cut-off varied (Fig. 4). For the respective choices of the knobs, we obtained 353, 340, and 369 biclusters for the TCGA, METABRIC, and GEO datasets, respectively (Supplementary Table 2). Permutation tests showed that no gene pairs had P -values less than the cut-offs (Fig. S1). Moreover, after adjustment for multiple hypothesis testing, none of the gene pair P -values were statistically significant.

Enrichment of bicluster samples in top (bottom) sample sets of bicluster genes

As pointed out earlier, the identification of largest cliques as seeds of our biclusters was based on the expectation that samples present in the final biclusters were enriched specifically in the up-regulated (down-regulated) samples for each gene making up the cliques. We tested the hypothesis that the subsets of samples making up the biclusters were enriched by the samples that make up the top (bottom) sets for each gene in the bicluster. For example, suppose a dataset consists of 1,000 samples, wherein on application of TuBA for high expression, 1 bicluster is identified to comprise 100 genes and 200 samples. For each of those 100 genes, we identify their top 200 samples and test whether these 200 samples are enriched in the 200 samples making up the bicluster. The null hypothesis is that these 2 sets of samples are independent and therefore we should not expect to see statistically significant associations between them. We applied this test for each of TuBA's biclusters in the TCGA, METABRIC, and GEO datasets. For high expression, we observed that all genes in all 353 biclusters from TCGA showed significant enrichment (hypergeometric test $FDR < 0.001$). In the case of METABRIC we observed 2 biclusters (biclusters 2 and 269) out of 340 biclusters with only 1% of their constituent genes not exhibiting enrichment, while in the case of GEO we observed only 1 bicluster (bicluster 22) out of 369 with 1% of its constituent genes not exhibiting enrichment. For the low-expression analysis of TCGA, we observed 2 biclusters (biclusters 14 and 165) out of 203 biclusters that comprised a few

genes that did not exhibit enrichment. A crucial observation across all the datasets was that even in the few biclusters that included a few genes with enrichment $FDR > 0.001$, none of these genes were constituents of the seeds of those biclusters. We therefore found it justified to rely on the subsets of genes that make up the seeds for future gene ontological enrichment tests. This enables us to identify the core functional signatures of biclusters.

We performed a similar analysis for the bicluster samples. Following the previous example of a high-expression bicluster with 100 genes and 200 samples, we aimed to identify the genes in the bicluster that had a given sample in their top 200 samples (based on the expression levels of the genes). Therefore, for each sample, we evaluated whether their corresponding subset of genes had significant overlaps with the complete set of genes making up the bicluster. So, we tested the null hypothesis that overlaps between them were not statistically significant. For high expression, we observed that 95% of biclusters (336 of 353) from the TCGA, 97% of biclusters (329 of 340) in the METABRIC, and 89% of biclusters (328 of 369) in the GEO databases had $>95%$ of samples enriched (hypergeometric test $FDR < 0.001$). For the low-expression analysis of TCGA, we observed that 98% of biclusters (199 of 203) had 95% of their samples enriched in the bottom 200 samples for the corresponding genes in the biclusters. Based on these analyses, the FDR values for each gene (sample) in any given bicluster can be viewed as their scores—the closer the value of the FDR is to zero for a gene (sample), the stronger is the association of the gene (sample) to the bicluster. We used the FDR values for the genes and samples within bicluster i to evaluate its overall quality, $Q(B_i)$, defined as the minimum of the fraction of the genes in the bicluster with $FDR < 0.05$ or the fraction of samples in the bicluster with $FDR < 0.05$. $Q(B_i)$ takes values between 0 and 1, where values close to zero indicate weak associations between the constituent genes and samples within the bicluster, while values close to 1 would indicate strong associations.

Consistency of TuBA within a dataset

To investigate whether TuBA could consistently discover biclusters within the TCGA RFS cohort, the 908 samples were divided randomly into 2 groups of 454 samples each. This was done 5 times to generate 5 pairs of datasets. TuBA was applied to each dataset pair using a percentile set size of 5% and an overlap significance cut-off of $FDR \leq 1e-08$. Pairwise comparisons (between sets of genes) of biclusters from the 5 trials showed that on average 73% biclusters from 1 dataset in each pair were enriched ($FDR < 0.001$) in ≥ 1 bicluster from the other (Supplementary Table 3). We found a significant difference (Mann-Whitney U -test $P < 1e-05$) in the number of genes contained in biclusters that matched among trials, compared to the number of genes in biclusters that did not; while the median size of biclusters that matched was 20 (range: 3–840), the median size of biclusters that did not match was 3 (range: 3–18) (note that 3 is the smallest sized bicluster generated by TuBA). Overall, TuBA was able to consistently identify matching sets of co-expressed genes from randomly sampled subsets of data within a dataset.

Consistency of TuBA's biclusters among independent datasets

Using common sets of genes, we compared the biclusters obtained from (i) TCGA and METABRIC, (ii) TCGA and GEO, and (iii) METABRIC and GEO. Pairwise comparisons of biclusters

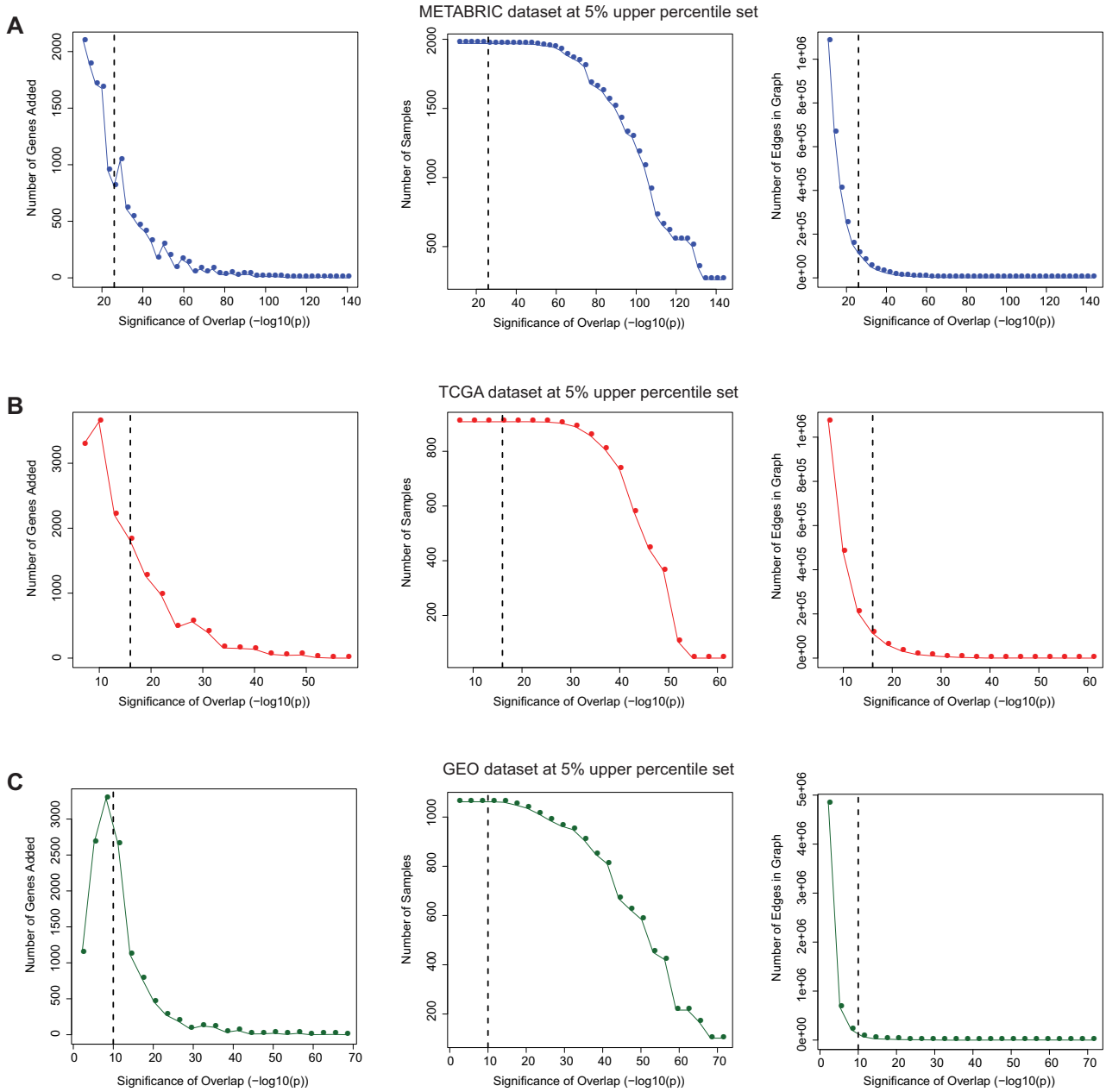


Figure 4: The effect of TuBA's parameters on the number of genes and samples in the graph. Plots for the number of genes added to the graph for every incremental decrease in the significance level for overlap $[-\log_{10}(P)]$, the number of samples in the graph at different significance levels of overlap, and the total number of edges in the graph at different significance levels of overlap corresponding to a percentile set size of 5% for (A) METABRIC, (B) TCGA, and (C) GEO datasets.

obtained from the 2 datasets were used to identify the biclusters that shared a significant proportion of their genes ($FDR < 0.001$). In the TCGA vs METABRIC comparison, 64% of biclusters obtained in 1 dataset were enriched in ≥ 1 bicluster in the other. In the TCGA vs GEO comparison, 69% of biclusters obtained in 1 dataset were enriched in ≥ 1 bicluster in the other. Finally, in the METABRIC vs GEO comparison, 76% of the biclusters obtained in 1 dataset were enriched in ≥ 1 bicluster in the other. Once again, we found that the biclusters that did not match were significantly smaller (median number of genes: 3–5) than the biclusters that matched (median number of genes: 20–25) between the datasets (Mann-Whitney U -test $P < 0.001$).

TuBA identifies subtype-specific biclusters

We classified BRCA samples based on the expression levels of the ESR1 (ER) and ERBB2 (human epidermal growth factor receptor 2 [HER2]) genes into 4 subtypes: (i) ER-/HER2-, (ii) ER+/HER2-, (iii) ER-/HER2+, and (iv) ER+/HER2+ (where the plus sign indicates over-expressed and minus sign indicates under-expressed). A substantial proportion of biclusters were enriched in the ER-/HER2- subtype: 53% for METABRIC (Fig. 5A and B), 54% for TCGA (Fig. 5C and D), and 40% for GEO (Fig. S6) (Supplementary Table 5).

