

# UC San Diego

## UC San Diego Previously Published Works

### Title

HiDeF: identifying persistent structures in multiscale omics data.

### Permalink

<https://escholarship.org/uc/item/5h11h6sx>

### Journal

Genome Biology, 22(1)

### Authors

Zheng, Fan

Zhang, She

Churas, Christopher

et al.

### Publication Date

2021-01-07

### DOI

10.1186/s13059-020-02228-4

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>


Peer reviewed

SHORT REPORT

Open Access

# HiDeF: identifying persistent structures in multiscale 'omics data



Fan Zheng<sup>1\*†</sup> , She Zhang<sup>2†</sup>, Christopher Churas<sup>1</sup>, Dexter Pratt<sup>1</sup>, Ivet Bahar<sup>2</sup> and Trey Ideker<sup>1\*</sup>

\* Correspondence: [fzheng@health.ucsd.edu](mailto:fzheng@health.ucsd.edu); [tideker@health.ucsd.edu](mailto:tideker@health.ucsd.edu)

<sup>†</sup>Fan Zheng and She Zhang contributed equally to this work.

<sup>1</sup>Division of Genetics, Department of Medicine, University of California San Diego, La Jolla, CA 92093, USA  
Full list of author information is available at the end of the article

## Abstract

In any 'omics study, the scale of analysis can dramatically affect the outcome. For instance, when clustering single-cell transcriptomes, is the analysis tuned to discover broad or specific cell types? Likewise, protein communities revealed from protein networks can vary widely in sizes depending on the method. Here, we use the concept of persistent homology, drawn from mathematical topology, to identify robust structures in data at all scales simultaneously. Application to mouse single-cell transcriptomes significantly expands the catalog of identified cell types, while analysis of SARS-COV-2 protein interactions suggests hijacking of WNT. The method, HiDeF, is available via Python and Cytoscape.

**Keywords:** Systems biology, Multiscale, Persistent homology, Community detection, Resolution, Single-cell clustering, Protein-protein interaction network

## Background

Significant patterns in data often become apparent only when looking at the right scale. For example, single-cell RNA sequencing data can be clustered coarsely to identify broad categories of cells (e.g., mesoderm, ectoderm), or analyzed more sharply to delineate highly specific subtypes (e.g., pancreas islet  $\beta$ -cells, thymus epithelium) [1–3]. Likewise, protein-protein interaction networks can inform groups of proteins spanning a wide range of spatial dimensions, from protein dimers (e.g., leucine zippers) to larger complexes of dozens or hundreds of subunits (e.g., proteasome, nuclear pore) to entire organelles (e.g., centriole, mitochondria) [4]. Many different approaches have been devised or applied to detect structures in biological data, including standard clustering, network community detection, and low-dimensional data projection [5–7], some of which can be tuned for sensitivity to objects of a certain size or scale (the so-called resolution parameters) [8, 9]. Even tunable algorithms, however, face the dilemma that the particular scales at which the significant biological structures arise are usually unknown in advance.

Guidelines for detecting robust patterns across scales come from the field of topological data analysis, which studies the geometric “shape” of data using tools from algebraic topology and pure mathematics [10]. A fundamental concept in this field is

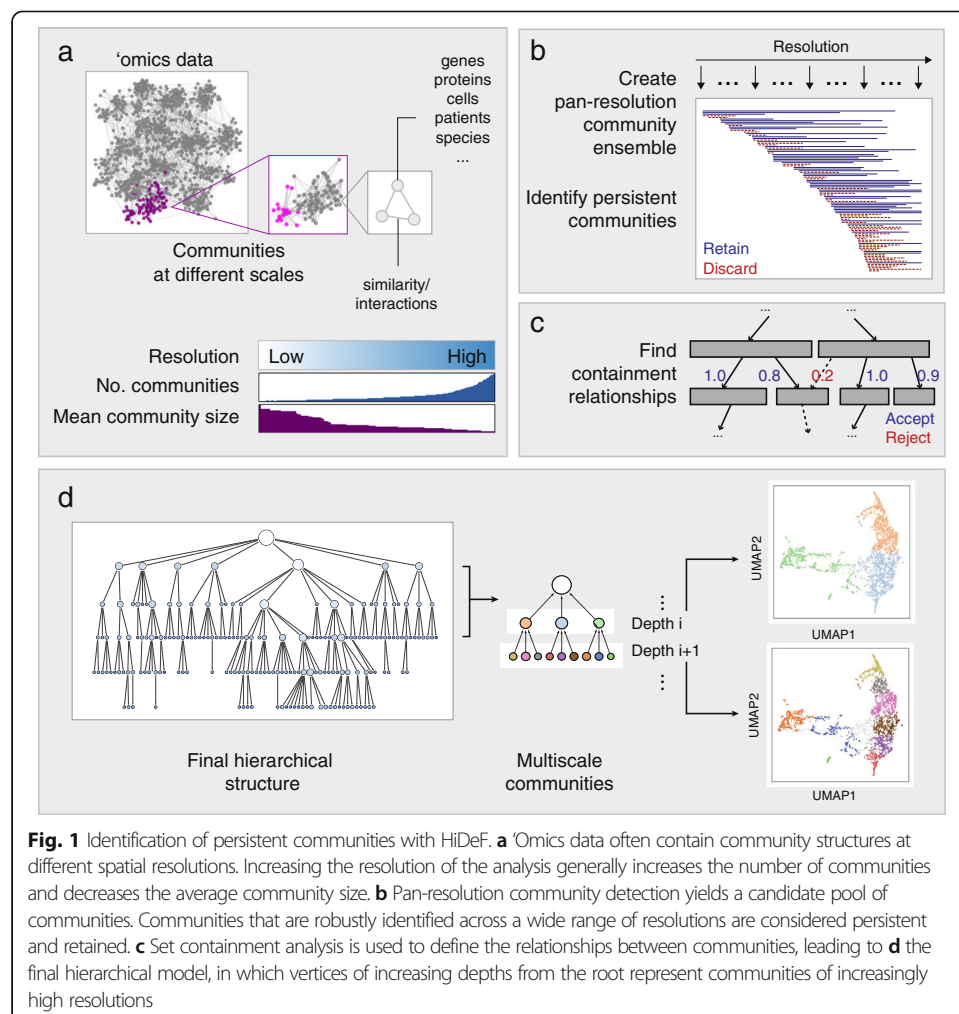


© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

“persistent homology” [11], the idea that the core structures intrinsic to a dataset are those that persist across different scales. Recently, this concept has begun to be applied to analysis of ‘omics data and particularly biological networks [12, 13]. Here, we sought to integrate concepts from persistent homology with existing algorithms for network community detection, resulting in a fast and practical multiscale approach we call the *Hierarchical community Decoding Framework* (HiDeF).

