
Internal Report num. 13369

Evaluation of the SDM-Open framework using benchmark datasets

Jožef Stefan Institute

Version FINAL

Abstract: Within Task 1.3 we aim to evaluate the developed SDM-Open framework using benchmark datasets. Using extensive evaluation on benchmark datasets (described below), we selected the best network mining and the best SDM methods to be included in the SDM-Open methodology.

Document administrative information	
Project acronym:	SDM-Open
Project number:	N2-0078
Deliverable number:	D1.3
Deliverable full title:	Evaluation of the SDM-Open framework using benchmark datasets
Document identifier:	SDM-Open D1.3-v0
Lead partner short name:	JSI
Report version:	FINAL
Report preparation date:	20/11/2020
Lead author:	Blaž Škrlič
Co-authors:	Matej Martinc, Jan Kralj, Nada Lavrač, Senja Pollak
Status:	Final

Introduction

The goal of WP1 of the SDM-OPEN project is to develop semantic data mining (SDM) methods to learn patterns from linked data fusing network analysis with data mining into a single, flexible framework to discover explanations from linked best describing input data. Moreover, the aim is to develop the SDM-Open framework to support experts using the combined methodology in support of Linked Open Data Science analytic processes.

Within Task 1.3 we aimed to evaluate the developed SDM-Open framework using benchmark datasets. Using extensive evaluation on benchmark datasets, we selected the best network mining and the best SDM methods to be included in the SDM-Open methodology.

We have developed Python library Py3plex for the visualization and analysis of multilayer networks, and performed benchmarking experiments. This work is published in the paper entitled "*Py3plex toolkit for visualization and analysis of multilayer networks*" and the benchmarking experiments are covered in the Section Experimental evaluation. The full publication text is added to this report.

We have also revisited the notion of layer entanglement and extended it to coupled multilayer networks and temporal networks and presented it in the paper "*Layer entanglement in multiplex, temporal multiplex, and coupled multilayer networks*". To investigate entanglement, we have proposed a random generator for coupled multilayer networks, and generated a large set of synthetic ones. We have evaluated entanglement intensity and homogeneity in all cases, and compared to static and temporal real world networks.

We hypothesised that different properties, which can be related to the very nature of the data they model (or to events in time-dependent data), may be reflected in the way layers are

intertwined. In the paper, we investigated these through the prism of layer entanglement in coupled multilayer networks. We tested over 30 real-life networks in six different disciplines (social, genetic, transport, co-authorship, trade, and neuronal networks). The full publication text is added to this report.

RESEARCH

Open Access

Py3plex toolkit for visualization and analysis of multilayer networks



Blaž Škrlić^{1,2*} , Jan Kralj² and Nada Lavrač^{1,2}

*Correspondence: blaz.skrlic@ijs.si

¹Jožef Stefan International Postgraduate School, Jamova 39, 1000 Ljubljana, Slovenia
²Jožef Stefan Institute, Jamova 39, 1000 Ljubljana, Slovenia

Abstract

Complex networks are used as means for representing multimodal, real-life systems. With increasing amounts of data that lead to large multilayer networks consisting of different node and edge types, that can also be subject to temporal change, there is an increasing need for versatile visualization and analysis software. This work presents a lightweight Python library, Py3plex, which focuses on the visualization and analysis of multilayer networks. The library implements a set of simple graphical primitives supporting intra- as well as inter-layer visualization. It also supports many common operations on multilayer networks, such as aggregation, slicing, indexing, traversal, and more. The paper also focuses on how node embeddings can be used to speed up contemporary (multilayer) layout computation. The library's functionality is showcased on both real and synthetic networks.

Keywords: Multilayer networks, Network visualization, Complex systems, Network embedding

Introduction

Analysis and visualization of complex networks offers novel opportunities to study intractable systems, such as protein interaction networks, transportation networks or social networks (Boccaletti et al. 2014; Wang et al. 2015; Pavlopoulos et al. 2008). As this vibrant research field offers novel tools at an increasing pace, development of freely available, scalable software resources is becoming a relevant research direction. Despite many existing tools for analysis and visualization of homogeneous networks, i.e. networks with only a single node type and non-annotated edges, tools that consider multilayer networks—also referred to as heterogeneous networks—are an active research area. In multilayer networks, many different types of annotations are taken into account, including relation-labeled edges and multiple node types. Such networks are considered, for example, when multiple layers of biological information (e.g., protein-protein, gene-gene interactions etc.) are available and need to be taken into account when studying diseases. We consider networks as multilayer when they contain at least two types of nodes or at least two types of edges.

This work is inspired by previous studies on multilayer networks formalized by Kiela et al. (2014). The paper first presents Py3plex, a Python library for analysis and visualization of heterogeneous (and homogeneous) networks. The visualization suite simplifies displaying of multilayered networks and network communities as well as network

embeddings (Grover and Leskovec 2016). Py3plex also implements a state-of-the-art procedure for converting heterogeneous networks to homogeneous networks (Kralj et al. 2018), a set of subroutines for partition enrichment (i.e. the process of learning qualitative explanations relating individual partitions, e.g., communities) using expert-curated domain knowledge Škrlić et al. (2018, 2019), and offers intuitive visualization of network-topological properties, such as the community structure.

In this paper, we also demonstrate the performance of embedding based network layout, where state-of-the-art network embedding algorithms—i.e., algorithms which map a given network to a predefined vector space—are used to obtain node coordinates in two dimensions. We show that this method is notably faster on larger networks, and can serve as a relevant improvement of contemporary force-directed layout computation. Finally, we explore how multiplex dynamic social networks can be visualized using the presented visualization technique. Such networks consist of the same set of nodes projected across different contexts. For example, behavior of users can be monitored across the social media they participate in (e.g., Twitter, Facebook etc.). We conclude with a discussion regarding both positive and negative aspects of the presented library along with suggestions for the further work.

This paper significantly extends the work on multilayer network visualization and analysis presented in our previous conference publication (Škrlić et al. 2019). First, the related work section (“[Related work](#)” section) has been extended by including tools Pajek (Batagelj and Mrvar 2001) and Tulip (Auber 2004; Auber et al. 2017) that inspired the creation of Py3plex. Next, we extended the description of the library’s features (“[Key features](#)” section).

We discuss that—apart from simple statistical analysis—the library also offers the functionality to learn from networks using some of the recently introduced machine learning approaches. Further, we explore in more detail the functionality related to community enrichment, as it offers the functionality to explain a given network’s partitioning using symbolic learning. Next, we added a section which explores how network embeddings can be used to construct network layout (“[Embedding-based network layout](#)” section). In this section, we first discuss relevant machine learning methodology and how it relates to two distinct steps in the layout construction. We next demonstrate how the presented layout compares to efficient Barnes-Hut-based minimization. As a demonstration of novel functionality, we added an example where a multiplex dynamic social network is visualized (see “[Experimental evaluation](#)” section). Finally, we extended the discussion and conclusions (“[Conclusions and further work](#)” section), where we discuss some of the current limitations, and the existing uses of the library.

Related work

This section presents the state-of-the-art network analysis libraries, relevant to the development of Py3plex. The most common approaches to network analysis can be split into two groups: GUI-based solutions and API-based solutions.

The most used GUI-based solutions include Cytoscape (Shannon 2003), Gephi (Bastian et al. 2009), and Pajek (Batagelj and Mrvar 2001). Cytoscape (Shannon 2003) is one of the largest network analysis projects to date. It supports custom manipulation of the loaded network, and is hence flexible both in terms of network visualization as well as analysis. The Gephi (Bastian et al. 2009) suite offers a similar set of functionalities, but it is known

for better visualization capabilities. Similarly, Pajek (Batagelj and Mrvar 2001) is used to analyse simple graphs, and was shown to scale well to large social networks. These solutions are mostly used in the final step of a network analysis project, where pre-computed node properties are used as part of the input. The tools are constantly updated with novel functionality, offering many graph analysis algorithms out-of-the-box.

The API-based solutions, which provide programmatic access from popular languages (such as Python, R, JavaScript, C++), are preferred when the entire network analysis and visualization is performed in the same environment. The NetworkX library (Hagberg et al. 2008) is prevalent for the Python environment, while the igraph library (Nepusz and Csárdi 2006) is commonly used among the R users. C, and C++ alternatives include SNAP (Leskovec and Sosič 2016), Boost Graph Library (BGL) (The Boost Graph Library 2002), and Tulip (Auber 2004; Auber et al. 2017). Compared to GUI-based analysis, API-based approaches result in a series of high-level function calls, which generate the desired output. Such approaches are commonly used for analysis when either the considered networks are large or the number of networks under consideration is high.

The above approaches focus on homogeneous networks consisting of single node (and edge) types. However, recent advances in multilayer network analysis have proven that the additional information associated with node and edge type provides insights regarding network structure and dynamics. Multilayer networks can be represented as higher order tensors (De Domenico et al. 2013), encoding inter- as well as intra-layer connections of varying intensity. In this paper, we focus on state-of-the-art implementations of this formalism, their functionality, and some of their drawbacks. The currently used libraries for visualization and analysis of multilayer networks include (1) libraries for the Python environment that include Pymnet¹ (Kivelä et al. 2014) and MultinetX² (Amato et al. 2017), as well as (2) a library for the R environment Muxviz³ (De Domenico et al. 2015). Detailed analysis of these approaches is presented in “[Comparison of multilayer network libraries](#)” section, where they are compared to the presented Py3plex library.

Py3plex library architecture

This section explains the presented Py3plex library’s architecture, followed by the description of individual components and subroutines.

Module organization

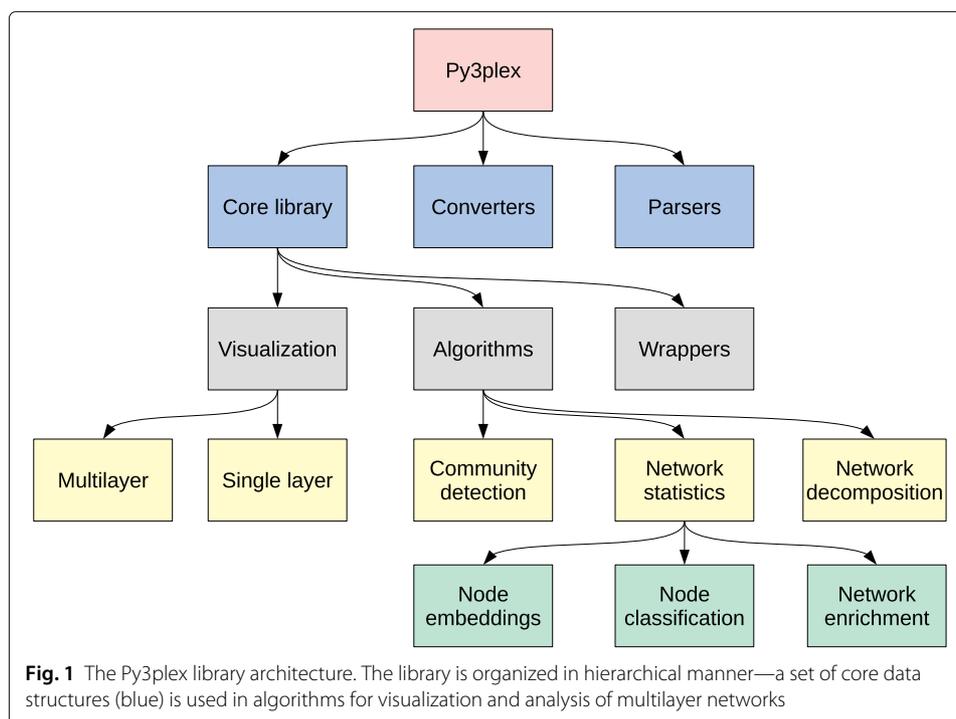
A high-level organization of the Py3plex library is shown in Fig. 1. The library includes methods for parsing and converting graphs from and to various formats, while the core library includes three modules:

- **The visualization module** consists of different subroutines used for network visualization of multilayer and single layer networks. Detailed description of the in-house developed visualization is given in “[Py3plex multilayer network visualization](#)” section.
- **The wrappers module** As many procedures are given as standalone executables, Py3plex offers a set of wrapper subroutines, useful for calling external state-of-the-art algorithms, for example, highly optimized network embedding routines (Grover and

¹<http://www.mkivela.com/pymnet/>

²<https://github.com/nkoub/multinetx>

³<http://muxviz.net/>



Leskovec 2016) as well as the InfoMap community detection algorithm (Rosvall et al. 2009). The methods offer easy access to multiplex community detection capabilities of Infomap, as well as community computation directly on the supra-adjacency matrix of a given network.

- **The algorithms module** includes implementations of many commonly used algorithms, such as Louvain community detection (Blondel et al. 2008), node ranking, network statistics, the recently introduced network topology enrichment (Škrlić et al. 2018), and network node classification and network embeddings.

All the modules are built in an extensible manner, allowing for new routines to be easily added. As Py3plex is built on top of the NetworkX (Hagberg et al. 2008), network operations can include any of the methods primarily designed for homogeneous networks. This functionality provides flexibility when implementing novel multilayer network analysis algorithms.

Key features

As Py3plex consists of multiple building blocks, we next describe key building blocks comprising the library. The set of data structures underlying all aforementioned modules offers intuitive access to manipulation, creation, and visualization of multilayer networks. Additionally, users of the library can use this core set of functionalities to implement their own algorithms. Other modules comprising Py3plex also build on top of the “multilayer network” object, a data structure we introduce for easier manipulation and consistency throughout the library.

The visualization module includes graphical primitives, needed to construct multi- and single-layered visualizations of considered networks. The module is built so that it accepts the core network structure (described in the previous section) as input and subsequently

outputs a visualization of choice. Detailed formulation of the supported visualizations is given in the following sections.

A substantial portion of Py3plex's functionality is devoted to analysis of multilayer networks. Here, functionality such as network aggregation, centrality computation, and similar, offers straightforward analysis of the network's properties. This part of the library also includes some of the state-of-the-art approaches for network decomposition and enrichments described in the following section.

Comparison of multilayer network libraries

In this section, we compare the functionality of different multilayer network analysis and visualization libraries. Where relevant, we emphasize the novelties Py3plex offers and how they compare to existing alternatives⁴. Each library offers some functionality that is not available in other libraries, hence no library is exhaustive. We evaluate different aspects of functionality, ranging from their computational efficiency to the number of various analysis functions provided. Table 1 presents the comparison of multilayer network analysis and visualization libraries, of which functionality is discussed in this section.

To ensure sufficient speed of execution, Pymnet (Kivelä et al. 2014), Py3plex and MultinetX (Amato et al. 2017) include subroutine implementations in C (via Cython), but also Numpy (Walt et al. 2011) and Scipy (Jones et al. 2001) libraries for optimized linear algebra-based computation. The R-based MuxViz (De Domenico et al. 2015) uses the Octave library for efficient array-based operations.