According to the PAM50 classification, there are 5 subtypes of BRCA: (i) basal-like, (ii) HER2-enriched, (iii) luminal A, (iv) lumi-

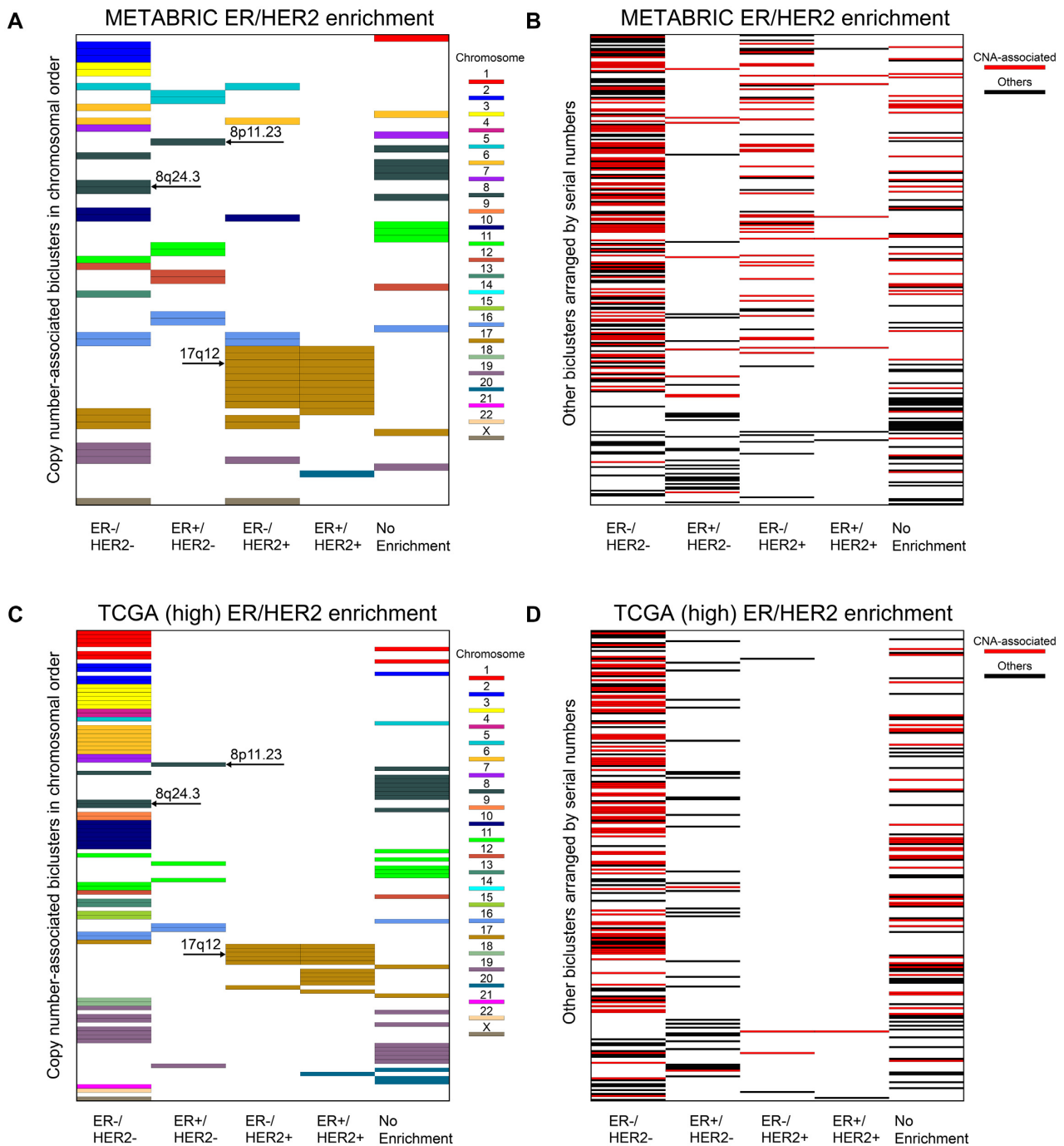


Figure 5: Subtype enrichment of CNV-associated and non-CNV biclusters. Enrichment of biclusters consisting of proximally located genes with copy number gains in the 4 subtypes based on ER/HER2 status for (A) METABRIC and (C) TCGA. The biclusters are represented by horizontal bars in each panel, color-coded according to the chromosome number of their constituent genes. Panels (B) and (D) show the remaining biclusters arranged according to their serial numbers in Supplementary Table 4 for METABRIC and TCGA, respectively. The ones that are associated with copy number gains of genes located at distant chromosomal sites are shown in red, while the rest are shown in black. Note that the thickness of the bar in each panel depends on the total number of biclusters displayed in that panel and so does not represent its chromosomal extent.

nal B, and (v) normal-like [41]. We observed a significant fraction of biclusters enriched in the basal-like subtype: 52% for METABRIC (Fig. S4A and B) and 55% for TCGA PAM50 (Fig. S4C and D). Although tumors of the basal-like or triple-negative subtypes accounted for only ~15% of all BRCA in the population, most of the altered expression profiles captured by our biclusters were in tumors of this subtype.

TuBA identifies down-regulated subtype-specific biclusters in RNA-Seq data

RNA-Seq offers a significant advantage over microarray assays. Theoretically, only the depth of sequencing limits the dynamic range of RNA-Seq data [42, 43]. Given that TCGA's RNA-Seq data have adequate sequencing depth, we expected a reliable quan-

tification of even low-expression transcripts. We therefore applied TuBA to the TCGA datasets to explore transcriptional profiles associated with low expression. We found that 46% of biclusters from TCGA were enriched in the ER-/HER2- subtype (Fig. S5A and B), while 48% of biclusters from the TCGA PAM50 dataset were enriched in the basal-like subtype (Fig. S4E and F). Thus, biclusters associated with low expression were predominantly enriched in the ER-/HER2- or basal-like subtypes. This further underscores the tremendous heterogeneity of altered transcriptional profiles within tumors of this subtype.

TuBA highlights biclusters with proximally located genes

We observed that several biclusters discovered by TuBA across the 3 datasets comprise genes that are proximally located on the chromosomes, suggesting copy number amplification (CNA) as an underlying mechanism. Copy number data were used to calculate the significance of the proportion of samples present in each bicluster that exhibited copy number gains. For each gene in a given bicluster, we computed a *P*-value for the significance of the proportion of samples with CNA present in the bicluster. These *P*-values were then combined using Fisher's method to yield a single *P*-value for the bicluster. This showed that 56% and 64% of biclusters from the METABRIC and TCGA datasets, respectively, were enriched for CNA ($FDR < 0.001$). Closer scrutiny revealed that only 60 (18%) biclusters from METABRIC were associated exclusively with CNA of proximally located genes (Fig. 4A); the remaining biclusters associated with CNA were enriched in genes from distant chromosomal locations (Fig. 4B). Similarly, 112 (32%) biclusters from TCGA were associated with CNA of proximally located genes (Fig. 4C). Many of these biclusters were associated with loci previously identified to exhibit copy number gains in BRCA [44, 45]. In order to explore the association between the biclusters obtained from the low-expression analysis and loss of copy number, we repeated the copy number analysis described above. We observed that 52% of biclusters from the TCGA dataset were enriched in copy number losses. However, only 21 biclusters contained genes located on the same chromosome (Fig. S5A); the remaining biclusters associated with copy number loss were enriched in genes from distant chromosomal locations. Similar analyses for PAM50 subtype enrichment for METABRIC and TCGA are summarized in Fig. S4.

To compare CNA-associated biclusters between TCGA and METABRIC, we prepared 2 datasets that contained genes that were common in the 2 cohorts (17,209 genes). Pairwise comparison of the set of genes in the CNA-enriched biclusters between the 2 datasets revealed that 61% of biclusters from TCGA matched ≥ 1 CNA-associated bicluster from METABRIC. On the other hand, 91% of biclusters from METABRIC were enriched in ≥ 1 CNA-associated bicluster from TCGA. This suggests that most of the CNA-enriched biclusters identified in the METABRIC microarray dataset were independently identified in the RNA-Seq dataset of TCGA.

We also observed some biclusters with proximally located genes that were not associated with gain in copy number. For TCGA, 14 biclusters out of 353 consisted of genes located proximally, while 18 biclusters out of 340 for METABRIC consisted of genes located near each other. Details of the genes and subtype-specific enrichments for some of these biclusters are summarized in Supplementary Table 6. Examples of biclusters from this category include the biclusters consisting of genes from the Cancer-Testis antigens family—

MAGEA2, MAGEA3, MAGEA6, MAGEA10, CSAG1, CSAG2, CSAG3 (Xq28)/CT45A3, CT45A5, CT45A6 (Xq26.3). These genes are known to be aberrantly expressed in triple-negative breast tumors [46], as well as in a few other tumor types [47].

TuBA identifies biclusters associated with non-tumor expression signatures

We also discovered biclusters that seemed to be associated with non-tumor cells. For instance, biclusters associated with immune response were among the largest identified independently in all 3 datasets. The top 5 Gene Ontology—Biological Processes (GO-BP) terms for the bicluster associated with immune response were T-cell co-stimulation, T-cell receptor signaling pathway, T-cell activation, regulation of immune response, and positive regulation of T-cell proliferation (Supplementary Table 7). This indicates immune cell infiltration in a significant number of tumor samples. To corroborate this, we stratified TCGA samples on the basis of their ESTIMATE [48] scores for the infiltration level of immune cells in tumor tissues into 3 groups—(i) top 25 percentile, (ii) intermediate 50 percentile, and (iii) bottom 25 percentile—and verified that samples in these biclusters associated with immune response were enriched in samples with the highest levels of immune infiltration ($FDR < 0.001$).

For all 3 datasets, we also observed a bicluster associated with the stromal adipose tissue. The top 5 GO-BP terms for this bicluster were response to glucose, triglyceride biosynthetic process, triglyceride catabolic process, retinoid metabolic process, and retinol metabolic process. An analysis based on the ESTIMATE scores for the level of stromal cells present in tumor tissue of TCGA samples confirmed that this bicluster was enriched within the top 25 percentile samples for stromal cell level. Subtype enrichment revealed that the bicluster was enriched in ER-/HER2-, basal-like (PAM50), and normal-like (PAM50) subtypes.

