

ARTICLE

DOI: 10.1038/s41467-018-04948-5

OPEN

# Prioritizing network communities

Marinka Zitnik<sup>1</sup>, Rok Sosič<sup>1</sup> & Jure Leskovec<sup>1,2</sup>

Uncovering modular structure in networks is fundamental for systems in biology, physics, and engineering. Community detection identifies candidate modules as hypotheses, which then need to be validated through experiments, such as mutagenesis in a biological laboratory. Only a few communities can typically be validated, and it is thus important to prioritize which communities to select for downstream experimentation. Here we develop CRANK, a mathematically principled approach for prioritizing network communities. CRANK efficiently evaluates robustness and magnitude of structural features of each community and then combines these features into the community prioritization. CRANK can be used with any community detection method. It needs only information provided by the network structure and does not require any additional metadata or labels. However, when available, CRANK can incorporate domain-specific information to further boost performance. Experiments on many large networks show that CRANK effectively prioritizes communities, yielding a nearly 50-fold improvement in community prioritization.

<sup>1</sup>Computer Science Department, Stanford University, 353 Serra Mall, Stanford, CA 94305, USA. <sup>2</sup>Chan Zuckerberg Biohub, 499 Illinois St., San Francisco, CA 94158, USA. Correspondence and requests for materials should be addressed to J.L. (email: [jure@cs.stanford.edu](mailto:jure@cs.stanford.edu))

Networks exhibit modular structure<sup>1</sup> and uncovering it is fundamental for advancing the understanding of complex systems across sciences<sup>2,3</sup>. Methods for community detection<sup>4</sup>, also called node clustering or graph partitioning, allow for computational detection of modular structure by identifying a division of network's nodes into groups, also called communities<sup>5–10</sup>. Such communities provide predictions/hypotheses about potential modules of the network, which then need to be experimentally validated and confirmed. However, in large networks, community detection methods typically identify many thousands of communities<sup>6,7</sup> and only a small fraction can be rigorously tested and validated by follow-up experiments. For example, gene communities detected in a gene interaction network<sup>11</sup> provide predictions/hypotheses about disease pathways<sup>2,3</sup>, but to confirm these predictions scientists have to test every detected community by performing experiments in a wet laboratory<sup>3,8</sup>. Because experimental validation of detected communities is resource-intensive and generally only a small number of communities can be investigated, one must prioritize the communities in order to choose which ones to investigate experimentally.

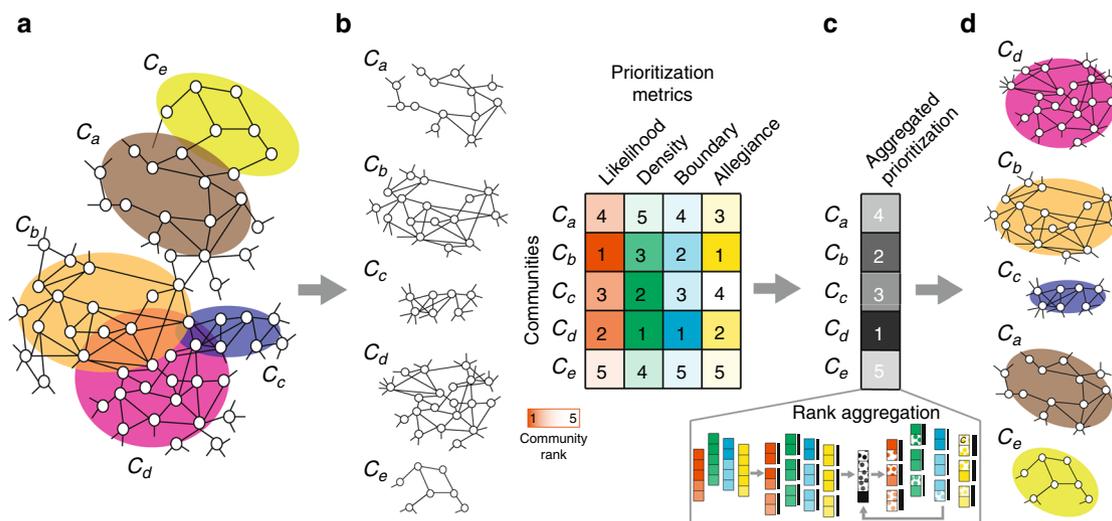
In the context of biological networks, several methods for community or cluster analysis have been developed<sup>2,3,12–15</sup>. However, these methods crucially rely and depend on knowledge in external databases, such as Gene Ontology (GO) annotations<sup>16</sup>, protein domain databases, gene expression data, patient clinical profiles, and sequence information, in order to calculate the quality of communities derived from networks. Furthermore, they require this information to be available for all communities. This means that if genes in a given community are not present in a gene knowledge database then it is not possible for existing methods to even consider that community. This issue is exacerbated because knowledge databases are incomplete and biased toward better-studied genes<sup>11</sup>. Furthermore, these methods do

not apply in domains at the frontier of science where domain-specific knowledge is scarce or non-existent, such as in the case of cell-cell similarity networks<sup>17</sup>, microbiome networks<sup>18,19</sup>, and chemical interaction networks<sup>20</sup>. Thus, there is a need for a general solution to prioritize communities based on network information only.

Here, we present CRANK, a general approach that takes a network and detected communities as its input and produces a ranked list of communities, where high-ranking communities represent promising candidates for downstream experiments. CRANK can be applied in conjunction with any community detection method (Supplementary Notes 2 and 5) and needs only the network structure, requiring no domain-specific meta or label information about the network. However, when domain-specific supervised information is available, CRANK can integrate this extra information to boost performance (Supplementary Notes 9 and 10). CRANK can thus prioritize communities that are well characterized in knowledge bases, such as GO annotations, as well as poorly characterized communities with limited or no annotations. Furthermore, CRANK is based on rigorous statistical methods to provide an overall rank for each detected community.

## Results

**Overview of CRANK prioritization approach.** CRANK community prioritization approach consists of the following steps (Fig. 1). First, CRANK finds communities using an existing, preferred community detection method (Fig. 1a). It then computes for each community four CRANK defined community prioritization metrics, which capture key structural features of the community (Fig. 1b), and then it combines the community metrics via an aggregation method into a single overall score for each community (Fig. 1c). Finally, CRANK prioritizes communities by ranking them by their decreasing overall score (Fig. 1d).



