

 Open access • Posted Content • DOI:10.1101/142281

## Phenotype-Driven Transitions In Regulatory Network Structure — [Source link](#)

Megha Padi, John Quackenbush

**Institutions:** Harvard University

**Published on:** 25 May 2017 - bioRxiv (Cold Spring Harbor Labs Journals)

**Topics:** Biological network

Related papers:

- [Evolutionary conservation and network structure characterize genes of phenotypic relevance for mitosis in human](#)
- [Complexity in cancer biology: is systems biology the answer?](#)
- [DENSE: efficient and prior knowledge-driven discovery of phenotype-associated protein functional modules](#)
- [ModuleBlast: identifying activated sub-networks within and across species](#)
- [Understanding tissue-specificity with human tissue-specific regulatory networks](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/phenotype-driven-transitions-in-regulatory-network-structure-50wnrw750>

## ARTICLE OPEN

## Detecting phenotype-driven transitions in regulatory network structure

Megha Padi<sup>1</sup> and John Quackenbush<sup>2,3</sup>

Complex traits and diseases like human height or cancer are often not caused by a single mutation or genetic variant, but instead arise from functional changes in the underlying molecular network. Biological networks are known to be highly modular and contain dense “communities” of genes that carry out cellular processes, but these structures change between tissues, during development, and in disease. While many methods exist for inferring networks and analyzing their topologies separately, there is a lack of robust methods for quantifying differences in network structure. Here, we describe ALPACA (ALtered Partitions Across Community Architectures), a method for comparing two genome-scale networks derived from different phenotypic states to identify condition-specific modules. In simulations, ALPACA leads to more nuanced, sensitive, and robust module discovery than currently available network comparison methods. As an application, we use ALPACA to compare transcriptional networks in three contexts: angiogenic and non-angiogenic subtypes of ovarian cancer, human fibroblasts expressing transforming viral oncogenes, and sexual dimorphism in human breast tissue. In each case, ALPACA identifies modules enriched for processes relevant to the phenotype. For example, modules specific to angiogenic ovarian tumors are enriched for genes associated with blood vessel development, and modules found in female breast tissue are enriched for genes involved in estrogen receptor and ERK signaling. The functional relevance of these new modules suggests that not only can ALPACA identify structural changes in complex networks, but also that these changes may be relevant for characterizing biological phenotypes.

*npj Systems Biology and Applications* (2018)4:16; doi:10.1038/s41540-018-0052-5

## INTRODUCTION

We tend to think of phenotypes as being characterized by differentially expressed genes or mutations in particular genes. However, the individual genes that show the greatest changes in expression in a phenotype do not tend to be drivers of that phenotype.<sup>1,2</sup> Despite the increasing power and depth of sequencing studies, identifying the causal mutations and single-nucleotide polymorphisms (SNPs) that are responsible for determining heritable traits and disease susceptibility remains challenging. Indeed, many studies have found thousands of genetic variants of small effect size contribute to common traits.<sup>3–5</sup> It has become apparent that complex regulatory interactions between multiple genes and variants can contribute to defining the state of the cell. Modeling such phenotypes requires that we have a clearer picture of how genes and proteins work together to perform normal cellular functions, and how remodeling the interactions between genes can cause changes in phenotype including disease.

In this context, it is useful to make a subtle shift and think of a phenotype as being defined by a network of interacting genes and gene products. It has been shown that analyzing the mathematical properties of such networks can provide important biological insight into phenotypic properties. For example, high-degree “hubs” in protein–protein interaction (PPI) networks are enriched for genes essential to growth.<sup>6</sup> Biological networks are known to have modular structure and contain closely interacting groups of nodes, or “communities”, that work together to carry

out cellular functions.<sup>7–9</sup> There are many analytical and experimental methods for inferring network models associated with different phenotypic states, and for computing topological properties like centrality and community structure.<sup>10–13</sup> However, the most significant questions we can ask of biological networks—how networks differ from each other, and how these differences in network structure drive functional changes—remain largely unanswered. A significant challenge in this area is the lack of computational approaches for finding meaningful changes in the structure of large complex networks.

Previous work on comparative analysis of biological networks has focused on the so-called “differential network”, the set of edges that are altered relative to a reference network.<sup>14</sup> While the advantage of this approach is its simplicity, there are several issues that arise in such an edge-based analysis. First, biological network inference has a relatively high rate of false negatives due to noise in both the experimental data that are used and in the network inference methods themselves. Consequently, it can be difficult to determine whether the appearance or disappearance of a single edge is “real”. The uncertainty in the estimate of the difference between two edge weights is the sum of the uncertainties in each individual edge, which inflates noise in the final differential network. Second, the perturbed network will in general contain both positive and negative changes in edge weight relative to the reference network, and it is challenging to analyze and interpret a differential network with mixed signs. If we only consider the new edges associated with a phenotype, we would miss the functional

<sup>1</sup>Department of Molecular and Cellular Biology, University of Arizona, Tucson, AZ, USA; <sup>2</sup>Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, MA, USA and <sup>3</sup>Department of Biostatistics, Harvard School of Public Health, Boston, MA, USA  
Correspondence: Megha Padi (mpadi@email.arizona.edu)

Received: 14 August 2017 Revised: 29 March 2018 Accepted: 2 April 2018  
Published online: 19 April 2018

effects of decreases in edge activity. Third, by focusing only on the altered edges and discarding common edges, the differential interactions are taken out of their functional context, making it difficult to connect them to global cellular changes. For example, adding or deleting ten scattered edges in a network may have very different consequences on the phenotype than would the same number of changes concentrated in a local functional neighborhood of the network.

One way to address these issues and find more robust differences between networks is to identify changes in groups of nodes, rather than in individual edges. Computational methods that have been developed to do this fall into several categories. First there are methods that evaluate differences in pre-specified network features, like user-defined gene sets, small regulatory motifs, or global topological characteristics. For example, Gamberdella et al. evaluated the statistical significance of differences in co-expression of a user-defined gene set between two conditions.<sup>15</sup> Similarly, the coXpress method defines clusters using co-expression in the reference condition, and tests for significant changes in each cluster under a new condition.<sup>16</sup> Landeghem et al. developed a method for inferring the best differential network that contrasts two datasets, and new measures have been devised to test whether global modular structure and degree characteristics are different between two networks.<sup>17–20</sup> However, these methods are limited to examining pre-defined gene modules and network features, and fail to take full advantage of the network structure. As such, they lack the ability to discover new pathways and network modules that functionally distinguish different phenotypes.

Other methods have been developed to discover de novo gene modules that differ between conditions. The DiffCoEx algorithm iteratively groups genes that are differentially co-expressed to find new modules.<sup>21,22</sup> Valcarcel et al. compared metabolite correlation networks to discover groups of metabolites that changed their correlation pattern between normal weight and obese mice.<sup>23</sup> These methods are based on first computing the most differential edges and then grouping them together, which increases the uncertainty of each edge estimate and does not incorporate functional edges that are present in both conditions,<sup>14,24</sup> thus losing network context.

Another class of methods attempts to identify “active modules”, which are groups of genes that are differentially expressed in a particular disease or condition and also highly connected in a reference network, such as the PPI network.<sup>25</sup> However, the “active modules” framework only uses differential gene expression and so focuses on the nodes rather than accounting for changes in the strength of regulatory edges.

We present a new graph-based approach called ALtered Partitions Across Community Architectures (ALPACA) that compares two networks and identifies de novo the gene modules that best distinguish the networks. ALPACA is based on modularity maximization, a technique commonly used to find communities in a single graph. As applied previously, modularity is a measure of the observed edge density of the communities as compared to their expected density in a degree-matched random graph. Although this technique is powerful, it has a “resolution limit” because communities can only be identified if they are larger than the typical cluster size in random graph configurations.<sup>26</sup> This lack of resolution is especially disadvantageous when studying transcriptional networks, which tend to have a dense and hierarchical structure, and whose functional units only become evident under different environmental conditions.<sup>27</sup> A framework based on modularity maximization has been created to find common community structure among multiple networks,<sup>28</sup> but the only way to detect differences is to apply modularity maximization to each network separately, followed by brute-force comparison of the two resulting community structures.

In ALPACA, we adapt the modularity framework to compare condition-specific networks to each other rather than to a random graph null model. We define a score called the “differential modularity” that compares the density of modules in the “perturbed” network to the expected density in a matched “baseline” network, allowing us to contrast, for example, networks from disease and healthy tissue samples and partition the nodes into optimal differential modules, without relying on predefined gene sets or pathways. In contrast to methods that simply cluster the most differential edges, ALPACA compares the full network structures active in each condition and reduces the noise from individual edges by estimating an aggregated null model. And because the null model is based on the community structure of a known reference network rather than on a random graph, the “resolution limit” is substantially smaller, and ALPACA can detect small disease modules otherwise hidden within larger regulatory programs associated with normal cellular functions.

To demonstrate the utility of ALPACA, we show that it can identify changes in the modular structure of simulated networks, and that it exhibits higher resolution and robustness than other network approaches. We then apply it to compare transcriptional networks derived from non-angiogenic and angiogenic subtypes of ovarian cancer, normal human fibroblasts and fibroblasts expressing tumor virus oncogenes, and male and female breast tissue from the Genotype-Tissue Expression (GTEx) project. In each case, we find that ALPACA identifies modules enriched in biological processes relevant to the phenotypes we are comparing.