TuBA's proximity measure was applied to gene expression data from 214 normal breast tissue samples from the Genotype-Tissue Expression (GTEx) public dataset. We observed that only 6.75% of biclusters obtained for the TCGA vs GTEx comparison were enriched in gene pair associations identified in the GTEx dataset. The bicluster associated with the adipose tissue signature was one of the biclusters found enriched in GTEx. Another group of biclusters enriched in the 3 cancer datasets, as well as in GTEx, were those associated with translation and ribosomal assembly. The top 5 GO-BP terms for these biclusters were translation, ribosomal RNA (rRNA) processing, ribosomal small subunit biogenesis, ribosomal large subunit assembly, and ribosomal large subunit biogenesis. These biclusters were enriched in the ER-/HER2- subtype ($FDR < 0.001$).

TuBA identifies clinically pertinent biclusters

We performed a Kaplan-Meier analysis of RFS, comparing the patients present in each bicluster to the rest for METABRIC and GEO. (The number of patients with incidence of recurrence in TCGA was insufficient for this kind of survival analysis to be statistically robust.) As expected for METABRIC, patients in the bicluster (bicluster 25) associated with the HER2 amplicon (17q12) had significantly shorter RFS time compared to the rest (Fig. S7). This is because patients in the METABRIC study were enrolled before the general availability of trastuzumab [49].

We also observed biclusters associated with CNA at the 8q24.3 locus in all 3 datasets (TCGA—biclusters 39 and 113, METABRIC—biclusters 26, 56, and 167, GEO—biclusters 16, 24, 37, 55, 74, 118, and 302). These patients also had significantly

shorter RFS times compared with those patients whose tumors did not have amplification of this locus (Fig. 6A–C). A similar result was obtained when we restricted the samples to ER+/HER2– tumors, validating an earlier observation that copy number gain of the 8q24.3 locus may confer resistance to ER targeted therapy [50]. We note, however, that biclusters with amplification of the 8q24.3 locus were enriched in the ER–/HER2– subtype ($P < 0.001$). Hence, amplification of this locus may be even more relevant in determining treatment for patients with ER–/HER2– breast cancers assigned into an intermediate (ambiguous) risk class by Oncotype DX [50]. Genes at 8q24.3 that may be considered promising candidates on the basis of their degrees in the biclusters include PUF60, EXOSC4, COMMD5, and HSF1. Specifically, PUF60 is an RNA-binding protein known to contribute to tumor progression by enabling increased MYC expression and greater resistance to apoptosis [51].

For both METABRIC and GEO, patients in biclusters associated with copy number gains of the 8p11.21–p11.23 loci (METABRIC—bicluster 289, GEO—bicluster 25) had significantly shorter RFS times compared with patients without amplification of this locus (Fig. 6D–F). We found that patients in this bicluster were enriched in the luminal B subtype, which has poorer prognosis than the luminal A subtype among ER+/HER2– tumors [52]. This suggested that amplification of the 8p11.21–p11.23 loci may be another marker of potential failure of ER targeted therapy.

Similarly, we found that patients whose tumors have copy number gains of the 17q22–q23.3 locus (METABRIC—biclusters 33 and 119, GEO—biclusters 15 and 160) had significantly shorter RFS times compared with patients whose tumors do not exhibit such a copy number gain (Fig. 6G–I). For METABRIC, this cohort was enriched in the luminal B (PAM50), ER+/HER2+, and ER–/HER2+ subtypes (FDR < 0.001). For GEO, this cohort was enriched in the ER+/HER2+ and ER–/HER2+ subtypes (FDR < 0.05). This suggests that amplification of this locus may confer additional risk of recurrence in HER2+ breast cancers.

Note that the biclusters discussed above were not the only ones that exhibited differential relapse outcomes. For METABRIC, 61 biclusters out of 340 were found to exhibit differential relapse outcomes for the patients present in the biclusters. Of these 61 biclusters, 69% were enriched in the ER–/HER2– subtype (64% for basal-like), with a significant proportion (67%) of these associated with copy number gains. For GEO, there were 48 such biclusters (13%) that exhibited differential relapse outcomes; 25% of these were enriched in the ER–/HER2– subtype.

Tests for enrichment of biclusters in tumors of higher grades revealed that 8 biclusters from TCGA were enriched in tumors of grade 3C. Some of these biclusters were associated with GO-BP terms related to angiogenesis, vasculogenesis, blood vessel maturation, and so forth. For METABRIC, 4 biclusters were enriched in tumors of grade 3, out of which 2 were associated with the HER2 amplicon (17q12). For GEO, 68 biclusters were enriched in tumors of grade 3, including biclusters associated with CNA at the HER2 amplicon.

We also looked at the lymph node status of patients and observed that 4 biclusters in TCGA were enriched in samples with positive lymph node status in the corresponding patients. One was associated with the HER2 amplicon, while the others were associated with CNA at the 8q22.1–q22.3 loci, 17q23.1–q23.3 loci, and the 19q13.43 locus, respectively. Similarly in METABRIC, we observed 4 biclusters enriched in samples with positive lymph node status in the corresponding patients—2 of them were associated with copy number gains at the HER2 amplicon, and the other 2 were associated with copy number gains at 19q13.11–q13.12 and 1q21.3–q25.1, respectively. Interestingly, biclusters

associated with CNA at 8q24.3, 8p11.21–p11.23, and 17q22–q23.3 that exhibited poor RFS outcomes were not enriched in tumors of higher grades or in patients with positive lymph node status in any of the 3 datasets. In the case of METABRIC, we additionally confirmed that none of these biclusters (8q24.3, 8p11.21–p11.23, 17q23.1–q23.3) were among the 36 biclusters enriched in samples with the poorest expected 5-year survival outcome (Nottingham Prognostic Index > 5.4) [53, 54]. This highlights the importance of these altered transcriptomic signatures for reclassification of patients into the category with higher risk of recurrence.

Hierarchical clustering of biclusters reveals shared mechanisms

Sample membership-based hierarchical clustering of biclusters revealed distinct groups of biclusters that presumably share common functional mechanisms (Fig. 7). These included clusters associated with cell cycle and proliferation, immune response, cell adhesion (extracellular matrix), translation, mitochondrial translation, and ribosomal RNA processing pathways. Because a significant fraction of our biclusters were associated with copy number alterations, we also found distinct groups of biclusters associated with significant copy number changes such as the ones associated with the HER2 amplicon, the 8p11.21–p11.23 loci, or the 8q24.3 locus.

Similarly, we used hierarchical clustering to group samples that were enriched in similar sets of biclusters, highlighting differential clinical outcomes. In particular, we observed 2 sets of samples enriched in biclusters associated with CNA at the 8q24.3 locus. In one group, the samples were enriched in biclusters related to immune response; this group showed significantly lower incidence of recurrence compared to those without enrichment in immune response-related biclusters. Both of these sets of samples were enriched in biclusters associated with cell division and proliferation. In contrast, we observed a cluster of samples enriched in biclusters associated with 8q24.3 copy number gain and a number of other loci; however, these were not enriched in biclusters associated with cell division and proliferation. This group exhibited low incidence of recurrence. We also observed a cluster of samples with significantly poor RFS that were enriched in biclusters associated with CNA at 17q25.1–q25.3, and in biclusters associated with cell division and proliferation.

TuBA compared to other biclustering methods

TuBA's proximity measure distinguishes its biclusters from those identified by other algorithms by leveraging the size of the datasets to identify subsets of tumor samples that co-express subsets of genes at their most extreme levels (high or low) relative to other samples. We emphasize that TuBA is designed to identify biclusters with samples that correspond to the extremals for the corresponding sets of genes and does not consider other subsets of conditions for the same sets of genes for biclustering. In contrast, most biclustering methods seek submatrices with constant, or coherent gene expression patterns. Given this key difference, only those biclusters that exhibit such expression patterns in the extremal (top or bottom) subsets of samples for some subsets of genes are expected to have agreement with the biclusters identified by TuBA. Therefore, a direct comparison between the biclusters discovered by other algorithms and those identified by TuBA would necessarily be limited.