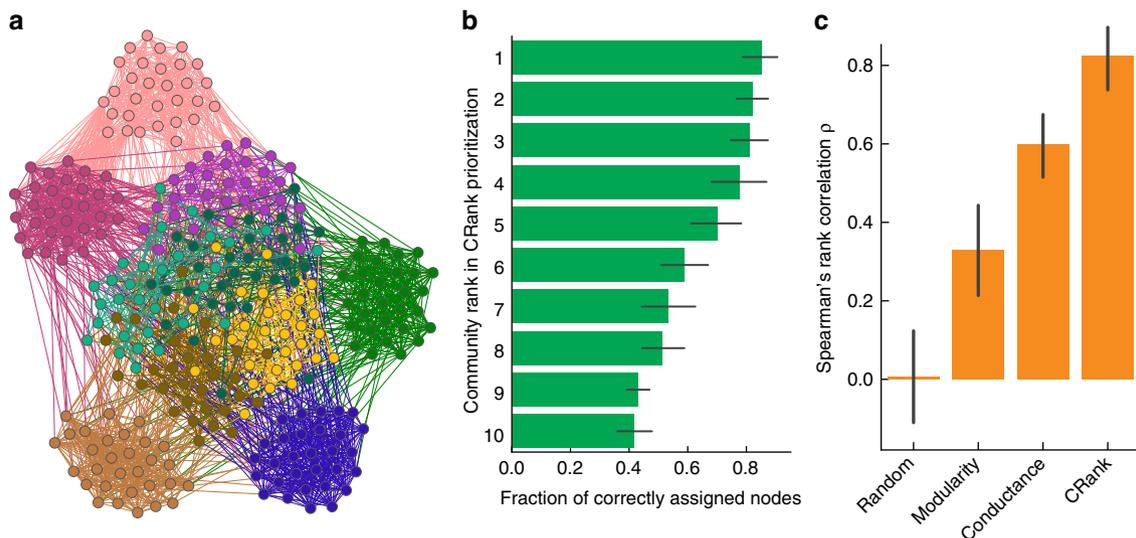
**Fig. 1** Prioritizing network communities. **a** Community detection methods take as input a network and output a grouping of nodes into communities. Highlighted are five communities, ( $C_a, \dots, C_e$ ), that are detected in the illustrative network. **b** After communities are detected, the goal of community prioritization is to identify communities that are most promising targets for follow-up investigations. Promising targets are communities that are most associated with external network functions, such as cellular functions in protein-protein interaction networks, or cell types in cell-cell similarity networks. CRANK is a community prioritization approach that ranks the detected communities using only information captured by the network structure and does not require any external data about the nodes or edges of the network. However, when external information about communities is available, CRANK can make advantage of it to further improve performance (Supplementary Notes 9 and 10). CRANK starts by evaluating four different structural features of each community: the overall likelihood of the edges in the community (likelihood), internal connectivity (density), external connectivity (boundary), and relationship with the rest of the network (allegiance). CRANK can also integrate any number of additional user-defined metrics into the prioritization without any further changes to the method. **c** CRANK then applies a rank aggregation method to combine the metrics and **d** produce the final ranking of communities

CRANK uses four different metrics to characterize network connectivity features for each detected community (Methods section). These metrics evaluate the magnitude of structural features as well as their robustness against noise in the network structure. The rationale here is that high-priority communities have high values of metrics and are also stable with respect to network perturbations. If a small change in the network structure—an edge added here, another deleted there—significantly changes the value of a prioritization metric then the community will not be considered high priority. We derive analytical expressions for calculating these metrics, which make CRANK computationally efficient and applicable to large networks (Supplementary Note 2). Because individual metrics may have different importance in different networks, a key element of CRANK is a rank aggregation method. This method combines the values of the four metrics into a single score for each community, which then determines the community's rank (see Methods and Supplementary Note 4). CRANK's aggregation method adjusts the impact of each metric on the ranking in a principled manner across different networks and also across different communities within a network, leading to robust rankings and a high-quality prioritization of communities (Supplementary Note 5).

**Synthetic networks.** We first demonstrate CRANK by applying it to synthetic networks with planted community structure (Fig. 2a). The goal of community prioritization is to identify communities that are most promising candidates for follow-up investigations.

Since communities provide predictions about the modular structure of the network, promising candidates are communities that best correspond to the underlying modules. Thus, in this synthetic example, the aim of community prioritization can be seen as to rank communities based on how well they represent the underlying planted communities, while only utilizing information about network structure and without any additional information about the planted communities. We quantify prioritization quality by measuring the agreement between a ranked list of communities produced by CRANK and the gold standard ranking (Supplementary Note 7). In the gold standard ranking, communities are ordered in the decreasing order of how accurately each community reconstructs its corresponding planted community.

We experiment with random synthetic networks with planted community structure (Fig. 2a), where we use a generic community detection method<sup>7</sup> to identify communities and then prioritize them using CRANK. We observe that CRANK produces correct prioritization—using only the unlabeled network structure, CRANK places communities that better correspond to planted communities towards the top of the ranking (Fig. 2b), which indicates that CRANK can identify accurately detected communities by using the network structure alone and having no other data about planted community structure. Comparing the performance of CRANK to alternative ranking techniques, such as modularity<sup>5</sup> and conductance<sup>21</sup> (Supplementary Note 7), we observe that CRANK performs 149 and 37% better than modularity and conductance, respectively, in terms of the Spearman's rank correlation between the generated ranking and



**Fig. 2** Synthetic networks with planted community structure. **a–c** In networks with known modular structure we can evaluate community prioritization by quantifying the correspondence between detected communities and the planted communities. **a** Benchmark networks on  $N = 300$  nodes are created using a stochastic block model with 10 planted communities<sup>10</sup>. Each planted community has 30 nodes, which are colored by their planted community assignment. Planted communities use different values for within-community edge probability  $p_{in}$ , five use  $p_{in} = 0.6$  and five use  $p_{in} = 0.2$ . As a result, planted communities with smaller within-community probability  $p_{in}$  are harder to detect. For each benchmark network we apply a community detection method<sup>6</sup> to detect communities and then use CRANK to prioritize them. CRANK produces a ranked list of detected communities. The gold standard rank of each community is determined by how accurately it corresponds to its planted counterpart. **b** Each bar represents one detected community and the bars are ordered by CRANK's ranking with the highest-ranked community located at the top and the lowest ranked community located at the bottom. As a form of validation, the width of each bar corresponds to the fraction of nodes in a community that are correctly classified into a corresponding planted community, with error bars showing the 95% confidence intervals over 500 benchmark networks. A perfect prioritization ranks the bars by decreasing width. Notice that CRANK perfectly prioritizes the communities even though it only uses information about the network structure, and has no access to information about the planted communities. **c** Prioritization performance is measured using Spearman's rank correlation  $\rho$  between the generated ranking and the gold standard ranking of communities. A larger value of  $\rho$  indicates a better performance. Across all benchmark networks, CRANK achieved average Spearman's rank correlation of  $\rho = 0.82$ . Alternative approaches resulted in poorer average performance: ranking based on modularity and conductance achieved  $\rho = 0.33$  and  $\rho = 0.60$ , respectively, whereas random prioritization obtained  $\rho = 0.00$