## RESULTS

### Modularity maximization and detecting community structure

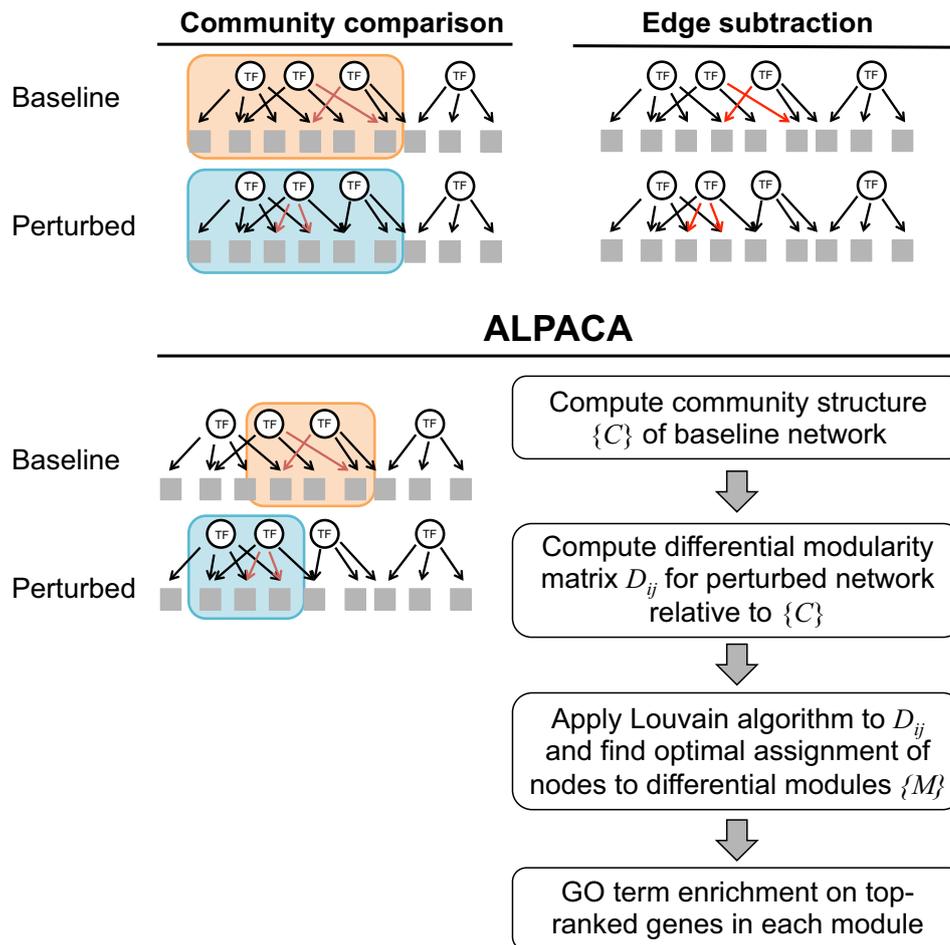
Many methods for determining the community structure of a network are based on maximizing the modularity.<sup>13</sup>

$$Q = \frac{1}{2m} \sum_{i,j} \left( A_{ij} - \frac{d_i d_j}{2m} \right) \delta(C_i, C_j). \quad (1)$$

Here,  $A_{ij}$  indicates the adjacency matrix of the network,  $m$  is the number of edges,  $d_i$  is the degree of node  $i$ , and  $C_i$  is the community assignment of node  $i$ . The modularity represents to what extent the proposed communities have more edges within them than expected in a randomly connected graph with the same degree properties; this null expectation is represented in the second term of the equation above. The modularity is optimized over the space of all possible partitions  $\{C\}$  and the value of  $C_i$  corresponding to the maximum modularity then determines the community structure of the network. An exhaustive search is not possible for large networks, but many methods have been developed to find locally optimal community structure, including ones based on edge betweenness, label propagation, and random walks.<sup>13,29,30</sup> The Louvain algorithm is a particularly efficient way to find high-quality local optima of the modularity function.<sup>31</sup>

### Community comparison and edge subtraction

Having arrived at a pair of inferred networks corresponding to different phenotypic states, there are two straightforward ways to compare the community structures based on the modularity metric (Fig. 1). One method, which we will call “community comparison”, consists of using modularity maximization to find the community structure for each network individually, and then finding the nodes that alter their community membership between the two networks. Another method, which we will call “edge subtraction”, is to compute the differences in the edge weights between the two networks, and then apply modularity maximization to the resulting subtracted weights.



**Fig. 1** Methods to compare networks and find changes in modular structure. “Community comparison” identifies communities separately in each network and looks for nodes that change their community membership. “Edge subtraction” finds communities by subtracting the networks and finding communities in the resulting differential edges (red arrows). ALPACA looks for groups of genes that are more interconnected in the perturbed network than expected given the community structure of the baseline network. Flowchart shows the major steps in the implementation of ALPACA

Both methods can detect large, dramatic changes in network structure. However, there are important differences in these methods. “Community comparison” is limited in its ability to detect structural changes smaller than the average community size in each individual network. In contrast, “edge subtraction” acts on the difference of the edge weights, which reduces the density of the network and increases the resolution, but this method is also more strongly affected by noise in the individual edges. Further, only positive edge weight differences can be used to run modularity maximization in the subtracted network, so edges that are lost are not appropriately accounted for; incorporating both positive and negative edge weight differences requires more complex techniques.<sup>32,33</sup>

**ALPACA:** a new method for detecting changes in community structure

To overcome some of the limitations of the community comparison and edge subtraction methods, we developed ALPACA, a new algorithm based on modularity maximization. The unique aspect of ALPACA is that, rather than comparing edge distributions to a random null model, we compare edges of the “perturbed” network to a null model based on the “baseline” network to find differential gene modules between the two networks (Fig. 1). ALPACA optimizes a new quantity called

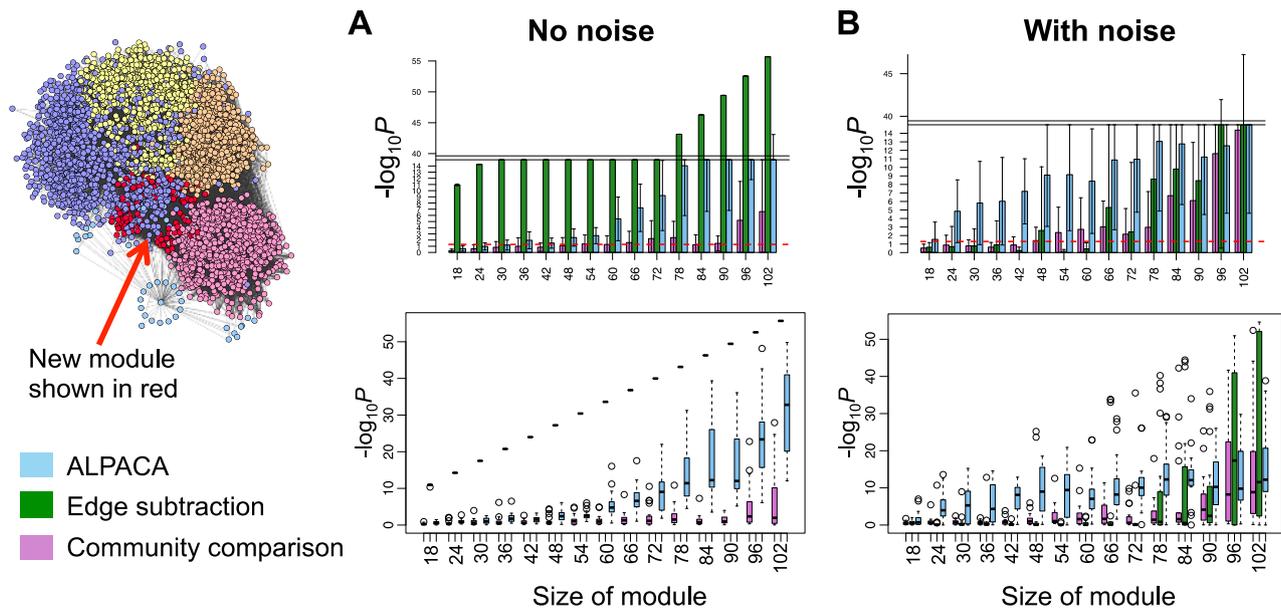
“differential modularity”, which we define as:

$$D = \frac{1}{m_p} \sum_{ij} D_{ij} \delta(M_i, M_j) = \frac{1}{m_p} \sum_{ij} (A_{ij}^p - N_{ij}) \delta(M_i, M_j). \quad (2)$$