## Results and discussion

HiDeF works in the three phases to analyze the structure of a biological dataset (see the “Methods” section). To begin, the dataset is formulated as a similarity network, depicting a set of biological entities (e.g., genes, proteins, cells, patients, or species) and pairwise connections among these entities (representing similarities in their data profiles). The goal of the first phase is to detect network communities, i.e., groups of densely connected biological entities. Communities are identified continually as the spatial resolution is scanned, producing a comprehensive pool of candidates across all scales of analysis (Fig. 1a). In the second phase, candidate communities arising at different resolutions are pairwise aligned to identify those that have been redundantly identified and are thus *persistent* (Fig. 1b). In the third phase, persistent communities are

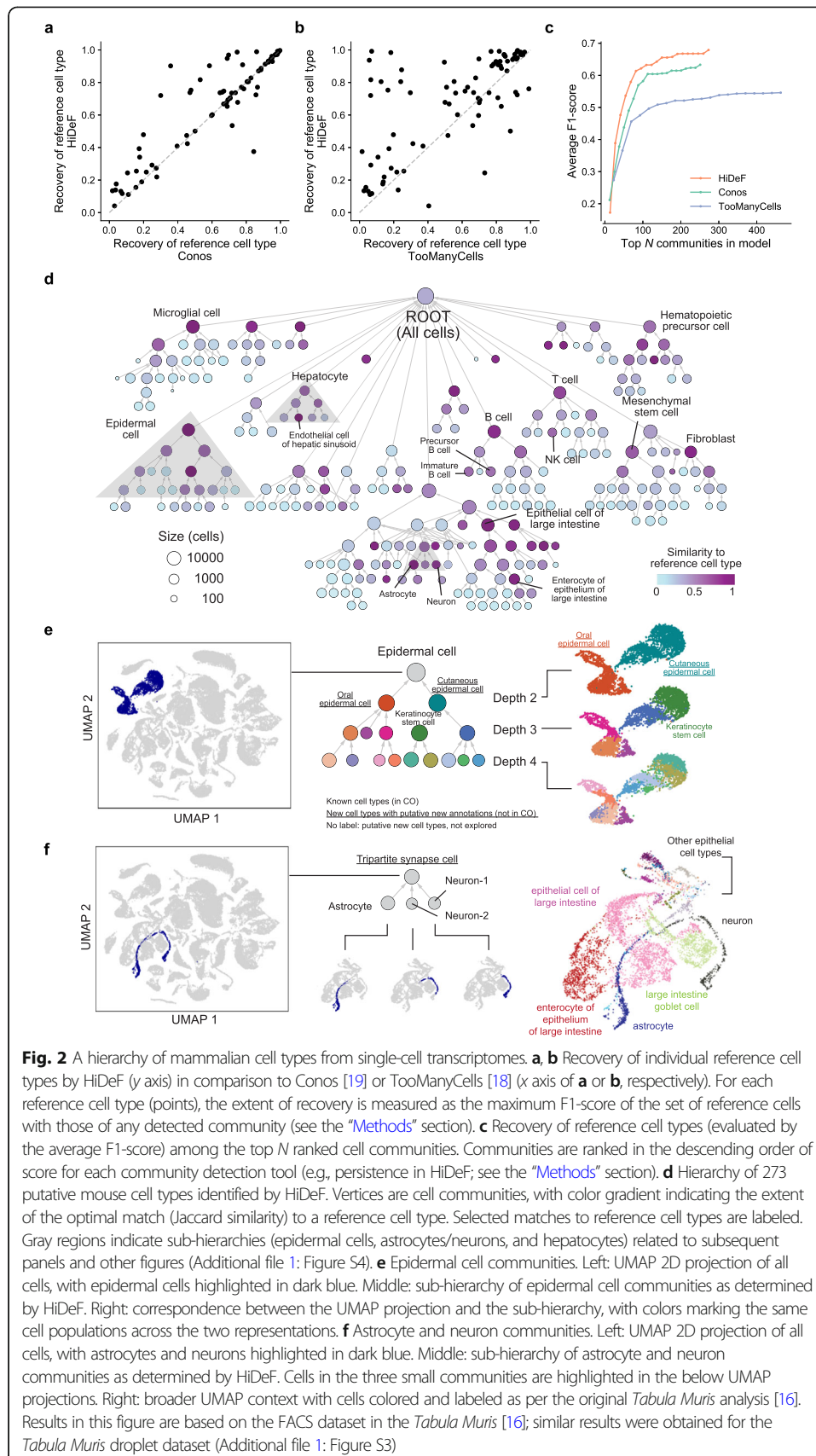


analyzed to identify cases where a community is fully or partially contained within another (typically larger) community, resulting in a hierarchical assembly of nested and overlapping biological structures (Fig. 1c, d). HiDeF is implemented as a Python package and can be accessed interactively in the Cytoscape network analysis and visualization environment [14] (Availability of data and materials).

We first explored the idea of measuring community persistence via analysis of synthetic datasets [15] in which communities were simulated and embedded in the similarity network at two different scales (Additional file 1: Figure S1a; see the “Methods” section). Notably, the communities determined to be most persistent by HiDeF were found to accurately recapitulate the simulated communities at the two scales (Additional file 1: Figure S1b-g). In contrast, applying community detection algorithms at a fixed resolution had limited capability to capture both scales of simulated structures simultaneously (Additional file 1: Figure S2; see the “Methods” section).

We next evaluated whether persistent community detection improves the characterization of cell types. We applied HiDeF to detect robust nested communities within cell-cell similarity networks based on the mRNA expression profiles of 100,605 single cells gathered across the organs and tissues of mice (obtained from two datasets in the *Tabula Muris* project [16]; see the “Methods” section). These cells had been annotated with a controlled vocabulary of cell types from the Cell Ontology (CO) [17], via analyses of cell-type-specific expression markers [16]. We used groups of cells sharing the same annotations to define a panel of 136 reference cell types and measured the degree to which each reference cell type could be recapitulated by a HiDeF community of cells (see the “Methods” section). We compared these results to TooManyCells [18] and Conos [19], two recently developed methods that generate nested communities of single cells in divisive and agglomerative manners, respectively (see the “Methods” section). Reference cell types tended to better match communities generated by HiDeF than those of other approaches, with 65% (89/136) having a highly overlapping community (Jaccard index  $> 0.5$ ) in the HiDeF hierarchy (Fig. 2a, b; Additional file 1: Figure S3a, b). This favorable performance was observed consistently when adjusting HiDeF parameters to formulate a simple hierarchy, containing only the strongest structures, or a more complex hierarchy including additional communities that are less persistent but still significant (Fig. 2c; Additional file 1: Figure S3c).