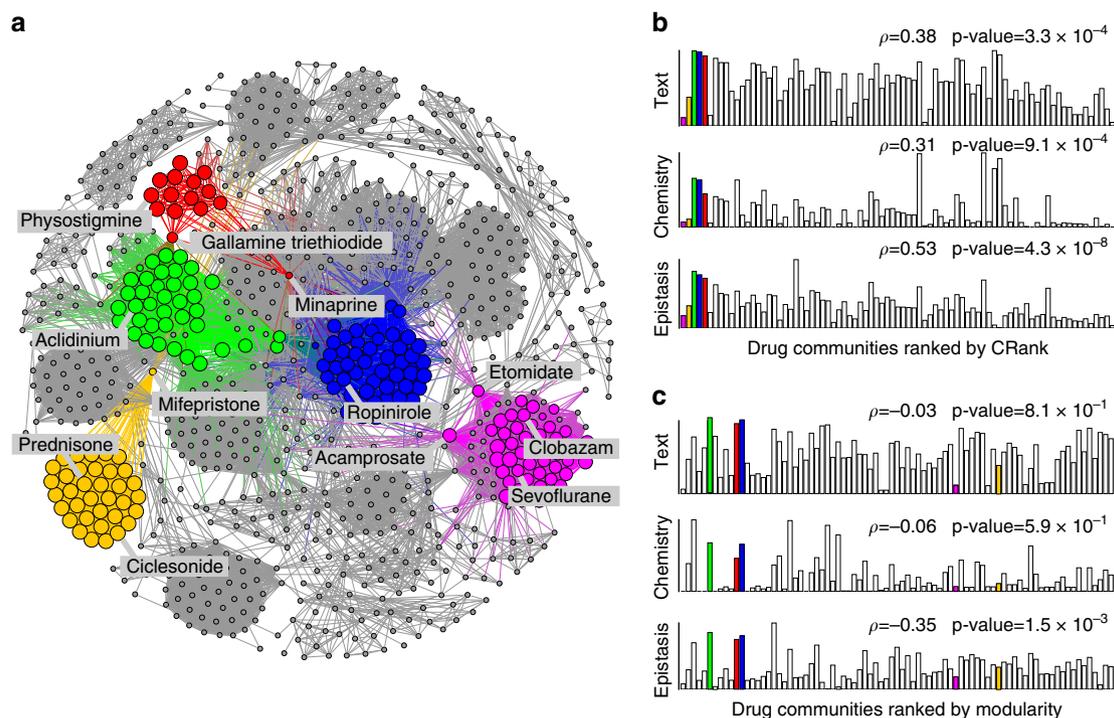
the gold standard community ranking (Fig. 2c). Moreover, we observed no correlation with the gold standard ranking when randomly ordering the detected communities. Although zero correlation is expected, poor performance of random ordering is especially illuminating because prioritization of communities is typically ignored in current network community studies.

**Networks of medical drugs with shared target proteins.** Community rankings obtained by CRANK provide a rich source of testable hypotheses. For example, we consider a network of medical drugs where two drugs are connected if they share at least one target protein (Fig. 3a). Because drugs that are used to treat closely related diseases tend to share target proteins<sup>22</sup>, we expect that drugs belonging to the same community in the network will be rich in chemicals with similar therapeutic effects. Identification of these drug communities hence provides an attractive opportunity for finding new uses of drugs as well as for studying drugs' adverse effects<sup>22</sup>.

After detecting drug communities using a standard community detection method<sup>7</sup>, CRANK relies only on the network structure to prioritize the communities. We evaluate ranking performance by

comparing it to metadata captured in external chemical databases and not used by the ranking method. We find that CRANK assigns higher priority to communities whose drugs are pharmacogenomically more similar (Fig. 3b), indicating that higher-ranked communities contain drugs with more abundant drug-drug interactions, more similar chemical structure, and stronger textual associations. In contrast, ranking communities by modularity score gives a poor correspondence with information in the external chemical databases (Fig. 3c).

We observe that the top-ranked communities are composed from an unusual set of drugs (Fig. 3a and Supplementary Data set), yet drugs with unforeseen community assignment may represent novel candidates for drug repurposing<sup>22</sup>. Examining the highest-ranked communities, we do not expect mifepristone, an abortifacient used in the first months of pregnancy, to appear together with a group of drugs used to treat inflammatory diseases. Another drug with unanticipated community assignment is minaprine, a psychotropic drug that is effective in the treatment of various depressive states<sup>23</sup>. Minaprine is an antidepressant that antagonizes behavioral despair; however, it shares target proteins with several cholinesterase inhibitors. Two examples of such inhibitors are physostigmine, used to treat glaucoma, and galantamine, a drug investigated for the treatment



**Fig. 3** Prioritizing network communities in the network of medical drugs. **a** The network of medical drugs connects two drugs if they share at least one target protein. Communities were detected by a community detection method<sup>7</sup>, and then prioritized by CRANK. Highlighted are five highest-ranked communities as determined by CRANK. Nodes of the highlighted communities are sized by their score of the Likelihood prioritization metric (Supplementary Note 3). Investigation reveals that these communities contain drugs used to: treat asthma and allergies (e.g., prednisone, ciclesonide; yellow nodes), induce anesthesia or sedation (e.g., clobazam, etomidate, sevoflurane, acamprosate; magenta nodes), block neurotransmitters in central and peripheral nervous systems (e.g., physostigmine, minaprine, gallamine triethiodide; red nodes), block the activity of muscarinic receptors (e.g., acclidinium; green nodes), and activate dopamine receptors (e.g., ropinirole; blue nodes). **b, c** We evaluate community prioritization against three external databases (Supplementary Note 6) that were not used during community detection or prioritization. For each community we measure: (1) drug-drug interactions between the drugs (“Epistasis”), (2) chemical structure similarity of the drugs (“Chemistry”), and (3) associations between drugs derived from text data (“Text”). We expect that a true high-priority community will have more drug-drug interactions, higher similarity of chemical structure, and stronger textual associations between the drugs it contains. Taking this into consideration, the external chemical databases define three gold standard rankings of communities against which CRANK is evaluated. Bars represent communities; bar height denotes similarity of drugs in a community with regard to the gold standard based on external chemical databases. In a perfect prioritization, bars would be ordered such that the heights would decrease from left to right. **b** CRANK ranking of drug communities outperforms ranking by modularity **c** across all three chemical databases (as measured by Spearman’s rank correlation  $\rho$  with the gold standard ranking). CRANK ranking achieves  $\rho = 0.38, 0.31, 0.53$ , while modularity obtains  $\rho = -0.03, -0.06, -0.35$

of moderate Alzheimer's disease<sup>24</sup>. In the case of minaprine, an antidepressant, it was just recently shown that this drug is also a cognitive enhancer that may halt the progression of Alzheimer's disease<sup>25</sup>. It is thus attractive that CRANK identified minaprine as a member of a community of primarily cholinesterase inhibitors, which suggests minaprine's potential for drug repurposing for Alzheimer's disease.

The analysis here was restricted to drugs approved for medical use by the U.S. Food and Drug Administration, because these drugs are accompanied by rich metadata that was used for evaluating community prioritization. We find that when CRANK integrates drug metadata into its prioritization model, CRANK can generate up to 55% better community rankings, even when the amount of additional information about drugs is small (Supplementary Note 10). However, approved medical drugs represent less than one percent of all small molecules with recorded interactions. Many of the remaining 99% of these molecules might be candidates for medical usage or drug repurposing but currently have little or no metadata in the chemical databases. This fact further emphasizes the need for methods such as CRANK that can prioritize communities based on network structure alone while not relying on any metadata in external chemical databases.