This score compares the number of edges in a module  $M$  in the perturbed network—whose adjacency matrix is given by  $A_{ij}^p$  and total edge weight is  $m_p$ —to the expected number of edges  $N_{ij}$  based on the pre-computed community structure  $\{C\}$  of the baseline network. Here,  $N_{ij}$  is defined as:

$$N_{ij} = \frac{\left( \sum_{b \in C_j} \tilde{w}_{ib} \right) \left( \sum_{a \in C_i} \tilde{w}_{aj} \right)}{\sum_{a \in C_i, b \in C_j} \tilde{w}_{ab}}, \quad (3)$$

where  $C_i$  is the community assignment of node  $i$  in the baseline network, and  $\tilde{w}_{ab}$  is the normalized weight of the edge between node  $a$  and node  $b$  in the baseline network:  $\tilde{w}_{ab} = \left( \frac{m_p}{m} \right) w_{ab}$ . For the normalization, we have chosen to globally scale the edge weights of the baseline network so that the total matches  $m_p$ , the sum of the edge weights in the perturbed network. This allows a fair comparison between two networks that could be derived from two datasets of differing quality or sample size and may have different global sensitivity properties. To identify the modules  $\{M\}$  that maximize the differential modularity, we use the following two-step procedure. First, we determine the community structure



**Fig. 2** Performance of three methods on simulated networks with added module. Network at left visualizes the regulatory network derived from normal human fibroblasts, with purple, yellow, orange, pink, and blue denoting the pre-existing community structure, and red nodes depicting the synthetically added module. Bar graphs show performance of each method—ALPACA, edge subtraction or community comparison—on network simulations with (a) or without (b) resampling of edges among the pre-existing communities.  $P$ -values were computed using a one-sided Wilcoxon test. Bar graphs show mean of  $-\log_{10}P$  over 20 network simulations, and error bars depict the corresponding standard deviation. Boxplots represent same data as the bar plots. Boxplot elements are defined as follows: center line, median; box limits, upper and lower quartiles; whiskers, 1.5 $\times$  interquartile range; points, outliers. Note that the color of the boxplots for the edge subtraction method in a is not visible because the distribution is very narrow

of the baseline network using established methods.<sup>9,31</sup> Second, we compute the differential modularity matrix  $D_{ij}$  and apply the Louvain optimization algorithm to iteratively aggregate the nodes into modules.<sup>31</sup>

Note that the equation above is presented in a form that applies to weighted bipartite networks, as we will be applying it to analyze transcription factor (TF)–gene interactions. It can be easily adapted to analyze other types of networks. More details about the implementation of all three methods—community comparison, edge subtraction, and differential modularity—are presented in the Materials and Methods section.

#### Evaluating the performance of ALPACA on simulated networks

We reasoned that ALPACA would be more sensitive to small changes in modular structure than methods based on standard community detection, because the null model is computed using detailed properties of the baseline network rather than relying on random graphs. We also believed that ALPACA would be less sensitive to noise in individual edge weights than edge subtraction, because the null model is estimated by averaging over communities in the baseline network. We set out to test these properties in a setting that resembles real biological networks as much as possible, but where we have control over the changes in modularity.

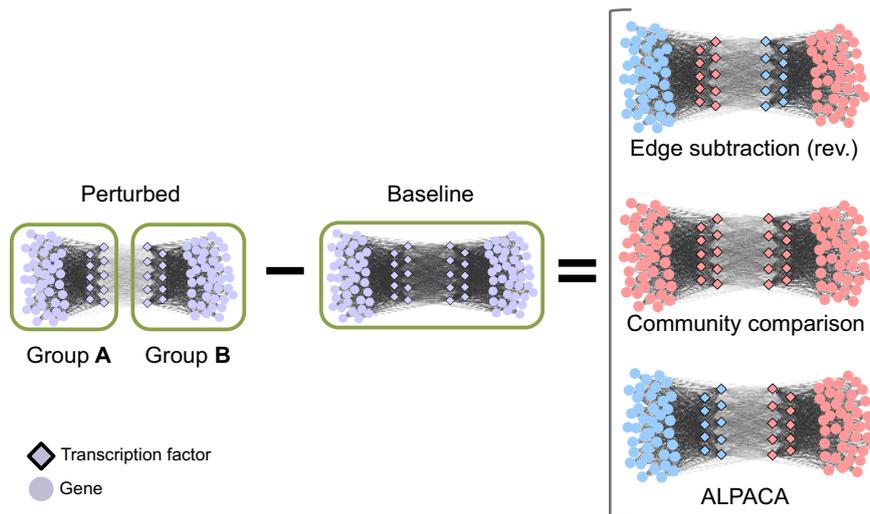
To do this, we constructed a baseline network and then created new modules through the “addition” of new edges, resulting in a perturbed network. For the noiseless version of this simulation, we inferred a regulatory network by integrating known human TF-binding sites with gene expression data in normal human fibroblasts using the algorithm PANDA<sup>34</sup> (see Materials and Methods for further details). After thresholding the edge weights and applying CONDOR,<sup>9</sup> a method for community detection in bipartite networks, we found that the baseline network had five communities of varying sizes. Next, we simulated a set of

perturbed networks by choosing a random subset of TFs and genes and adding new edges between them, thus artificially creating a new module. The new module consisted of between 3 and 21 TFs, and five times as many genes as TFs.

To these simulated networks, we applied three differential community detection methods—community comparison, edge subtraction, and ALPACA—and ranked the nodes by their contribution to the final score for each method. We then used Kolmogorov–Smirnov and Wilcoxon tests to evaluate whether the “true” module ranked higher than expected by chance in each ranked list. The edge subtraction method demonstrated superior performance for recovering modules of all sizes (Fig. 2a); this is to be expected, since the only new edges added to the networks were within the new modules. Examining the results from the other two methods, we observed that ALPACA is substantially better than community comparison at detecting smaller modules ranging down to a size of 50 nodes.

We then introduced edge noise into the “addition” simulation while retaining the modular structure of the underlying network. To do this, we made another series of perturbed networks, where, in addition to introducing the new module as described above, we also randomly resampled the edges from the baseline network while retaining the inter-community and intra-community edge density. In this more realistic set of simulations, we found that ALPACA outperformed the other methods on modules in the range of 18–90 nodes (Fig. 2b).

To check that these results are independent of the particular optimization algorithm used, we repeated the analysis using the Louvain method instead of CONDOR for initial community detection in the community comparison and edge subtraction methods. The results were similar in both cases, and in particular, ALPACA still outperformed the other methods on modules in the range of 18–54 nodes (Supplementary Fig. 1). This indicates that the superior performance of ALPACA is not due to the



**Fig. 3** Performance of three methods on perturbations that decrease edge density. Left-hand side shows a network transition involving a decrease in edge weights between nodes in groups A and B. All other edges remain the same. Right-hand side shows the results of three methods when comparing these two networks. Each method identified up to two differential modules, which are distinguished by their light blue and light pink colors in each case. Note that the “edge subtraction” method needs to be applied in the reverse manner, comparing the baseline network against the perturbed network, in order to have positive differential edge weights

optimization method used, but rather arises directly from the definition of the differential modularity.

While the edge subtraction method works well to detect “added” modules under low noise conditions, it becomes problematic if edges are deleted or if their weights decrease in the perturbed state relative to the control, because most network clustering methods are only formulated for positive edge weights. One might suggest transformation of edge weights, but any simple transformation of negative edge weights to make them positive (for example, by exponentiation or a linear shift) would bias the results. Algorithms that directly incorporate negative edge weights are complex and involve multiple steps and assumptions.<sup>32,33</sup> In contrast, ALPACA’s differential modularity matrix  $D_{ij}$  contains both negative and positive values, corresponding to areas of decreasing and increasing edge density relative to the baseline network and its community structure. By optimizing over the sum of  $D_{ij}$ , ALPACA incorporates positive and negative changes in edge density in a symmetric fashion.

As a simple demonstration of ALPACA’s ability to detect community structure changes with negative weights, we created “subtracted” simulations in which selected edges in a baseline network are reduced in weight to produce a substantially different perturbed network structure (Fig. 3 and Supplementary Fig. 2; see Materials and Methods for more details). In Fig. 3, for example, the network consists of two dense node groups, A and B, which are more strongly connected together in the baseline condition (edge weight 0.8) than in the perturbed condition (edge weight 0.2). Therefore, the perturbation causes groups A and B to separate and perform distinct functions; intuitively, this means groups A and B characterize the change in modular structure between the two networks. Because the only change in edge weights is the decrease in edges between A and B, the edge subtraction method results in a network with negative edge weights.