The top-level communities in the HiDeF hierarchy corresponded to broad cell lineages such as “T cell,” “B cell,” and “epidermal cell.” Finer-grained communities mapped to more specific known subtypes (Fig. 2d) or, more frequently, putative new subtypes within a lineage. For example, “epidermal cell” was split into two distinct epidermal tissue locations, skin and tongue; further splits suggested the presence of still more specific uncharacterized cell types (Fig. 2e). HiDeF communities also captured known cell types that were not apparent from 2D visual embeddings (Additional file 1: Figure S4a, b), and also suggested new cell-type combinations. For example, astrocytes were joined with two communities of neuronal cells to create a distinct cell type not observed in the hierarchies of TooManyCells [18], Conos [19], or a two-dimensional data projection with UMAP [20] (Fig. 2f; Additional file 1: Figure S4c). This community may correspond to the grouping of a presynaptic neuron, postsynaptic neuron, and a surrounding astrocyte within a so-called tripartite synapse [21].

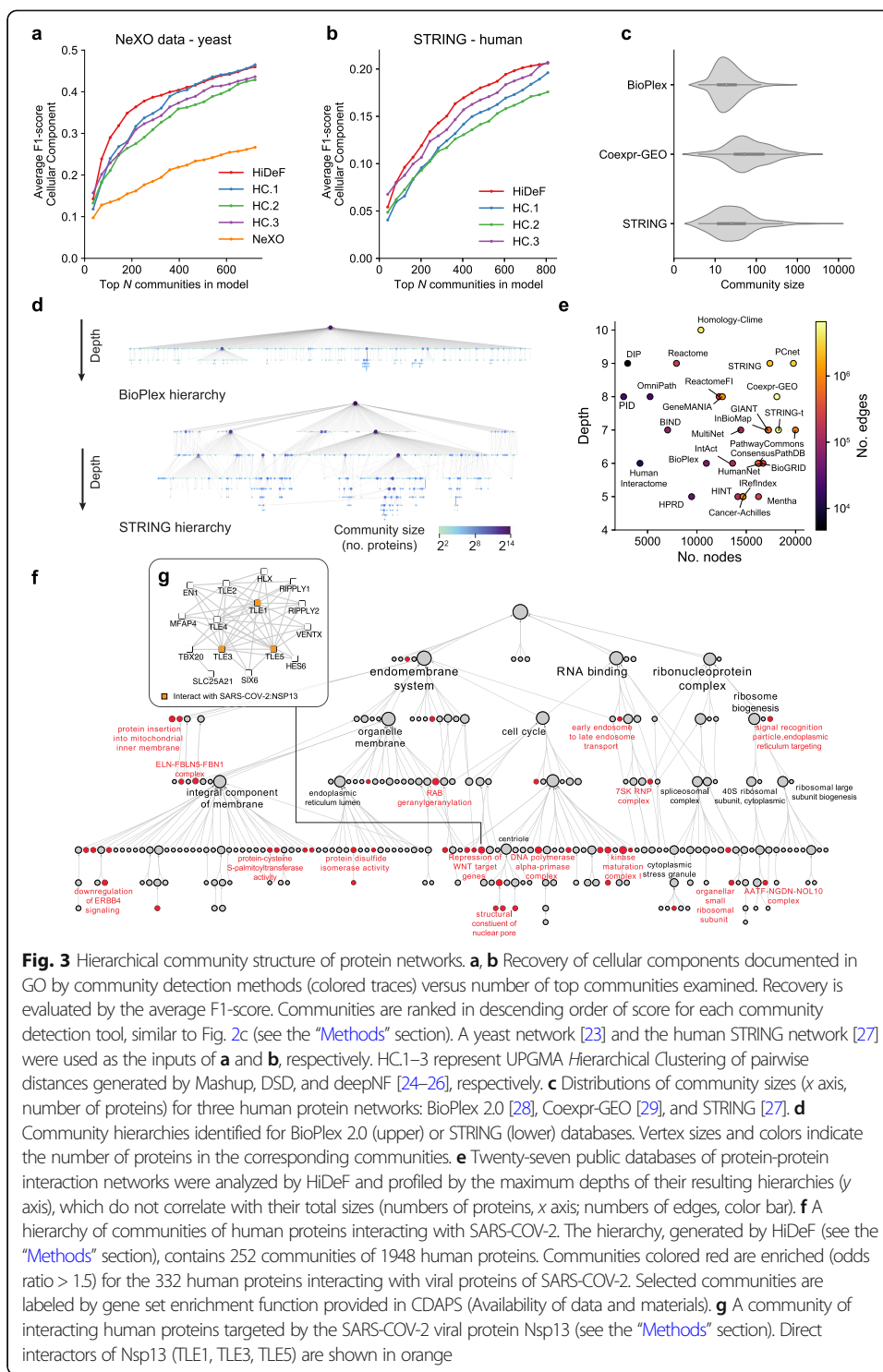


Next, we applied HiDeF to analyze protein-protein interaction networks, with the goal of characterizing protein complexes and higher-order protein assemblies spanning spatial scales. We benchmarked this task by the agreement between HiDeF communities and the Gene Ontology (GO) [22], a database that manually assigns proteins to cellular components, processes, or functions based on curation of literature (see the “Methods” section). Application to protein-protein interaction networks from budding yeast and human found that HiDeF captured knowledge in GO more significantly than previous pipelines proposed for this task, including the NeXO approach to hierarchical community detection [23] and standard hierarchical clustering of pairwise protein distances calculated by three recent network embedding approaches [24–26] (Fig. 3a, b; Additional file 1: Figures S5, S6; see the “Methods” section). HiDeF could be directly applied to the original interaction networks or to network embedded versions to further improve the performance and robustness (Additional file 1: Figure S7).

We also applied HiDeF to analyze a collection of 27 human protein interaction networks [30, 31]. We found significant differences in the distributions of community sizes across these networks, loosely correlating with the different measurement approaches used to generate each network. For example, BioPlex 2.0, a network characterizing biophysical protein-protein interactions by affinity-purification mass-spectrometry (AP-MS) [28], was dominated by small communities of 10–50 proteins, whereas a network based on mRNA co-expression [29] tended towards larger-scale communities of > 50 proteins. In the middle of this spectrum, the STRING network, which integrated biophysical protein interactions and gene co-expression with a variety of other features [27], contained both small and large communities (Fig. 3c). In agreement with the observation above, the hierarchy of BioPlex had a relatively shallow shape in comparison to that of STRING (and other integrated networks including GIANT and PCNet [30, 32]), in which communities across many scales formed a deep hierarchy (Fig. 3d, e; Availability of data and materials).