**Gene and protein interaction networks.** CRANK can also prioritize communities in molecular biology networks, covering a spectrum of physical, genetic, and regulatory gene interactions<sup>11</sup>. In such networks, community detection is widely used because gene communities tend to correlate with cellular functions and thus provide hypotheses about biological pathways and protein complexes<sup>2,3</sup>.

CRANK takes a network and communities detected in that network, and produces a rank-ordered list of communities. As before, while CRANK ranks the communities purely based on network structure, the external metadata about molecular functions, cellular components, and biological processes is used to assess the quality of the community ranking (Supplementary Note 7).

Considering highest-ranked gene communities, CRANK's ranking contains on an average five times more communities whose genes are significantly enriched for cellular functions, components, and processes<sup>16</sup> than random prioritization, and 13% more significantly enriched communities than modularity-based or conductance-based ranking (Supplementary Note 11). For example, in the human protein–protein interaction network, the highest-ranked community by CRANK is composed of 20 genes, including *PORCN*, *AQP5*, *FZD6*, *WNT1*, *WNT2*, *WNT3*, and other members of the Wnt signaling protein family<sup>26</sup> (Supplementary Note 11). Genes in that community form a biologically meaningful group that is functionally enriched in the Wnt signaling pathway processes ( $p$ -value =  $6.4 \times 10^{-23}$ ), neuron differentiation ( $p$ -value =  $1.6 \times 10^{-15}$ ), cellular response to retinoic acid ( $p$ -value =  $2.9 \times 10^{-14}$ ), and in developmental processes ( $p$ -value =  $9.2 \times 10^{-10}$ ).

Functional annotation of molecular networks is largely unavailable and incomplete, especially when studied objects are not genes but rather other entities, e.g., miRNAs, mutations, single-nucleotide variants, or genomic regions outside protein-coding loci<sup>27</sup>. Thus it is often not possible to simply rank the communities by their functional enrichment scores. In such scenarios, CRANK can prioritize communities reliably and accurately using only network structure without necessitating any external databases. Gene communities that rank at the top according to CRANK represent predictions that could guide scientists to prioritize resource-intensive laboratory experiments.

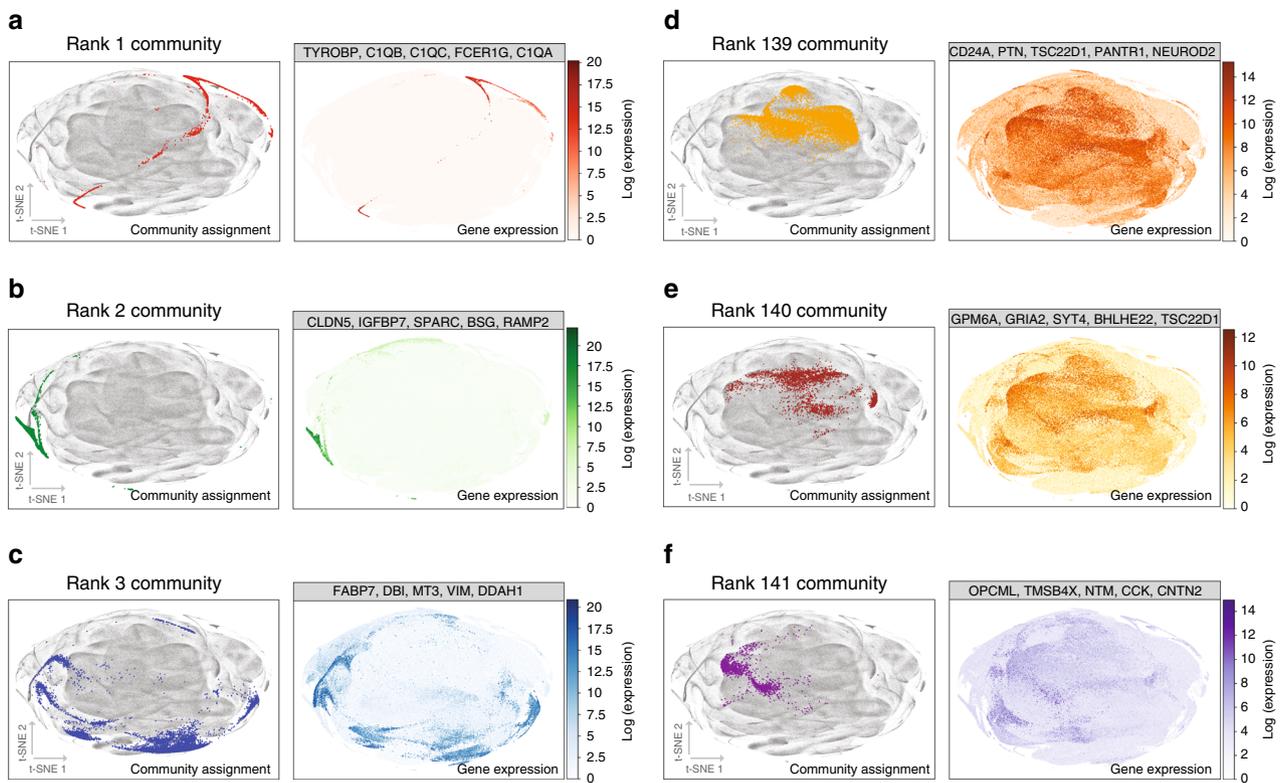
**Megascale cell–cell similarity networks.** Single-cell RNA sequencing has transformed our understanding of complex cell populations<sup>28</sup>. While many types of questions can be answered using single-cell RNA-sequencing, a central focus is the ability to survey the diversity of cell types and composition of tissues within a sample of cells.

To demonstrate that CRANK scales to large networks, we used the single-cell RNA-seq data set containing 1,306,127 embryonic mouse brain cells<sup>29</sup> for which no cell types are known. The data set was preprocessed using standard procedures to select and filter the cells based on quality-control metrics, normalize and scale the data, detect highly variable genes, and remove unwanted sources of variation<sup>9</sup>. The data set was represented as a weighted graph of nearest neighbor relations (edges) among cells (nodes), where relations indicated cells with similar gene expression patterns calculated using diffusion pseudotime analysis<sup>30</sup>. To partition this graph into highly interconnected communities we apply a community detection method proposed for single-cell data<sup>8</sup>. The method separates the cells into 141 fine-grained communities, the largest containing 18,788 (1.8% of) and the smallest only 203 (0.02% of) cells. After detecting the communities, CRANK takes the cell–cell similarity network and the detected communities, and generates a rank-ordered list of communities, assigning a priority to each community. CRANK's prioritization of communities derived from the cell–cell similarity network takes <2 min on a personal computer.