Table 1 Comparison of multilayer network analysis and visualization libraries

	Py3plex	Pymnet (Kivelä et al. 2014)	MuxViz (De Domenico et al. 2015)	MultinetX (Amato et al. 2017)
Core features				
Programming language	Python 3 and C (via Numpy and Cython)	Python 3 and C (via Numpy)	R	Python 3 and C (via Numpy)
Basic statistics	✓	✓	✓	✓
Visualization of large networks	✓	-	✓	-
Visualization in 3D	-	✓	-	✓
Aggregation/decomposition	✓	✓	✓	-
Random graph generators	✓	✓	✓	✓
Adjacency matrix manipulation	✓	✓	✓	✓
[3pt] Additional features				
[3pt] Node classification	✓	-	-	-
Isomorphisms	✓	✓	-	-
Community detection	✓	-	✓	-
GUI version	-	-	✓	-
Tensor manipulation	✓	✓	✓	✓
Node ranking	✓	✓	✓	✓
Semantic topology enrichment	✓	-	-	-
Temporal networks	-	-	✓	✓
Network embedding	✓	-	-	-

⁴For a comprehensive overview of visualization tools, we refer the reader to the recent survey (McGee et al. 2019).

In terms of visualization, all packages compared, include at least some form of network visualization. Apart from 2D layouts, The MuxViz and Pymnet libraries also support 3D layouts. As discussed in “[Py3plex multilayer network visualization](#)” section, Py3plex offers a new approach to multilayer network visualization, whereas MultinetX, for example, is capable of displaying the supra-adjacency matrix (De Domenico et al. 2013) of the network. The multitude of different multilayer network visualization applications indicates that no single type of visualization is optimal for a given task, as each visualization is optimal for a particular range of network size and density.

Apart from MultinetX, all other libraries compared above include different methods for network aggregation and decomposition. A similar compendium of functionality is offered by Pymnet and MuxViz. All packages offer intuitive access to a network’s supra-adjacency matrix, and hence provide many opportunities for implementation of computationally efficient multilayer network analysis algorithms. However, Py3plex also leverages full functionality of the recently developed HINMINE (Kralj et al. 2018) methodology for heterogeneous network decomposition and network node classification. As standard aggregation schemes (De Domenico et al. 2015) yield a homogeneous network, for example, by summing over edges in individual layers, HINMINE-based aggregation is based on frequency of relation-annotated directed paths (of length two) between the target node type. Compared to standard aggregation, HINMINE is thus suitable for situations where domain knowledge was used for annotating the network edges. Additionally, Py3plex also includes basic aggregation sub-routines, as well as supports classification using Personalized PageRank-based feature vectors is also supported (Kralj et al. 2018). An interested reader can find the details in Appendix A.

All of the libraries listed above include methods for obtaining quick insights into a network’s structure. The Pymnet and Py3plex are currently the only libraries supporting different isomorphism algorithms for evaluating network similarity.

In terms of other functionalities, we observe the following. MuxViz also supports GUI-based analysis solutions, as well as API-based ones. Once installed, it can be executed locally within the browser as standalone software.

We observe Pymnet offers one of the most intuitive API interfaces for manipulation of tensor representations of multilayer networks, such as slicing, indexing etc. In comparison, MultinetX and MuxViz offer similar functionality, whereas Py3plex operates on attribute-rich list-based representations of a network and is not necessarily suitable for all multilayer slicing tasks.

Aside from network analysis and visualization capabilities, Py3plex is the only library which also supports various forms of learning from multilayer networks. It provides the methodology for semantic enrichment of complex networks. The main objective in this emerging field is to associate previously known domain knowledge with topological structures of a network. For example, the functional characterization of a network’s communities can only be assessed using curated domain knowledge. We refer to the use of such knowledge to understand network-topological properties as network enrichment.

The supported network enrichment, following the Community-Based Semantic Subgroup Discovery (CBSSD) methodology (Škrlić et al. 2018), can be described in two steps:

- 1 Partition detection. The network is partitioned into separate subnetworks, which commonly represent some form of functional similarity.
- 2 Enrichment. Individual communities are compared to the remainder of the network by using domain knowledge in the form of ontologies. Should a given ontology term (concept) be over- or under-represented in a considered community, the community shall be considered enriched with respect to this term.

The end result of such analysis are thus community-term pairs, sorted by e.g., term significance. Currently, Py3plex is the only library that supports both rule-based enrichment as well as standard Fisher's exact test-based enrichment, both introduced as part of the CBSSD methodology.

The current implementation of Py3plex contains wrapper routines for widely used network embedding algorithms, such as Node2vec and similar (Grover and Leskovec 2016). Such functionality is currently not offered in any other libraries.

Overall, Py3plex offers novel algorithms for understanding network structure, an intuitive interface for manipulation of multilayer networks as well as network decomposition and aggregation. However, the primary features of the current version of Py3plex are primarily network visualization and spatially efficient network manipulation (linear in terms of nodes and edges). This functionality offers additional state-of-the-art performance, not observed in the other libraries. Below, we first present the proposed network visualization ("Py3plex multilayer network visualization" section), followed by embedding-based node layout computation ("Embedding-based network layout" section), and finally on empirical evaluation of the proposed approaches on a set of artificial random multilayer networks, as well as real world biological and social networks ("Experimental evaluation" section).

Py3plex multilayer network visualization

One of the key contributions of this paper is a novel method for visualization of multilayer networks. In this section, we showcase the proposed visualization method and its implementation. Next, we evaluate the method's performance on a series of random networks, where we subsequently compare the results with Pymnet's network visualization capabilities.

Visualization methodology and implementation

The proposed set of visualization techniques aims to address the issue of visualizing multilayer networks. The implementation of the presented layout algorithm operates in three main steps described below.

Intra-layer layout computation. First, subnetworks consisting of same-typed nodes and edges between them are considered. For each node type, we compute a fast, force-directed layout on the single-type subnetwork. This results in (x, y) coordinate pairs for individual nodes that are scaled so that both coordinates are between 0 and 1. Individual, single-layer networks are derived from the existing NetworkX graph library (Hagberg et al. 2008), consequently offering full NetworkX functionality. Additionally, Py3plex provides the means to customize the majority of network properties, including edge shapes and colors, individual layer layouts, node sizes and colors, and overall network organization. Computationally expensive operations are vectorized using the Numpy numeric library (Walt et al. 2011). The Force Atlas 2

algorithm (based on Barnes-Hut n -body minimization), which was transpiled from Python to C, is used as the default option during visualization (Jacomy et al. 2014). Throughout the work, we refer to such layouts as spring layout plots.

Multilayer network drawing. In the second step, nodes are drawn along a diagonal line based on their type (i.e. layer they belong to). This is achieved by calculating the actual coordinates of each point by separating each consecutive layer from the preceding layer by exactly a single coordinate unit on both axes. The result is that, for a node in the i -th layer, where we calculated coordinates (x, y) in step 1, we now update the coordinates to $(x + i, y + i)$. As nodes in different layers are now separated by at least one coordinate unit, each subnetwork corresponding to different node type is drawn in a separate region with no overlap between different layers. Each layer includes a network with its own local organization, independent of the inter-layer connections.

Inter-layer edge drawing. The final step involves drawing of inter-layer edges, which are represented as arcs connecting different nodes. One of the main problems we faced during the implementation of this step was edge positioning and parameterization, as there are many different ways of drawing an arc between two nodes. We considered three possible scenarios in terms of inter-layer edge drawing:

- 1 The edges are only on the upper part of the layer diagonal;
- 2 The edges are only on the bottom part of the layer diagonal;
- 3 The edges are on both sides of the diagonal projection.

An inter-layer edge between nodes n_1 and n_2 of the network (represented by points (x_1, y_1) and (x_2, y_2) , respectively) is drawn as follows. First, an artificially introduced point (x_3, y_3) is calculated by taking the midpoint between (x_1, y_1) and (x_2, y_2) and scaling its y coordinate by a factor of τ (a parameter of the method). Scaling the y coordinate by a fixed amount (τ) instead of shifting it ensures that each edge is, in effect, shifted up by a different value, allowing all edges to be visible in the final image. Changing the value of parameter τ offers simple manipulation of inter-layer edges as follows:

- 1 when $\tau > 1$, the arc will be located above the diagonal;
- 2 when $\tau = 1$, the arc will be a line between the two nodes (a third point on $f(x)$);
- 3 when $\tau < 1$, the arc will be located below the diagonal.

Once n_1 , n_2 and n_3 are obtained, points n_1 and n_2 are connected by a parabolic arc that passes through all three points. Even though the current implementation constructs inter-layer edges using interpolation over three points, Py3plex supports adding an arbitrary number of intermediary points, in which case cubic arc interpolation is used to produce the line between n_1 and n_2 . In our experiments, however, we show that even a single intermediary point allows for clear visualization.

We further propose a heuristic for automatically determining an arc's position, i.e. whether an arc should be located above or below the diagonal. We achieve this as follows. Given n_1 and n_2 as defined in the previous paragraph, the arc is located above the diagonal if and only if $\text{int}(n_3) > n_3$, where $\text{int}(n_3)$ represents the integer representation of the coordinates of n_3 . Integer-based rounding works, as layers are separated by exactly a single coordinate unit. All operations for arc computation, scaling and transformation are vectorized for better performance. Should the network

appear incomprehensible, natural logarithms are applied to node representations—filled circles based on individual node degrees—which simplifies visualization of denser networks.

Py3plex can easily visualize more than ten layers with tens of thousands of nodes. Compared to existing solutions, diagonal projection of multiple layers enables visualization using standard layout algorithms, with additional specification of inter-layer edges. An example visualization using the presented Py3plex library is shown in Fig. 2, with an alternative visualization using a single layer force-directed layout of the whole network is included for comparison.

Strengths and potential drawbacks

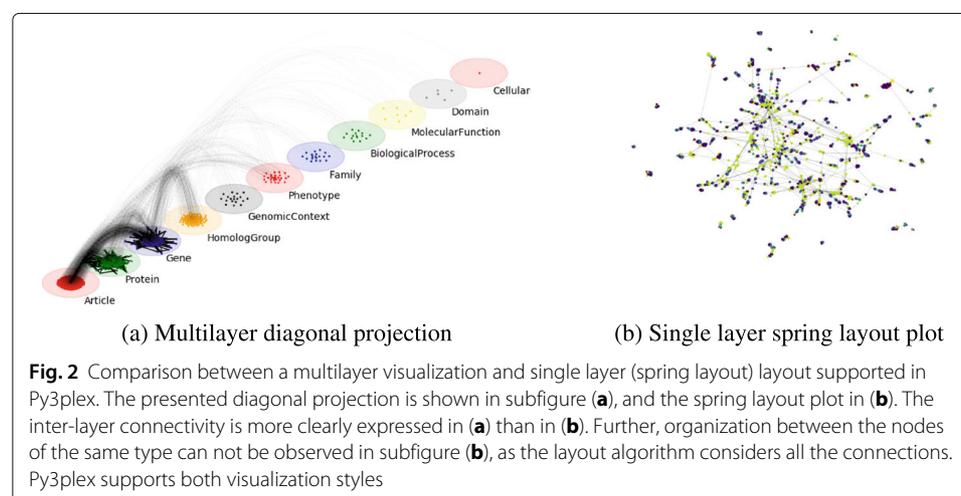
In this section, we discuss the strengths and weaknesses of the proposed visualization. We begin by describing the strengths and continue with the discussion of the cognitive load and other potential drawbacks.

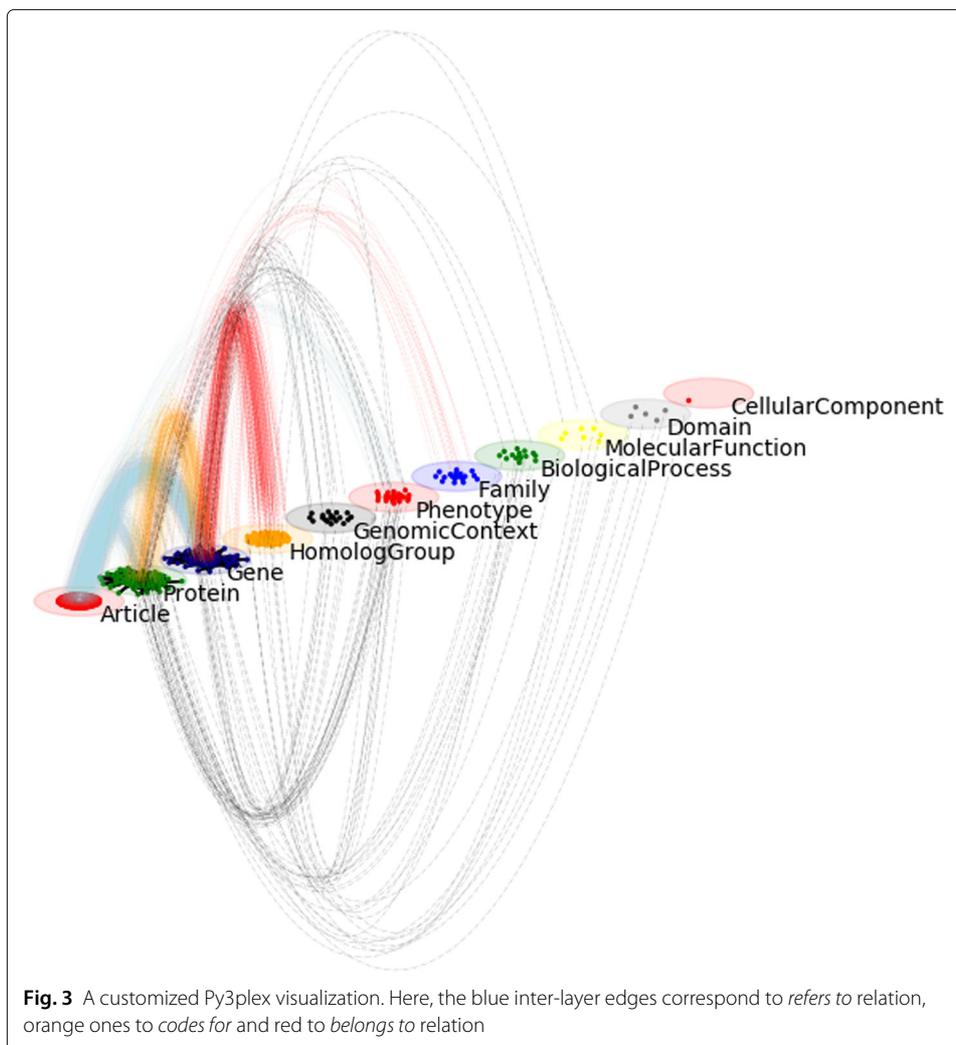
Strengths

Many aspects of the visualization presented in this paper can be customized to emphasize either the node or layer properties. For example, in Fig. 3 we colored differently the inter-layer edges corresponding to specific relations. Additionally, such edges can be plotted either on the upper or the lower side of the diagonal containing networks, making it possible for the user to emphasize only the selected inter-layer edges. The height, types of lines, colors and transparency can be fine-tuned to the user's preferences. The colors of intra-layer edges and nodes can also be customized—for example, special nodes or sets of nodes can be colored to emphasize the intra-layer structure.