In contrast to clustering frameworks, HiDeF recognizes when a community is contained by multiple parent communities, which in the context of protein-protein networks suggests that the community participates in diverse pleiotropic biological functions. For example, a community corresponding to the MAPK (ERK) pathway participated in multiple larger communities, including RAS and RSK pathways, sodium channels, and actin capping, consistent with the central roles of MAPK signaling in these distinct biological processes [33] (Additional file 1: Figure S8). The hierarchies of protein communities identified from each of these networks have been made available as a resource in the NDEx database [34] (Availability of data and materials).

To explore multiscale data analysis in the context of an urgent public health issue, we considered a recent application of AP-MS that characterized interactions between the 27 SARS-COV-2 viral subunits and 332 human host proteins [35]. We used network propagation to select a subnetwork of the BioPlex 3.0 human protein interactome [36] proximal to these 332 proteins (1948 proteins and 22,835 interactions) and applied HiDeF to identify its community structure (see the “Methods” section). Among the 251 persistent communities identified (Fig. 3f), we noted one consisting of human Transducin-Like Enhancer (TLE) family proteins, TLE1, TLE3, and TLE5, which interacted with SARS-COV2 Nsp13, a highly conserved RNA synthesis protein in corona and other nidoviruses (Fig. 3g) [37]. TLE proteins are well-known inhibitors of the Wnt



signaling pathway [38]. Inhibition of WNT, in turn, has been shown to reduce coronavirus replication [39] and recently proposed as a COVID-19 treatment [40]. If interactions between Nsp13 and TLE proteins can be shown to facilitate activation of WNT, TLEs may be of potential interest as drug targets.

## Conclusions

Community persistence provides a basic metric for distilling biological structure from data, which can be tuned to select only the strongest structures or to include weaker patterns that are less persistent but still significant. This concept applies to diverse biological subfields, as demonstrated here for single-cell transcriptomics and protein interaction mapping. While these subfields currently employ very different analysis tools which largely evolve separately, it is perhaps high time to seek out core concepts and broader fundamentals around which to unify some of the ongoing development efforts. To that effect, the methods explored here have wide applicability to analyze the multi-scale organization of many other biological systems, including those related to chromosome organization, the microbiome and the brain.

## Methods

### Overview of the approach

Consider an undirected network graph  $G$ , representing a set of biological *objects* (vertices) and a set of *similarity relations* between these objects (edges). Examples of interest include networks of cells, where edges represent pairwise cell-cell similarity in transcriptional profiles characterized by single-cell RNA-seq, or networks of proteins, where edges represent pairwise protein-protein biophysical interactions. We seek to group these objects into *communities* (subsets of objects) that appear at different scales and identify approximate *containment relationships* among these communities, so as to obtain a hierarchical representation of the network structure. The workflow is implemented in three phases. Phase I identifies communities in  $G$  at each of a series of spatial resolutions  $\gamma$ . Phase II identifies which of these communities are *persistent* by way of a *pan-resolution community graph*  $G_C$ , in which vertices represent communities, including those identified at each resolution, and each edge links pairs of similar communities arising at different resolutions. Persistent communities correspond to large components in  $G_C$ . Phase III constructs a final hierarchical structure  $H$  that represents containment and partial containment relationships (directed edges) among the persistent communities (vertices).

### Phase I: Pan-resolution community detection

Community detection methods generally seek to maximize a quantity known as the *network modularity*, as a function of community assignment of all objects [41]. A *resolution parameter* integrated into the modularity function can be used to tune the scale of the communities identified [9, 42, 43], with larger/smaller scale communities having more/fewer vertices on average (Fig. 1a). Of the several types of resolution parameter that have been proposed, we adopted that of the Reichardt-Bornholdt configuration model [42], which defines the generalized modularity as:

$$Q(\vec{G}, \gamma) = \sum_{ij} \left( A_{ij} - \gamma \frac{k_i k_j}{2m} \right) \delta(i, j)$$

where  $\vec{G}$  defines a mapping from objects in  $G$  to community labels,  $k_i$  is the degree of vertex  $i$ ,  $m$  is the total number of edges in  $G$ ,  $\gamma$  is the resolution parameter,  $\delta(i, j)$  indicates that vertices  $i$  and  $j$  are assigned to the same community by  $\vec{G}$ , and  $A$  is the adjac-



community matrix of  $G$ . To determine  $\vec{G}$ , we use the extended Louvain algorithm implemented in the Python package `louvain-igraph` (<http://github.com/vtraag/louvain-igraph>; version 0.6.1). Values of  $\gamma$  are sampled logarithmically between lower and upper bounds  $\gamma_{\min}$  and  $\gamma_{\max}$  at a minimum density such that for all  $\gamma$  there exist at least 10 nearby  $\gamma'$  satisfying:

$$\left| \log_{10}(\gamma') - \log_{10}(\gamma) \right| < 0.1$$

Two  $\gamma$  values satisfying the above formula are defined as  $\gamma$ -proximal. The sampling step was practically set to 0.1 to sufficiently capture the interesting structures in the data; it is conceptually similar to the Nyquist sampling frequency in signal processing [44]. We used  $\gamma_{\min} = 0.001$ , which we found always resulted in the theoretical minimum number of communities, equal to the number of connected components in  $G$ . We used  $\gamma_{\max} = 20$  for single-cell data (Fig. 2; Additional file 1: Figures S3, S4) and  $\gamma_{\max} = 50$  for simulated networks (Additional file 1: Figures S1, S2) and protein interaction networks (Fig. 3; Additional file 1: Figures S5-8). Performing Louvain community detection at each  $\gamma$  over this defined progression of values resulted in a set of communities  $\vec{G}$  at each  $\gamma$ .

### Phase II: Identification of persistent communities

To identify persistent communities, we define the pairwise similarity between any two communities  $a$  and  $b$  as the Jaccard similarity of their sets of objects,  $s(a)$  and  $s(b)$ :

$$J(a, b) = \frac{s(a) \cap s(b)}{s(a) \cup s(b)}$$

Pairwise community similarity is computed only for pairs of communities identified at two different  $\gamma$ -proximal resolution values, as communities within a resolution do not overlap. To represent these similarities, we define a *pan-resolution community graph*  $G_C$ , in which vertices are communities identified at any resolution and edges connect pairs of similar communities having  $J(a, b) > \tau$ . Each component of  $G_C$  defines a family of similar communities spanning resolutions, for which the *persistence* can be naturally defined by the number of distinct  $\gamma$  values covered by the component. For each component in  $G_C$  larger than a persistence threshold  $\chi$ , the biological objects participating in more than  $p\%$  of communities represented by the vertices of that component define a *persistent community*.