In the cell–cell similarity network, one could assume that top-ranked communities contain highly distinct marker genes<sup>31</sup>, while low-ranked communities contain marker genes whose expression levels are spread out beyond cells in the community. To test this hypothesis, we identify marker genes for each detected community. In particular, for each community we find genes that are differentially expressed in the cells within the community<sup>9</sup> relative to all cells that are not in the community.

We find that high-ranked communities in CRANK contain cells with distinct marker genes, confirming the above hypothesis (average  $z$ -score of marker genes with respect to the bulk mean gene expression was above 200 and never smaller than 150) (Fig. 4a, b). In contrast, cells in low-ranked communities show a weak expression activity diffused across the entire network and no community-specific expression activity (Fig. 4c, d). Examining cells assigned to the highest-ranked community (rank 1 community) in CRANK, we find that most differentially expressed genes are *TYROBP*, *CIQB*, *CIQC*, *FCER1G*, and *CIQA* (at least a 200-fold difference in normalized expression with respect to the bulk mean expression<sup>9</sup>). It is known that these are immunoregulatory genes and that they play important roles in signal transduction in dendritic cells, osteoclasts, macrophages, and microglia<sup>32</sup>. In contrast, low-ranked communities (Fig. 4 visualizes rank 139, rank 140, and rank 141 communities) contain predominantly cells in which genes show no community-specific expression. Genes in communities ranked lower by CRANK hence do not have localized mRNA expression levels, suggesting there are no good marker genes that define those communities<sup>28</sup>. Since the expression levels of mRNA are linked to cellular function and can be used to define cell types<sup>28</sup>, the analysis here points to the potential of using highest-ranked communities in CRANK as candidates to characterize cells at the molecular level, even in data sets where no cells are yet classified into cell types.

**Analysis of CRANK prioritization approach.** The CRANK approach can be applied with any community detection method and can operate on directed, undirected, and weighted networks.



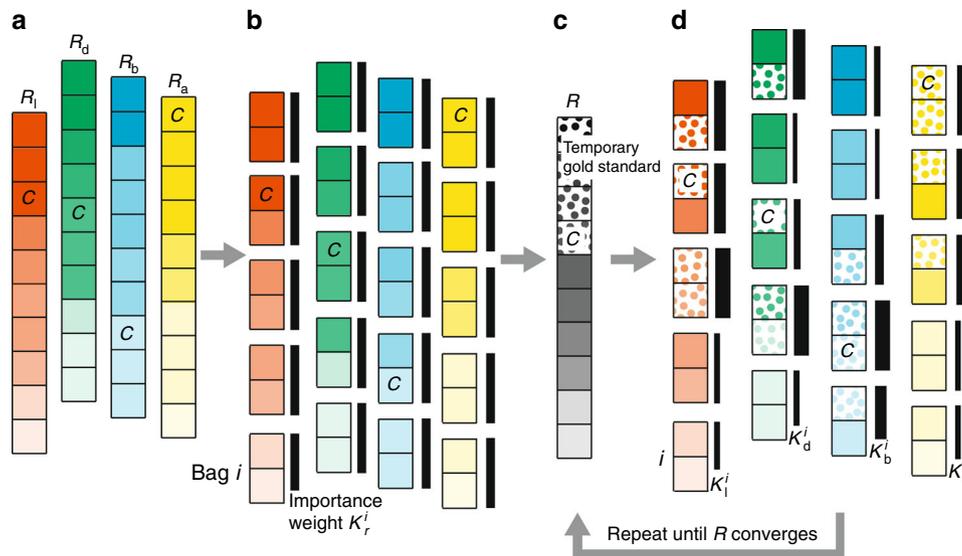
**Fig. 4** Prioritizing network communities in the megascale cell-cell similarity network. The network of embryonic mouse brain has 1,306,127 nodes representing brain cells<sup>29</sup>. Communities are detected using a community detection method developed for single-cell RNA-seq data<sup>8</sup> and prioritized using CRANK, generating a rank-ordered list of detected communities. **a–c** Shown are three communities that are ranked high by CRANK; **a** rank 1, **b** rank 2, and **c** rank 3 community. t-SNE projections<sup>39</sup> show cells assigned to each community. t-SNE is a dimensionality reduction technique that is particularly well suited for visualization of high-dimensional data. Cells assigned to each community are distinguished by color, and all other cells are shown in gray. We investigate the quality of community ranking by examining gene markers for cells in each community<sup>28</sup>. We use the single-cell RNA-seq data set to obtain a gene expression profile for each cell, indicating the activity of genes in the cell. For each community we then identify marker genes, i.e., genes with the strongest differential expression between cells assigned to the community and all other cells<sup>9</sup>. In the t-SNE projection we then color the cells by how active the marker genes are. This investigation reveals that communities ranked high by CRANK are represented by clusters of cells whose marker genes have a highly localized expression. For example, marker genes for rank 1 community in **a** (the highest community in CRANK ranking) are *TYROBP*, *C1QB*, *C1QC*, *FCER1G*, and *C1QA*. Expression of these genes is concentrated in cells that belong to the rank 1 community. Similarly, marker genes for rank 2 and rank 3 communities are specifically active in cell populations that match well the boundary of each community. **d–f** t-SNE projections show cells assigned to 3 low-ranked communities; **d** rank 139, **e** rank 140, and **f** rank 141 community. t-SNE projections are produced using the same differential analysis as in **a–c**. Although these communities correspond to clusters of cells in the t-SNE projections, their marker genes have diluted gene expression that is spread out over the entire network, indicating that CRANK has correctly considered these communities to be low priority. For example, marker genes for rank 141 community in **f** are *OPCML*, *TMSB4X*, *NTM*, *CCK*, and *CNTN2*, which show a weak expression pattern that is diffused across the entire network

Furthermore, CRANK can also use external domain-specific information to further boost prioritization performance (Supplementary Note 10). Results on diverse biological, information, and technological networks and on different community detection methods show that the second best performing approach changes considerably across networks, while CRANK always produces the best result, suggesting that it can effectively harness the network structure for community prioritization (Supplementary Note 8). CRANK automatically adjusts weights of the community metrics in the prioritization, resulting in each metric participating with different intensity across different networks (Supplementary Fig. 6). This is in sharp contrast with deterministic approaches, which are negatively impacted by heterogeneity of network structures and network community models employed by different community detection methods. The four CRANK community prioritization metrics are essential and complementary. CRANK metrics considered together perform on average 45% better than the best single CRANK metric, and 26% better than any subset of three CRANK metrics (Supplementary Note 8). CRANK performs

on average 38% better than approaches that combine alternative community metrics (Supplementary Note 8). Furthermore, CRANK can easily integrate any number of additional and domain-specific community metrics<sup>2,12–15</sup>, and performs well in the presence of low-signal and noisy metrics (Supplementary Note 9). Furthermore, CRANK outperforms alternative approaches that combine the metrics by approximating NP-hard rank aggregation objectives (Supplementary Note 8).