Potential drawbacks

In this section, we also discuss some of the drawbacks the presented approach introduces, especially when considering larger networks. One of the main problems with such complex visualizations is the amount of overlapping edges, which was shown to be problematic for the user experience by Purchase (1997). We split the following discussion into two main parts: intra-layer overdraw and inter-layer overdraw.





Intra-layer edge overdraw. Non-planar, real-world networks are impossible to draw without some edge overlap. The problem with drawing networks within layers is thus inherent to all other libraries which visualize such network. Py3plex tackles this issue by enhancing transparency of edges based on density, edge thresholding, as well as controlling edge widths. Further, node sizes also take up substantial amounts of layout space, thus need to be adapted accordingly. Even though the space in which individual intra-layer networks are drawn is smaller than the whole canvas, the Force Atlas 2 algorithm which is used for intra-layer layout computation, disperses the nodes based on their connectivity patterns. This way, it partially separates the densely connected clusters (some arc overlap is reduced this way). To preserve topological properties of intra-layer networks, some overlap (e.g., within functional clusters) is inevitable.

Inter-layer edge overdraw. In our visualization, the majority of the inter-layer overlaps are noticed at the upper-most (or bottom-most) parts of the parabolic arcs. Techniques for emphasizing edges (e.g., transparency-based filtering etc.) that can be adopted to further emphasize individual edges are discussed in “Strengths” section. Additionally, as the presented inter-layer edges span between

layers on the upper and the lower part of the main diagonal, some of the edge overlap can be reduced by redistributing the edges accordingly across the empty regions of the canvas. Such positioning can be automatically determined by Py3plex. Next, heights of individual arcs can be manually configured, enabling definition of custom, less overlapping groups of arcs. While presented solutions do not entirely solve the problem, we believe that for non-planar graphs (especially in 2D), if the user knows what aspect to emphasize, the proposed solution can provide sensible visualizations. While separating the layers may incur more overall edge overdrawing, we believe that this cost is outweighed by the benefits (i.e. added visual clarity offered by visually separated layers) of our visualization.

Finally, we discuss how intra- as well as inter-layer edges contribute to understanding of the plots. Should the number of inter-layer edges increase in layers that are very close together, such setting is prone to cognitive overload and can be addressed by specifying a different layer order. Very sparse networks with many layers are also harder to visualize using the presented visualization, as with many layers, information regarding connectivity can become harder to comprehend—the inter-layer edges can span across larger regions of canvas and are harder to follow.

Cognitive load

In this section we discuss the potential implications of the presented methodology with respect to cognitive load, as this aspect of visualization can be critical in determining the usability of a given visualization method. This section follows guidelines from Huang et al. (2009), who investigated how complex networks remain understandable to a non-expert human observer. Via a variety of cognitive tasks (such as triangle counting), Huang et al. showed that networks with only tens of nodes and hundreds of edges can already pose a problem when it comes to their interpretation and understanding. While the work by Huang et al. is not focused on multilayer networks, visualization properties they recognize as relevant are also applicable when visualizing multilayer networks.

Huang et al. (2009) identify several key factors of cognitive load presented by a given network visualization. The factors that are most relevant for the following discussion regarding cognitive load of Py3plex plots are the following:

Domain complexity. Not all domains are equally complex. Huang et al. point out that visualizations of biological networks should differ from the ones used for displaying social networks. Py3plex is adapted in line with these findings. The layer-level diagonal visualization introduced in this work offers intuitive segmentation of e.g., biological information (e.g., DNA, RNA, protein etc.); however, such levels are not necessarily present in social networks. In order to address this issue, Py3plex offers functionality to aggregate layers, which can be adapted to specific use-cases.

Data complexity. This aspect is closely related to the studied domain. We observed for example, that biological networks contain more node and edge types, than the social networks, requiring different visualization strategies. The internal data structure used for visualization and manipulation is a heterogeneous information network with (optional) attributes assigned to nodes. This structure was expressive enough for the examples shown in this work, consisting of multiple node and edge types. It could be further adapted to e.g., hypergraphs, should the need arise.

Visual complexity. Visualizations can have varying degrees of complexity. Here, aspects such as edge or node overlaps and the number of different elements visualized need to be considered. The complexity of the visualizations obtained using Py3plex can be high, as it displays multiple layers along with inter- and intra-layer edges. We refer the reader to “[Potential drawbacks](#)” section for detailed discussion of the visualization aspects which influence the final output the most.

Interactivity

Because one of the output options supported by Py3plex is a Matplotlib canvas (Hunter 2007), the resulting visualizations can easily be:

- Zoomed-into. A square region of the visualization is selected and zoomed-into. This functionality offers e.g., a way to emphasize only certain layers of interest.
- Stretched. The interactive viewer offers simple functionality for adapting the shape of the resulting visualization, offering fine-tuning with respect to e.g., overlapping text.
- Animated. Matplotlib offers animation functionality, making possible the construction of e.g., dynamic visualizations, consisting of multiple (e.g., time-dependent) frames. An example of such an animation is available online in .gif format⁵.

Embedding-based network layout

Visualization of large networks commonly results in long computation times and incomprehensible layouts. Recent advancements in the field of machine learning on graphs can be leveraged to facilitate the process of network visualization. Even though embedding-based data visualization is becoming commonplace in contemporary machine learning, such techniques span back to Harel and Koren (2002), who initially investigated how network embeddings can be used for visualization. They use principal component analysis (PCA) as the main embedding mechanism. Even though PCA can offer valuable insights when projecting the data into orthogonal space of lower dimension, it does not necessarily maintain all network-topological properties which represent key parts of the considered network. The initially considered network embedding (PCA) does not take into account higher-order node neighborhoods, missing out on e.g., densely connected parts of networks that can only be accessible when considering longer random walks.

This section is structured as follows. First, we describe the role of network embedding algorithms in a standard machine learning setting. Next, we show how methods for non-linear dimensionality reduction can be combined with scalable node embeddings for network layout construction. Finally, we present a simple benchmark of the presented layout algorithm compared with a generic, force-directed layout.

Network embedding

Recent advancements in learning from complex networks commonly consist of two main steps: network embedding and learning. Recent approaches for embedding construction include DeepWalk (Perozzi et al. 2014), Node2vec (Grover and Leskovec 2016), Struc2vec (Ribeiro et al. 2017), and similar approaches, all of which attempt to capture node information and encode it in the form of d -dimensional vectors. In this work we are interested

⁵https://github.com/SkBlaz/Py3plex/raw/master/example_images/animation.gif

in the embedding phase of these algorithms, as well as the use of resulting node embeddings as the first step in the presented layout calculation algorithm. The feature matrix (table) can be used with less additional preprocessing compared to sophisticated relational nature of a graph. The presented visualization approach first constructs such an embedding, and subsequently projects it to a two-dimensional vector space; this space of (x, y) pairs represents initial node coordinates. Schematic representation of this idea is shown in Fig. 4.

We next describe some of the key steps of `node2vec`, as recently given in (Kralj et al. 2019). We believe understanding of how `node2vec` operates shall offer the reader intuition as to why use the selected embedding method as a building block of the proposed visualization.

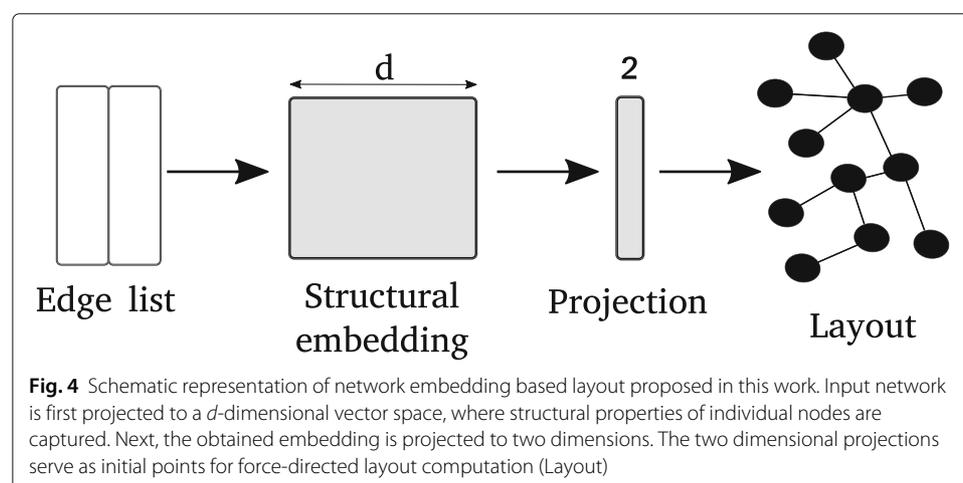
Network embedding using `node2vec`

A recently developed approach to vectorizing network nodes is the `node2vec` algorithm (Grover and Leskovec 2016), which uses the random walks to calculate features that express similarities between pairs of nodes.

The `node2vec` algorithm takes as input a network of n nodes, represented as a graph $G = (V, E)$ where V is the set of network nodes and E is the set of connections, or edges, in the network. The algorithm returns a matrix $f \in \mathbb{R}^{|V| \times d}$ with a pre-defined number of columns d . Matrix f is interpreted as a collection of d -dimensional feature vectors with the i -th row of the matrix corresponding to the feature vector of the i -th node in the network. We write $f(u)$ to mean the row of matrix f , corresponding to node u . The goal of the algorithm is to construct feature vectors $f(u)$ in such a way that the feature vectors of all nodes that share a certain neighborhood will be similar. Matrix f is calculated as follows:

$$\mathcal{E}(G) = \arg \text{Max}_{f \in \mathbb{R}^{|V| \times d}} \sum_{u \in V} \left(-\log \left(\sum_{v \in V} e^{f(u) \cdot f(v)} \right) + \sum_{n_i \in N_S(u)} f(n_i) \cdot f(u) \right) \quad (1)$$

where $N(u)$ denotes the network neighborhood of node u given a *sampling strategy*, and \mathcal{E} the embedding constructor (`node2vec` in this case). In the above expression the inner sum calculates the similarities between a node and all nodes in its neighborhood. This



sum is large if the feature vectors of nodes in the same neighborhood are collinear, however it also increases if feature vectors of nodes have a large norm. The first value of each summand decreases when the norms of feature vectors increase, thereby penalizing collections of feature vectors with large norms.

Expression (1) has a probabilistic interpretation which models a process of randomly selecting nodes from the network. The probability $P(n|u)$ of node n following node u in the selection process is proportional to $e^{f(n) \cdot f(u)}$. Assuming that selecting one node is independent from selecting any other node, we can calculate the probability of selecting all nodes from a given set A as $P(A|u) = \prod_{n \in A} P(n|u)$, and Eq. 1 can then be rewritten as follows:

$$\mathcal{E}(G) = \arg \text{Max}_{f \in \mathbb{R}^{|V| \times d}} \sum_{u \in V} \log(P(N_S(u)|f(u))). \quad (2)$$

Term $N_S(u)$ in Eqs. (1) and (2) denotes the neighborhood of u given a sampling strategy S and is calculated by simulating a random walker traversing the network starting at node u . The transition probabilities for traversing from node n_1 to node n_2 depends on the node n_0 the walker visited before node n_1 , making the process of traversing the network a second order random walk. The unnormalized transition probabilities are set using two parameters, p and q , and are equal to:

$$P(n_2 | \text{previous step moved from node } n_0 \text{ to } n_1) = \begin{cases} \frac{1}{p} & \text{if } n_2 = n_0 \\ 1 & \text{if } n_2 \text{ can be reached from } n_1 \\ \frac{1}{q} & \text{otherwise} \end{cases}$$

Parameters p and q are referred to as the *return* parameter and the *in-out* parameter, respectively. A low value of the return parameter p means that the random walker is more likely to backtrack its steps, meaning the random walk will be closer to a breadth first search. On the other hand, a low value of the parameter q encourages the walker to move away from the starting node and the random walk resembles a depth first search of the network. To calculate the maximizing vector f , a set of random walks of limited size is simulated starting from each node in the network to generate several samples of the set $N_S(u)$.

The function maximizing expression (1) is calculated using stochastic gradient descent. The value of (1) is estimated at each generated sampling of the neighborhoods $N_S(u)$ for all nodes in the network to discover the vector f that maximizes the expression for the simulated neighborhood set.

Extensions to multilayer networks

In this section, we discuss how the node2vec embedding algorithm relates to the considered multilayer network visualization. The original implementation of node2vec operates only on homogeneous, weighted networks. As such, we primarily use it to obtain intra-layer layout, which is also the computationally more expensive part of the layout computation. However, as Py3plex can easily return the supra-adjacency matrix, node2vec could also be applied to such matrix directly in order to obtain global node representations. The considered node2vec was also recently extended to multilayer tissue networks indicating such extensions are possible (Zitnik and Leskovec 2017). We test a similar idea with dynamic networks in “[Experiment one: benchmark of layout computation time](#)” section.

Reducing embedding dimensionality using t-SNE

Note that even though the nodes of a network can be embedded in 2 dimensional space directly, the obtained representations are normally not representative of the network's structure. If $\mathcal{E}(G)$ represents the network embedding function, we introduce an additional operation, $\mathcal{P}(\mathcal{E}(G))$, i.e. a projection of the obtained network embedding (see previous section) to a 2-dimensional, real-valued vector space. Embedding-based graph drawing was previously proven to scale to very large networks (Hachul and Jünger 2006), thus we believe similar ideas implemented with more recent approaches could offer similarly good performance on large multilayer networks with many layers.

In the experiments shown in the following sections, we use t-SNE projections (Maaten and Hinton 2008) for obtaining the final set of node coordinates. The t-SNE algorithm takes as input a set of high dimensional vectors and projects them to a low-dimensional vector space while maintaining (as much as possible) the similarities between the vectors. For use in data visualization, the low-dimensional space has dimension 2 or 3. The algorithm works in two steps.

- 1 Similarities between pairs of input vectors (in our case, the d -dimensional node embeddings) are calculated.
- 2 Low-dimensional vectors are calculated such that the vector-pair-similarities of the low dimensional vectors re-create the original similarities as closely as possible.