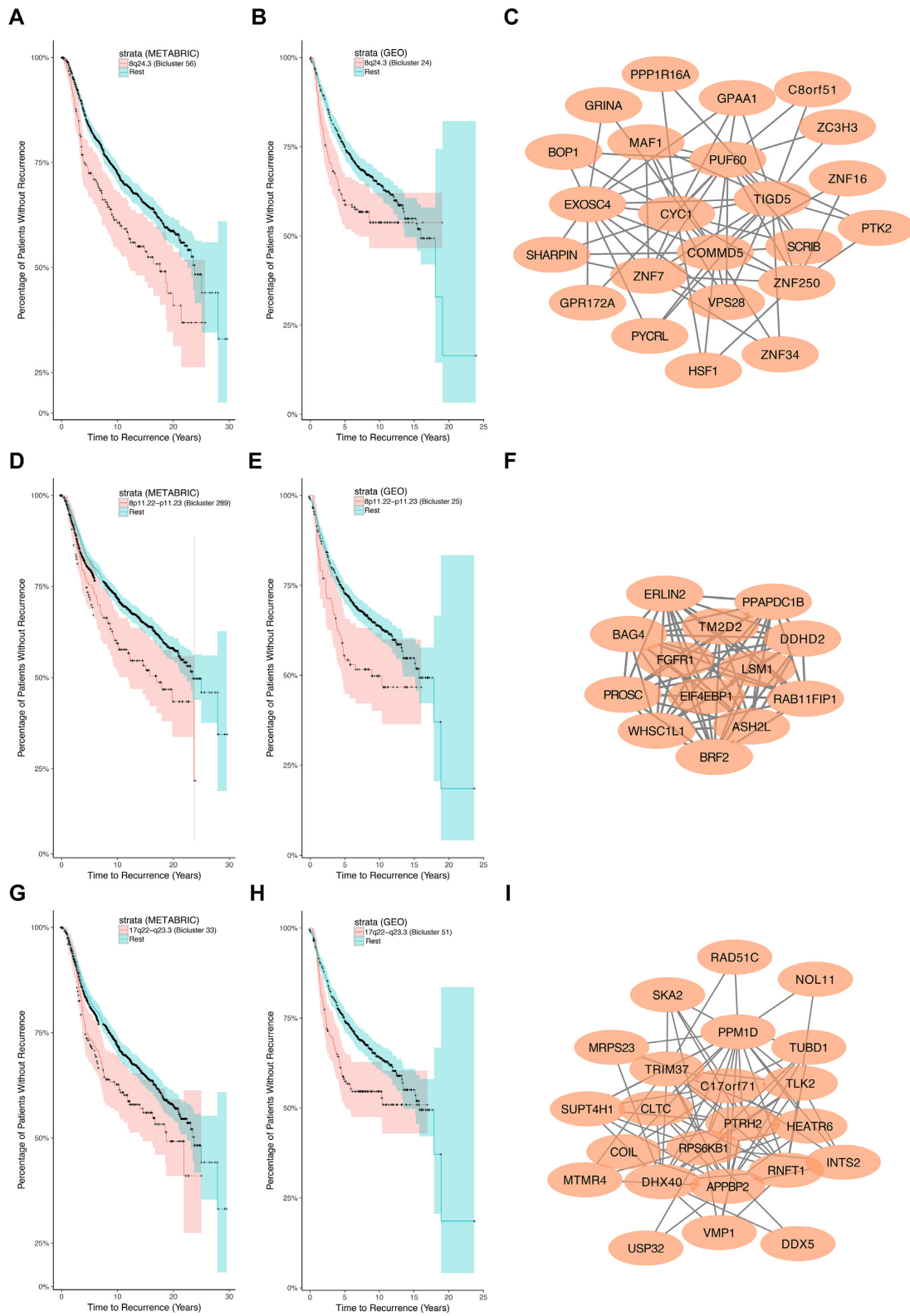


Figure 6: Clinically pertinent biclusters. Kaplan-Meier survival curves for the set of patients in the bicluster (red) compared to the remaining set of patients (blue) for METABRIC and GEO datasets, together with the graphs corresponding to the biclusters for 8q24.3 (A–C), 8p11.22–p11.23 (D–F), and 17q22–q23.3 (G–I).

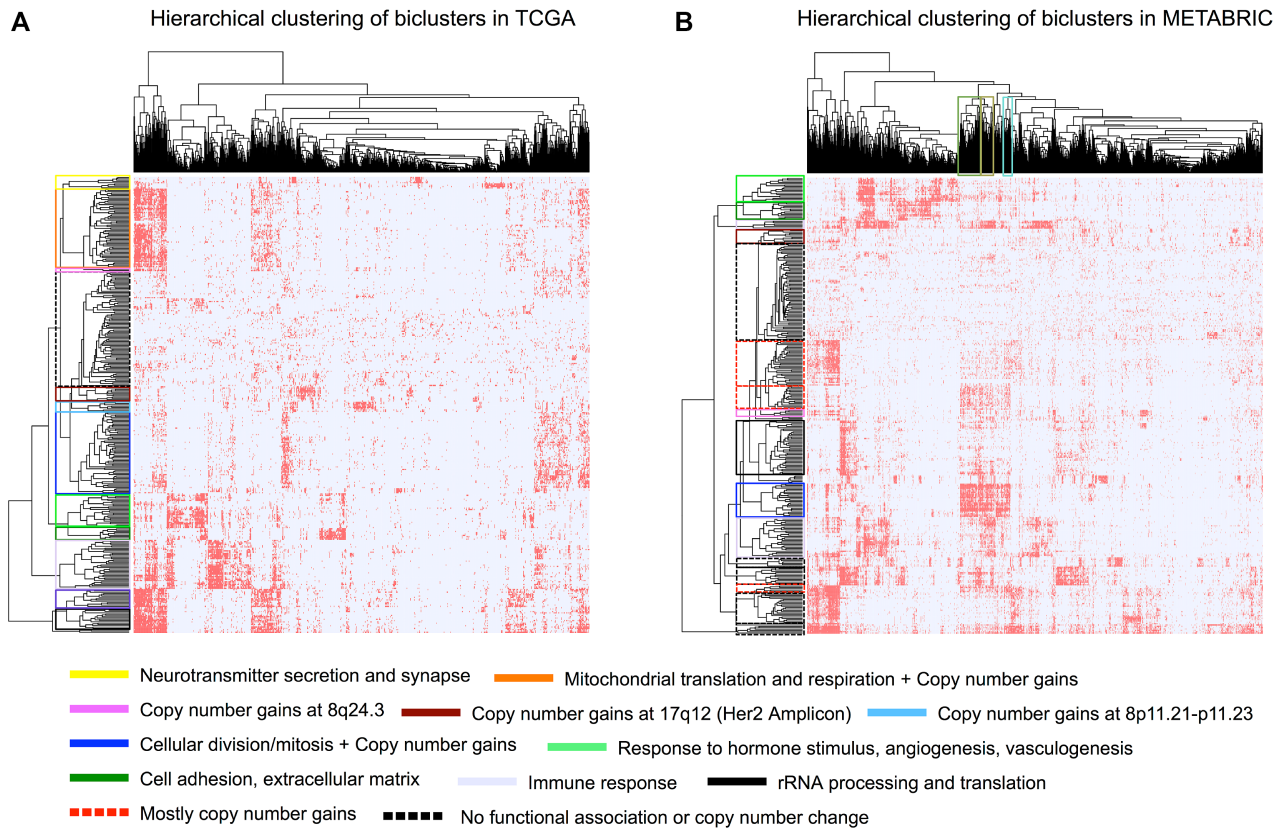


Figure 7: Hierarchical clustering of biclusters-samples. Hamming distance was applied to the biclusters-samples binary matrix for (A) TCGA and (B) METABRIC RFS datasets. The clusters of samples marked by green, brown, and cyan on top in panel (B) exhibit poor recurrence-free survival. The green and brown clusters are associated with copy number gains at 8q24.3, while the cyan cluster is associated with copy number gains at 17q25.1–q25.3. Additionally, all 3 of them were enriched in gene signatures associated with cellular division and proliferation.

In their paper on DeBi [55], a novel biclustering method that identifies differentially expressed biclusters based on a frequent itemset approach, Serin and Vingron applied their method to both synthetic and real gene expression datasets, including diffuse large B-cell lymphoma (DLBCL) data comprising 661 genes and 180 samples [56]. Apart from DeBi, they applied ISA [57], OPSM [58], QUBIC [59], and SAMBA [60] to this dataset; we used these results to evaluate TuBA against these methods. To ensure a uniform and unbiased comparison between the enrichment results for biclusters from different algorithms, we used GeneSCF [39] to perform GO-BP enrichment on the biclusters obtained by all methods. Fig. 8A shows the proportions of GO-BP-enriched biclusters for 5 different significance levels (FDRs): 0.001%, 0.1%, 0.5%, 1%, and 5%. For the FDR cut-off of <5%, almost all the biclusters for every algorithm were enriched in ≥ 1 GO-BP term. TuBA had 2 non-enriched biclusters out of 94, SAMBA had 2 non-enriched biclusters out of 128, and QUBIC had 1 bicluster out of 100 that were not enriched in a GO-BP term (Supplementary Table 7). As the FDR cut-off was lowered, TuBA had lower proportions of enriched biclusters compared to other algorithms for the corresponding FDR cut-offs. This can be partly attributed to the fact that the other algorithms discover biclusters that can have arbitrary overlaps between their genes. Because most of the biclusters discovered by other algorithms shared genes with other biclusters, we could expect a certain amount of redundancy in enriched GO-BP terms. In contrast, TuBA precludes any overlap between the genes of the seeds of its biclusters. Moreover, its biclusters often include proximally

located genes with aberrant expression due to copy number changes, which may not show enrichment in GO-BP terms.

For a closer examination of the redundancy in the GO terms enrichment, we identified the top 5 GO-BP terms for every bicluster obtained by each algorithm (not every bicluster was enriched in 5 distinct GO-BP terms; some had <5, while others were not enriched in any term). For each algorithm, we prepared lists of all the unique GO-BP terms for the entire set of biclusters. The ratios of the number of elements in these lists to the total number of biclusters for each algorithm at 5 different significance levels show that TuBA-identified biclusters were enriched in a more extensive array of biological process terms (Fig. 8C).

In addition to the DLBCL dataset, we also analyzed the TCGA dataset with the following biclustering algorithms: (i) BIMAX [61], (ii) ISA, (iii) QUBIC, and (iv) SAMBA, using their respective default parameters (Supplementary Tables 7 and 8). For succinct descriptions of each of these algorithms we refer the reader to Prelic et al [61] and Pontes et al [15]. We used the *biclust* package in R for BIMAX [62], the *isa2* package in R for ISA [63], the *QUBIC* package in R for QUBIC [64], and the *Expander* software for running SAMBA [65]. Fig. 8B shows the proportion of GO-BP terms enriched in biclusters of each algorithm for 5 different significance levels. TuBA compared favorably with other algorithms, especially when we accounted for the redundancy of the GO terms that were found enriched in the biclusters. We observed again that TuBA's biclusters were enriched in a larger set of distinct biological process terms (Fig. 8D). We observed similar re-