## Discussion

The task of community prioritization is to rank-order communities detected by a community detection method such that communities with best prospects in downstream analysis are ranked towards the top. We demonstrated that prioritizing communities in biological, information, and technological networks is important for maximizing the yield of downstream analyses and experiments. Prior efforts crucially depend on external meta information to calculate the quality of communities



**Fig. 5** Combining community prioritization metrics without an external gold standard. **a** The rank aggregation algorithm starts with four ranked lists of communities,  $R_r$ , each one arising from the values of a different community prioritization metric  $r$  (where  $r$  is one of “l”—likelihood, “d”—density, “b”—boundary, “a”—allegiance). Communities are ordered by the decreasing value of the metric. We use  $C$  to indicate the rank of an illustrative community by the community prioritization metrics and at different stages of the algorithm. **b** Each ranked list is partitioned into equally sized groups, called bags. Each bag  $i$  in ranked list  $R_r$  has attached importance weight  $K_r^i$  whose initial values are all equal (represented by black bars all of same width). CRANK uses the importance weights  $K_r^i$  to initialize aggregate prioritization  $R$  as a weighted average of community ranks  $R_l, R_d, R_b, R_a$ . **c** The top-ranked communities (denoted as dotted cells) in the aggregated prioritization  $R$  serve as a temporary gold standard, which is then used to iteratively update the importance weights  $K_r^i$ . **d** In each iteration, CRANK updates importance weights using the Bayes factor calculation<sup>36</sup> (Supplementary Note 4). Given bag  $i$  and ranked list  $R_r$ , CRANK updates importance weight  $K_r^i$ , based on how many communities from the temporary gold standard appear in bag  $i$ . Updated importance weights then revise the aggregated prioritization in which the new rank  $R(C)$  of community  $C$  is expressed as:  $R(C) = \sum_r \log K_r^{i_r(C)} R_r(C)$ , where  $K_r^{i_r(C)}$  indicates the importance weight of bag  $i_r(C)$  of community  $C$  for metric  $r$ , and  $R_r(C)$  is the rank of  $C$  according to  $r$ . By using an iterative approach, CRANK allows for the importance of a metric not to be predetermined and to vary across communities

with an additional constraint that this information has to be available for all communities. We devised a principled approach for the task of community prioritization. Although the approach does not need any meta information, it can utilize such information if it is available. Furthermore, CRANK is applicable even when the meta information is noisy, incomplete, or available only for a subset of communities.

The CRANK community ranking is based on the premise that high-priority communities produce high values of community prioritization metrics and that these metrics are stable with respect to small perturbations of the network structure. Our findings support this premise and suggest that both the magnitude of the metrics and the robustness of underlying structural features have an important role in the performance of CRANK across a wide range of networks (Supplementary Note 8). CRANK can easily be extended using existing network metrics and can also consider new domain-specific scoring metrics (Supplementary Notes 9 and 10). Thus, it would be especially interesting to apply it to networks, where rich meta information exists and interesting domain-specific scoring metrics can be developed, such as protein interaction networks with disease pathway meta information<sup>33</sup>, and molecular networks with genome-wide associations<sup>34</sup>. We believe that the CRANK approach opens the door to principled methods for prioritizing communities in large networks and, when coupled with experimental validation, can help us to speed-up scientific discovery process.

**Methods**

**Community prioritization model.** CRANK prioritizes communities based on the robustness and magnitude of multiple structural features of each community. For each feature  $f$ , we specify a corresponding prioritization metric  $r_f$ , which captures the magnitude and the robustness of  $f$ . Robustness of  $f$  is defined as the change in the value of  $f$  between the original network and its randomly perturbed version.

The intent here is that high-quality communities will have high values of  $f$  and will also be robust to perturbations of the network structure. We define and discuss specific prioritization metrics later. Here, we first present the overall prioritization model.

Random perturbations of the network are based on rewiring of  $\alpha$  fraction of the edges in a degree preserving manner<sup>35</sup> (Supplementary Note 2). Parameter  $\alpha$  measures perturbation intensity; a value close to zero indicates that the network has only a few edges rewired whereas a value close to one corresponds to a maximally perturbed network, which is a random graph with the same degree distribution as the original network.

Even though the prioritization model is framed conceptually in terms of perturbing the network by rewiring its edges, CRANK never actually rewires the network when calculating the prioritization metrics. Network rewiring is a computationally expensive operation. Instead, we derive analytical expressions that evaluate the metrics in a closed form without physically perturbing the network (Supplementary Note 2), which leads to a substantial increase in scalability of CRANK.

Given structural feature  $f$ , we define prioritization metric  $r_f$  to quantify the change in the value of  $f$  between the original and the perturbed network. We want  $r_f$  to capture the magnitude of feature  $f$  in the original network as well as the change in the value of  $f$  between the network and its perturbed version. We define prioritization metric  $r_f$  for community  $C$  as:

$$r_f(C; \alpha) = \frac{f(C)}{1 + d_f(C, \alpha)}, \tag{1}$$

where  $f(C)$  is the feature value of community  $C$  in the original network,  $\alpha$  measures perturbation intensity,  $d_f(C, \alpha) = |f(C) - f(C|\alpha)|$  is the change of the feature value for community  $C$  between the network and its  $\alpha$ -perturbed version, and  $f(C|\alpha)$  is the value of feature  $f$  in the  $\alpha$ -perturbed version of the network.

Generally, higher priority communities will have higher values of  $r_f$ . In particular, as  $f$  can take values between zero and one, then  $r_f$  also takes values between zero and one.  $r_f$  attains value of zero for community  $C$  whose value of  $f(C)$  is zero. When  $f(C)$  is nonzero, then  $r_f(C; \alpha)$  down-weights it according to the sensitivity of community  $C$  to network rewiring.  $f(C)$  is down-weighted by the largest amount when it changes as much as possible under the network perturbation (i.e.,  $d_f(C, \alpha) = 1$ ). And,  $f(C)$  remains unchanged when community  $C$  is maximally robust to network perturbation (i.e.,  $d_f(C, \alpha) = 0$ ).

**Community prioritization metrics.** Prioritization metric  $r_f(C)$  captures the magnitude as well as the robustness of structural feature  $f$  of community  $C$ . We define four different community prioritization metrics  $r_f$ . Through empirical analysis we show that these metrics holistically and non-redundantly quantify different features of network community structure (Supplementary Note 8). Each metric is necessary and contributes positively to the performance of CRANK. We combine these metrics into a global ranking of communities using a rank aggregation method that we describe later.

Given a network  $G(\mathcal{V}, \mathcal{E}, C)$  with nodes  $\mathcal{V}$ , edges  $\mathcal{E}$ , and detected communities  $C$ , CRANK can be applied in conjunction with any statistical community detection method that allows for computing the following three quantities: (1) the probability of node  $u$  belonging to a given community  $C$ ,  $p_C(u) = p(u \in C)$ , (2) the probability of an edge  $p(u, v) = p((u, v) \in \mathcal{E})$ , and (3) a contribution of community  $C$  towards the existence of an edge  $(u, v)$ ,  $p_C(u, v) = p((u, v) \in \mathcal{E} | u, v \in C)$ . Many commonly used community detection methods allow for computing the above three quantities (Supplementary Note 5).