Similarities between input vectors $\{x_1, \dots, x_n\}$ are modeled as follows. Let x_i and x_j denote two points in the input d -dimensional embedding. The similarity between data point x_j and x_i is modeled as the probability that x_i would pick x_j as its neighbor if neighbors were picked in proportion to their probability density under a Gaussian centered at x_i and is calculated as follows:

$$p_{ij} = \frac{e^{-\|x_i - x_j\|^2 / 2\sigma_i^2}}{\sum_{g \neq h} e^{-\|x_g - x_h\|^2 / 2\sigma_i^2}},$$

where σ_i is the variance of the Gaussian, centered on data point x_i . In t-SNE, σ_i is determined so that the *perplexity* of the Gaussian equals a fixed value, specified by the user. The desired perplexity parameter can be interpreted as the number of neighbors of each data point. From p_{ij} , the joint probability p_{ij} is calculated as $p_{ij} = \frac{p_{ji} + p_{ij}}{2n}$.

In finding the low-dimensional representation $\{y_1, \dots, y_n\}$, t-SNE models similarities between pairs of representation vectors using the Cauchy distribution, calculated as follows:

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{g \neq h} (1 + \|y_g - y_h\|^2)^{-1}}.$$

Using this model, t-SNE uses stochastic gradient descent to minimize the Kullback-Leibler divergence (Kullback and Leibler 1951) between the original probabilities p_{ij} and the new probabilities q_{ij} . We refer the interested reader to the original paper (Maaten and Hinton 2008) for technical details of this optimization.

Final formulation

Embedding \mathcal{E} , described in “[Network embedding using node2vec](#)” section and projection \mathcal{P} , described in “[Reducing embedding dimensionality using t-SNE](#)” section, represent

two fundamentally different mappings. First, \mathcal{E} attempts to capture a given node's neighbourhood information, whilst \mathcal{P} attempts to construct a low-dimensional embedding by preserving the input's dimension (distances between input feature vectors).

The mappings \mathcal{E} and \mathcal{P} are applied sequentially to obtain coordinate tuples $t_{xy} \in \mathbb{R}^2$, where t_{xy} can already serve as the coordinates for plotting individual nodes, or can be used as the initialization of the force-directed layout. The presented algorithm can thus be described in three simple steps:

- 1 Network embedding into d -dimensional vector space.
- 2 Projection of d -dimensional embeddings into two dimensions.
- 3 (Optional) few iterations of distance minimization.

For the experimental evaluation discussed in the next section we used `node2vec` for the network embedding part, and t-SNE for projections. Note that the idea discussed in this section is both embedding, as well as projection-agnostic—arbitrary embedding algorithm's output can be projected to two dimensions using an arbitrary projection method. We recognize as relevant related work the body of literature focusing on node embedding learning, summarized in (Goyal and Ferrara 2018), a survey in which `node2vec` proved to be one of the best performing algorithms. Similarly, the recently introduced Embedding Projector (Smilkov et al. 2016) offers visualization of embeddings projected via PCA or other non-linear projections.

Experimental evaluation

In the following sections, we discuss the performance of `Py3plex` with respect to visualization, as well as analysis tasks. We begin by comparing the proposed embedding-based layout to some of the contemporary layout algorithms. Next, we demonstrate the scalability of the library and conclude with an analysis of a dynamic multiplex social network.

Experiment one: benchmark of layout computation time

We next present a simple benchmark, where we tested the speed of the presented method in comparison with the Force Atlas 2 algorithm (FA2) (which uses Barnes-Hut approximation for the n-body problem for faster minimization). We used the FA2 algorithm as it is widely used in software such as Gephi (Bastian et al. 2009), representing a relevant baseline for this task.

We compared the computation time of the two algorithms on a real-life protein-protein interaction network described below. The IntAct protein-protein interaction network is currently one of the largest resources for mining the human proteome. To perform the experiments, we first downloaded the current version of protein-protein interaction network from the IntAct database (Orchard et al. 2013), which at the time of writing consists of more than 350,000 nodes and approximately 3.8 million edges. In IntAct, the nodes represent individual proteins, and the (undirected) edges represent their interactions. The edges are weighted, where the edge weights correspond to experimental reliability of the interactions between the corresponding proteins, and take values between 0 and 1. This data base consists of protein-protein interaction pairs, scored with a real value representing the confidence of a given interaction. For comparing layouts, we used all edges with score of at least 0.2, yielding a network with more than 100,000 nodes

and 400,000 edges. We computed layouts for each network five times, and averaged the computation times.

The obtained benchmark times along with computed layouts are summarized in Table 2. It can be observed that the embedding based layout looks notably different to any of the force-directed ones. After many rounds of minimization, the embedding based layout starts to resemble the force-directed one. We believe this experiment demonstrates the power of using network embeddings for qualitative analysis. We additionally discuss the obtained results in “Conclusions and further work” section. We continue with a benchmark study, where we compare the visualization time of Py3plex, compared with Pymnet, an alternative Python based library for visualization of such networks.

Experiment two: overall performance

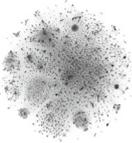
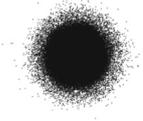
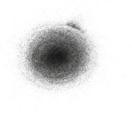
In this section we first present the result of comparing visualization times of Py3plex and Pymnet libraries. Next we discuss a practical case study where we visualize a multiplex dynamic social network. Multiplex networks are comprised of the same set of nodes projected across layers. Here, no physical inter-layer edges are commonly present, as they correspond to *is a* relation.

We present the results for the times needed to obtain a visualization of networks of different complexities; here, we compare Py3plex to Pymnet library. Even though Pymnet does not support the diagonal projection and Py3plex does not support the default 3D linear projections, both methods are useful for visualization of different aspects of multilayer networks.

The main aim of this section is to demonstrate that Py3plex offers better support for visualization of larger networks.

The experimental setting for this task was designed as follows. Random multilayer Erdős-Rényi (ER) networks were generated using the Pymnet (Kivelä et al. 2014) library. We used this network model purely for its simplicity and control over the node and edge space. While ER networks do not exhibit such real-world properties as, for example, Stochastic Block Models (Holland et al. 1983), we believe the larger ER networks considered exhibit enough complexity for comparisons to be meaningful for the considered

Table 2 Comparison of force directed layout with the presented embedding based layout

Iterations	Time (BH)	Visualization (BH)	Time (Embedding)	Visualization (Embedding)
0	1min		23 min	
10	24min		61 min	
100	426min		480 min	

The BH denotes the Barnes-Hut-based Force Atlas2 layout computation. Note that apart from force minimization of the non-embedded network we also present results of minimizing a network's coordinates, initialized using embedding projections

comparisons — in this section, we are only interested in comparing the computation time needed to visualize the networks. As the computation and visualization times are mostly dependent on the numbers of nodes and edges, if Py3plex outperforms Pymnet on these networks, we also expect it to outperform Pymnet on other networks of comparable size.

Such random networks are parameterized using parameter N corresponding to the total number of nodes, parameter L corresponding to the number of layers, and parameter p corresponding to the re-wiring probability. We generated the networks in the following parameter ranges:

$$p \in \{0.05, 0.1, 0.2, 0.3\},$$

$$N \in \{5, 10, 20, 50, 80, 100\},$$

$$L \in \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}.$$

Once generated, the time needed for visualization was recorded. We compared visualization times of Pymnet and Py3plex, as the remaining Python-based alternative, MultinetX, does not support drawing of coupled edges.

We show the final results of our experiment in the form of box plots, where $|N|$ or $|E|$ is plotted on the x -axis and the time needed is plotted on the y -axis (Figs. 5 and 6). Networks with more than 100 nodes were not considered for this benchmark, as it took Pymnet more than two hours for visualization. Nonetheless, Py3plex was able to visualize a network with $|N| = 4,000$ and $|E| = 18,600$ under two hours, even though the obtained network is not necessarily useful for visualization purposes. The machine used for benchmark testing was an off-the-shelf Lenovo y510p laptop.

Visualization of dynamic multiplex networks

Real-world multilayer or multiplex networks are commonly subject to either edge or node dynamics; for example, friendships form over time, or biological phenotypes emerge and disappear (Secrier et al. 2012). Here, the number of e.g., edges can be subject to notable

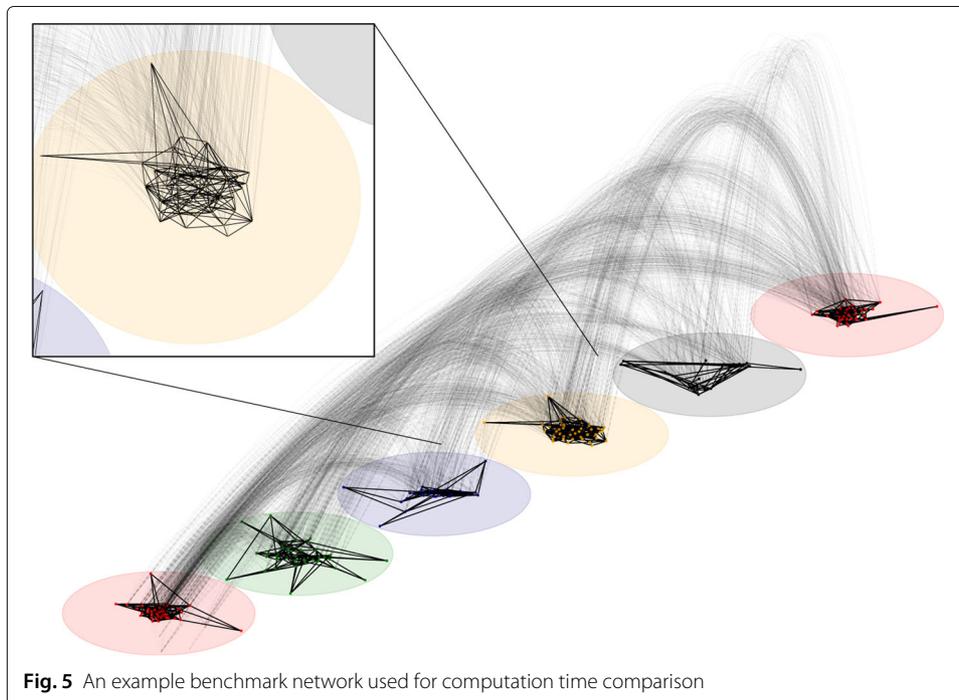
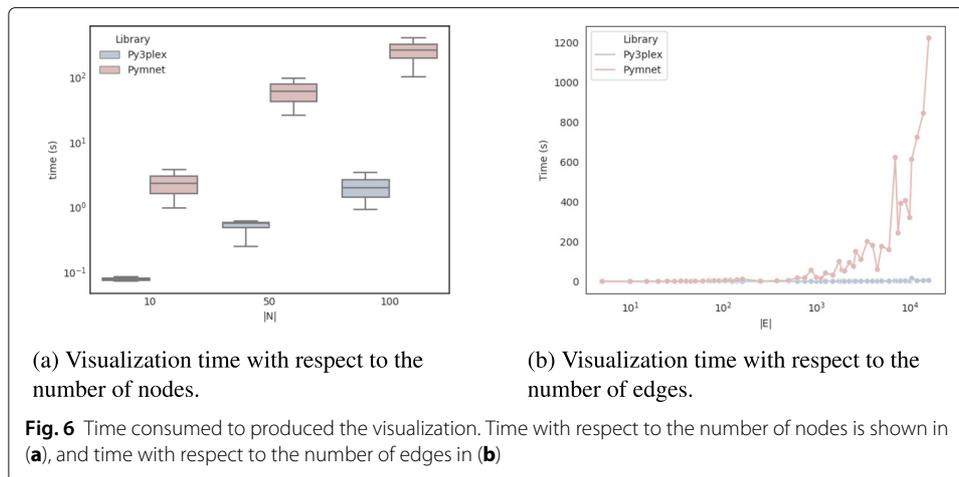


Fig. 5 An example benchmark network used for computation time comparison



change, hence the temporal component of a given network can not be neglected. In this section we showcase some of Py3plex's rather simple extensions to dynamic networks. A natural representation for a dynamic network is some form of video output, where individual e.g., time windows or slices are used as frames. We present Py3plex's functionality for production of such frames.

We consider a dynamic multiplex social network initially introduced in Omodei et al. (2015). Here, the same set of users is observed in terms of retweets, mentions and comments. We represented the network by considering each *type of interaction* as a separate layer, as the users remain the same. The network consists of 392,000 users, we consider more than 100,000 time points, each representing an event (edge) between the users. The edges represent either retweets, mentions or comments.

Visualization preparation

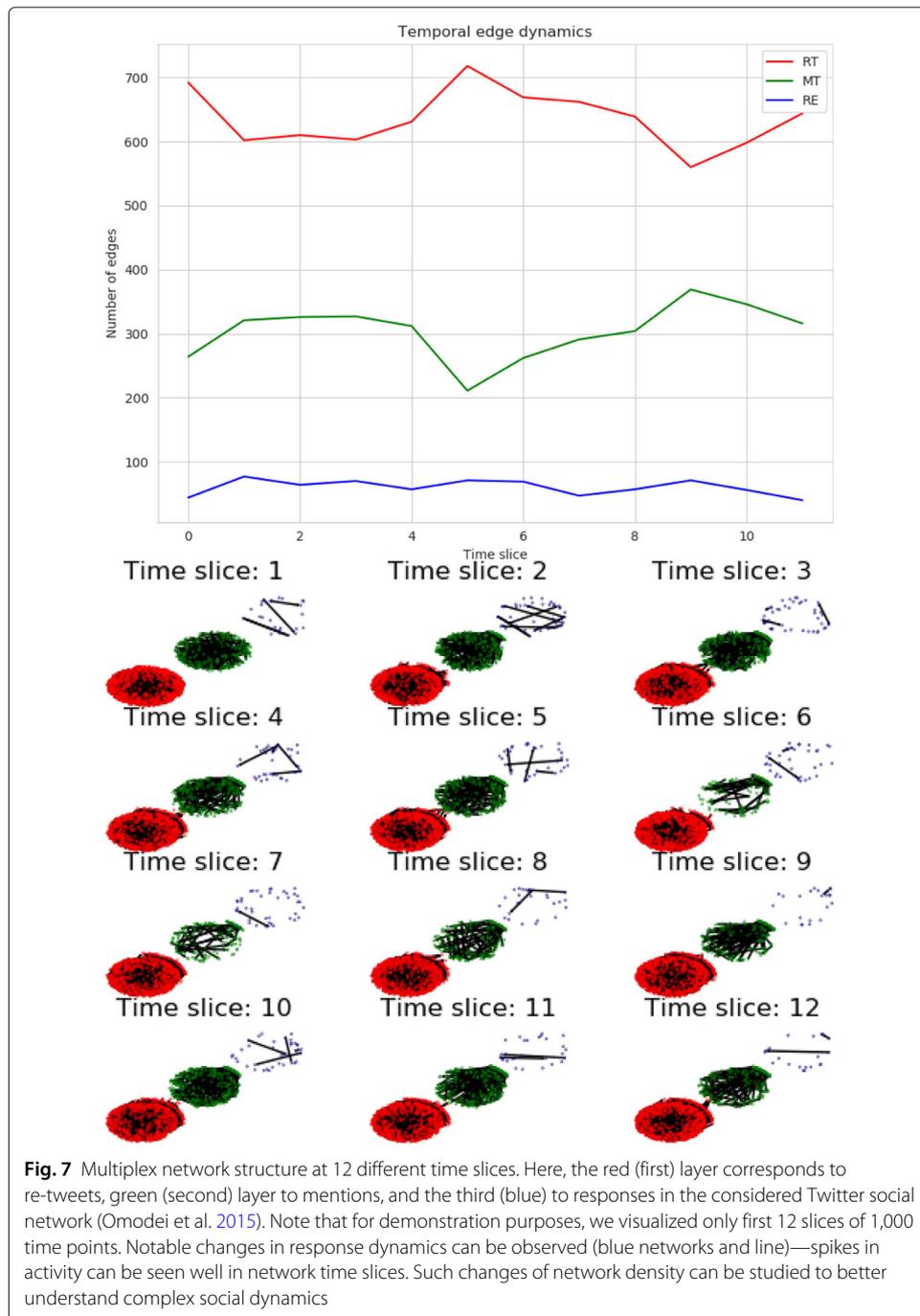
As the users in separate layers are the same, we do not plot any inter-layer edges, as in the presented visualization they do not contain any added value and only result in cognitive overload. Next, we compute layouts for individual nodes as follows. We first create a simpler, homogeneous network consisting of all users. There, we compute the layout for each node. We use the same layout for each layer, hence, for example, a node's position in e.g., the *retweet* layer directly corresponds to the one in e.g., *mentions* layer.