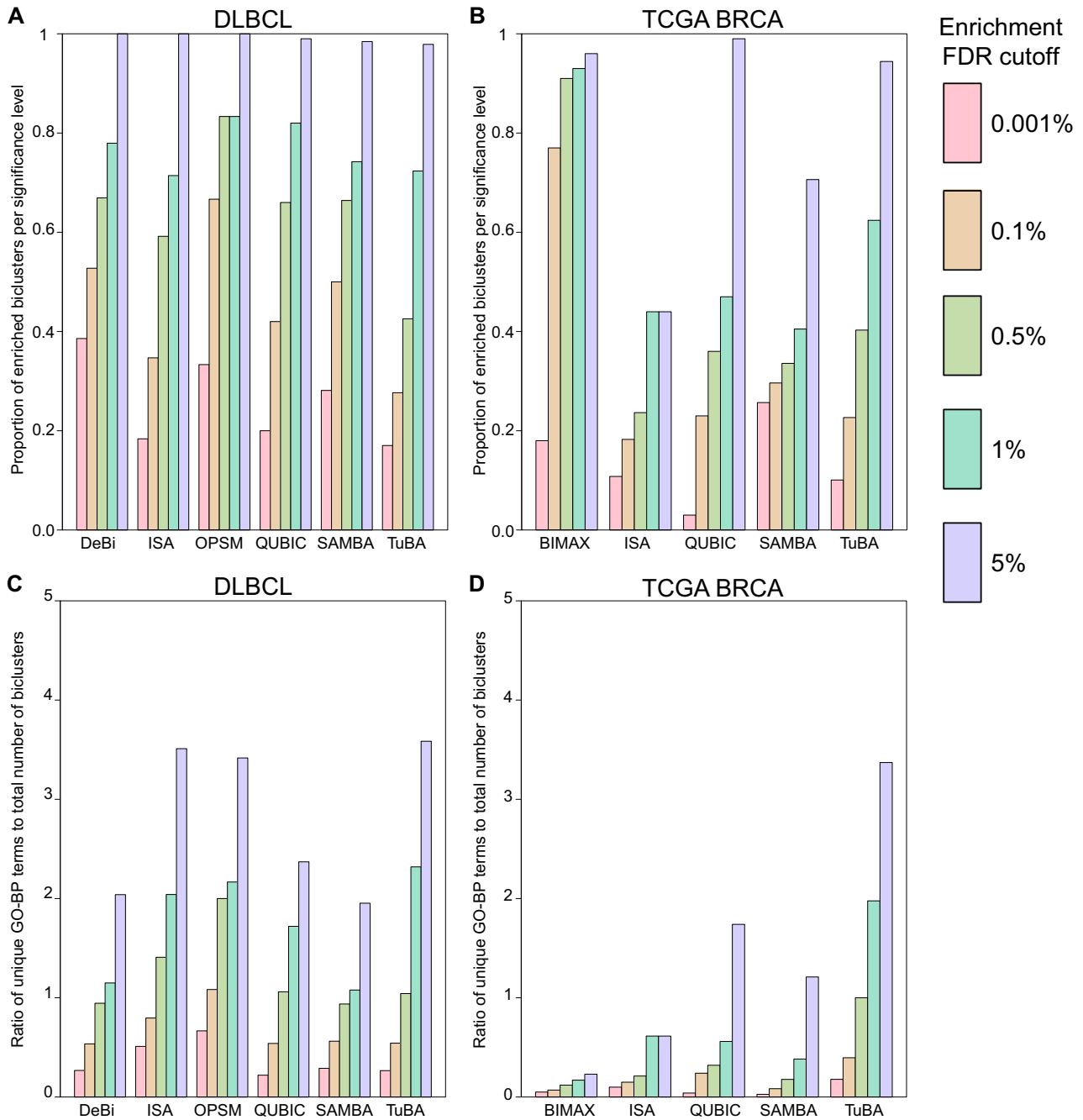


Figure 8: TuBA compared to other biclustering methods. Proportions of GO-BP terms–enriched biclusters for each biclustering method at 5 different significance levels for (A) the DLBCL dataset, (B) the TCGA BRCA dataset. Ratios of number of unique GO-BP terms and total number of biclusters at 5 different significance levels for (C) the DLBCL dataset, (D) the TCGA BRCA dataset.

results for METABRIC (Fig. S8A and B). The results of the comparative analysis are summarized in the Supplementary Methods.

In these analyses, the choice of the parameters is a crucial factor in determining the performance of each biclustering method. It is possible that different results could be obtained by more prudent choices of parameters for the other algorithms; however, a detailed analysis of optimal parameter choices for each of these algorithms is beyond the scope of this study. We must point out for TuBA that, for any given dataset, there is no

optimal (or default) choice of its 2 parameters; the biclusters obtained for any given choice of the parameters simply satisfy the basic requirements laid down by those choices. We looked at GO-BP term enrichments for TuBA’s biclusters for 5 different choices of the overlap cut-off (Supplementary Methods). Although the total number of biclusters obtained differed for each choice, the proportion of enriched biclusters at different significance levels remained similar irrespective of the parameter choice (Fig. S9A). Similarly, the ratio of the number of unique GO-BP terms and the

total number of biclusters was consistent across all 5 choices of overlap cut-offs (Fig. S9B).

In earlier studies comparing biclustering algorithms [8, 61], synthetic datasets were generated with constant, shifting, and/or scaling patterns of expression for subsets of conditions and genes. The algorithms were evaluated on the basis of how well they were able to identify the known biclusters implanted in these synthetic datasets. Because TuBA is not based on a mathematical model of the data in its biclusters, a comparison based on synthetic datasets is not feasible. However, in the case of tumor datasets we have the benefit of complementary genomic data that could provide us with truth-known scenarios for validation. For example, alterations at the genomic level can directly influence the expression levels of genes; it is well known that a significant proportion of tumors across multiple tumor types frequently exhibit genomic alterations such as gains or losses in the copy numbers of genes. Quite often, these alterations are not limited to a single gene but include multiple genes located at neighboring chromosomal locations. If such alterations are located at transcriptionally active sites, then co-expression of the neighboring genes that are affected by it will be observed. In BRCA for instance, ~15–20% of tumors possess extensive gains in copy numbers of genes at the 17q12 cytoband locus (includes, e.g., *ERBB2* [Her2], *STARD3*, *GRB7*, *PNMT*, *PGAP3*, *MED1*). Identification of co-expression of genes at this locus in the subset of samples that are histologically HER2-positive (HER2+) represents a simple truth-known scenario that can be used to verify whether a given biclustering algorithm identifies the co-expression of these genes in the subset of samples that exhibit this alteration. We identified HER2+ samples in the TCGA dataset and, for each biclustering algorithm, selected those biclusters that were enriched in these samples (hypergeometric test $FDR < 0.001$). BICMAX and SAMBA did not discover any, but ISA identified 2 biclusters (biclusters 71 and 72) enriched in HER2+ samples. Although the genes from the 17q12 amplicon, such as *ERBB2*, *STARD3*, *GRB7*, *PNMT*, *PGAP3*, and *MED1*, were present in ISA's enriched biclusters, they made up a small subset within the genes in them—bicluster 71 had 639 genes, while bicluster 72 had 539 genes. QUBIC also had 4 biclusters that were enriched in HER2+ samples; however, they did not contain any genes from the HER2 amplicon (including *ERBB2*). In contrast, not only did TuBA identify a bicluster (bicluster 256) exclusively associated with the HER2 amplicon, it identified many other biclusters associated exclusively with CNA of genes located near each other.

In summary, apart from TuBA, only ISA identified co-expression of the genes located at the HER2 amplicon. However, ISA's co-expression module corresponding to the amplicon was embedded within much larger sets of genes. In the absence of information about copy number gain of the *ERBB2* gene, it would be a challenge to explicitly identify the co-expression module corresponding to the amplicon, and in turn infer the underlying mechanism for their co-expression. TuBA successfully uncovers those co-expressed sets of genes that are associated with CNA of neighboring sites on the chromosome and is particularly efficient at identifying transcriptionally active copy number gains, as compared to other algorithms.

The nature of our proximity measure allows us to determine differential co-expression signatures without the need to specify subsets of samples in advance. Gao et al. [66] proposed a biclustering method, Bicmix, based on a Bayesian statistical model to infer subsets of co-regulated genes that covary in all samples, or in only a subset of samples. They also developed a principled method to recover context-specific gene co-expression networks from the sparse biclustering matrices obtained by Bicmix.

They applied Bicmix to the breastCancerNKI dataset and identified 432 genes that were differentially co-expressed in ER+ and ER- samples. Of these 432 genes, 430 were up-regulated in ER- samples and down-regulated in ER+ samples, while 2 genes were down-regulated in ER- samples and up-regulated in ER+ samples. We applied TuBA (for high expression) to the same dataset with the following choice of parameters: (i) percentile set size: 10%, and (ii) overlap significance cut-off: $FDR \leq 1e-08$. We obtained 549 biclusters, several of which consisted solely of probes that correspond to the same gene (Supplementary Table 5). This is reasonable because probes corresponding to the same gene are expected to demonstrate higher expression levels in the same set of samples. We inquired whether some of the biclusters discovered by TuBA corroborated the differential co-expression signature between ER+ and ER- samples identified by Bicmix. Using Fisher's exact test, we determined that the set of 430 genes up-regulated in ER- samples and down-regulated in ER+ samples were enriched in 30 biclusters discovered by TuBA—bicluster 5 shows the maximum enrichment ($FDR < 1e-165$). In fact, the genes that had the highest degrees in the co-expression network discovered by Bicmix—*CD247*, *CD53*, *IL10RA*, and *CXCR3*—were among the ones with highest degrees in bicluster 5 discovered by TuBA. The 2 genes (*SFRP2* and *COL12A1*) that were up-regulated in ER+ samples and down-regulated in ER- samples were also found to be co-expressed in a TuBA bicluster (bicluster 115). TuBA also identified biclusters corresponding to amplicons at 17q12 (HER2), enriched in ER- samples ($FDR = 0.02$); 8q24.3, enriched in ER- ($FDR = 0.003$) samples; and 17q25-q25.3, enriched in ER- samples ($FDR = 7.09e-05$). Thus, in addition to the differential co-expression network identified by Bicmix, TuBA recovered biclusters associated with genomic alterations such as CNA, several of which are differentially expressed between ER+ and ER- samples. Overall, TuBA recovered 144 biclusters enriched in ER- samples and 31 biclusters enriched in ER+ samples ($FDR < 0.05$). This is consistent with our earlier observation that a significant proportion of biclusters discovered independently in the TCGA, METABRIC, and GEO datasets were enriched in the ER-/HER2- subtype.

Runtime analysis

TuBA's graph-based algorithm relies on the identification of largest cliques, which is a computationally hard problem. Large graphs (both in terms of the number of genes, and edges) can potentially lead to long computation times. The size of our graphs is principally determined by the choice of the cut-off for the second parameter—the significance level of overlap between percentile sets. We varied the cut-offs for the TCGA and METABRIC datasets such that the total number of edges in the resultant graphs ranged between 10,000 and 250,000.

For choices of overlap cut-offs consistent with our suggested heuristic, we recorded TuBA's computation time to generate final biclusters for each dataset (Fig. 9A). The computation time for TCGA increased dramatically as the number of edges in the graphs exceeded 150,000. In particular, while the computation time for a graph with 200,000 edges for METABRIC was ~70 minutes, the computation time for a graph of similar size for TCGA was ~42 hours. Thus, although METABRIC is the larger dataset with 24,368 genes and 1,970 samples compared to TCGA's 20,241 genes and 908 samples, more iterations were required to identify all the largest cliques in the graphs for TCGA given its respective choices of parameters. We therefore conclude that TuBA's computation time depends on the nature and complexity of the graphs themselves.