Our rationale in defining the prioritization metrics is to measure properties that determine a high-quality community, which is also robust and stable with respect to small perturbations of the network. For example, a genuine high-quality community should provide good support for the existence of edges between its members in the original network as well as in the perturbed version. If a small change in the network structure—an edge added here, another deleted there—can completely change the value of the prioritization metric then the community should not be considered high quality. Analogously, a high-quality community should have low confidence for edges pointing outside of the community both in the original as well as in the perturbed network.

**Community likelihood.** The community likelihood metric quantifies the overall connectivity of a given community. It measures the likelihood of the network structure induced by the nodes in the community. Note that the metric does not simply count the edges but considers them in a probabilistic way. As such it quantifies how well the observed edges can be explained by the community  $C$ . The intuition is that high-quality community will contribute a large amount of likelihood to explain the observed edges. We formalize the community likelihood for a given community  $C$  as follows:

$$f_l(C|\alpha) = \prod_{u \in C} p_C(u) \prod_{v \in C} s_C(u, v|\alpha), \quad (2)$$

where  $s_C(u, v|\alpha)$  is defined as follows:

$$s_C(u, v|\alpha) = \begin{cases} p_C(v)p_C(u, v|\alpha) & \text{if } (u, v) \in \mathcal{E} \\ p_C(v)(1 - p_C(u, v|\alpha)) & \text{if } (u, v) \notin \mathcal{E}. \end{cases}$$

Here,  $p_C(u, v|\alpha)$  is a contribution of community  $C$  towards the creation of edge  $(u, v)$  under network perturbation intensity  $\alpha$ . We derive analytical expressions for  $p_C(u, v|\alpha)$  which allows us to compute their values without ever actually perturbing the network (Supplementary Note 2).

Here (and for the other three prioritization metrics) we evaluate the feature in the original network ( $f_l(C) = f_l(C|\alpha = 0)$ ) as well as in the slightly perturbed version of the network ( $f_l(C|\alpha = 0.15)$ ). We then combine the two scores using the prioritization metric formula in Eq. (1).

**Community density.** In contrast to community likelihood, which quantifies the contribution of a community to the overall edge likelihood, community density simply measures the overall strength of connections within the community. By considering edge probabilities that are not conditioned on the community  $C$ , density implicitly takes into consideration potentially hierarchical and overlapping community structures. When a community is nested inside other communities, these enclosing communities contribute to the increased density of community's internal edges. Formally, we define the density of a community as the joint probability of the edges between community members. Assuming network perturbation intensity  $\alpha$ , density of community  $C$  is defined as:

$$f_d(C|\alpha) = \prod_{(u, v) \in \mathcal{E} \cap C, v \in C} p(u, v|\alpha), \quad (3)$$

where  $p(u, v|\alpha)$  is the probability of edge  $(u, v)$  under network perturbation intensity  $\alpha$ . We derive analytical expression for  $p(u, v|\alpha)$  which allows us to compute their values without ever actually perturbing the network (Supplementary Note 2).

**Community boundary.** To complement the internal connectivity measured by community density, community boundary considers the strength of edges leaving the community. A structural feature of a high-quality community is its good separation from the surrounding parts of the network. In other words, a high-quality community should have sharp edge boundary, i.e.,  $B_C = \{(u, v) \in \mathcal{E}; u \in C, v \notin C\}$ <sup>4</sup>. This intuition is captured by accumulating the

likelihood against edges connecting the community with the rest of the network:

$$f_b(C|\alpha) = \prod_{u \in C, v \in \mathcal{V} \setminus C} (1 - p(u, v|\alpha)). \quad (4)$$

The evaluation of Eq. (4) takes computational time linear in the size of the network, which is impractical for large networks with many detected communities. To speed up the calculations, we use negative sampling (Supplementary Note 2) to calculate the value of Eq. (4), and thereby reduce the computational complexity of the boundary metric to time that depends linearly on the number of edges leaving the community.

**Community allegiance.** Last we introduce community allegiance. We define community allegiance as the preference for nodes to attach to other nodes that belong to the same community. Allegiance measures the fraction of nodes in a community for which the total probability of edges pointing inside the community is larger than probability of edges that point to the outside of the community. For a given community  $C$  and network perturbation intensity  $\alpha$ , community allegiance is defined as:

$$f_a(C|\alpha) = \frac{1}{|C|} \sum_{u \in C} \delta \left( \sum_{v \in N_u \cap C} p(u, v|\alpha) \geq \sum_{v \in N_u \setminus C} p(u, v|\alpha) \right), \quad (5)$$

where  $N_u$  is a set of network neighbors of  $u$  and  $\delta$  is the indicator function,  $\delta(x) = 1$  if  $x$  is true, and  $\delta(x) = 0$ , otherwise.

Community has high allegiance if nodes in the community tend to be more strongly connected to other members of the community than to the rest of the network. In a community with no significant allegiance this metric takes a value that is close to zero or changes substantially when the network is only slightly perturbed. However, in the presence of substantial community allegiance, the metric takes large values and is not sensitive to edge perturbation.

**Combining community prioritization metrics.** We just defined four community prioritization metrics: likelihood, density, boundary, and allegiance. Each metric on its own provides a useful signal for prioritizing communities (Supplementary Note 8). However, scores of each metric might be biased, have high variance, and behave differently across different networks (Supplementary Fig. 6). It is thus essential to combine the values of individual metrics into a single aggregated score.

We develop an iterative unsupervised rank aggregation method that, without requiring an external gold standard, combines the prioritization metrics into a single aggregated prioritization of communities. The method is outlined in Fig. 5. It naturally takes into consideration the fact that importance of individual prioritization metrics varies across networks and across community detection methods. The aggregation method starts by representing the values of each prioritization metric with a ranked list. In each ranked list, communities are ordered by the decreasing value of the metric. The method then determines the contribution of each ranked list to the aggregate prioritization by calculating importance weights. The calculation is based on Bayes factors<sup>36–38</sup>, an established tool in statistics. Each ranked list has associated a set of importance weights. Importance weights can vary with rank in the list. The method then calculates the aggregated prioritization of communities in an iterative manner by taking into account uncertainty that is present across different ranked lists and within each ranked list.

To calculate the weights without requiring gold standard, the method uses a two-stage iterative procedure. After initializing the aggregated prioritization, the method alternates between the following two stages until no changes in the aggregated prioritization are observed: (1) use the aggregated prioritization to calculate the importance weights for each ranked list, and (2) re-aggregate the ranked lists based on the importance weights calculated in the previous stage.

The model for aggregating community prioritization metrics, the algorithm, and the analysis of its computational time complexity are detailed in Supplementary Notes 4 and 5 (Supplementary Note 1). The complete algorithm of CRANK approach is provided in Supplementary Note 5.

**Data availability.** All relevant data are public and available from the authors of original publications. The project website is at: <http://snap.stanford.edu/crank>. The website contains preprocessed data used in the paper and additional examples of CRANK's use.