We take into account the network dynamics (temporal links) by considering time slices of size 12,000. Although there exist more intricate methods for taking into account the temporal information, we believe the presented visualization serves the demonstration purpose, as it can easily be extended to e.g., a moving window scenario. For clarity, we do not display isolated nodes.

Finally, we fix the order of layers, and plot each frame (slice) as a block in a grid of 4×3 images, showing the evolution of the network. Note that the obtained frames can easily be used to construct a video, should the need arise. Apart from the grid of networks we also attach simple line plots, which show the overall layer activity with respect to a given time slice. This simpler representation of temporal evolution is used to assess whether the presented visualization captured any distinct events in the evolution of the considered network.

Visualization result

The resulting visualization is shown in Fig. 7. The considered time granularity (1000 points) yields re-tweet and mention networks (red and green) with constant activity. The response network (blue) is subject to the highest variability in terms of link quantity. At e.g., slices 3, 10 and 22 the network is subject to high activity — we believe this indicates the spread of viral news. Other slices, such as for example 8 and 9 indicate there are also periods, when only a few mentions appear.



Conclusions and further work

We have developed and implemented Py3plex, a network analysis and visualization library focused on multilayer networks. In addition to most of the functionality offered by other state-of-the-art libraries, Py3plex introduces a novel visualization, suitable for larger multilayer networks, as to our knowledge, there currently do not exist any methods that can visualize networks composed of thousands of nodes along with their intra- as well as inter-layer organization. Py3plex is suitable for the development of algorithms, as it offers fast lookup and indexing routines, which scale to networks consisting of multiple layers with hundreds of thousands of edges.

One of the new concepts introduced in this paper is a novel multilayer network visualization layout. Py3plex plots intra-layer connections separately to the inter-layer ones, yielding more comprehensive visualizations. We believe the presented diagonal projection could also be extended into an additional dimension (orthogonal), hence offering more efficient space consumption. The presented library currently only includes primitive tools for *animating* multilayer networks—here, the temporal component is split into small slices, which are animated (network snapshots are merged into a single animation). We believe this aspect of Py3plex could be further improved, as many real-world multilayer networks also have a dynamic component.

Another novelty of this work is embedding based visualization. We believe the experiments conducted serve as a proof of concept, yet we believe many other embedding algorithms could also be used to improve the layout's quality. Furthermore, the recently introduced UMAP algorithm (McInnes et al. 2018) could also be used instead of t-SNE, making the layout computation even more scalable. We believe this embedding could be especially useful in situations where larger networks with distinct topological structure are considered. Even though we demonstrated the embedding based layout on a biological network we believe it could serve as a viable visualization alternative to visualize other types of networks, such as for example the transportation networks.

In this work, we compared the visualization capabilities of Py3plex to these of the Pymnet library. Even though we demonstrated that Py3plex layout scales better, Pymnet's 3D visualization is better suited for small demonstrative purposes, where the key idea of a studied system needs to be conveyed as clearly as possible. We believe that the two libraries are in this sense complementary, as we found Py3plex to be the preferable tool used for visualization of actual, larger multilayer networks.

In terms of scalability, we attempted to re-implement some of the common bottlenecks, such as the layout computation, aggregation, and indexing using efficient vectorized operations. However, we have yet to improve the traversal algorithms, as they are not the key focus of this work. We will re-implement random graph generation and crawling first in C, and eventually on GPUs for maximum performance.

Finally, we believe that Py3plex will serve as a testbed for the development of novel algorithms, especially the ones focusing on multilayer dynamics and structure. An example which proves this functionality is the recently implemented variation of layer

entanglement (Renoust et al. 2014), which was constructed using only the data structures offered by Py3plex. By offering fast prototyping, this library can serve to exchange the ideas related to multilayer network analysis.

One of the main drawbacks of Py3plex visualizations is the cognitive load. While the user can currently observe high-level organization of a multilayer network, it can be hard to interpret the structure of individual layers and their contributions to the structure of the whole network. Even though Py3plex can plot directly to an interactive canvas, we believe this issue represents an interesting research direction and will be addressed in follow up work.

As the Py3plex analysis suite is by no means exhaustive, further work includes the implementation of network dynamics analysis methods as well as the more efficient, GPU-based random samplers. Py3plex aims to bridge the gap between machine learning approaches and complex networks, and it does not yet include extensive tensor manipulation, unfolding, and construction routines offered in the Pymnet library.

Appendix A: Network decomposition functions

In this section, we give an overview of the heterogeneous network decomposition functions introduced in HINMINE (Kralj et al. 2018), supported by Py3plex.

Given a heuristic function f , a weight w of an edge between the two nodes u and v is computed as follows:

Let B denote the set of all nodes of the base type. We use the following notations: $f(t, d)$ denotes the number of times a term t appears in the document d and D denotes the corpus (a set of documents). We assume that the documents in the set are labeled, each document belonging to a class c from a set of all classes C . We use the notation $t \in d$ to describe that a term t appears in document d . Where used, the term $P(t)$ is the probability that a randomly selected document contains the term t , and $P(c)$ is the probability that a randomly selected document belongs to class c . We use $|d|$ to denote the length (in words) of a document, and $avgdl$ denotes the average document length in the corpus. All weight functions (heuristics) supported by Py3plex are summarized in Table 3.

Table 3 Term weighing schemes, taken from (Kralj et al. 2018), tested for decomposition of heterogeneous networks and their corresponding formulas

Scheme	Formula
tf	$f(t, d)$
if-idf	$f(t, d) \cdot \log \left(\frac{ D }{ \{d' \in D: t \in d'\} } \right)$
chi ²	$f(t, d) \cdot \sum_{c \in C} \frac{(P(t \wedge c)P(\neg t \wedge \neg c) - P(t \wedge \neg c)P(\neg t \wedge c))^2}{P(t)P(\neg t)P(c)P(\neg c)}$
ig	$f(t, d) \cdot \sum_{c \in C, c' \in \{c, \neg c\}} \sum_{t' \in \{t, \neg t\}} \left(P(t', c') \cdot \log \frac{P(t' \wedge c')}{P(t')P(c')} \right)$
gr	$f(t, d) \cdot \sum_{c \in C} \frac{\sum_{c' \in \{c, \neg c\}} \sum_{t' \in \{t, \neg t\}} (P(t', c') \cdot \log \frac{P(t' \wedge c')}{P(t')P(c')})}{-\sum_{c' \in \{c, \neg c\}} P(c') \cdot \log P(c')}$
delta-idf	$f(t, d) \cdot \sum_{c \in C} \left(\log \frac{ c }{ \{d' \in D: d' \in c \wedge t \in d'\} } - \log \frac{ c }{ \{d' \in D: d' \notin c \wedge t \notin d'\} } \right)$
rf	$f(t, d) \cdot \sum_{c \in C} \log \left(2 + \frac{ \{d' \in D: d' \in c \wedge t \in d'\} }{ \{d' \in D: d' \notin c \wedge t \notin d'\} } \right)$
bm25	$f(t, d) \cdot \log \left(\frac{ D }{ \{d' \in D: t \in d'\} } \right) \cdot \frac{k+1}{f(t, d) + k} \cdot \frac{1}{(1-b + b \cdot \frac{ d }{avgdl})}$

Acknowledgements

We would like to thank to Lucija Luetič for proofreading the manuscript, which notably improved the work's quality. The work of the first author was funded by the Slovenian Research Agency through a young researcher grant. The work of other authors was supported by the Slovenian Research Agency (ARRS) core research programme *Knowledge Technologies* (P2-0103) and ARRS funded research project *Semantic Data Mining for Linked Open Data* (financed under the ERC Complementary Scheme, N2-0078). The work was supported also by European Union's Horizon 2020 research and innovation programme under grant agreement No 825153, project EMBEDDIA (Cross-Lingual Embeddings for Less-Represented Languages in European News Media).

Authors' contributions

The authors' contributions are as follows. BŠ and JK implemented the library and conducted the experiments. BŠ, JK and NL designed the experiments and conducted theoretical overviews and analysis. All authors read and approved the final manuscript.

Availability of data and materials

The Py3plex library as well as the datasets used in this paper are freely accessible at web site <https://github.com/SkBlaz/Py3plex>, where many working examples of the Py3plex functionality are also offered.

Competing interests

The authors declare that they have no competing interests.

Received: 9 March 2019 Accepted: 30 August 2019

Published online: 29 October 2019

References

- Amato R, Kouvaris NE, San Miguel M, Díaz-Guilera A (2017) Opinion competition dynamics on multiplex networks. *New J Phys* 19(12)
- Auber D, Archambault D, Bourqui R, Delest M, Dubois J, Lambert A, Mary P, Mathiaut M, Mélançon G, Pinaud B, et al. (2017) TULIP 5. Springer
- Auber D (2004). In: Jünger M, Mutzel P. (eds). *Tulip — A Huge Graph Visualization Framework*. Springer, Berlin. pp 105–126
- Bastian M, Heymann S, Jacomy M (2009) Gephi: an open source software for exploring and manipulating networks. In: *Third International AAAI Conference on Weblogs and Social Media*
- Batagelj V, Mrvar A (2001) Pajek—analysis and visualization of large networks. In: *International Symposium on Graph Drawing*. Springer. pp 477–478
- Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E (2008) Fast unfolding of communities in large networks. *J Stat Mech Theory Exp* 2008(10):10008
- Boccaletti S, Bianconi G, Criado R, del Genio CI, Gómez-Gardeñes J, Romance M, Sendiña-Nadal I, Wang Z, Zanin M (2014) The structure and dynamics of multilayer networks. *Phys Rep* 544(1):1–122
- De Domenico M, Nicosia V, Arenas A, Latora V (2015) Structural reducibility of multilayer networks. *Nat Commun* 6:6864
- De Domenico M, Porter MA, Arenas A (2015) MuxViz: A tool for multilayer analysis and visualization of networks. *J Complex Netw* 3(2):159–176
- De Domenico M, Solé-Ribalta A, Cozzo E, Kivela M, Moreno Y, Porter MA, Gómez S, Arenas A (2013) Mathematical formulation of multilayer networks. *Phys Rev X* 3(4). <https://doi.org/10.1103/physrevx.3.041022>
- Goyal P, Ferrara E (2018) Graph embedding techniques, applications, and performance: A survey. *Knowl-Based Syst* 151:78–94
- Grover A, Leskovec J (2016) Node2vec: Scalable feature learning for networks. In: *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*. ACM, New York. pp 855–864
- Hachul S, Jünger M (2006) An experimental comparison of fast algorithms for drawing general large graphs. In: Healy P, Nikolov NS (eds). *Graph Drawing*. Springer, Berlin, Heidelberg. pp 235–250
- Hagberg A, Swart P, S Chult D (2008) Exploring network structure, dynamics, and function using NetworkX. In: *Proceedings of the 7th Python in Science Conference (SciPy)*
- Harel D, Koren Y (2002) Graph drawing by high-dimensional embedding. In: *International Symposium on Graph Drawing*. Springer. pp 207–219
- Holland PW, Laskey KB, Leinhardt S (1983) Stochastic blockmodels: First steps. *Soc Networks* 5(2):109–137
- Huang W, Eades P, Hong S-H (2009) Measuring effectiveness of graph visualizations: A cognitive load perspective. *Inf Vis* 8(3):139–152
- Hunter JD (2007) Matplotlib: A 2d graphics environment. *Comput Sci Eng* 9(3):90
- Jacomy M, Venturini T, Heymann S, Bastian M (2014) ForceAtlas2, A continuous graph algorithm for handy network visualization designed for the Gephi software. *PLoS ONE* 9(6):98679
- Jones E, Oliphant T, Peterson P, et al (2001) SciPy: Open source scientific tools for Python
- Kivela M, Arenas A, Barthelemy M, Gleeson JP, Moreno Y, Porter MA (2014) Multilayer networks. *J Complex Netw* 2(3):203–271
- Kullback S, Leibler RA (1951) On information and sufficiency. *Ann Math Stat* 22(1):79–86
- Kralj J, Robnik-Šikonja M, Lavrač N (2018) HINMINE: Heterogeneous Information Network Mining with Information Retrieval Heuristics. *J Intell Inf Syst* 50(1):29–61
- Kralj J, Robnik-Šikonja M, Lavrač N (2019) Netsdm: Semantic data mining with network analysis. *J Mach Learn Res* 20(32):1–50
- Leskovec J, Sosič R (2016) Snap: A general-purpose network analysis and graph-mining library. *ACM Trans Intell Syst Technol* 8(1):1–1120
- Maaten Lvd, Hinton G (2008) Visualizing data using t-sne. *J Mach Learn Res* 9(Nov):2579–2605