If instead we reverse the process and subtract the perturbed network from the baseline network, the resulting positive edge weight network produces two modules, one consisting of TFs in group A linked with genes in group B, the other consisting of TFs in group B linked with genes in group A. This does not match the intuitive result we are looking for. The community comparison method detects no change because both the baseline and perturbed networks are composed of the same two node

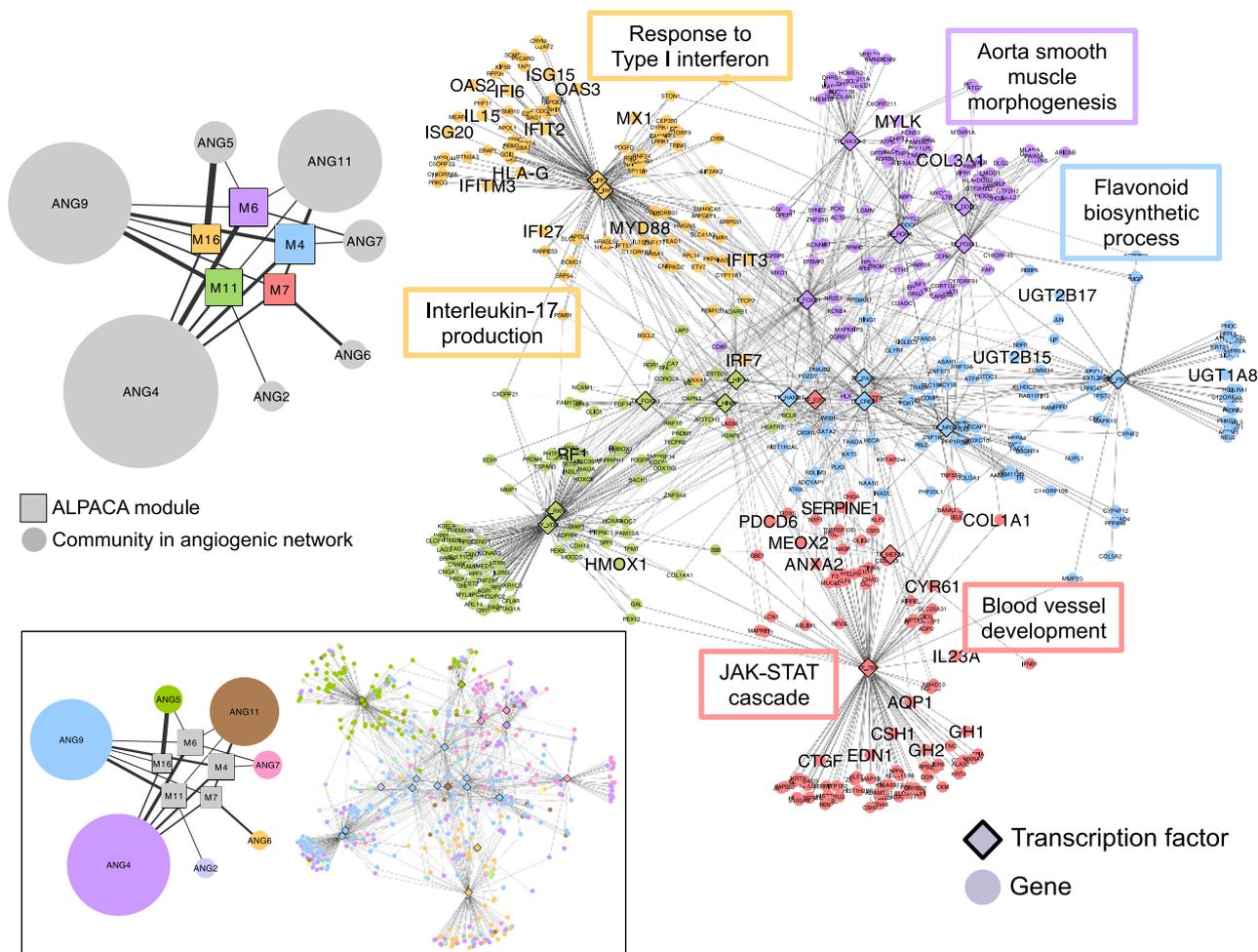
communities. However, ALPACA correctly identifies groups A and B as the differential modules characterizing this transition.

An example with three node groups is shown in Supplementary Fig. 2. Again, we find that ALPACA identifies the key change in modular structure and edge subtraction cannot. Although these examples are simple, such areas of decreased edge density will be locally embedded in any realistic biological network and will strongly influence the identification of neighboring modules.

#### Angiogenic vs. non-angiogenic ovarian cancer tumors

Ovarian cancer is the second most common cause of cancer death among women in the developed world. Available treatment options for ovarian cancer, such as platinum-based therapies, often lead to chemoresistance and recurrence. Ovarian cancer tumors can be stratified by gene expression profile, tissue of origin, or other characteristics, in order to better understand heterogeneity and predict patient-specific therapeutic strategies. We previously found that a gene signature associated with angiogenesis is able to classify ovarian cancer patients into a poor-prognosis subtype.<sup>35</sup>

We classified 510 ovarian cancer patients from The Cancer Genome Atlas into 188 angiogenic and 322 non-angiogenic tumors and used PANDA to infer separate gene regulatory networks for the two subtypes, as previously described.<sup>36</sup> We then applied a variety of methods to look for changes in community structure associated with the angiogenic tumors, ranked the nodes by their contribution to the total score for each method (see Materials and Methods), and evaluated the core genes in each set for functional enrichment. In order to evaluate the unique contributions of ALPACA, we first applied standard community detection techniques to identify communities in each subtype-specific network, using both the Louvain method and CONDOR, and we looked for GO terms that were statistically enriched in the angiogenic network but not in the non-angiogenic network. Next, we applied edge subtraction, community comparison, and ALPACA to directly identify differential modules associated with angiogenic tumors. Finally, we also computed the differentially expressed genes between the non-angiogenic and angiogenic cancer subtypes. The GO term enrichment with  $P_{adj} < 0.05$  for each method is presented in full in Supplementary Table 1.



**Fig. 4** ALPACA modules associated with angiogenic ovarian tumors. Right-hand side shows five of the modules, with nodes colored by their membership. Edge opacity is proportional to its contribution to the differential modularity. Network is annotated with representative enriched GO terms with  $P_{adj} < 0.05$ , and the genes annotated by the shown GO terms are labeled in larger font. Left-hand side shows the relationship between the ALPACA modules (denoted by M) and the community structure of the angiogenic network (denoted by ANG). Edge thickness depicts the fraction of genes in that differential module that are present in a particular angiogenic network community. The size of each node is proportional to the number of genes in that module or community. Bottom inset: Same networks as above, but colored by community membership in the angiogenic network rather than by membership in the ALPACA modules

Consistent with what we observed in the simulated networks, ALPACA had higher resolution than the other methods and identified 25 modules specific to the angiogenic network. Strikingly, ALPACA was the only network method that identified a gene module enriched in “blood vessel development”, the pathway that we know drives the phenotypic difference between these two ovarian cancer subtypes. Standard community detection methods did not find such a cluster. The non-angiogenic network communities were enriched for histone methylation, embryo development, G-protein-coupled receptor signaling, interferon signaling, and chromatin assembly, whereas the angiogenic communities were enriched for cAMP biosynthetic process, response to fibroblast growth factor, MAPKK activity, and interferon signaling (Supplementary Table 1). The community comparison method did not yield any enriched GO terms. The edge subtraction method resulted in four large modules enriched for general processes like regulation of cell shape, extracellular matrix organization, nucleosome assembly, and immune response (Supplementary Table 1). The differentially expressed genes did exhibit enrichment for “blood vessel development”, but that is to be expected given that the two subtypes of tumors were defined using a gene signature associated with angiogenesis. The

remainder of the differentially expressed genes was not enriched for functional groups that overlapped with those we identified using network-based methods (see Supplementary Information for more details).

ALPACA led to more specific GO term enrichment than the other methods, suggesting that it was able to more carefully refine differential module structure. For example, instead of general GO terms like “immune response”, the ALPACA modules were enriched for particular immune-related pathways like Type I interferon response, interleukin production, regulation of the NF $\kappa$ B pathway, and inflammation. Other enriched pathways included JAK-STAT and growth hormone signaling, urogenital development, triglyceride homeostasis, flavonoid glucuronidation, and cell migration (Fig. 4). Some of these pathways, like JAK-STAT and cell migration, have well-established associations with ovarian tumor progression, while others like flavonoids and triglycerides have only tentative connections with risk of ovarian cancer. Nevertheless, there is substantial support for the biological relevance of these pathways with disease etiology. We provide a detailed discussion of the modules and their biological functions in the Supplementary Information.

Most of the ALPACA GO term results could not be found by running community detection on the angiogenic network alone, which shows that ALPACA partitions nodes in a novel manner that does not merely reflect the underlying community structure of the disease network but instead highlights the changes in modular structure between conditions (Fig. 4, inset). We note that running ALPACA in reverse, to find modules present in the non-angiogenic network as compared to the angiogenic network, results in a substantially smaller set of enriched GO terms, which fall mostly into the metabolic and immune categories, with no enrichment in blood vessel development (Supplementary Table 1). ALPACA therefore selectively identifies network modules associated with the specific phenotype under study. We also computed the correlation in expression among the genes in each ALPACA module. The average absolute value of the Pearson correlation among all gene pairs in each module ranged only from 0.06 to 0.13, suggesting that network module detection does not merely reflect correlation-based clustering (Supplementary Fig. 3). Finally, we ranked the genes by their contribution to the differential modularity and used Gene Set Enrichment Analysis (GSEA) to evaluate enrichment for GO terms across the whole network (see Materials and Methods). The results included some of the biological processes found through module-level enrichment, like “blood vessel remodeling” and “response to Type I interferon”, but not others, likely due to decreased resolution caused by combining genes from different ALPACA modules (Supplementary Table 4).