### Phase III: A hierarchy of nested and overlapping communities

We initialize a hierarchical structure represented by  $H$ , a directed acyclic graph (DAG) in which each vertex represents a persistent community. A *root vertex* is added to represent the community of all objects. The containment relationship between two vertices,  $v$  and  $w$ , is quantified by the *containment index* (CI):

$$CI(v, w) = \frac{|s(v) \cap s(w)|}{|s(w)|}$$

which measures the fraction of objects in  $w$  shared with  $v$ . An edge is added from  $v$  to  $w$  in  $H$  if  $CI(v, w)$  is larger than a threshold  $\sigma$  ( $w$  is  $\sigma$ -contained by  $v$ ). Since  $J(v, w) < \tau$

for all  $v, w$  (a property established by the procedure for connecting similar communities in phase II), setting  $\sigma \geq 2\tau/(1 + \tau)$  guarantees  $H$  to be acyclic. In practice, we used a relaxed threshold  $\sigma = \tau$ , which we found generally maintains the acyclic property but includes additional containment relations. In the (in our experience rare) event that cycles are generated in  $H$ , i.e.,  $CI(v, w) \geq \tau$  and  $CI(w, v) \geq \tau$ , we add a new community to  $H$ , the union of  $v$  and  $w$ , and remove  $v$  and  $w$  from  $H$ .

Finally, redundant relations are removed by obtaining a transitive reduction [45] of  $H$ , which represents the hierarchy returned by HiDeF describing the organization of communities. The biological objects assigned to each community are expanded to include all objects assigned to its descendants. Throughout this study, we used the parameters  $\tau = 0.75$ ,  $\chi = 5$ ,  $p = 75$ . Note that since  $\chi$  is a threshold of minimum persistence, the results under a larger value of  $\chi'$  can be produced by simply removing communities with persistence lower than  $\chi'$  (Figs. 2c and 3a, b; Additional file 1: Figures S2, S3c, S5). Generally, we observed that the conclusions drawn in this study were robust to this choice of parameters. The persistence of communities only moderately correlates with community sizes, with the consequence that different choices of persistent threshold  $\chi$  do not strongly favor structures at particular scales (Additional file 1: Figure S9). Different combinations of parameters  $\tau$  and  $p$  typically do not significantly change the performance of HiDeF in the benchmark tests on protein-protein interaction networks (Additional file 1: Figure S6), except that certain parameters (e.g.,  $\tau = 0.9$ ) are less robust to network perturbation (i.e., randomly deleting edges from networks). We found that combining HiDeF with node embedding resolved this issue and further improved the performance and robustness (Additional file 1: Figure S7; see sections below).

### Simulated benchmark networks

Simulated network data were generated using the Lancichinetti-Fortunato-Radicchi (LFR) method [15] (Additional file 1: Figures S1, S2). We used an available implementation (LFR benchmark graphs package 5 at <http://www.santofortunato.net/resources>) to generate benchmark networks with two levels of embedded communities, a coarse-grained (macro) level and a fine-grained (micro) level. Within each level, a vertex was exclusively assigned to one community. Two parameters,  $\mu_c$  and  $\mu_f$ , were used to define the fractions of edges violating the simulated community structures at the two levels. All other edges were restricted to occur between vertices assigned to the same community (Additional file 1: Figure S1a). We fixed other parameters of the LFR method to values explored by previous studies [9]. In particular,  $N = 1000$  (number of vertices),  $k = 10$  (average degree),  $maxk = 40$  (maximum degree),  $minc = 5$  (minimum number of vertices for a micro-community),  $maxc = 20$  (maximum number of vertices for a micro-community),  $minC = 50$  (minimum number of vertices for a macro-community),  $maxC = 100$  (maximum number of vertices for a macro-community),  $t_1 = 2$  (minus exponent for the degree sequence), and  $t_2 = 1$  (minus exponent for the community size distribution). The numbers of coarse-grained communities and fine-grained communities in each simulated network were approximately bounded by  $minC$ ,  $maxC$ ,  $minc$ , and  $maxc$  (10–20 and 50–200, respectively), and the sizes of communities within each level were set to be close to each other (as  $t_2 = 1$ ).

Some community detection algorithms include iterations of local optimization and vertex aggregation, a process that, like HiDeF, also defines a hierarchy of communities, albeit as a tree rather than a DAG. We demonstrated that without scanning multiple resolutions, this process alone was insufficient to detect the simulated communities at all scales (Additional file 1: Figure S2). We used Louvain and Infomap [46, 47], which have stable implementations and have shown strong performance in previous community detection studies [48]. For Louvain, we optimized the standard Newman-Girvan modularity (equivalent to  $\gamma = 1$ , see above) using the implementation at <http://github.com/vtraag/louvain-igraph> (version 0.6.1). For Infomap, we used the version 1.0.0-beta.47 from <https://www.mapequation.org/>, and set “Markov time” (the “resolution parameter” of Infomap) to 1 and other parameters to default. In general, these settings generated trees with two levels of communities. Note that Infomap sometimes determined that the input network was non-hierarchical, in which cases the coarse- and fine-grained communities were identical by definition.

### Single-cell RNA-seq data

Mouse single-cell RNA-seq data (Fig. 2; Additional file 1: Figure S3) were obtained from the Tabula Muris project [16] (<https://tabula-muris.ds.czbiohub.org/>), which contains two datasets generated with different experimental methods of separating bulk tissues into individual cells: FACS and microfluidic droplet. We applied HiDeF to the shared nearest neighbor graph parsed from the data files (R objects; accessible at <https://doi.org/10.6084/m9.figshare.5821263.v2>) provided in that study. All data normalization and pre-processing procedures have been described in the *Tabula Muris* paper [16]. Briefly, counts were log-normalized using the natural logarithm of 1 + counts per million (for FACS) or 1 + counts per ten thousand (for droplet). A threshold (0.5) for the standardized log dispersion was used to select variable genes. A shared nearest neighbor (SNN) graph was then created by the Seurat *FindNeighbors* function [3] using the first 30 principal components of each dataset. Note that the input of HiDeF is a plain text file representing this graph, and the HiDeF codebase does not depend on the installation of a particular software suite for single-cell analysis. The shared nearest neighbor graph of cells could be generated from scratch or with other parameter combinations of the *FindNeighbors* function in Seurat.

Identical analyses were applied to the FACS and the droplet datasets, respectively, yielding a hierarchy of 273 and 279 communities, respectively (Fig. 2d). ScanPy 1.4.5 [49] was used to create tSNE or UMAP embeddings and associated two-dimensional visualizations [20] as baselines for comparison (Fig. 2e, f; Additional file 1: Figure S3a, b). Through previous analysis of the single-cell RNA data, all cells in these datasets had been annotated with matching cell-type classes in the Cell Ontology (CO) [17]. Before comparing these annotations with the communities detected by HiDeF, we expanded the set of annotations of each cell according to the CO structure, to ensure the set also included all of the ancestor cell types of the type that was annotated. For example, CO has the relationship “[keratinocyte] (is\_a) [epidermal\_cell],” and thus, all cells annotated as “keratinocyte” are also annotated as “epidermal cell.” The CO was obtained from <http://www.obofoundry.org/ontology/cl.html> and processed by the Data Driven Ontology Toolkit (DDOT) [50] retaining “is\_a” relationships only.