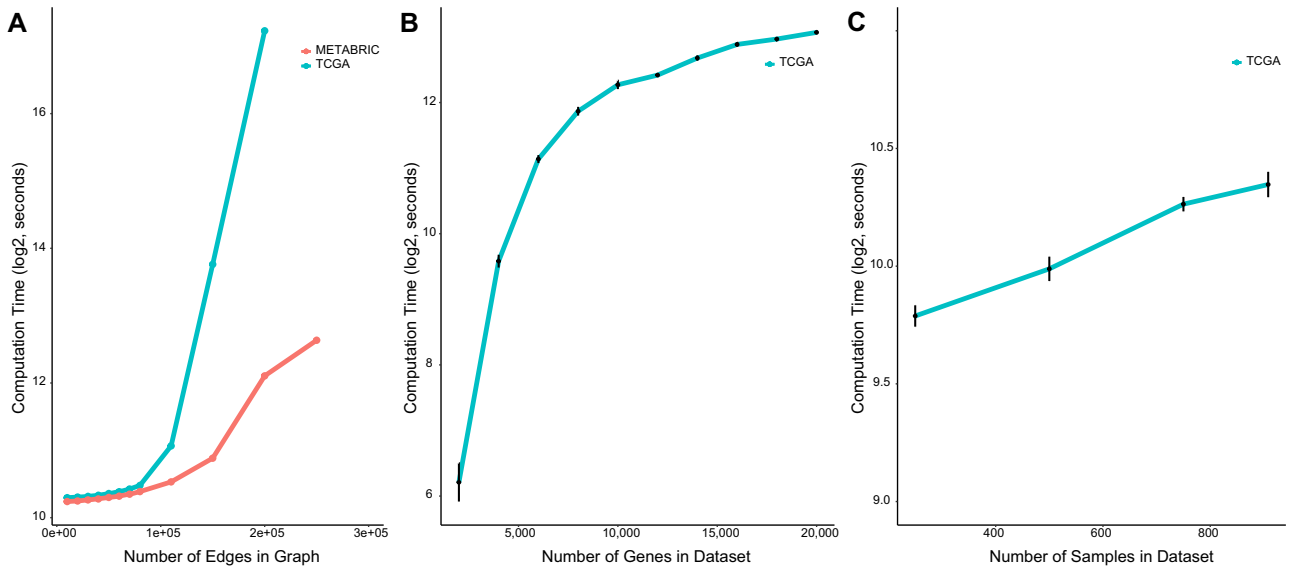


Figure 9: Runtime analysis of TuBA. (A) Computation time taken by TuBA to discover biclusters for different sizes (number of edges) of the graphs based on different choices of overlap cut-offs for the TCGA (blue) and METABRIC (red) datasets. (B) Dependence of computation time on the number of rows (genes) for choices of overlap cut-offs consistent with our suggested heuristic. Here, we chose overlap cut-offs consistent with our suggested heuristic and ensured that comparable numbers of edges were generated for different datasets; this explains why the computation times for datasets with upwards of 14,000 genes were quite similar to each other. (C) Dependence of computation time on the number of columns (samples) for choices of overlap cut-offs consistent with our suggested heuristic.

We also investigated the impact of the size of datasets on computation time. We created new subsets from the TCGA dataset by randomly sampling a fixed number of genes. We varied the number of genes from 2,000 to 20,000 and created 10 randomly sampled datasets for each gene number (Fig. 9B). For each individual run, we chose an overlap cut-off consistent with our suggested heuristic and ensured that comparable numbers of edges were generated for different datasets.

We also investigated the impact of the number of samples in a dataset. For this, we created 5 randomly selected subsets from the TCGA dataset each with 250, 500, 750, and 908 samples (Fig. 9C). As expected, TuBA's computation time did not depend strongly on the number of samples in the datasets.

In its current implementation, using a 2.7-GHz Intel Xeon processor and 48 GB of RAM, TuBA has longer runtime than most other existing algorithms. Depending on the choice of the overlap cut-off, the runtimes can vary between 15 and 120 minutes for datasets with $\sim 20,000$ genes and 1,000 samples.

Discussion

Global clustering approaches have successfully unveiled distinct disease subtypes in tumors, prompting the community to look beyond traditional clinico-pathological signatures to identify relevant disease processes. However, the extensive heterogeneity, even within tumors of a given subtype, confounds the identification of many altered transcriptional programs by such unsupervised clustering methods.

In this paper, we introduce an algorithm called TuBA based on a proximity measure specifically designed to extract gene co-expression signatures that correspond to the extremes of expression (both high and low for RNA-Seq data, and high for array-based platform). This enables us to preferentially identify co-aberrant gene signatures associated with the disease states of tumors. The identification of altered transcriptional profiles can be particularly relevant for those tumors that have so far eluded targeted drug development for therapy. This is exempli-

fied by tumors of the basal-like or triple-negative subtypes for BRCA. Although these tumors account for only $\sim 15\%$ of all BRCAAs in the population, a significant fraction of biclusters identified by TuBA corresponded to alterations associated with tumors of these subtypes (Fig. S10). For each dataset, a simple estimation of enrichment of samples in a given bicluster within any other bicluster revealed that the samples in the biclusters corresponding to CNA at 8p11.21–p11.23 or 17q12 were enriched ($FDR < 0.001$) independently in $\sim 5\%$ of all biclusters for both TCGA and METABRIC, respectively. In sharp contrast, 30–40% of all biclusters were enriched in samples with copy number gains at the 8q24.3 locus ($FDR < 0.001$). Additionally, 51% of all biclusters obtained from the low-expression analysis of TCGA were enriched in the samples corresponding to the 8q24.3 bicluster. Previous studies have also identified the amplicon at 8q24.3 by representational difference analysis as a location of oncogenic alterations in breast cancer that can occur independent of neighboring MYC amplifications [67]. Although the 8q24.3 bicluster itself is enriched with ER-/HER2- samples, these observations, together with poor RFS outcome observed independently in both METABRIC and GEO, highlight this locus as a promising prognostic marker for BRCA tumors, irrespective of subtype.

We must mention a notable exception in the biclusters discovered by TuBA for all 3 datasets—none of the biclusters contained the ESR1 gene, which codes for estrogen receptor (ER). Closer inspection revealed that ESR1 had statistically significant associations with several genes; however, the level of significance of overlap with these genes was much lower ($FDR > 1e-07$) than the chosen cut-off for all 3 datasets. While it is known that $\sim 70\%$ of BRCAAs exhibit elevated expression of ER, overexpression of ER may not be a sufficient condition to drive co-expression of genes involved in other pathways [68]. This may explain why co-expression of ESR1 with other genes was not as significant as the other associations that were extracted and summarized in TuBA's biclusters.

Apart from highlighting the heterogeneity of CNA-associated alterations in tumors of the ER-/HER2- subtype or basal-like sub-

type, TuBA offered a glimpse into the utility, the limitations, and the potential pitfalls with the current subtype classification approaches. In the ER/HER2-based subtype enrichment, we observed a significant proportion of biclusters that were not specifically enriched in any of the 4 subtypes. For instance, several CNA-associated biclusters from chromosome 8 were not subtype-enriched. In the case of PAM50 subtype classification, however, we observed that most of these biclusters were enriched in the luminal B subtype for METABRIC (and to a limited extent for TCGA). While this appears to indicate that PAM50 offers an improvement on the traditional clinico-pathological approach to subtype classification, it unfortunately fails to classify several samples associated with overexpression of ERBB2 as HER2-positive. As a consequence, several of our biclusters associated with the HER2 amplicon and copy number gains in the neighboring locations on chromosome 17 (17q.21.1–q21.2 and 17q21.32–q21.33), for both METABRIC and TCGA, were observed to be enriched in the luminal B subtype. This corroborates the modest level of agreement with PAM50 classification reported by Parker et al. [41], as well as disagreements in later studies [69]. Given that trastuzumab is a clinically proven therapeutic drug for HER2+ tumors, misclassification of these patients into any other subtype can be highly disadvantageous.

Change in copy number is often not a sufficient condition for elevated (or suppressed) expression levels of transcripts because there are multiple layers of regulation of transcription in cells [70, 71]. TuBA specifically identifies sets of genes with copy number changes that are transcriptionally active (or inactive), filtering out the ones that are unlikely to influence disease progression. Moreover, the graph-based approach allows us to infer the relative importance of each gene within a bicluster, based on its degree. In the case of high-expression analysis, the degree of each gene is an indicator of how frequently it is expressed aberrantly at high levels by the subset of samples that make up any given bicluster. As an example, consider the CNA-associated bicluster from TCGA corresponding to gains at the 8q22.1–q22.3 loci. The bicluster exhibited enrichment in patients with lymph node-positive disease (the corresponding bicluster in METABRIC has a significance level of $FDR = 0.052$ for patients with positive lymph node status). The gene with the highest degree in the bicluster was *MTDH* (metadherin), which has been shown to be associated with increased chemotherapy resistance and metastasis in BRCA [72–74].

Clustering analysis of biclusters and samples based on the membership of samples within biclusters allowed us to identify the sites that were altered concomitantly within the same subsets of samples. Moreover, we improved our perspective on the tumor microenvironment in the subsets of samples that exhibit non-tumor-associated signatures (e.g., immune, extracellular matrix). Differences in disease progression due to distinct microenvironments in tumors with similar transcriptional alterations can help us better understand the potential role of the microenvironment within the context of tumors harboring these specific alterations. For instance, we noticed a difference in RFS outcomes between 2 groups of patients who exhibit copy number gains at 8q24.3; the group that was additionally associated with an immune response signature was observed to have better RFS outcomes compared to the group that did not exhibit a strong association with the immune response.

Unlike most biclustering methods, TuBA does not allow arbitrary overlaps between its biclusters. This is because it is designed to discover biclusters with samples that correspond to the extremals for the corresponding gene set; biclusters with other conditions are not permitted for the same gene set. How-

ever, our biclusters are not exclusive, and some overlap between their genes and samples is permitted. For example, in the case of an ER-/HER2- BRCA sample that exhibits CNA at 8q24.3, because of high immune-cell infiltration in the tumor, the same sample may also be present in the biclusters enriched in the sets of genes associated with immune response.