#### Tumor virus perturbations in primary human cells

DNA viruses hijack the host cell cycle to jumpstart viral genome replication. Tumor viruses can do this so effectively that they lead to aberrant cell proliferation and tumorigenesis, and studying tumor viruses can shed light on the molecular mechanisms behind cancer. Previously, we expressed a panel of 63 proteins from four families of DNA tumor viruses—Epstein-Barr virus, human papillomaviruses, polyomaviruses, and adenovirus—in IMR90 primary human fibroblasts and generated gene expression profiles for each cell line.<sup>37</sup> To construct regulatory networks, we divided the gene expression data into two groups, the first corresponding to the 37 viral proteins classified as “transforming” due to their tumorigenic properties, and the second corresponding to all the control cell lines that contain either empty vectors or GFP. We used PANDA to infer networks by combining gene expression from each sample group with a prior map of cell-type-specific DNase-I-hypersensitive TF-binding sites.<sup>34</sup>

We first ran standard community detection on each network, using the Louvain method for modularity maximization. The control network contained communities enriched in cell migration, axon guidance, and wound response (Supplementary Table 2). The communities in the transforming viral oncogene network were enriched for epithelial–mesenchymal transition, cell migration, axon guidance, and wound response. Since an important function of fibroblasts is to migrate and heal wounds, many of the results from standard community detection appear to be cell-type-specific processes that are not specific to viral oncogenes. The genes with the biggest changes in community assignment were enriched in BMP response and natural killer cell development (Supplementary Table 2). Applying the edge subtraction method using Louvain or CONDOR optimization methods resulted in enrichment for chromatin modification, the Toll-like Receptor pathway, and immune response.

We then applied ALPACA to compare the two networks. Like the edge subtraction method, ALPACA also revealed changes in immune response and chromatin modification but, importantly, it also found significant enrichment for “mitotic cell cycle”, which is the main process we expect to be perturbed by tumor viruses (Fig. 5). Consistent with this, we had previously found that

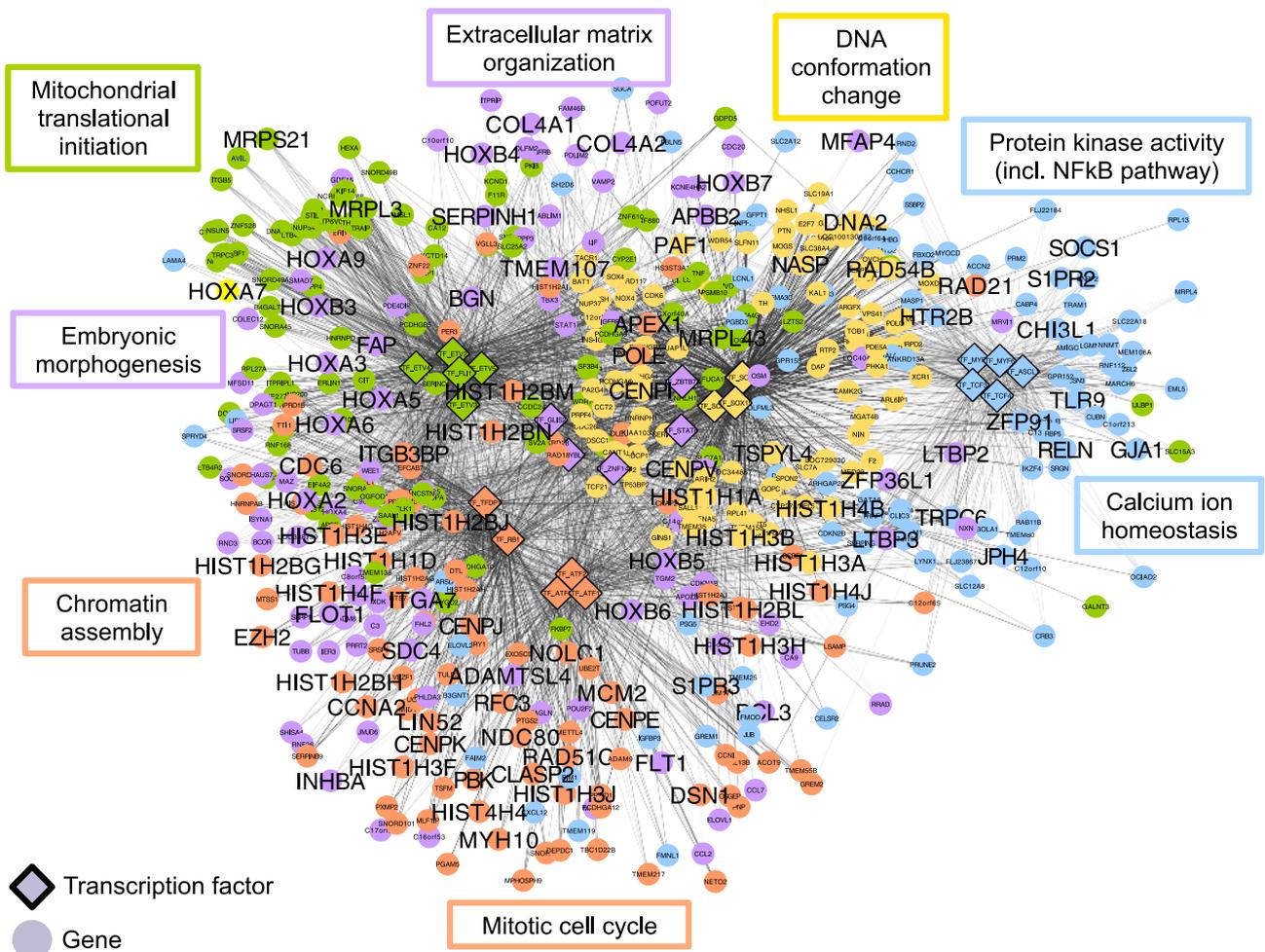
fibroblast cell lines expressing transforming viral oncogenes have significantly altered growth rates.<sup>37</sup> The mitotic cell cycle GO term was not found using any other network-based method, nor was it enriched among the differentially expressed genes (Supplementary Table 2). The TFs regulating this module—RB1, TFDP1, and ATF family members—were not differentially expressed either and were only found using ALPACA. The average absolute value of the Pearson correlation coefficient among genes in the ALPACA modules ranged from 0.22 to 0.34 (Supplementary Fig. 3). The functional significance of the ALPACA modules is further described in the Supplementary Information. As was true with ovarian cancer, there is substantial support for the relevance of these modules in viral transformation.

#### Sexual dimorphism in normal breast tissue

The GTEx consortium has generated gene expression data using tissue collected from 51 body sites and in nearly 600 individuals. Not surprisingly, the tissue with the greatest difference between males and females in autosomal gene expression is the breast.<sup>38</sup> We used PANDA to create tissue-specific regulatory networks to study the effect of sex on regulatory networks in breast tissue.<sup>38</sup> We first applied the Louvain method to detect communities separately in the networks derived from male and female breast tissue and tested for functional enrichment of GO terms in the male and female communities. We found that both the networks were enriched for the same biological processes: GTPase-mediated signal transduction and protein catabolic process (Supplementary Table 3). Therefore, despite what one might expect to be substantially different, the global structure of the male and female networks failed to identify sex-specific patterns of regulation. We also used the edge subtraction method to search for modular differences between the sexes and tested modules for GO term enrichment, but this too failed to identify any significant GO biological processes. In contrast, ALPACA detected several functional modules that differentiate male and female breast tissue (Fig. 6 and Supplementary Information). These modules were enriched in developmental and signaling pathways that are relevant to breast tissue and are often dysregulated in breast cancer. Notably, ALPACA uniquely identified a module associated with female breast tissue that was enriched for “intracellular estrogen receptor signaling pathway”, the hormonal process we expect to be critical for female breast development and overall function. Most of the biological processes identified by ALPACA were not enriched among genes differentially expressed between male and female breast tissue (Supplementary Table 3 and Supplementary Information). The average absolute value of the Pearson correlation coefficient among genes in the ALPACA modules ranged from 0.16 to 0.2 (Supplementary Fig. 3).

## DISCUSSION

Biological networks have complex modular and hierarchical topologies that allow organisms to carry out the functions necessary for survival. Various perturbations, such as diseases, environmental conditions, or mutations, can lead to changes in the phenotype of the organism. Techniques such as differential expression analysis can be used to characterize the transition between different cellular states, but changes in gene expression are ultimately driven by changes in regulatory pathways. If we are to fully understand the basis of complex phenotypes and diseases, we need computational methods that can analyze how regulatory networks change with phenotype. To address this challenge, we developed ALPACA, an algorithm for comparing the topology of two large networks, using a metric we call the “differential modularity”, to find groups of nodes that characterize differences in network structure. ALPACA differs from other community



**Fig. 5** ALPACA modules associated with transforming viral oncogenes. Network shows five modules, with nodes colored by membership in differential modules. Edge opacity is proportional to its contribution to the differential modularity. Network is annotated with representative enriched GO terms with  $P_{\text{adj}} < 0.05$ . Genes annotated by the shown GO terms are labeled in large font