Received: 10 January 2018 Accepted: 5 June 2018

Published online: 29 June 2018

## References

- Benson, A. R., Gleich, D. F. & Leskovec, J. Higher-order organization of complex networks. *Science* **353**, 163–166 (2016).

2. Menche, J. et al. Uncovering disease-disease relationships through the incomplete interactome. *Science* **347**, 1257601 (2015).
3. Costanzo, M. et al. A global genetic interaction network maps a wiring diagram of cellular function. *Science* **353**, 1381 (2016).
4. Fortunato, S. Community detection in graphs. *Phys. Rep.* **486**, 75–174 (2010).
5. Newman, M. E. Modularity and community structure in networks. *Proc. Natl Acad. Sci. USA* **103**, 8577–8582 (2006).
6. Gopalan, P. K. & Blei, D. M. Efficient discovery of overlapping communities in massive networks. *Proc. Natl Acad. Sci. USA* **110**, 14534–14539 (2013).
7. Yang, J., McAuley, J. & Leskovec, J. Detecting cohesive and 2-mode communities in directed and undirected networks. In *Proceedings of the 7th ACM International Conference on Web Search and Data Mining*, 323–332 (2014).
8. Levine, J. H. et al. Data-driven phenotypic dissection of aml reveals progenitor-like cells that correlate with prognosis. *Cell* **162**, 184–197 (2015).
9. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.* **33**, 495–502 (2015).
10. Newman, M. & Clauset, A. Structure and inference in annotated networks. *Nat. Commun.* **7**, 11863 (2016).
11. Rolland, T. et al. A proteome-scale map of the human interactome network. *Cell* **159**, 1212–1226 (2014).
12. Subramanian, A. et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl Acad. Sci. USA* **102**, 15545–15550 (2005).
13. Glaab, E., Baudot, A., Krasnogor, N., Schneider, R. & Valencia, A. EnrichNet: network-based gene set enrichment analysis. *Bioinformatics* **28**, 451–457 (2012).
14. Baryshnikova, A. Systematic functional annotation and visualization of biological networks. *Cell Syst.* **2**, 412–421 (2016).
15. Hofree, M. et al. Network-based stratification of tumor mutations. *Nat. Methods* **10**, 1108–1115 (2013).
16. Ashburner, M. et al. Gene Ontology: tool for the unification of biology. *Nat. Genet.* **25**, 25–29 (2000).
17. Regev, A. et al. Science forum: the human cell atlas. *eLife* **6**, e27041 (2017).
18. da Silva, R. R., Dorrestein, P. C. & Quinn, R. A. Illuminating the dark matter in metabolomics. *Proc. Natl Acad. Sci. USA* **112**, 12549–12550 (2015).
19. Layeghifard, M., Hwang, D. M. & Guttman, D. S. Disentangling interactions in the microbiome: a network perspective. *Trends Microbiol.* **25**, 217–228 (2017).
20. Aksenov, A. A., da Silva, R., Knight, R., Lopes, N. P. & Dorrestein, P. C. Global chemical analysis of biology by mass spectrometry. *Nat. Rev. Chem.* **1**, s41570–017 (2017).
21. Schaeffer, S. E. Graph clustering. *Comput. Sci. Rev.* **1**, 27–64 (2007).
22. Guney, E., Menche, J., Vidal, M. & Barabasi, A.-L. Network-based in silico drug efficacy screening. *Nat. Commun.* **7**, 10331 (2016).
23. Law, V. et al. DrugBank 4.0: shedding new light on drug metabolism. *Nucleic Acids Res.* **42**, 1091–1097 (2014).
24. Erkinjuntti, T. et al. Efficacy of Galantamine in probable vascular dementia and Alzheimer’s disease combined with cerebrovascular disease: a randomised trial. *Lancet* **359**, 1283–1290 (2002).
25. Cecilia Rodrigues Simoes, M. et al. Donepezil: an important prototype to the design of new drug candidates for Alzheimer’s disease. *Mini Rev. Med. Chem.* **14**, 2–19 (2014).
26. Kessenbrock, K. et al. Diverse regulation of mammary epithelial growth and branching morphogenesis through noncanonical Wnt signaling. *Proc. Natl Acad. Sci. USA* **114**, 3121–3126 (2017).
27. Kellis, M. et al. Defining functional DNA elements in the human genome. *Proc. Natl Acad. Sci. USA* **111**, 6131–6138 (2014).
28. Trapnell, C. Defining cell types and states with single-cell genomics. *Genome Res.* **25**, 1491–1498 (2015).
29. Zheng, G. X. et al. Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.* **8**, 14049 (2017).
30. Haghverdi, L., Buettner, M., Wolf, F. A., Buettner, F. & Theis, F. J. Diffusion pseudotime robustly reconstructs lineage branching. *Nat. Methods* **13**, 845–848 (2016).
31. Zeisel, A. et al. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* **347**, 1138–1142 (2015).
32. Ma, J., Jiang, T., Tan, L. & Yu, J.-T. TYROBP in Alzheimer’s disease. *Mol. Neurobiol.* **51**, 820–826 (2015).
33. Agrawal, M., Zitnik, M. & Leskovec, J. Large-scale analysis of disease pathways in the human interactome. In *Pacific Symposium on Biocomputing*, vol. **23**, 111 (World Scientific, Singapore, 2018).
34. Choobdar, S. et al. Open community challenge reveals molecular network modules with key roles in diseases. Preprint at *bioRxiv* <https://www.biorxiv.org/content/early/2018/02/15/265553> (2018).
35. Karrer, B., Levina, E. & Newman, M. E. Robustness of community structure in networks. *Phys. Rev. E* **77**, 046119 (2008).
36. Kass, R. E. & Raftery, A. E. Bayes factors. *J. Am. Stat. Assoc.* **90**, 773–795 (1995).
37. Berger, J. O. & Pericchi, L. R. The intrinsic bayes factor for model selection and prediction. *J. Am. Stat. Assoc.* **91**, 109–122 (1996).
38. Casella, G. & Moreno, E. Assessing robustness of intrinsic tests of independence in two-way contingency tables. *J. Am. Stat. Assoc.* **104**, 1261–1271 (2012).
39. Maaten, L.v.d. & Hinton, G. Visualizing data using t-SNE. *J. Mach. Learn. Res.* **9**, 2579–2605 (2008).

### Acknowledgements

M.Z., R.S., and J.L. were supported by NSF, NIH BD2K, DARPA SIMPLEX, Stanford Data Science Initiative, and Chan Zuckerberg Biohub. We thank Austin R. Benson and William L. Hamilton for their valuable feedback.

### Author contributions

M.Z., R.S., and J.L. designed and performed research, contributed new analytic tools, analyzed data, and wrote the paper.

### Additional information

**Supplementary Information** accompanies this paper at <https://doi.org/10.1038/s41467-018-04948-5>.

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

**Reprints and permission** information is available online at <http://npg.nature.com/reprintsandpermissions/>



**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