Another limitation of TuBA is that it can only be applied reliably for large datasets that contain ≥ 100 samples. Depending on cohort heterogeneity, some of the overlaps between percentile sets may not be significant in smaller datasets. However, the deliberate design of our proximity measure leveraging the size of the datasets offers a significant benefit—it not only enables the identification of the plethora of gene co-aberrations associated with the tumors but also enables the estimation of the extent or prevalence of the identified alterations in the population. This is where the tunable aspect of TuBA becomes relevant—the 2 “knobs” should be viewed as valuable aids that help estimate the extents of the prevalence of various alterations in the tumor population and their clinical relevance. Although transcriptomic changes are not the ultimate determinants of progression, our algorithm holds the promise to improve therapeutic selection and design by identifying significantly altered transcriptional patterns associated with tumors.

Conclusions

TuBA is quite distinct from other biclustering algorithms, in that it is designed to identify biclusters with samples that correspond to the extremals for the corresponding sets of genes. Most biclustering algorithms are designed to identify nearly constant or coherent gene expression levels in subsets of genes across subsets of samples. However, we were able to show that TuBA outperforms other algorithms in identification of co-expressed genes located in transcriptionally active copy number-altered sites. Moreover, from a differential co-expression perspective, TuBA offers an advantage over other methods because no prior specification of subsets of samples (context) is necessary; the nature of our proximity measure ensures that such differential co-expression signatures are preferentially identified. Given these considerations, TuBA offers great promise as a biclustering method that can identify biologically relevant gene co-expression signatures that are not successfully captured by other unsupervised clustering or biclustering approaches. These signatures, along with the ones identified by other biclustering methods, would enable a comprehensive understanding of the underlying alterations and shared mechanisms in subsets of tumors.

Availability of data and materials

TuBA is open-sourced and available in R scripts at <https://github.com/KhiabanianLab>. TCGA dataset was obtained from UCSC Xena Portal (<http://xena.ucsc.edu>). METABRIC dataset was obtained from the cBioPortal (<http://www.cbioportal.org>). GEO and breastCancerNKI datasets were obtained from Gyorffy & Schafer [26], and van't Veer et al. [27] and van de Vijver et al. [28], respectively. Supporting data and materials are available in the GigaScience GigaDB database [75].

Availability of supporting source code and requirements

Project Name: TuBA: Tunable Biclustering Algorithm

Project home page: <https://github.com/KhiabanianLab/TuBA>
 Operating System(s): Platform Independent
 Programming language: R
 Other requirements: R 3.3.0 or higher
 License: GNU GPL v3
 RRID:SCR.017121

Additional files

Fig. S1. (A) Histogram for overlap significance values (in $-\log_{10}$ scale) based on a permutation test on the METABRIC dataset. The P-values were not corrected for multiple hypothesis testing. (B) Histogram for overlap significance values (in $-\log_{10}$ scale) after correcting for multiple hypothesis testing.

Fig. S2. Impact of the size of top percentile set for the bicluster corresponding to the HER2 amplicon (17q12). (A) Number of samples in bicluster as the overlap significance is lowered from 10^{-20} to 10^{-16} for percentile set size of 5%, (B) number of samples in bicluster as the overlap significance is lowered from 10^{-35} to 10^{-23} for percentile set size of 10%.

Fig. S3. Enrichment of biclusters consisting of proximally located genes with copy number gains in the PAM50 subtypes for (A) METABRIC and (C), (E) TCGA. The biclusters are represented by horizontal bars in each panel, color-coded according to the chromosome number of their constituent genes. Panels (B), (D), and (F) show the remaining biclusters arranged according to their serial numbers in Supplementary Table 3 for METABRIC and TCGA, respectively. The ones that are associated with copy number (CN) gains of genes located at distant chromosomal sites are shown in red while those associated with loss are shown in green. The rest are shown in black. Note, the thickness of the bar in each panel depends on the total number of biclusters displayed in that panel and so does not represent its chromosomal extent.

Fig. S4. (A) Enrichment of biclusters from TCGA consisting of proximally located genes with copy number loss in the ER/HER2 subtypes. The biclusters are represented by horizontal bars in each panel, color-coded according to the chromosome number of their constituent genes. Panel (B) shows the remaining biclusters arranged according to their serial numbers in Supplementary Table 3. The ones that exhibit copy number loss of genes located at distant chromosomal sites are shown in green, while the rest are shown in black. Note, the thickness of the bar in each panel depends on the total number of biclusters displayed in that panel and so does not represent its chromosomal extent.

Fig. S5. Enrichment of biclusters from GEO in the subtypes based on ER and HER2 status. The biclusters have been arranged according to their serial numbers in Supplementary Table 3 (bicluster 1 on top) and are represented by horizontal black bars.

Fig. S6. Degrees of genes in the bicluster corresponding to the HER2 amplicon (17q12) for (A) TCGA, (B) METABRIC, and (C) GEO.

Fig. S7. (A) Kaplan-Meier survival curve for the set of patients in the Her2 (17q12) bicluster (red) compared to the remaining set of patients (blue) for the METABRIC dataset, and (B) the graph corresponding to the bicluster.

Fig. S8. (A) Proportions of GO-BP term-enriched biclusters for each biclustering method at 5 different significance levels for the METABRIC dataset. (B) Ratios of number of unique GO-BP terms and total number of biclusters at 5 different significance levels for the METABRIC dataset.

Fig. S9. (A) Proportions of GO-BP term-enriched biclusters obtained by TuBA at 5 different significance levels for 5 choices of the overlap significance cut-off for the TCGA dataset. (B) Ratios

of number of unique GO-BP terms to the total number of biclusters at 5 different significance levels for 5 choices of the overlap significance cut-off for the TCGA dataset.

Fig. S10. Summary of the results from primary analysis of breast cancer datasets.

Abbreviations

BRCA: breast invasive carcinoma; CNA: copy number amplification; CNV: copy number variation; DeBi: Differentially Expressed Biclusters; DLBCL: diffuse large B-cell lymphoma; ER: estrogen receptor; FDR: false discovery rate; GB: gigabytes; GEO: Gene Expression Omnibus; GISTIC2: Genomic Identification of Significant Targets in Cancer; GO-BP: Gene Ontology-Biological Processes; GTEx: Genotype-Tissue Expression; HER2+: human epidermal growth factor receptor 2 positive; ISA: Iterative Signature Algorithm; METABRIC: Molecular Taxonomy of Breast Cancer International Consortium; NIH: National Institutes of Health; OPSP: Order-Preserving Submatrix; PAM50: Prosigna Breast Cancer Prognostic Gene Signature Assay; QUBIC: Qualitative Biclustering; RAM: random access memory; RFS: recurrence-free survival; RNA-Seq: RNA sequencing; rRNA: ribosomal RNA; TCGA: The Cancer Genome Atlas; TuBA: Tunable Biclustering Algorithm; UCSC: University of California Santa Cruz.

Competing interests

The authors declare that they have no competing interests.

Funding

H.K. acknowledges support from the American Cancer Society (IRG-15-168-01) and a grant from the National Cancer Institute (R01CA233662). The funding agencies had no role in the design of the study and collection, analysis, and interpretation of data or in writing of the manuscript.

Authors' contributions

A.S., G.B., and H.K. conceived the study and designed the algorithm. A.S. implemented the algorithm and performed the statistical analyses. All authors contributed to the drafting of the manuscript and critical discussion of the results. All authors read and approved the final manuscript.

Acknowledgments

This research was partially supported by the Biomedical Informatics Shared Resource at Rutgers Cancer Institute of New Jersey (P30CA072720) as well as Rutgers Office of Advanced Research Computing (NIH 1S100D012346-01A1).

References

1. Eisen MB, Spellman PT, Brown PO, et al. Cluster analysis and display of genome-wide expression patterns. *Proc Natl Acad Sci U S A* 1998;95(25):14863–8.
2. Alizadeh AA, Eisen MB, Davis RE, et al. Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling. *Nature* 2000;403(6769):503–11.
3. Roth FP, Hughes JD, Estep PW, et al. Finding DNA regulatory motifs within unaligned noncoding sequences clustered by whole-genome mRNA quantitation. *Nat Biotechnol*