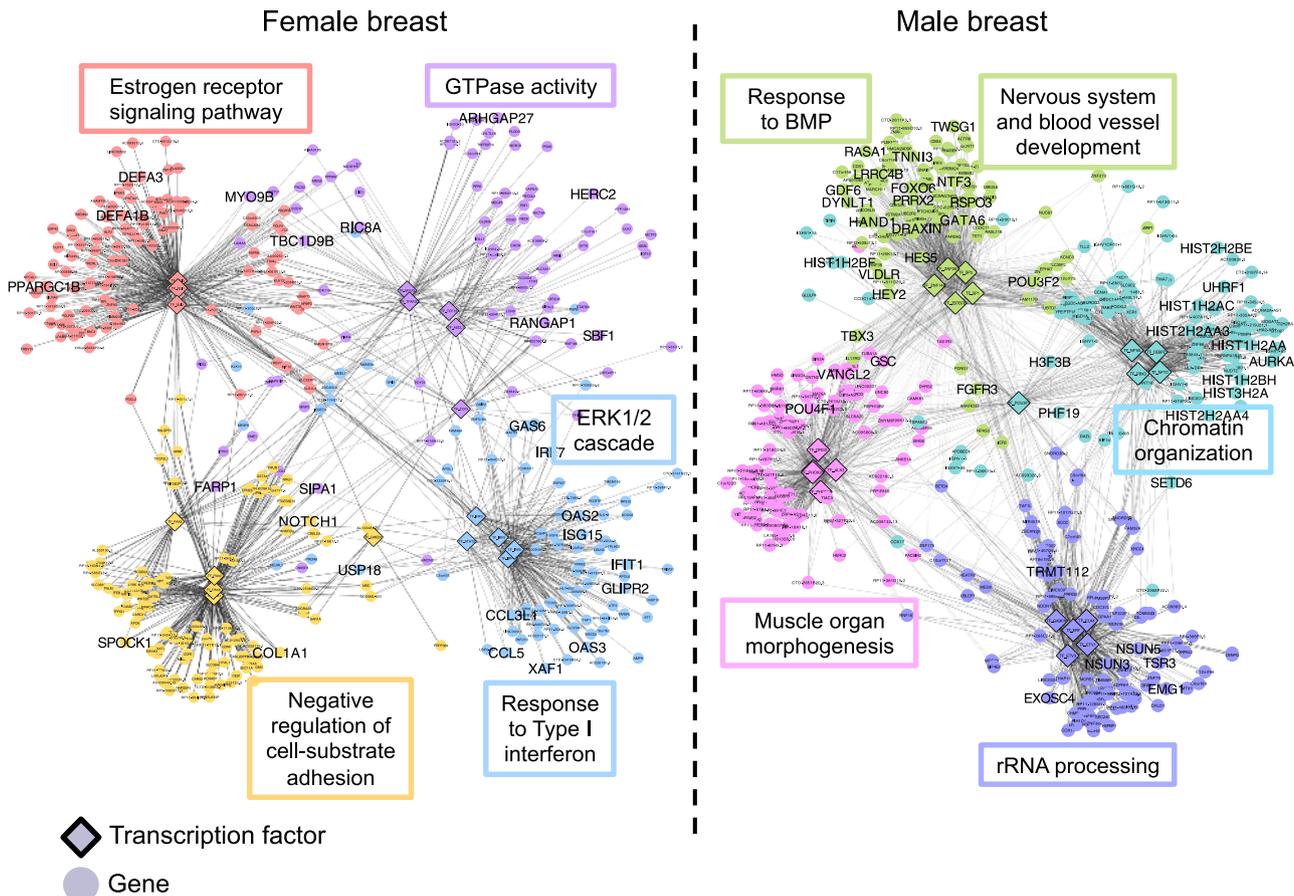
detection methods in that it compares the structure of networks to each other rather than to a random background network and is thus better able to detect subtle differences in network modular structure. This can potentially allow detection of small modules that function together in particular conditions or in disease.

We evaluated the performance of ALPACA on simulated networks and compared it to two available approaches for detecting changes in network modular structure: (i) “community comparison”, where one applies standard community detection to the baseline and perturbed networks separately and contrasts the resulting communities, and (ii) “edge subtraction”, which involves subtracting the two networks edge by edge, and clustering the resulting differential network. ALPACA was able to resolve smaller differential modules than the community comparison method. Intuitively, this is because modularity maximization in its standard form penalizes the splitting of a large dense community into smaller ones, whereas the differential modularity score used in ALPACA penalizes the formation of large communities similar to those present in the baseline network. In addition, ALPACA was more robust to noise in individual network edges and better at detecting small modules than the edge subtraction method. In the edge subtraction method, the uncertainty of the edges in the “differential” network is the sum of the uncertainties in the corresponding edges of the original networks. In contrast, ALPACA aggregates the signals coming from multiple edges in the

baseline network communities to derive a null model for edge density, so it is less sensitive to the uncertainty in individual edges.

ALPACA’s differential modularity metric directly compares the edges that one sees within a community to what you would expect based on the topology of a corresponding reference network. This adapts the well-established modularity maximization method to infer subtle changes in the community structure that arise when comparing distinct complex phenotypes. Unlike other methods that simply subtract networks, ALPACA preserves those secondary interactions that exist in both networks but allows them to shift their functional context as the edges around them change, which can capture new modular structures. The differential modularity also incorporates increased and decreased edge weights across the entire network into a single, simple framework for module detection. And unlike community comparison, ALPACA can detect new modules that form within the context of globally active regulatory programs that are present in both the baseline and perturbed networks.

We applied ALPACA to compare transcriptional networks that were inferred from a combination of gene expression and TF-binding data using the PANDA network inference algorithm. PANDA does not explicitly use the expression correlation between regulators and the target genes, and can therefore model TFs that are not changing in mRNA expression but whose activity is controlled through other mechanisms, like post-translational



**Fig. 6** Sexually dimorphic ALPACA modules in human breast tissue. Networks show four modules specific to either female (left-hand side) or male (right-hand side) breast tissue. Nodes are colored by membership in differential modules. Edge opacity is proportional to its contribution to the differential modularity. Networks are annotated with representative enriched GO terms with  $P_{adj} < 0.05$ . Genes annotated by the shown GO terms are labeled in large font

modification. PANDA also incorporates changes in promoter activity that could alter regulatory targeting patterns. Comparing angiogenic to non-angiogenic subtypes of ovarian cancer, we found functional modules that were enriched in expected disease pathways like blood vessel development, interleukin production, and JAK-STAT signaling. We also found enrichment for less expected processes like flavonoid biosynthesis and triglyceride homeostasis, which have been speculated to be relevant for ovarian cancer, but for which the underlying molecular pathways are not known.<sup>39–42</sup> All these modules were specific to the angiogenic subtype and uniquely revealed by ALPACA; they could not be found through standard community detection in the individual angiogenic and non-angiogenic networks or in an edge-subtracted network, or by running ALPACA in reverse on the non-angiogenic network compared to the angiogenic network.

In another application of the method, we compared normal male and female breast tissue to find sex-specific patterns of regulation. Many of the modules we found were enriched in known processes related to breast development and breast cancer, like ERK and Rho GTPase signaling. Perhaps most strikingly, the female breast network contained a differential module enriched for estrogen receptor signaling, which is one of the main sex-specific pathways known to be active in breast tissue. Once again, these results could not be found using other community detection and network comparison methods.

ALPACA builds on our growing understanding of how networks define phenotype. Differential expression is driven by changes in the activity and structure of gene regulatory networks. But adding

or subtracting edges does more than change individual regulatory interactions. With enough individual changes occurring in the right places in the starting network, changes in edges can lead to the creation or destruction of functional communities of genes and their regulators. While the global structure of the network may be largely unchanged, these new functional communities provide insight into coherent processes that differentiate one phenotype from another.

Consistent with this, we found that standard differential expression analysis was unable to detect enrichment in many of the biological functions found using ALPACA. This is because network-level analysis, and ALPACA in particular, helps organize both strongly and weakly differentially expressed genes into new modules that are under common regulatory control, identifying signaling pathways that could not have been distinguished if genes were ranked purely by differential expression. Moreover, ALPACA can also find the TFs that likely regulate these pathways but are not themselves differentially expressed. In addition, we found that the average magnitude of the Pearson correlation coefficient among genes in an ALPACA module ranged from 0.06 to 0.34, suggesting that ALPACA provides insights that go beyond correlation-based clustering of gene expression profiles.

ALPACA requires a minimum input of two graphs and could be applied to many types of biological networks, including metabolic, PPI, and expression Quantitative Trait Loci (eQTL) networks, all of which exhibit highly functional modular structures.<sup>9,43,44</sup> For example, we could imagine applying ALPACA to compare community structure in PPI networks with mutation-driven

“edgetic” perturbations, in order to discover functional changes in protein complexes and signaling associated with disease.<sup>45</sup> ALPACA could also be applied to compare eQTL networks in patient cohorts with differing pathologies to prioritize sets of SNPs and genes that influence complex traits. Any of these ALPACA modules could be interrogated for disease association using data from genome-wide association studies.<sup>46</sup>

As is true of all methods, some of the general difficulties in community detection will affect the performance of ALPACA. Many real-world networks, including those found in biology, do not have one clearly superior community partition, but instead exhibit a complex landscape of near-optimal community structures.<sup>47</sup> A number of studies have found that ensemble methods can help identify robust communities in this landscape.<sup>48,49</sup> In ALPACA, we use CONDOR to detect communities in the baseline network, and this is a deterministic method that finds one of the high-scoring partitions of the network. We then optimize the differential modularity using the Louvain method, which is stochastic and can access many near-optimal solutions. The performance of ALPACA could potentially be improved by employing a consensus approach in both these steps.

Another challenge in community detection is the resolution limit of various community detection methods.<sup>26</sup> ALPACA is based on modularity maximization and therefore still has a resolution limit, even though it is smaller than usual due to the fact that the null model is defined at the community level rather than at the global network level. But in principle, phenotypes and diseases may be driven by modules of any size. To identify all such modules, we could adapt ALPACA to operate at a variety of granularity levels by incorporating resistance parameters that modify the modularity function.<sup>50,51</sup>

As more genome-wide studies of molecular interactions and multi-omics data are generated, better statistical models for network analysis will be critical to making differential network biology a robust and reproducible platform for studying complex diseases.<sup>14</sup> ALPACA establishes a rigorous framework for comparing complex networks and identifying changes in modular structure, and is an important step forward in creating methodological platforms for predictive analysis of biological networks.