- McGee F, Ghoniem M, Melançon G, Otjacques B, Pinaud B (2019) The state of the art in multilayer network visualization. *Comput Graph Forum* 0(0). <https://doi.org/10.1111/cgf.13610>
- McInnes L, Healy J, Melville J (2018) Umap: Uniform manifold approximation and projection for dimension reduction. arXiv preprint arXiv:1802.03426
- Nepusz G, Csárdi G (2006) The igraph software package for complex network research. *Complex Syst* 1695(5):1–9
- Omodei E, De Domenico MD, Arenas A (2015) Characterizing interactions in online social networks during exceptional events. *Front Phys* 3:59
- Orchard S, Ammari M, Aranda B, Breuza L, Briganti L, Broackes-Carter F, Campbell NH, Chavali G, Chen C, Del-Toro N, et al (2013) The mintact project—intact as a common curation platform for 11 molecular interaction databases. *Nucleic Acids Res* 42(D1):358–363
- Pavlopoulos GA, O'Donoghue SI, Satagopam VP, Soldatos TG, Pafilis E, Schneider R (2008) Arena3d: visualization of biological networks in 3d. *BMC Syst Biol* 2(1):104
- Perozzi B, Al-Rfou R, Skiena S (2014) Deepwalk: Online learning of social representations. In: *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14*. ACM, New York. pp 701–710
- Purchase H (1997) Which aesthetic has the greatest effect on human understanding? In: *International Symposium on Graph Drawing*. Springer. pp 248–261
- Renoust B, Melançon G, Viaud M-L (2014). In: Missaoui R, Sarr I (eds). *Entanglement in Multiplex Networks: Understanding Group Cohesion in Homophily Networks*. Springer, Cham. pp 89–117
- Ribeiro LFR, Saverese PHP, Figueiredo DR (2017) Struc2vec: Learning node representations from structural identity. In: *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '17*. ACM, New York. pp 385–394
- Rosvall M, Axelsson D, Bergstrom CT (2009) The map equation. *Eur Phys J Spec Top* 178(1):13–23
- Secrier M, Pavlopoulos GA, Aerts J, Schneider R (2012) Arena3d: visualizing time-driven phenotypic differences in biological systems. *BMC Bioinformatics* 13(1):45
- Shannon P (2003) Cytoscape: A software environment for integrated models of biomolecular interaction networks. *Genome Res* 13(11):2498–2504
- Škrlić B, Kralj J, Lavrač N (2019) Cbssd: community-based semantic subgroup discovery. *J Intell Inf Syst*. <https://doi.org/10.1007/s10844-019-00545-0>
- Škrlić B, Kralj J, Vavpetič A, Lavrač N (2018) Community-based semantic subgroup discovery. In: Appice A, Loglisci C, Manco G, Masciari E, Ras ZW (eds). *New Frontiers in Mining Complex Patterns*. Springer, Cham. pp 182–196
- Škrlić B, Kralj J, Lavrač N (2019) Py3plex: A library for scalable multilayer network analysis and visualization. In: Aiello LM, Cherifi C, Cherifi H, Lambiotte R, Lió P, Rocha LM (eds). *Complex Networks and Their Applications VII*. Springer, Cham. pp 757–768
- Smilkov D, Thorat N, Nicholson C, Reif E, Viégas FB, Wattenberg M (2016) Embedding projector: Interactive visualization and interpretation of embeddings. arXiv preprint arXiv:1611.05469
- The Boost Graph Library (2002) *User Guide and Reference Manual*. Addison-Wesley Longman Publishing Co., Inc., Boston
- Walt Svd, Colbert SC, Varoquaux G (2011) The numpy array: a structure for efficient numerical computation. *Comput Sci Eng* 13(2):22–30
- Wang Z, Wang L, Szolnoki A, Perc M (2015) Evolutionary games on multilayer networks: a colloquium. *Eur Phys J B* 88(5):124
- Zitnik M, Leskovec J (2017) Predicting multicellular function through multi-layer tissue networks. *Bioinformatics* 33(14):190–198

Publisher's Note

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

Submit your manuscript to a SpringerOpen® journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► [springeropen.com](https://www.springeropen.com)

RESEARCH

Open Access



Layer entanglement in multiplex, temporal multiplex, and coupled multilayer networks

Blaž Škrlić^{1*}  and Benjamin Renoust²

*Correspondence:

blaz.skrlic@ijs.si

¹ Jožef Stefan International Postgraduate School and Jožef Stefan Institute, Ljubljana, Slovenia

Full list of author information is available at the end of the article

Abstract

Complex networks, such as transportation networks, social networks, or biological networks, capture the complex system they model by often representing only one type of interactions. In real world systems, there may be many different aspects that connect entities together. These can be captured using multilayer networks, which combine different modalities of interactions in a single model. Coupling in multilayer networks may exhibit different properties which can be related to the very nature of the data they model (or to events in time-dependent data). We hypothesise that such properties may be reflected in the way layers are intertwined. In this paper, we investigated these through the prism of layer entanglement in coupled multilayer networks. We test over 30 real-life networks in 6 different disciplines (social, genetic, transport, co-authorship, trade, and neuronal networks). We further propose a random generator, displaying comparable patterns of elementary layer entanglement and transition coupling entanglement across 1,329,696 synthetic coupled multilayer networks. Our experiments demonstrate difference of layer entanglement across disciplines, and even suggest a link between entanglement intensity and homophily. We additionally study entanglement in 3 real world temporal datasets displaying a potential rise in entanglement activity prior to other network activity.

Keywords: Multiplex networks, Layer entanglement, Temporal network, Network topology, Network generator

Introduction

A real world complex system often counts multiple interactions between multiple different entities. When these interactions are regrouped under multiple families of entities, multilayer network modelling becomes a tool of choice to capture the key components of the system. The use of this model emerges in all fields of science from social sciences to finances, logistics, biology, and many more (Kivelä et al. 2014).

With multilayer networks, the study of *multiple viewpoints* [or aspects (Kivelä et al. 2019)] on the same network data becomes possible. This is critical for example in social network analysis, to study the role of users in different networks, and compare them (for example the same individual may behave differently on LinkedIn, Twitter, or Facebook). These different networks form different types of links that may be overlaid.

Motivated by their practical interest, *multilayer networks* also show interesting structures (Battiston et al. 2014) that could be exploited to mine *community structures* or study the roles of nodes and edges through centrality, for example. These are also possible in a traditional network analysis standpoint but often requires some kind of simplification (such as one-mode projection) but recent advances show that interesting structures can be obtained *directly* from the multilayer networks (Gomez et al. 2013; Chen et al. 2018; Škrlj et al. 2019).

The key concept in multilayer networks are the layers themselves. Since the structure of such networks is driven by the layers and their aspect (Kivelä et al. 2014), understanding how the layers organise can reveal properties unique to a given multilayer network model (Renoust et al. 2015; Škrlj and Renoust 2019). Particularly, the intertwining of edges, or *layer entanglement* (Renoust et al. 2014, 2013), shows how layers overlap to form coherent structures and substructures.

Although recent works have focused on multilayer network analysis and description (Wang et al. 2018; Omodei et al. 2015), not many have focused on a large scale analysis of multilayer networks of different nature—and produced in different disciplines, while comparing them to synthetic models. One comparative study of flow analysis (De Domenico et al. 2015a) has particularly influenced this paper where emerging structures are described, albeit not compared to synthetic models.

In their seminal work, McPherson et al. (2001) discuss how ties emerge in social systems. They investigate how people similarity, i.e. homophily, is a strong driver to the formation of ties, with the addition to make them more durable in a dynamic system. They investigate social ties in a multilayer manner, and argue for further research: “*in the impact of multiplex ties on the patterns of homophily; [and] the dynamic of network change over time [...]*”. Our original work (Škrlj and Renoust 2019)—that we extend in this paper—particularly resonates with the first point of McPherson et al., in that we displayed a link between homophily (McPherson et al. 2001; Borgatti et al. 2009) in social networks and high entanglement intensity networks.

This paper extends (Škrlj and Renoust 2019), which originally contributed with an open source implementation of entanglement homogeneity and intensity for multiplex networks, while evaluating them over 30 real world networks. We proposed also a synthetic multiplex network generator. A generation of over 10k synthetic networks, and their comparison with the real world networks, displayed common patterns of entanglement homogeneity and intensity that could be specific to the families of applications that *generated* the networks. In this extended work, we contribute with:

- The theoretical extension of the entanglement computation to a fully multiplex model that takes into account coupling edges;
- The extension of our synthetic generator accordingly;
- The computations on a wider range of real and synthetic networks (1,329,696 synthetic networks were considered);
- The study of entanglement in large, temporal multiplex networks;
- An open-source implementation of all conducted experiments.

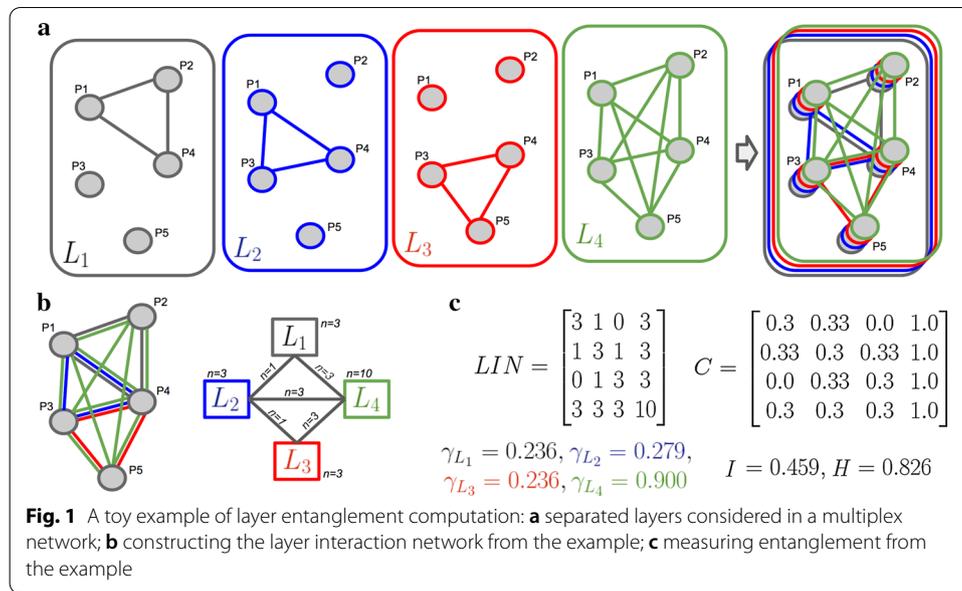
Coupled multilayer and multiplex networks

A multilayer network can be defined as a sequence $M = \{G_l\}_{l \in L} = \{(V_l, E_l)\}_{l \in L}$ where $E_l \subseteq V_l \times V_l$ is a set of edges in one network $l \in L$ of the sequence (Kivelä et al. 2014). Multilayer networks are commonly understood as layers comprised of interactions, where each layer corresponds to a specific aspect of the system. Coupling accounts for transitions between layers. Kivelä et al. (2014) consider a multiplex network as “*diagonally coupled multilayer networks in which each layer shares at least one node with some other layer in the network*”. They consider also *node-aligned multiplex networks*, which do not specifically address coupling of nodes, but assume that nodes are shared (and coupled) across all layers. In our context, we refer to *coupled multilayer networks* when we specifically consider networks with coupling between nodes across layers, and simply to *multiplex networks* when considering node-aligned multiplex networks. The difference between these two types of multiplex networks is only whether we consider or not the coupling between layers. In multiplex networks, nodes represent *the same* entity across all layers.

We represent a multiplex network as a structure $M' = (V_M, E_M)$, where V_M is the set of nodes and E_M the set of all edges (in *all* layers). \mathbb{V} denotes the super set of all nodes, and $\mathbb{E} = \mathbb{V} \times \mathbb{V}$ the super set of all edges, regardless of the layers. There may exist coupling edges connecting nodes through layers, forming *transition coupling*. This may concern, for example, coupled multilayer networks which are modelling transportation systems (Cozzo et al. 2015). In that case, we can differentiate the *elementary* layers (holding *inner-layer* edges) from the *transition* coupling (holding *coupling* edges). Each transition coupling $t = (l, l')$ between layer l and l' can be modelled similarly to a layer, with a set of nodes and edges. If $S \subset L$ represents the subset of all *elementary* layers, and $T \subset L$ the subset of all *transition* coupling, we may define our coupled multilayer network M as the union. It combines a multilayer network with elementary layers only, and another multilayer network with transition coupling only $M = \{G_l\}_{l \in L} = M_S \cup M_T = \{G_s\}_{s \in S} \cup \{G_t\}_{t \in T}$. The coupling can heavily influence the structural behaviour of multilayer networks (Cozzo and Moreno 2016). It can also influence the resilience of the network against failures (De Domenico et al. 2014) and naturally the diffusion phenomena (Tejedor et al. 2018) too.

Among other examples of coupled multilayer networks, a biological system can be studied at the protein, RNA, or gene level (Valdeolivas et al. 2018). Similarly, social networks can be studied by taking into account a person's presence on multiple platforms (Mittal and Bhatia 2019). For computational purposes, such networks are commonly represented in the form of supra-adjacency matrices, where *block-diagonal structures* connect the same node across individual layers emerges (Cozzo et al. 2015). Algorithms can operate on such matrices directly, and thus exploit additional information representing multiple aspects.

Algorithms for analysis of multilayer networks can also operate on sparse adjacency data structure of the multilayer network directly. Yet, they need to take into account that a given node is present in multiple layers. Such representation is suitable for this work, as we are focused primarily on how edges co-occur across *layers*. Hence, this work focuses primarily on the relations *between the layers* of a given multilayer network. We next discuss the two measures we consider throughout this work.



Entanglement in multiplex networks

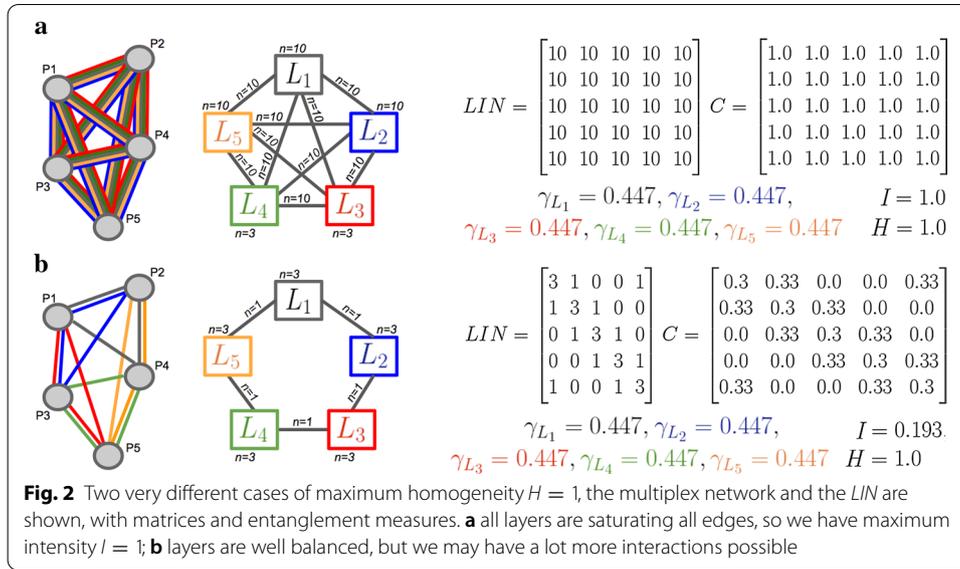
We briefly recall the definitions of entanglement measures from previous work (Renoust et al. 2014).