- 1998;16(10):939–45.
4. Golub TR, Slonim DK, Tamayo P, et al. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science* 1999;286(5439):531–7.
 5. Perou CM, Jeffrey SS, van de Rijn M, et al. Distinctive gene expression patterns in human mammary epithelial cells and breast cancers. *Proc Natl Acad Sci U S A* 1999;96(16):9212–7.
 6. Cheng Y, Church GM. Biclustering of expression data. *Proc Int Conf Intell Syst Mol Biol* 2000;8:93–103.
 7. Madeira SC, Oliveira AL. Biclustering algorithms for biological data analysis: a survey. *IEEE/ACM Trans Comput Biol Bioinform* 2004;1(1):24–45.
 8. Eren K, Deveci M, Kucuktunc O, et al. A comparative analysis of biclustering algorithms for gene expression data. *Brief Bioinform* 2013;14(3):279–92.
 9. Tanay A, Sharan R, Shamir R. Discovering statistically significant biclusters in gene expression data. *Bioinformatics* 2002;18(Suppl 1):S136–144.
 10. Oghabian A, Kilpinen S, Hautaniemi S, et al. Biclustering methods: biological relevance and application in gene expression analysis. *PLoS One* 2014;9(3):e90801.
 11. Yoon S, Nardini C, Benini L, et al. Discovering coherent biclusters from gene expression data using zero-suppressed binary decision diagrams. *IEEE/ACM Trans Comput Biol Bioinform* 2005;2(4):339–54.
 12. Getz G, Levine E, Domany E. Coupled two-way clustering analysis of gene microarray data. *Proc Natl Acad Sci U S A* 2000;97(22):12079–84.
 13. Sheng Q, Moreau Y, De Moor B. Biclustering microarray data by Gibbs sampling. *Bioinformatics* 2003;19(Suppl 2):ii196–205.
 14. Hochreiter S, Bodenhofer U, Heusel M, et al. FABIA: factor analysis for bicluster acquisition. *Bioinformatics* 2010;26(12):1520–7.
 15. Pontes B, Giraldez R, Aguilar-Ruiz JS. Biclustering on expression data: a review. *J Biomed Inform* 2015;57:163–80.
 16. Pontes B, Giraldez R, Aguilar-Ruiz JS. Quality measures for gene expression biclusters. *PLoS One* 2015;10(3):e0115497.
 17. van Dam S, Vosa U, van der Graaf A, et al. Gene co-expression analysis for functional classification and gene-disease predictions. *Brief Bioinform* 2018;19:575–92.
 18. Bron C, Kerbosch J. Algorithm 457: finding all cliques of an undirected graph. *Commun ACM* 1973;16(9):575–7.
 19. Karp RM. Reducibility among combinatorial problems. In: Miller RE, Thatcher JW, eds. *Proceedings of a Symposium on the Complexity of Computer Computations*, Yorktown Heights, NY, 1972. Plenum; 1972:85–103.
 20. Cancer Genome Atlas Network. Comprehensive molecular portraits of human breast tumours. *Nature* 2012;490(7418):61–70.
 21. Mermel CH, Schumacher SE, Hill B, et al. GISTIC2.0 facilitates sensitive and confident localization of the targets of focal somatic copy-number alteration in human cancers. *Genome Biol* 2011;12(4):R41.
 22. UCSC Xena. <http://xena.ucsc.edu>. Accessed April 15, 2019.
 23. cBioPortal for Cancer Genomics. <http://www.cbioportal.org>. Accessed April 15, 2019.
 24. Cerami E, Gao J, Dogrusoz U, et al. The cBio cancer genomics portal: an open platform for exploring multidimensional cancer genomics data. *Cancer Discov* 2012;2(5):401–4.
 25. Gao J, Aksoy BA, Dogrusoz U, et al. Integrative analysis of complex cancer genomics and clinical profiles using the cBioPortal. *Sci Signal* 2013;6(269):pl1.
 26. Györfy B, Schafer R. Meta-analysis of gene expression profiles related to relapse-free survival in 1,079 breast cancer patients. *Breast Cancer Res Treat* 2009;118(3):433–41.
 27. van 't Veer LJ, Dai H, van de Vijver MJ, et al. Gene expression profiling predicts clinical outcome of breast cancer. *Nature* 2002;415(6871):530–6.
 28. van de Vijver MJ, He YD, van't Veer LJ, et al. A gene-expression signature as a predictor of survival in breast cancer. *N Engl J Med* 2002;347(25):1999–2009.
 29. Schroeder M H-KB, Culhane A, Sotiriou C, et al. breastCancerNKI: Gene expression dataset. R package version 1.0.6. 2011. <http://bioconductor.org/packages/breastCancerNKI/>
 30. Hastie TTR, Sherlock G, Eisen M, et al. Imputing missing data for gene expression arrays. Technical Report, Division of Biostatistics, Stanford University. 1999. <http://www.web.stanford.edu/~hastie/Papers/missing.pdf>. Accessed 15 April, 2019.
 31. ESTIMATE. <http://bioinformatics.mdanderson.org/estimate>. Accessed April 15, 2019
 32. GTEx Portal. www.gtexportal.org/home/. Accessed April 15, 2019
 33. Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biol* 2010;11(10):R106.
 34. R Core Team. R. A Language and Environment for Statistical Computing. 2016.
 35. Csardi G, Nepusz T. The igraph software package for complex network research. *InterJournal Complex Syst* 2006;1695.
 36. Wickham H. The split-apply-combine strategy for data analysis. *J Stat Softw* 2011;40(1):1–29.
 37. Shannon P, Markiel A, Ozier O, et al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res* 2003;13(11):2498–504.
 38. Wickham H. ggplot2: Elegant Graphics for Data Analysis. New York: Springer; 2009.
 39. Subhash S, Kanduri C. GeneSCF: a real-time based functional enrichment tool with support for multiple organisms. *BMC Bioinformatics* 2016;17(1):365.
 40. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc Ser A* 1995;57(1):289–300.
 41. Parker JS, Mullins M, Cheang MC, et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *J Clin Oncol* 2009;27(8):1160–7.
 42. Marguerat S, Bahler J. RNA-seq: from technology to biology. *Cell Mol Life Sci* 2010;67(4):569–79.
 43. Mortazavi A, Williams BA, McCue K, et al. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods* 2008;5(7):621–8.
 44. Kallioniemi A, Kallioniemi OP, Piper J, et al. Detection and mapping of amplified DNA sequences in breast cancer by comparative genomic hybridization. *Proc Natl Acad Sci U S A* 1994;91(6):2156–60.
 45. Kao J, Salari K, Bocanegra M, et al. Molecular profiling of breast cancer cell lines defines relevant tumor models and provides a resource for cancer gene discovery. *PLoS One* 2009;4(7):e6146.
 46. Curigliano G, Viale G, Ghioni M, et al. Cancer-testis antigen expression in triple-negative breast cancer. *Ann Oncol* 2011;22(1):98–103.
 47. Simpson AJ, Caballero OL, Jungbluth A, et al. Cancer/testis antigens, gametogenesis and cancer. *Nat Rev Cancer* 2005;5(8):615–25.
 48. Yoshihara K, Shahmoradgoli M, Martinez E, et al. Inferring tumour purity and stromal and immune cell admixture from expression data. *Nat Commun* 2013;4:2612.

49. Curtis C, Shah SP, Chin SF, et al. The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature* 2012;**486**(7403):346–52.
50. Bilal E, Vassallo K, Toppmeyer D, et al. Amplified loci on chromosomes 8 and 17 predict early relapse in ER-positive breast cancers. *PLoS One* 2012;**7**(6):e38575.
51. Wang J, Liu Q, Shyr Y. Dysregulated transcription across diverse cancer types reveals the importance of RNA-binding protein in carcinogenesis. *BMC Genomics* 2015;**16**(Suppl 7):S5.
52. Inic Z, Zegarac M, Inic M, et al. Difference between luminal A and luminal B subtypes according to Ki-67, tumor size, and progesterone receptor negativity providing prognostic information. *Clin Med Insights Oncol* 2014;**8**:107–11.
53. Haybittle JL, Blamey RW, Elston CW, et al. A prognostic index in primary breast cancer. *Br J Cancer* 1982;**45**(3):361–6.
54. Galea MH, Blamey RW, Elston CE, et al. The Nottingham Prognostic Index in primary breast cancer. *Breast Cancer Res Treat* 1992;**22**(3):207–19.
55. Serin A, Vingron M. DeBi: Discovering Differentially Expressed Biclusters using a Frequent Itemset Approach. *Algorithms Mol Biol* 2011;**6**(1):18.
56. Rosenwald A, Wright G, Chan WC, et al. The use of molecular profiling to predict survival after chemotherapy for diffuse large-B-cell lymphoma. *N Engl J Med* 2002;**346**(25):1937–47.
57. Bergmann S, Ihmels J, Barkai N. Iterative signature algorithm for the analysis of large-scale gene expression data. *Phys Rev E Stat Nonlin Soft Matter Phys* 2003;**67**(3 Pt 1):031902.
58. Ben-Dor A, Chor B, Karp R, et al. Discovering local structure in gene expression data: the order-preserving submatrix problem. *J Comput Biol* 2003;**10**(3-4):373–84.
59. Li G, Ma Q, Tang H, et al. QUBIC: a qualitative biclustering algorithm for analyses of gene expression data. *Nucleic Acids Res* 2009;**37**(15):e101.
60. Tanay A, Sharan R, Kupiec M, et al. Revealing modularity and organization in the yeast molecular network by integrated analysis of highly heterogeneous genomewide data. *Proc Natl Acad Sci U S A* 2004;**101**(9):2981–6.
61. Prelic A, Bleuler S, Zimmermann P, et al. A systematic comparison and evaluation of biclustering methods for gene expression data. *Bioinformatics* 2006;**22**(9):1122–9.
62. Sebastian Kaiser RS, Khamiakova T, Sill M, et al. biclust: Bi-Cluster Algorithms. 2018. <https://cran.r-project.org/web/packages/biclust/index.html>. Accessed April 15, 2019.
63. Csardi G, Kutalik Z, Bergmann S. Modular analysis of gene expression data with R. *Bioinformatics* 2010;**26**(10):1376–7.
64. Zhang Y, Xie J, Yang J, et al. QUBIC: a bioconductor package for qualitative biclustering analysis of gene co-expression data. *Bioinformatics* 2017;**33**(3):450–2.
65. Shamir R, Maron-Katz A, Tanay A, et al. EXPANDER—an integrative program suite for microarray data analysis. *BMC Bioinformatics* 2005;**6**:232.
66. Gao C, McDowell IC, Zhao S, et al. Context specific and differential gene co-expression networks via bayesian biclustering. *PLoS Comput Biol* 2016;**12**(7):e1004791.
67. Mu D, Chen L, Zhang X, et al. Genomic amplification and oncogenic properties of the KCNK9 potassium channel gene. *Cancer Cell* 2003;**3**(3):297–302.
68. Planas-Silva MD, Donaher JL, Weinberg RA. Functional activity of ectopically expressed estrogen receptor is not sufficient for estrogen-mediated cyclin D1 expression. *Cancer Res* 1999;**59**(19):4788–92.
69. Guiu S, Michiels S, Andre F, et al. Molecular subclasses of breast cancer: how do we define them? The IMPAKT 2012 Working Group Statement. *Ann Oncol* 2012;**23**(12):2997–3006.
70. Lee TI, Young RA. Transcriptional regulation and its misregulation in disease. *Cell* 2013;**152**(6):1237–51.
71. Lelli KM, Slattery M, Mann RS. Disentangling the many layers of eukaryotic transcriptional regulation. *Annu Rev Genet* 2012;**46**:43–68.
72. Wan L, Kang Y. Pleiotropic roles of AEG-1/MTDH/LYRIC in breast cancer. *Adv Cancer Res* 2013;**120**:113–34.
73. Song Z, Wang Y, Li C, et al. Molecular modification of Metadherin/MTDH impacts the sensitivity of breast cancer to doxorubicin. *PLoS One* 2015;**10**(5):e0127599.
74. Shi X, Wang X. The role of MTDH/AEG-1 in the progression of cancer. *Int J Clin Exp Med* 2015;**8**(4):4795–807.
75. Singh A, Bhanot G, Khiabani H. Supporting data for “TuBA: Tunable biclustering algorithm reveals clinically relevant tumor transcriptional profiles in breast cancer” GigaScience Database 2019. <http://dx.doi.org/10.5524/100601>.