## METHODS

### ALPACA algorithm

ALPACA comprises the following two steps:

**Step 1:** The input network consists of edges between regulators and target genes. We first label the nodes that act as regulators and targets separately. In particular, a gene that encodes a TF becomes two separate nodes depending on whether we are modeling its mRNA expression level (target node) or protein activity (regulator node). For the weight of each edge, we use the final z-score output by the PANDA network inference algorithm. We then take the edges that have positive weight in the baseline condition, and run bipartite weighted network community detection using either CONDOR or the Louvain method.

**Step 2:** Compute  $D_{ij}$  for the perturbed network, using the definition in the main text and the baseline communities found in Step 1. It is possible that the numerator and denominator of  $N_{ij}$  are both zero, meaning that there were no edges between the communities  $C_i$  and  $C_j$ . This can happen if, for example, at least one of the nodes  $i$  or  $j$  were not connected to the baseline network to begin with. In this case, we define  $N_{ij}$  to be zero, since the “expected” number of edges between the two nodes is zero. We next apply a generalized Louvain procedure to assign nodes into communities based on  $D_{ij}$ .<sup>19</sup> Briefly, the Louvain method works as follows: (i) Start with every node in its own community, (ii) go through each node iteratively, and merge it with the node that produces the biggest increase in differential modularity, (iii) after reaching a local optimum, treat each of the resulting groups as “metanodes” in a new “metanetwork” and recalculate an effective adjacency matrix, and (iv) repeat steps (ii) and (iii) until convergence. For the purpose of reporting reproducible results, we iterate through the nodes in the same pre-determined order every time, and we break ties by selecting the first member of the set.

In an optional third step, we can evaluate the core genes in each module for enrichment in known biological pathways.

**Step 3:** The core genes are those that are most important to the integrity of the module and therefore potentially the most robust and essential members. To define the core genes, we score each node according to its contribution to the differential modularity of the module that it belongs to:

$$S_i = \frac{1}{m_p} \sum_j D_{ij} \delta(M_i, M_j). \quad (4)$$

We ranked the target genes in each module by their scores  $S_i$ . Since the size of typical modules found in ALPACA ranged from about 50 to 200 genes, we chose to use the top 50 core genes from each module to evaluate functional enrichment in an equitable manner across all the modules. We also repeated each analysis using the top 100 core genes in order to test the dependence of the enrichment on the cutoff. GO term enrichment was calculated using the GOSTATS package in R, with the following parameters: the gene universe is defined to be the set of all possible target genes in the initial networks, and the  $p$ -value calculation is conditioned on the GO hierarchy structure. In each module, the  $p$ -values were adjusted for multiple testing using the Benjamini–Hochberg method.

To run GSEA, we ranked all genes in the network by (i) their raw score  $S_i$  or by (ii) a version of the score  $\tilde{S}_i$ , normalized for the relative sizes of each module,<sup>9</sup> so that genes in smaller modules are not at a disadvantage simply because they have fewer potential connections to nodes within the module:

$$\tilde{S}_i = \frac{\sum_j D_{ij} \delta(M_i, M_j)}{\sum_{i,j} D_{ij} \delta(M_i, M_j)}. \quad (5)$$

GSEA pre-ranked was run against GO biological process gene sets using the desktop Java application (<http://software.broadinstitute.org/gsea>) with default parameters. An FDR threshold of 0.25 was used to identify significantly enriched GO terms.

### Edge subtraction method

For each edge, the edge weight of the baseline network was subtracted from the edge weight in the perturbed network to compute  $\Delta w_{ij}$ , and only edges with  $\Delta w_{ij} > 0$  were retained. We then used the  $\Delta w_{ij}$  values as new edge weights to perform community detection using CONDOR or Louvain optimization.<sup>9,31</sup>

### Community comparison method

We first used either CONDOR or Louvain method to find the community structure of the baseline and perturbed networks, in each case keeping only edges that had positive z-scores. We next aimed to efficiently map the two community structures to each other. To find the best approximation of a linear mapping, we computed  $R$  in the equation  $B = AR$ , where  $A$  is the  $N \times q$  matrix of node membership for the baseline community structure, and  $B$  is the corresponding matrix for the perturbed community structure (here  $N$  is the number of genes and  $q$  is the number of communities). To invert the matrix  $A$ , we used singular value decomposition to compute the pseudoinverse  $A^p = VD^{-1}U^T$ , where  $A = UDV^T$  and then computed  $R = A^p B$ . The entries of the  $q \times q$  matrix  $R$  represent an approximate linear transformation that maps the communities in the baseline network to the communities of the perturbed network. Finally, we scored each node according to how much its community membership remains the same between baseline and perturbed conditions, using the formula  $S_i^{(R)} = \sum_j A_{ij} R B_{ij}$ . Nodes were ranked from low to high values of  $S_i^{(R)}$  for further analysis. Low-scoring nodes represent the nodes that participate in altered community structure in the perturbed network.

### Creating simulated networks and evaluating differential community methods

To simulate “addition” networks, we started with the GFP-control network from the tumor virus dataset (see section on “Data preprocessing, differential expression, and network inference” for details on how this network was constructed) and thresholded the edges at a z-score of 2.7 (for noiseless simulation) or 2.9 (for noisy simulation). The threshold was chosen such that the resulting edges would form an unweighted network with a similar community structure as the full weighted network. We found that applying CONDOR to the GFP-control network at a threshold of 2.7 resulted in five communities containing 1336, 833, 781, 1018, and 44 nodes each. To add a module, we randomly chose a subset of these nodes

and added all possible new edges between them. To add noise in the second set of “addition” networks, we also resampled edges as follows: (i) start with an empty network with the same nodes as the GFP-control network, (ii) count the number of edges between TFs in community  $C_i$  and target genes in community  $C_j$ , for each pair  $i$  and  $j$ , in the GFP-control network, and (iii) add a matching number of edges randomly between the TFs in community  $C_i$  and target genes in community  $C_j$  in the new network.

We evaluated the results of each method on the simulated networks by comparing the ranks of true positives (the target genes in the added module) against a background consisting of target genes not in the added module. We used one-sided Kolmogorov–Smirnov and Wilcoxon tests to look for significant differences in the distribution of the ranks. Both tests gave similar results, and in the figures we present the one-sided Wilcoxon  $p$ -values.

To create the “subtracted” simulation with two node groups, we started with a fully connected network containing 100 nodes, with all edge weights set to a default value of 0.1. We then defined two node groups, A and B, each containing 10 TFs and 40 genes. Edges within each of these groups were set to edge weight 1.0. Next, to create the baseline network we set the weights of all edges between groups A and B to be 0.8. To create the perturbed network we set the weights of all edges between groups A and B to be 0.2. To create the three-group “subtracted” network, we first created a fully connected network containing 125 nodes, with all edge weights set to a default value of 0.1. We then defined three node groups A, B, and C containing 50, 25, and 50 nodes, respectively (of which 10, 5, and 10 were TFs). Edges within each group were set to weight 1.0, and all edges between groups B and C were set to weight 0.2. For the baseline network, the edges between groups A and B were set to weight 0.8 and for the perturbed network, the edges between groups A and B were set to weight 0.2.

## Data preprocessing, differential expression, and network inference

Preprocessing and network inference for ovarian cancer data was carried out as previously described.<sup>36</sup> Briefly, we ran the network inference algorithm PANDA (Passing Attributes between Networks for Data Assimilation) to integrate gene expression data with TF-binding sites to create regulatory networks for each subtype.<sup>34</sup> The prior network of binding sites for 111 TFs were defined as the occurrence of the corresponding motif in the promoter, defined as [−750, +250] base pairs around the transcription start site (TSS).

The viral oncogene gene expression data were normalized and batch-corrected, and a map of high-probability TF-binding sites was created by combining cell-type-specific DNase-I hypersensitivity data with motif occurrence in the promoters defined as [−25 kb, 25 kb] around each TSS, as previously described.<sup>37</sup> The binding sites and gene expression were combined to infer networks using PANDA with default parameters, as previously described.<sup>1</sup>

Sex-specific and tissue-specific transcriptional networks for the GTEx data were constructed as previously described.<sup>38,52</sup>

Differential expression analysis was carried out using the R package *limma*, and  $p$ -values were adjusted for multiple testing using the Benjamini–Hochberg method.<sup>53</sup>

## Code availability

ALPACA is implemented in R and is freely available for download through Github at <https://github.com/meghapadi/ALPACA>.

## Data availability

Ovarian cancer gene expression data are available from The Cancer Genome Atlas (TCGA) at <https://gdc.cancer.gov>. Tumor virus gene expression data are available from the Gene Expression Omnibus (GEO), accession number: GSE38467. Breast tissue data from the Genotype-Tissue Expression (GTEx) project can be found at <https://sites.google.com/a/channing.harvard.edu/kimberlyglass/tools/gtex-networks>.