We compared HiDeF to TooManyCells [18] and Conos [19] as baseline methods. The former is a divisive method which iteratively applies bipartite spectral clustering to the cell population until the modularity of the partition is below a threshold; the latter uses the Walktrap algorithm to agglomeratively construct the cell-type hierarchy [51]. We chose to compare with these methods because their ability to identify multiscale communities either was the main advertised feature or had been shown to be a major strength. TooManyCells (version 0.2.2.0) was run with the parameter “min-modularity” set to 0.025 as recommended in the original paper [18], with other settings set to default. This process generated dendrograms (binary trees) with 463 communities. The Walktrap algorithm was run from the Conos package (version 1.2.1) with the parameter “step” set to 20 as recommended in the original paper [19], yielding a dendrogram. The *greedyModularityCut* method in the Conos package was used to select  $N$  fusions in the original dendrogram, resulting in a reduced dendrogram with  $2N + 1$  communities (including  $N$  internal and  $N + 1$  leaf nodes). Here, we used  $N = 125$ , generating a hierarchy with 251 communities (Fig. 2c).

The communities in each hierarchy were ranked to analyze the relationships between cell-type recovery and model complexity (Fig. 2c, Additional file 1: Figure S3c). HiDeF communities were ranked by their persistence; Conos and TooManyCells communities were ranked according to the modularity scores those methods associate with each branch-point in their dendrograms. Conos/Walktrap uses a score based on the gain of modularity in merging two communities, whereas TooManyCells uses the modularity of each binary partition.

### Protein-protein interaction networks

We obtained a total of 27 human protein interaction networks gathered previously by survey studies [30, 31], along with one integrated network from budding yeast (*S. cerevisiae*) that had been used in a previous community detection pipeline, NeXO [23]. This collection contained two versions of the STRING interaction database, with the second removing edges from text mining (labeled STRING-t versus STRING, respectively; Fig. 3). Benchmark experiments for the recovery of the Gene Ontology (GO) were performed with STRING and the yeast network (Fig. 3a, b, Additional file 1: Figure S4). The reference GO for yeast proteins was obtained from <http://nexo.ucsd.edu/>. A reference GO for human proteins was downloaded from <http://geneontology.org/> via an API provided by the DDOT package [50].

HiDeF was directly applied to all of the above benchmark networks. The NeXO communities were obtained from <http://nexo.ucsd.edu/>, with a robustness score assigned to each community. To benchmark communities created by hierarchical clustering, we first calculated three versions of pairwise protein distances (HC.1–3; Fig. 3a, b; Additional file 1: Figure S4) using Mashup, DSD, and deepNF [24–26]. Mashup was used to embed each protein as a vector, with 500 and 800 dimensions for yeast and human, as recommended in the original paper. A pairwise distance was computed for each pair of proteins as the cosine distance between the two vectors. Similarly, deepNF was used to embed each protein into a 500-dimensional vector by default. DSD generates pairwise distances by default. Given these pairwise distances, UPGMA clustering was applied to generate binary hierarchical trees. Following the procedure given in the NeXO and Mashup papers [23, 24], communities with  $< 4$  proteins were discarded.

Since all methods had slight differences in the resulting number of communities, communities from each method were sorted in decreasing order of score, enabling comparison of results across the same numbers of top-ranked communities. HiDeF communities were ranked by persistence. NeXO communities were ranked by the robustness value assigned to each community in the original paper [23]. To rank each community  $c$  of hierarchical clustering (branch in the dendrogram), a one-way Mann-Whitney  $U$  test was used to test for significant differences between two sets of protein pairwise distances: (set 1) all pairs consisting of a protein in  $c$  and a protein in the sibling community of  $c$ ; (set 2) all pairs consisting of a protein in each of the two children communities of  $c$ . The communities were sorted by the one-sided  $p$  value of significance that distances in set 1 are greater than those in set 2.

#### Metric for evaluating the performance of multiscale structure identification

We adopted a metric average F1-score [52] to evaluate the overall performance of multiscale structure identification, focusing on the recovery of reference communities. Given a set of reference communities  $C^*$  and a set of computationally detected communities  $\vec{C}$ , the score was defined as:

$$\frac{1}{|C^*|} \sum_{C_i \in C^*} F1(C_i, \vec{C}_{g(i)})$$

where  $g(i)$  is the best match of  $C_i$  in  $\vec{C}$ , defined as follows:

$$g(i) = \operatorname{argmax}_j F1(C_i, \vec{C}_j)$$

and  $F1(C_i, \vec{C}_j)$  is the harmonic mean of  $\text{Precision}(C_i, \vec{C}_j)$  and  $\text{Recall}(C_i, \vec{C}_j)$ . The calculations were conducted by the *xmeasures* package (<https://github.com/eXascaleInfolab/xmeasures>) [53].

#### Combining HiDeF with network embedding

HiDeF was directly applied to the original networks in most of our analyses of protein-protein interaction networks and compared with the results of hierarchical clustering following the network embedding techniques [24, 26]. We sought to explore if we can combine the strength of network embedding and HiDeF to further improve the performance and robustness to parameter choices (Additional file 1: Figure S7). We borrowed the idea of shared nearest neighbor (SNN) graph that we had been using in the analyses of single-cell data. We made a customized script to use the 500-dimensional node embeddings of the STRING network as the input of the Seurat *FindNeighbors* function [3]. The parameters of this function remained as the default. The output SNN graph has  $1.65 \times 10^6$  edges, which is on the same magnitude as the original network ( $2.23 \times 10^6$  edges). We then applied HiDeF to this SNN graph with different combinations of parameters (Additional file 1: Figure S7).

#### Analysis of SARS-COV-2 viral-human protein network

Three hundred thirty-two human proteins identified to interact with SARS-COV-2 viral protein subunits were obtained from a recent study [35]. This list was expanded to include additional human proteins connected to two or more of the 332 virus-interacting

human proteins in the new BioPlex 3.0 network [36]. These operations resulted in a network of 1948 proteins and 22,835 interactions. HiDeF was applied to this network with the same parameter settings as for other protein-protein interaction networks (see the previous “[Methods](#)” sections), and enrichment analysis was performed via g:Profiler [54] (Fig. 3f, g).

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-020-02228-4>.

**Additional file 1: Supplementary figures 1-9.**

**Additional file 2.** Review history.

### Acknowledgements

We are grateful for the helpful discussions with Drs. Jianzhu Ma, Karen Mei, and Daniel Carlin.

### Review history

The review history is available as [Additional file 2](#).