Layer interaction network

Recall the notion of a multiplex network $M = (V_M, E_M) = \{G_l\}_{l \in L}$. As mentioned earlier, such a network really distinguishes itself from classical graphs through the use of different layers to connect nodes. These layers may have different patterns and may overlap together. There may even exist latent dependencies among these layers. To investigate this matter, each layer could be abstracted to one single node and form a new graph, the *Layer Interaction Network* (hereafter *LIN*) (Renoust et al. 2014). Visualizing the *LIN* is a key component for multiplex network visualization such as in Renoust et al. (2015). In the *LIN*, $LIN = (L, F)$, each node $u_l, u_{l'}, u_{l''}, \dots$ corresponds to a layer $l, l', l'', \dots \in L$ of the multiplex network M , and each edge $f \in F$ captures when two layers overlap through edges. More formally, there exists an edge $f = (u_l, u_{l'})$ whenever there exists at least two nodes $v, v' \in V_M$ with the condition that there exists at least one edge connecting these two nodes on each layer $e_M = (v, v') \in l$ and $e'_M = (v, v') \in l'$. The *LIN* can be interpreted as an edge-layer co-occurrence graph, and the weight of an edge $f = (u_l, u_{l'})$, denoted as $n_{l,l'}$ equals the number of times layers l and l' co-occur. By extension, $n_{l,l}$ is the number of edges on layer l . This process is illustrated in Fig. 1b.

Layer entanglement

The analysis of layer entanglement is inspired by the analysis of *relation content* in social networks (Burt and Schøtt 1985). The idea is to study the redundancy between



relation content, each forming in our formalism a different layer. The layer entanglement measures the “influence” of a layer in its neighbourhood.

This measure is recursively defined: the entanglement γ_l of a layer l is defined upon the entanglement of the layers it is entangled with. Similarly to the eigen centrality (Wasserman and Faust 1994), this translates into the recursive equation:

$$\gamma_l \cdot \lambda = \sum_{l' \in T} \frac{n_{l,l'}}{n_{l,l}} \gamma_{l'}$$

The entanglement of a layer γ_l can be retrieved from a vector γ which corresponds to the right eigenvector (associated to the maximum eigenvalue λ) of the layer overlap frequency matrix with corresponding overlap, defined as:

$$C = (c_{l,l'}), \quad \text{where} \quad c_{l,l'} = \frac{n_{l,l'}}{n_{l,l}} \quad \text{and} \quad c_{l,l} = \frac{n_{l,l}}{|\mathbb{E}|}$$

this metric was initially introduced in Burt and Schött (1985), then later constructed using the weights in the LIN (Renoust et al. 2014) (see Figs. 1 and 2).

Entanglement intensity and homogeneity

The layer entanglement γ_l measures the share of layer l overlapping with other layers. The more a group of layers interacts together, the more the nodes they connect will be cohesive in view of *these* layers, hence the more $\gamma_l \forall l \in L$ values will be similar (their share of entanglement will be similar). This is captured by the *entanglement homogeneity* (Renoust et al. 2014) which is then defined as the following cosine similarity:

$$H = \frac{\langle \mathbf{e}_L, \gamma \rangle}{\|\mathbf{e}_L\| \|\gamma\|} \in [0, 1].$$

With $\mathbf{e}_L = [1, 1, \dots, 1]_L$ the vector of size L all filled with 1’s. Optimal homogeneity is not necessarily reached only when all nodes are connected through all layers, but also when

all nodes are connected in a very balanced manner between all layers (see Fig. 2). Homogeneity thus permits various *symmetries* in a given *LIN*.

When a maximum overlap is reached through all layers in the network, the frequencies in the matrix C (of size $|L| \times |L|$) are saturated with $C_{i,j} = 1$. This gives us a theoretical limit to measure the amount of layer overlap through the *entanglement intensity* (Renoust et al. 2014), defined as:

$$I = \lambda/|L|.$$

In practice, both entanglement intensity and homogeneity have been used to measure the coherence of clusters of documents (Renoust et al. 2013).

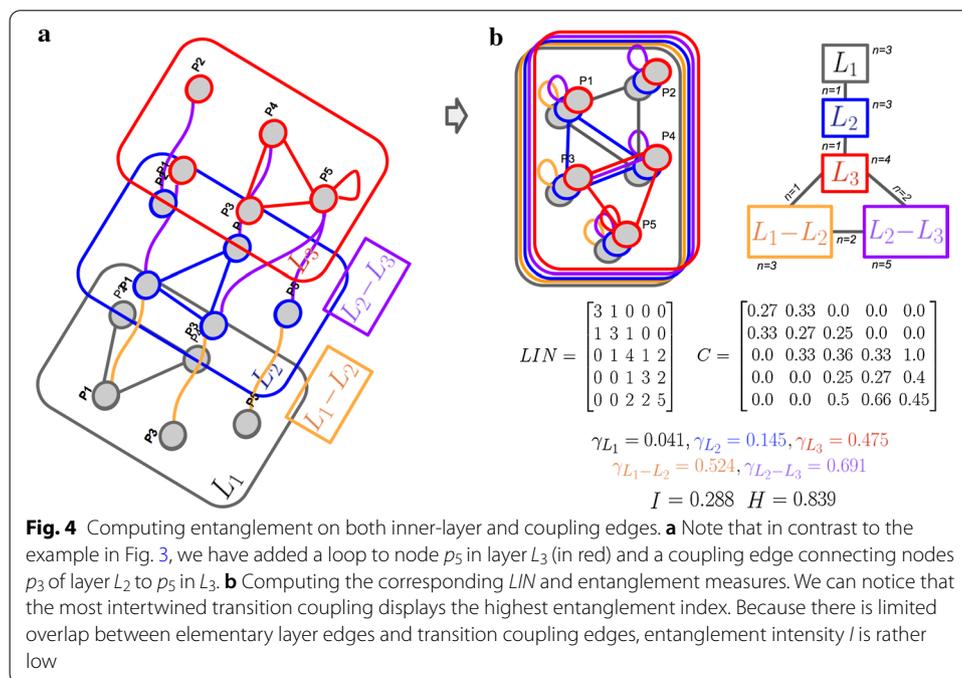
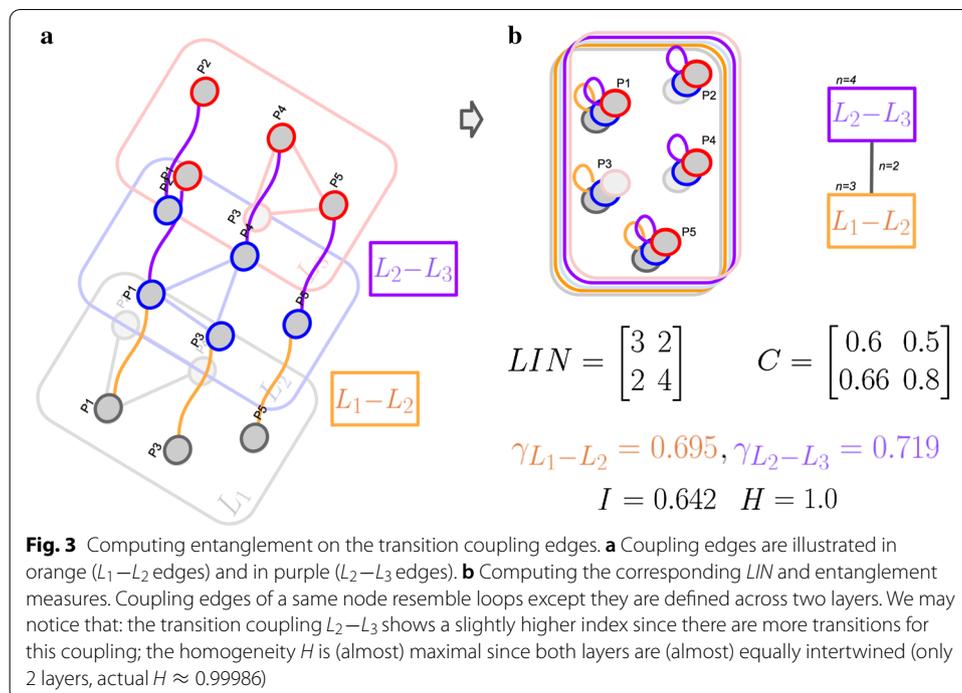
Transition coupling entanglement

We have defined the layer entanglement which measures overlap between layers of a multiplex network, but many multiplex networks include another critical parameter which is coupling edges (Battiston et al. 2014). The coupling often measures the transition of nodes *between* layers, hence the transitions of nodes are captured by edges connecting nodes *across* layers.

Recall our multiplex graph $M = (V_M, E_M)$. Suppose S is the set of *elementary* layers, we can then have transitions between any pair of elementary layers $l \in S$ and $l' \in S$. Let $u_l = (u, l)$, $u \in V_M$, $l \in S$, the connection of a node u within a layer l . A transition coupling edge e can be defined as follows: $e = (u_l, v_{l'}) \in E_M$ such that e connects nodes $\{u, v\} \subseteq V_M$ across layers $l \neq l'$, $\{l, l'\} \subseteq S$. Coupling edges often connect a same node across two layers and may be used to model a physical transition, such as a change from subway to train in a station of a transportation network. As a consequence, a pair of layers $(l, l') = t$ forms a transition coupling $t \in T$ when there exists at least one such edge $e = (u_l, v_{l'}) \in E_M$. Note that taken together, these elementary and transition coupling subsets form the set of all layers $S \cup T = L$, and that the size of T is bounded by the size of S such that $|T| \leq \frac{1}{2}|S|(|S| - 1)$.

Now, given this definition, nothing limits the computation of entanglement (introduced in “[Layer interaction network](#)”, “[Layer entanglement](#)” and “[Entanglement intensity and homogeneity](#)” section) only to the elementary layers part of M_S , as illustrated in Fig. 3. Entanglement can also be used to characterise the coupling between these *elementary* layers if applied only to the edges of the *transition* coupling M_T .

The nature of coupling often captures a very distinct characteristic of the network in comparison to its elementary layers. A transition coupling edge mostly connects the same node across layers, while elementary layers do not always display loops. These cases may happen on rare occasions, one example being an underground path connecting subway stations being modelled as a transition coupling, but the literature is very poor regarding such examples. It is however technically possible to consider both elementary layers and transition coupling in one multiplex network M to compute entanglement (as shown in Fig. 4), but we keep this discussion for the “[Appendix](#)”. In practice, the intensity and homogeneity greatly differ between them, and often result in clearly separated components of the *LIN*.



A coupled multilayer network generator

In this section, we describe an algorithm which generates synthetic coupled multilayer networks, i.e. multilayer networks which share some nodes across some layers, but do not guarantee that all nodes are being shared between all layers. These kinds of networks make the link between general multilayer networks and node-aligned

multiplex networks [for which the assumption is that all nodes are shared through all layers (Kivelä et al. 2014)].

The algorithm is based on the following observations. Let $M = (V_M, E_M)$ represent a coupled multilayer network with layer set L . Each node is associated to a random number of layers $\{l_1, l_2, \dots, l_i\} \subseteq L$. Now for each layer $l_i \in L$ there is a set of nodes $V_{l_i} \subseteq V_M$ which forms a potential set of edges of size $|E_{l_i}| = \frac{1}{2}|V_{l_i}|(|V_{l_i}| - 1)$. We introduce o , a parameter determining the probability of a node occurring at a given layer. We then introduce the probability p of an edge to be created between any pair of nodes belonging to a layer so we may avoid cliques to form on each layer. We referred in our previous work to the edge dropout (Škrlj and Renoust 2019), which is $d = 1 - p$ as the share of links we drop from the clique model. Intuitively, the more similar a given random multiplex is to a clique over each layer, the higher its elementary layer intensity should be. Hence, high intensity implies larger probability that two given nodes will have an edge between them on more than one layer. The generator also accounts for coupling by adding transition coupling edges. These coupling edges are connecting nodes across two layers. We introduce q , the probability for a same node to be connected across two layers. The higher q , the more nodes will be connected through layers. Note that in our initial work (Škrlj and Renoust 2019), neither o nor q were considered (o was in fact picked uniformly).

Algorithm 1: A coupled multilayer network generator.

Parameters : Number of nodes n , number of layers m , inner-layer edge probability p , coupling edge probability q ,
Result: A coupled multilayer network M

```

1  $M \leftarrow$  emptyMultilayerObject;
2 for node in  $[1 \dots n]$  do
3   layerNodes  $\leftarrow$  assignNodeToLayers(node,  $o$ ,  $m$ )
4   ;  $\triangleright$  Nodes are assigned to layers among  $m$  with probability  $o$ .
5   update( $M$ , layerNodes);  $\triangleright$  Update global network.
6 end
7 for layer  $l_i$  with corresponding node set  $V_{l_i}$  do
8   nodeClique  $\leftarrow$  generator of node pairs from  $V_{l_i}$ ;  $\triangleright$  With or without possible loops.
9   innerLayerEdges  $\leftarrow$  sampleWithProbability(nodeClique,  $p$ );  $\triangleright$  Sample via  $p$ .
10  update( $M$ , innerLayerEdges);  $\triangleright$  Update global network.
11 end
12 for layers  $l_i, l_j$  with shared node set  $V_{l_i, l_j}$  do
13   sameNodeTransitionCouplingEdges  $\leftarrow$  sampleWithProbability( $V_{l_i, l_j}$ ,  $q$ );  $\triangleright$  Sample via  $q$ .
14   update( $M$ , sameNodeTransitionCouplingEdges);  $\triangleright$  Update global network.
15 end
16 return  $M$ ;

```

The purpose of this generator is to offer a simple *testbed* for further exploration, as well as additional evidence of the relation between homogeneity and intensity on many random, synthetic networks. The Algorithm 1 represents the proposed procedure.