## ACKNOWLEDGEMENTS

This work was supported by NIH grants K25 HG006031 (M.P.) and R01 HL111759 and R35 CA197449 (J.Q.).

## AUTHOR CONTRIBUTIONS

M.P. conceived of the project, performed analysis, and wrote the paper. J.Q. helped refine the analysis and wrote the paper.

## ADDITIONAL INFORMATION

**Supplementary information** accompanies the paper on the *npj Systems Biology and Applications* website (<https://doi.org/10.1038/s41540-018-0052-5>).

**Competing interests:** The authors declare no competing financial interests.

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

## REFERENCES

- Padi, M. & Quackenbush, J. Integrating transcriptional and protein interaction networks to prioritize condition-specific master regulators. *BMC Syst. Biol.* **9**, 80 (2015).
- Giaever, G. et al. Functional profiling of the *Saccharomyces cerevisiae* genome. *Nature* **418**, 387–391 (2002).
- Locke, A. E. et al. Genetic studies of body mass index yield new insights for obesity biology. *Nature* **518**, 197–206 (2015).
- Wood, A. R. et al. Defining the role of common variation in the genomic and biological architecture of adult human height. *Nat. Genet.* **46**, 1173–1186 (2014).
- Fuchsberger, C. et al. The genetic architecture of type 2 diabetes. *Nature* **536**, 41–47 (2016).
- Jeong, H., Mason, S. P., Barabasi, A. L. & Oltvai, Z. N. Lethality and centrality in protein networks. *Nature* **411**, 41–42 (2001).
- Hartwell, L. H., Hopfield, J. J., Leibler, S. & Murray, A. W. From molecular to modular cell biology. *Nature* **402**, C47–C52 (1999).
- Menche, J. et al. Disease networks. Uncovering disease–disease relationships through the incomplete interactome. *Science* **347**, 1257601 (2015).
- Platig, J., Castaldi, P. J., DeMeo, D. & Quackenbush, J. Bipartite community structure of eQTLs. *PLoS Comput. Biol.* **12**, e1005033 (2016).
- ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- Marbach, D. et al. Wisdom of crowds for robust gene network inference. *Nat. Methods* **9**, 796–804 (2012).
- Marbach, D. et al. Tissue-specific regulatory circuits reveal variable modular perturbations across complex diseases. *Nat. Methods* **13**, 366–370 (2016).
- Newman, M. E. & Girvan, M. Finding and evaluating community structure in networks. *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.* **69**, 026113 (2004).
- Ideker, T. & Krogan, N. J. Differential network biology. *Mol. Syst. Biol.* **8**, 565 (2012).
- Gambardella, G. et al. Differential network analysis for the identification of condition-specific pathway activity and regulation. *Bioinformatics* **29**, 1776–1785 (2013).
- Watson, M. CoXpress: differential co-expression in gene expression data. *BMC Bioinformatics* **7**, 509 (2006).
- Van Landeghem, S., Van Parys, T., Dubois, M., Inze, D. & Van de Peer, Y. Diffany: an ontology-driven framework to infer, visualise and analyse differential molecular networks. *BMC Bioinformatics* **17**, 18 (2016).
- Gill, R., Datta, S. & Datta, S. A statistical framework for differential network analysis from microarray data. *BMC Bioinformatics* **11**, 95 (2010).
- Danon, L., Diaz-Guilera, A., Duch, J. & Arenas, A. Comparing community structure identification. *J. Stat. Mech. Theory Exp.* **9**, P09008 (2005).
- Perotti, J. I., Tessone, C. J. & Caldarelli, G. Hierarchical mutual information for the comparison of hierarchical community structures in complex networks. *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.* **92**, 062825 (2015).
- Tesson, B. M., Breitling, R. & Jansen, R. C. DiffCoEx: a simple and sensitive method to find differentially coexpressed gene modules. *BMC Bioinformatics* **11**, 497 (2010).
- Amar, D., Safer, H. & Shamir, R. Dissection of regulatory networks that are altered in disease via differential co-expression. *PLoS Comput. Biol.* **9**, e1002955 (2013).
- Valcarcel, B. et al. Genome metabolome integrated network analysis to uncover connections between genetic variants and complex traits: an application to obesity. *J. R. Soc. Interface* **11**, 20130908 (2014).
- Mitra, K., Carvunis, A. R., Ramesh, S. K. & Ideker, T. Integrative approaches for finding modular structure in biological networks. *Nat. Rev. Genet.* **14**, 719–732 (2013).
- Ideker, T., Ozier, O., Schwikowski, B. & Siegel, A. F. Discovering regulatory and signalling circuits in molecular interaction networks. *Bioinformatics* **18**, S233–S240 (2002).

26. Fortunato, S. & Barthelemy, M. Resolution limit in community detection. *Proc. Natl Acad. Sci. USA* **104**, 36–41 (2007).
27. Gerstein, M. B. et al. Architecture of the human regulatory network derived from ENCODE data. *Nature* **489**, 91–100 (2012).
28. Mucha, P. J., Richardson, T., Macon, K., Porter, M. A. & Onnela, J. P. Community structure in time-dependent, multiscale, and multiplex networks. *Science* **328**, 876–878 (2010).
29. Raghavan, U. N., Albert, R. & Kumara, S. Near linear time algorithm to detect community structures in large-scale networks. *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.* **76**, 036106 (2007).
30. Rosvall, M. & Bergstrom, C. T. Maps of random walks on complex networks reveal community structure. *Proc. Natl Acad. Sci. USA* **105**, 1118–1123 (2008).
31. Blondel, V., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. *J. Stat. Mech.*, P10008 (2008).
32. Esmailian, P. & Jalili, M. Community detection in signed networks: the role of negative ties in different scales. *Sci. Rep.* **5**, 14339 (2015).
33. Traag, V. A. & Bruggeman, J. Community detection in networks with positive and negative links. *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.* **80**, 036115 (2009).
34. Glass, K., Huttenhower, C., Quackenbush, J. & Yuan, G. C. Passing messages between biological networks to refine predicted interactions. *PLoS ONE* **8**, e64832 (2013).
35. Bentink, S. et al. Angiogenic mRNA and microRNA gene expression signature predicts a novel subtype of serous ovarian cancer. *PLoS ONE* **7**, e30269 (2012).
36. Glass, K., Quackenbush, J., Spentzos, D., Haibe-Kains, B. & Yuan, G. C. A network model for angiogenesis in ovarian cancer. *BMC Bioinformatics* **16**, 115 (2015).
37. Rozenblatt-Rosen, O. et al. Interpreting cancer genomes using systematic host network perturbations by tumour virus proteins. *Nature* **487**, 491–495 (2012).
38. Chen, C.-Y. et al. Sexual dimorphism in gene expression and regulatory networks across human tissues. *bioRxiv* <https://doi.org/10.1101/082289> (2016).
39. Cassidy, A., Huang, T., Rice, M. S., Rimm, E. B. & Tworoger, S. S. Intake of dietary flavonoids and risk of epithelial ovarian cancer. *Am. J. Clin. Nutr.* **100**, 1344–1351 (2014).
40. Gates, M. A. et al. Flavonoid intake and ovarian cancer risk in a population-based case-control study. *Int. J. Cancer* **124**, 1918–1925 (2009).
41. Hua, X. et al. Association among dietary flavonoids, flavonoid subclasses and ovarian cancer risk: a meta-analysis. *PLoS ONE* **11**, e0151134 (2016).
42. Tania, M., Khan, M. A. & Song, Y. Association of lipid metabolism with ovarian cancer. *Curr. Oncol.* **17**, 6–11 (2010).
43. Spirin, V. & Mirny, L. A. Protein complexes and functional modules in molecular networks. *Proc. Natl Acad. Sci. USA* **100**, 12123–12128 (2003).
44. Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N. & Barabasi, A. L. Hierarchical organization of modularity in metabolic networks. *Science* **297**, 1551–1555 (2002).
45. Sahni, N. et al. Widespread macromolecular interaction perturbations in human genetic disorders. *Cell* **161**, 647–660 (2015).
46. Lamparter, D., Marbach, D., Rueedi, R., Kutalik, Z. & Bergmann, S. Fast and rigorous computation of gene and pathway scores from SNP-based summary statistics. *PLoS Comput. Biol.* **12**, e1004714 (2016).
47. Good, B. H., de Montjoye, Y. A. & Clauset, A. Performance of modularity maximization in practical contexts. *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.* **81**, 046106 (2010).
48. Lancichinetti, A. & Fortunato, S. Consensus clustering in complex networks. *Sci. Rep.* **2**, 336 (2012).
49. Zhang, P. & Moore, C. Scalable detection of statistically significant communities and hierarchies, using message passing for modularity. *Proc. Natl Acad. Sci. USA* **111**, 18144–18149 (2014).
50. Arenas, A., Fernandez, A. & Gomez, S. Analysis of the structure of complex networks at different resolution levels. *N. J. Phys.* **10**, 053039 (2008).
51. Reichardt, J. & Bornholdt, S. Statistical mechanics of community detection. *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.* **74**, 016110 (2006).
52. Sonawane, A. R. et al. Understanding tissue-specific gene regulation. *Cell Rep.* **21**, 1077–1088 (2017).
53. Smyth, G. K. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. *Stat. Appl. Genet. Mol. Biol.* **3**, Article 3 (2004).



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

© The Author(s) 2018