### Peer review information

Andrew Cosgrove was the primary editor on this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

### Authors' contributions

F.Z. and T.I. designed the study and developed the conceptual ideas. F.Z. and S.Z. designed and implemented the main algorithm. F.Z. and D.P. collected the input data and conducted the analysis. S.Z. and F.Z. made the Python package. C.C. developed the server that provides the Cytoscape integration. F.Z. and T.I. wrote the manuscript with contributions from S.Z., D.P., and I.B. The authors read and approved the final manuscript.

### Funding

This work has been supported by the NIH grants to T.I. (R01 HG009979, U54 CA209891) and I.B. (P41 GM103712, P01 DK096990).

### Availability of data and materials

HiDeF is available through CDAPS (Community Detection APplication and Service) [55], which enables simultaneous visualization of the hierarchical model and the underlying network data and is integrated with the Cytoscape visualization and analysis environment. The Cytoscape App can be downloaded at: <http://apps.cytoscape.org/apps/cycommunitydetection>. HiDeF is separately available as a Python package: <https://github.com/fanzheng10/HiDeF>. HiDeF and CyCommunityDetection are licensed under the BSD-3 license. The version of source code (HiDeF v1.0.0) used in this manuscript is deposited in Zenodo [56].

The hierarchical models generated in this study can be obtained as a network collection within the Network Data Exchange (NDEX) database [34, 57]. These models include the hierarchy of murine cell types (Fig. 2), the hierarchies of yeast and human protein communities identified through protein network analysis, and the hierarchy of human protein complexes targeted by SARS-COV2 (Fig. 3).

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Competing interests

T.I. is cofounder of Data4Cure, is on the Scientific Advisory Board, and has an equity interest. T.I. is on the Scientific Advisory Board of Ideaya BioSciences and has an equity interest. The terms of these arrangements have been reviewed and approved by the University of California San Diego, in accordance with its conflict of interest policies.

### Author details

<sup>1</sup>Division of Genetics, Department of Medicine, University of California San Diego, La Jolla, CA 92093, USA.

<sup>2</sup>Department of Computational and Systems Biology, School of Medicine, University of Pittsburgh, Pittsburgh, PA 15213, USA.

Received: 25 June 2020 Accepted: 8 December 2020

Published online: 07 January 2021

### References

1. Regev A, Teichmann SA, Lander ES, Amit I, Benoist C, Birney E, Bodenmiller B, Campbell P, Carninci P, Clatworthy M, et al. The human cell atlas. *Elife*. 2017;6:e27041.

2. Levine JH, Simonds EF, Bendall SC, Davis KL, Amir E-AD, Tadmor MD, Litvin O, Fienberg HG, Jager A, Zunder ER, et al. Data-driven phenotypic dissection of AML reveals progenitor-like cells that correlate with prognosis. *Cell*. 2015;162:184–97.
3. Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol*. 2018;36:411–20.
4. Harold FM. Molecules into cells: specifying spatial architecture. *Microbiol Mol Biol Rev*. 2005;69:544–64.
5. Jain AK, Murty MN, Flynn PJ. Data clustering: a review. *ACM Comput Surv*. 1999;31:264–323.
6. Fortunato S, Hric D. Community detection in networks: a user guide. *Phys Rep*. 2016;659:1–44.
7. Lvd M, Hinton G. Visualizing data using t-SNE. *J Mach Learn Res*. 2008;9:2579–605.
8. Arenas A, Fernández A, Gómez S. Analysis of the structure of complex networks at different resolution levels. *New J Phys*. 2008;10:053039.
9. Traag VA, Krings G, Van Dooren P. Significant scales in community structure. *Sci Rep*. 2013;3:2930.
10. Carlsson G. Topology and data. *Bull Am Math Soc*. 2009;46:255–308.
11. Edelsbrunner H, Harer J. Persistent homology—a survey. *Contemp Math*. 2008;453:257–82.
12. Dabaghian Y, Mémoli F, Frank L, Carlsson G. A topological paradigm for hippocampal spatial map formation using persistent homology. *PLoS Comput Biol*. 2012;8:e1002581.
13. Petri G, Expert P, Turkheimer F, Carhart-Harris R, Nutt D, Hellyer PJ, Vaccarino F. Homological scaffolds of brain functional networks. *J R Soc Interface*. 2014;11:20140873.
14. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res*. 2003;13:2498–504.
15. Lancichinetti A, Fortunato S, Radicchi F. Benchmark graphs for testing community detection algorithms. *Phys Rev E Stat Nonlinear Soft Matter Phys*. 2008;78:046110.
16. Tabula Muris Consortium. Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris. *Nature*. 2018;562:367–72.
17. Diehl AD, Meehan TF, Bradford YM, Brush MH, Dahdul WM, Dougall DS, He Y, Osumi-Sutherland D, Ruttenberg A, Samtiovijai S, et al. The Cell Ontology 2016: enhanced content, modularization, and ontology interoperability. *J Biomed Semantics*. 2016;7:44.
18. Schwartz GW, Zhou Y, Petrovic J, Fasolino M, Xu L, Shaffer SM, Pear WS, Vahedi G, Faryabi RB. TooManyCells identifies and visualizes relationships of single-cell clades. *Nat Methods*. 2020;17:405–13.
19. Barkas N, Petukhov V, Nikolaeva D, Lozinsky Y, Demharter S, Khodosevich K, Kharchenko PV. Joint analysis of heterogeneous single-cell RNA-seq dataset collections. *Nat Methods*. 2019;16:695–98.
20. Becht E, McInnes L, Healy J, Dutertre C-A, IWH K, Ng LG, Ginhoux F, Newell EW. Dimensionality reduction for visualizing single-cell data using UMAP. *Nat Biotechnol*. 2018;37:38–44.
21. Perea G, Navarrete M, Araque A. Tripartite synapses: astrocytes process and control synaptic information. *Trends Neurosci*. 2009;32:421–31.
22. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT, et al. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet*. 2000;25:25–9.
23. Dutkowsky J, Kramer M, Surma MA, Balakrishnan R, Cherry JM, Krogan NJ, Ideker T. A gene ontology inferred from molecular networks. *Nat Biotechnol*. 2013;31:38–45.
24. Cho H, Berger B, Peng J. Compact integration of multi-network topology for functional analysis of genes. *Cell Syst*. 2016;3:540–548.e545.
25. Cao M, Zhang H, Park J, Daniels NM, Crovella ME, Cowen LJ, Hescott B. Going the distance for protein function prediction: a new distance metric for protein interaction networks. *PLoS One*. 2013;8:e76339.
26. Gligorijevic V, Barot M, Bonneau R. deepNF: deep network fusion for protein function prediction. *Bioinformatics*. 2018;34:3873–81.
27. Szklarczyk D, Gable AL, Lyon D, Junge A, Wyder S, Huerta-Cepas J, Simonovic M, Doncheva NT, Morris JH, Bork P, et al. STRING v11: protein–protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Res*. 2019;47:D607–13.
28. Huttlin EL, Bruckner RJ, Paulo JA, Cannon JR, Ting L, Baltier K, Colby G, Gebreab F, Gygi MP, Parzen H, et al. Architecture of the human interactome defines protein communities and disease networks. *Nature*. 2017;545:505–9.
29. Subramanian A, Narayan R, Corsello SM, Peck DD, Natoli TE, Lu X, Gould J, Davis JF, Tubelli AA, Asiedu JK, et al. A next generation connectivity map: L1000 platform and the first 1,000,000 profiles. *Cell*. 2017;171:1437–1452.e1417.
30. Huang JK, Carlin DE, Yu MK, Zhang W, Kreisberg JF, Tamayo P, Ideker T. Systematic evaluation of molecular networks for discovery of disease genes. *Cell Syst*. 2018;6:484–495.e485.
31. Choobdar S, Ahsen ME, Crawford J, Tomasoni M, Fang T, Lamparter D, Lin J, Hescott B, Hu X, Mercer J, et al. Assessment of network module identification across complex diseases. *Nat Methods*. 2019;16:843–52.
32. Greene CS, Krishnan A, Wong AK, Ricciotti E, Zelaya RA, Himmelstein DS, Zhang R, Hartmann BM, Zaslavsky E, Sealfon SC, et al. Understanding multicellular function and disease with human tissue-specific networks. *Nat Genet*. 2015;47:569–76.
33. Cargnello M, Roux PP. Activation and function of the MAPKs and their substrates, the MAPK-activated protein kinases. *Microbiol Mol Biol Rev*. 2012;76:496.
34. Pratt D, Chen J, Pillich R, Rynkov V, Gary A, Demchak B, Ideker T. NDEx 2.0: a clearinghouse for research on cancer pathways. *Cancer Res*. 2017;77:e58–61.
35. Gordon DE, Jang GM, Bouhaddou M, Xu J, Obernier K, White KM, O’Meara MJ, Rezelj WV, Guo JZ, Swaney DL, et al. A SARS-CoV-2 protein interaction map reveals targets for drug repurposing. *Nature*. 2020;583:459–68.
36. Huttlin EL, Bruckner RJ, Navarrete-Perea J, Cannon JR, Baltier K, Gebreab F, Gygi MP, Thornock A, Zarraga G, Tam S, et al. Dual proteome-scale networks reveal cell-specific remodeling of the human interactome. *bioRxiv*. 2020; <https://doi.org/10.1101/2020.01.19.905109>.
37. Snijder EJ, Decroly E, Ziebuhr J. The nonstructural proteins directing coronavirus RNA synthesis and processing. *Coronaviruses*. 2016;96:59–126.
38. Chodaparambil JV, Pate KT, Hepler MRD, Tsai BP, Muthurajan UM, Luger K, Waterman ML, Weis WI. Molecular functions of the TLE tetramerization domain in Wnt target gene repression. *EMBO J*. 2014;33:719–31.

39. Wu C-J, Jan J-T, Chen C-M, Hsieh H-P, Hwang D-R, Liu H-W, Liu C-Y, Huang H-W, Chen S-C, Hong C-F, et al. Inhibition of severe acute respiratory syndrome coronavirus replication by niclosamide. *Antimicrob Agents Chemother*. 2004;48:2693–6.
40. Xu J, Shi P-Y, Li H, Zhou J. Broad spectrum antiviral agent niclosamide and its therapeutic potential. *ACS Infect Dis*. 2020;6:909–15.
41. Newman MEJ, Girvan M. Finding and evaluating community structure in networks. *Phys Rev E Stat Nonlinear Soft Matter Phys*. 2004;69:026113.
42. Reichardt J, Bornholdt S. Statistical mechanics of community detection. *Phys Rev E Stat Nonlinear Soft Matter Phys*. 2006;74:016110.
43. Jeub LGS, Sporns O, Fortunato S. Multiresolution consensus clustering in networks. *Sci Rep*. 2018;8:3259.
44. Oshana R. Chapter 1 - Introduction to digital signal processing. In Robert Oshana eds. *DSP for embedded and real-time systems*; Newnes; 2012. p. 1–14. <https://doi.org/10.1016/B978-0-12-386535-9.00001-9>. <http://www.sciencedirect.com/science/article/pii/B9780123865359000019>.
45. Aho AV, Garey MR, Ullman JD. The transitive reduction of a directed graph. *SIAM J Comput*. 1972;1:131–7.
46. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. *J Stat Mech*. 2008;2008:P10008.
47. Rosvall M, Bergstrom CT. Maps of random walks on complex networks reveal community structure. *Proc Natl Acad Sci U S A*. 2008;105:1118–23.
48. Yang Z, Algesheimer R, Tessone CJ. A comparative analysis of community detection algorithms on artificial networks. *Sci Rep*. 2016;6:30750.
49. Wolf FA, Angerer P, Theis FJ. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol*. 2018;19:15.
50. Yu MK, Ma J, Ono K, Zheng F, Fong SH, Gary A, Chen J, Demchak B, Pratt D, Ideker T. DDOT: a Swiss army knife for investigating data-driven biological ontologies. *Cell Syst*. 2019;8:267–273.e263.
51. Pons P, Latapy M. Computing communities in large networks using random walks. *J Graph Algorithms Appl*. 2006;10:191–218.
52. Yang J, Leskovec J. Overlapping community detection at scale: a nonnegative matrix factorization approach. In: *Proceedings of the sixth ACM international conference on Web search and data mining*; 2013. p. 587–96.
53. Lutov A, Khayati M, Cudré-Mauroux P. Accuracy evaluation of overlapping and multi-resolution clustering algorithms on large datasets. *IEEE International Conference on Big Data and Smart Computing (BigComp)*, Kyoto, Japan, 2019, p. 1–8.
54. Raudvere U, Kolberg L, Kuzmin I, Arak T, Adler P, Peterson H, Vilo J. g:Profiler: a web server for functional enrichment analysis and conversions of gene lists (2019 update). *Nucleic Acids Res*. 2019;47:W191–8.
55. Singhal A, Cao S, Churas C, Pratt D, Fortunato S, Zheng F, Ideker T. Multiscale community detection in Cytoscape. *PLoS Comput Biol*. 2020;16:e1008239.
56. Zheng F, Zhang S: HiDeF v1.0.0. Zenodo. <https://doi.org/10.5281/zenodo.4059074>. 2020.
57. Zheng F: HiDeF collection on NDEx, <https://doi.org/10.18119/N9ZP58>. 2020.

## Publisher's Note

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

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