The generator first randomly assigns the same node index to the many layers (lines 2–5). Once assigned, the layers are processed by applying sampling on $\binom{|V_{l_i}|}{2}$ possible edges in layer l_i . Note that in line 7, this whole clique is virtually generated. The global multiplex is updated during this process (lines 6–10). These steps are then repeated for each transition coupling i.e. pairs of elementary layers (lines 11–14). The implementation thus uses a generator, for which *lazy evaluation* avoids potential combinatorial explosion when considering a large number of nodes and low edge probability.

Some theoretical properties of the generator

In this section we show two properties of the proposed generator. We denote $n = |V_M|$ the parameter setting the number of nodes of the network, $m = |L|$ the parameter setting the number of edge layers in the network, and p the inner-layer edge probability. Let $\phi \in \mathbb{N}^+$ represent the number of possible edges. Then $\phi \leq m \cdot \binom{n}{2}$. Let $o = 1$. Each layer can have at most n nodes. Assuming they form a clique, each layer is thus comprised of $\binom{n}{2}$ edges. As there are m layers, there can be at most $m \cdot \binom{n}{2}$ edges — a clique of n nodes in each layer (assuming $p = 1$). We refer to this bound as $\phi \leq m \cdot \binom{n}{2}$.

In the limit, as $p \rightarrow 1$, a full clique needs to be constructed, assuming each node is projected across all layers. The complexity *w.r.t.* the number of layers and edges is: $\mathcal{O}(m \cdot \binom{n}{2}) = \mathcal{O}(|E_M|)$. Note that, even though theoretically, the proposed generator creates a clique and then samples from it, current, lazy implementation only *generates* the edges needed to satisfy a given p percentage. In practice, only when $p \approx 1$, the generator needs larger portions of space (and time). As such, fully connected networks do not represent real systems, we were able to generate a multitude of very diverse networks. This generator-based implementation does not imply that large spatial overheads are not possible: such situations occur when very dense networks are considered.

We next discuss the impacts of q parameter. The number of coupling edges has a worst case complexity of $\mathcal{O}\left(\binom{m}{2} \cdot n\right)$ since q directly depends on the number of layers available. Let l_a and l_b represent a given pair of layers, where each layer consists of all n possible nodes. As each node couples only to itself, there are at most n edges between l_a and l_b . As there are $\binom{m}{2}$ possible layer pairs, if nodes are in each pair fully coupled, the network can have at most $\binom{m}{2} \cdot n$ coupling edges.

However, is that also the case when considering only transition coupling? Consider the following example of a multiplex network without the coupling edges. No matter what p is employed, if $q \approx 0$, coupling intensity will be low—very few coupling edges are introduced, the observed *LIN* will be very sparse. Hence, we posit that the distribution of intensity shall be *constant* with respect to a given p . The proof of this claim is by contradiction. We assume that p would indeed influence coupling entanglement intensity. Since transition coupling intensity is defined solely based on the coupling edges, this claim would imply a dependency between p and q , which is by the definition (and design) not the case. Even if the nodes are *isolated* in each layer, transition coupling intensity can be high. Note also that the node positioning, governed by o , directly impacts both elementary and transition coupling entanglement, since there is higher possibilities for edges to overlap when nodes belong to many layers. These points are illustrated in our “[Empirical evaluation](#)” section and further in the “[Appendix](#)”.

Layer entanglement in temporal multiplex networks

Analysis of temporal multiplex networks has shown promising results in multiple fields of science, such as for example healthcare and transportation (Sannino et al. 2017).

Since patterns of layer interaction networks result in typical entanglement values, considering temporal entanglement is related to sizing particular topologies of a temporal multiplex network. For example, a high intensity among members in a multiplex social network communicating through different social media corresponds to a synchronization of communications between them. When such a synchronization corresponds to the preparation of a particular event, understanding such synchronization could help forecast the event.

In this section, we first discuss how we define temporal multiplex networks and entanglement time series. We limit the following discussion to the consideration of entanglement between elementary layers only, i.e. only inner-layer edges.

Temporal multiplex networks and entanglement

Real-life networks often evolve over time, making them behave differently at different points. In our current setting, we define the temporal aspect of our network such as each edge e_t is defined at a specific time point t . A multiplex network M_d can then be defined for a given time window d . A time window $d = [t_0, t_f]$ covers a time frame (beginning at t_0 and ending t_f), and the multiplex network M_d is defined such as each edge exists within the time window:

$$M_d = (V_M, \{e_t \in E_M\}_{t \in d}).$$

The second scenario we considered is that of *moving time windows*. Here, edges from the *past* windows are considered when constructing a given network M , i.e.,

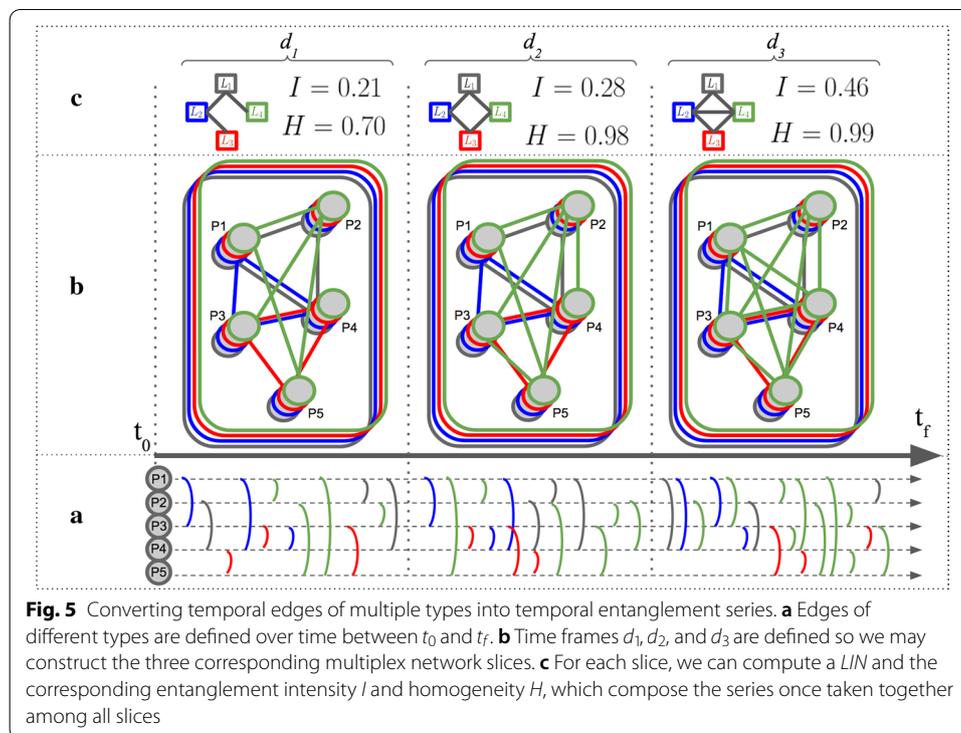
$$M_f = (V_M, \{e_t \in E_M\}_{t \in \{d-f, \dots, d-1\}}).$$

Our intuition is to compare the shape of a network at different moving time windows. For example, we could compare political social networks under different rulers of a country (Renoust et al. 2016a, b). To do so, we can simply compute entanglement homogeneity and intensity for each time window and compare them. Since our computation only focuses on edge, we consider the network as multiplex, the nodes are shared across all time frames.

Slicing the time windows is a very different topic and many options are open (Gomez et al. 2013; Beck et al. 2014). For example, it could be achieved manually, with equal time slices, moving window, or with *volume of changes*. In our context, we consider the identification of time window through slices of equal duration in time, but the principle can be extended. We refer to the duration r in time of the slices as *time resolution*.

We may now investigate entanglement homogeneity and intensity properties with respect to time resolution (r), and verify if patterns of intensity/homogeneity variation can be predicted. Note that one challenge of slice-based modelling of temporal multiplex networks is the problem of selecting the correct resolution r , i.e. how coarse (or fine)-grained the intervals must be in order to capture desired dynamics.

In a system covering a global period of D , once a slicing resolution is chosen, we can observe values of homogeneity and intensity at the time series level, i.e. for each slice $d \in D$, and define the intensity time series $\mathbb{S}_I = \{I_{M_d}\}, \forall d \in D$ and the homogeneity time series as $\mathbb{S}_H = \{H_{M_d}\}, \forall d \in D$. These intensity and homogeneity time series can now feed further processing. Note that \mathbb{S}_{I_f} and \mathbb{S}_{H_f} are defined analogously (entanglement for



the past f slices, moving in the increments of one slice). The whole processing from temporal edges to time series is illustrated in Fig. 5.

In our following evaluation (“Entanglement in temporal multiplex networks” section), we explore \mathbb{S}_I and \mathbb{S}_H when also considering a moving window of previous f time slices. The rationale for considering past f slices up to the considered time point is that such information only includes past data, and could indicate whether entanglement can be also used for *forecasting purposes*. The second option considered, where only the current time slice was plotted, can shed insight on whether online monitoring based on I or H is a sensible option.

Empirical evaluation

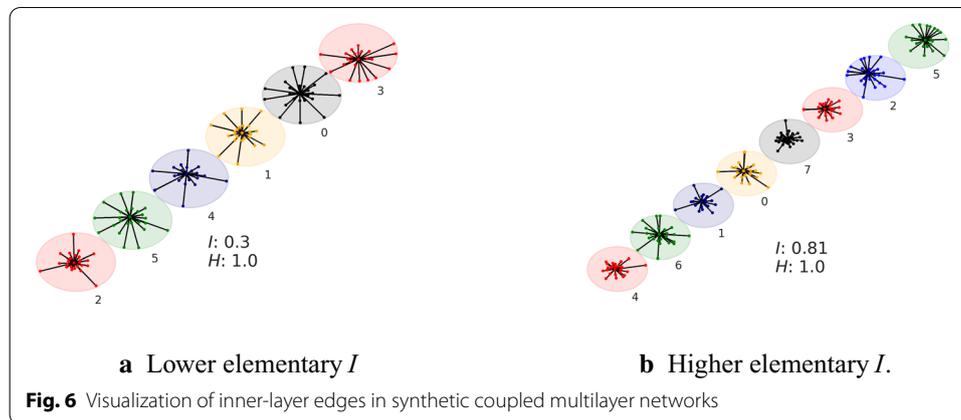
We now study entanglement intensity and homogeneity across different series of networks. We first investigate entanglement measures across different parameters of synthetic settings. We follow with investigations on a large panel of real world networks. We finish our study with the study of entanglement in temporal multiplex networks.

Entanglement in synthetic networks

In this first study, we compare entanglement measures over a series of synthetic multiplex networks, using our proposed generator.

We consider for all our generations, the following key parameters:

- Number of nodes (n) from 10 to 200 in increments of 10.
- Number of layers (m) in 1, 2, 3, 4, 6, 7, 9, 10.
- Layer assignment probability (o), from 0 to 1 in increments of 0.05



- Edge probability (p) from 0 to 1 in increments of 0.05.
- Transition coupling edge probability (q) from 0 to 1 in increments of 0.05.

Multiplex networks without transition coupling

A first generation concerns multiplex networks settings in which transition coupling is not specified (for example, friendship over different social platforms), so we do not consider parameter q here.

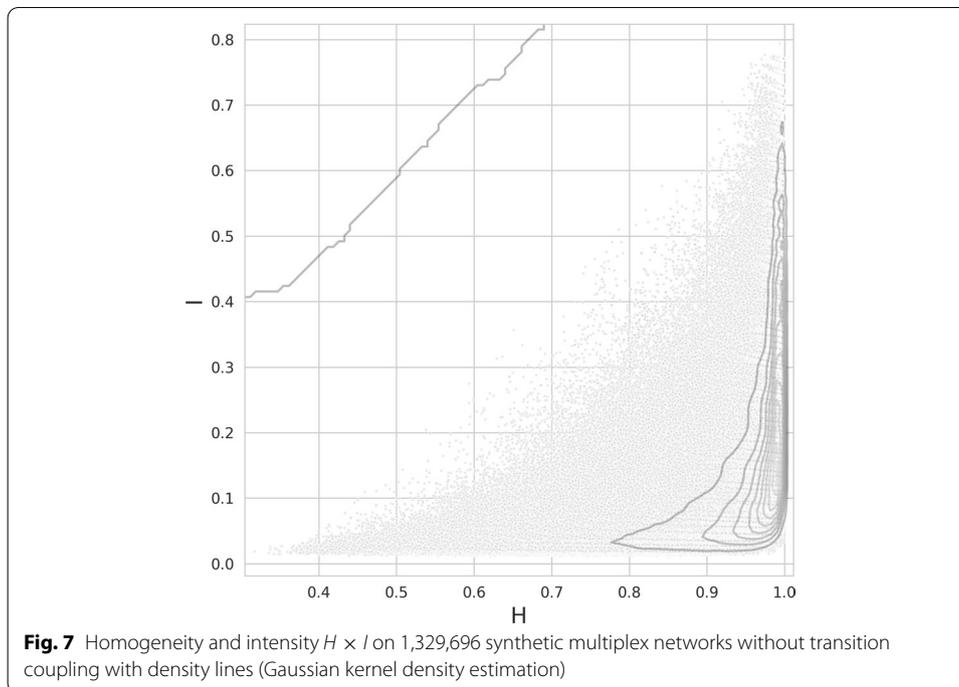
We have generated in total 1,329,696 synthetic networks (a couple are illustrated in Fig. 6).

We measure entanglement intensity I and homogeneity H on each generated network (averaged over all connected components of the layer overlap frequency matrix). We investigate the role of the different parameters over the entanglement measures, as illustrated in Figs. 7, 8 and 9.

There is an obvious dependency between entanglement intensity and homogeneity since we cannot obtain low homogeneity with high intensity values (Fig. 7). This is due to the nature of both measures. With a high intensity, most of the layers are overlapping over most of the network. As a consequence, there is little space for permutations in the way layers overlap, this means the entanglement of all individual layers γ_l tends to align, hence resulting in high values of homogeneity. This leads to a denser production of high homogeneity networks as illustrated by the density lines in Fig. 7.

The number of nodes n and edges m do not show a strong dependency with homogeneity, but a slight one on intensity. Higher values of n and m make it easier to obtain sparser networks, with the consequence of resulting lower values of intensity. We further illustrate these in Fig. 8. This effect mitigates quickly with higher numbers of nodes and layers.

We further explore the layer assignment probability of a node o , and the inner-layer edge probability p in Fig. 9. There is a first dependency appearing on the layer assignment probability o , for which higher values tend to produce higher homogeneity (Fig. 9b). Higher homogeneity is reached when all layers contribute equally, meaning that a higher o shows more chances for each layer to contain most of the nodes. We may also observe apparent linear trend between the edge probability p (sparseness) and



entanglement intensity (Fig. 9d). This trend confirms that sparser networks (i.e. lower p) are less “intensely” overlapping over edges. As intensity directly measures this property, this result outlines one of the *desired* properties of the proposed network generator.

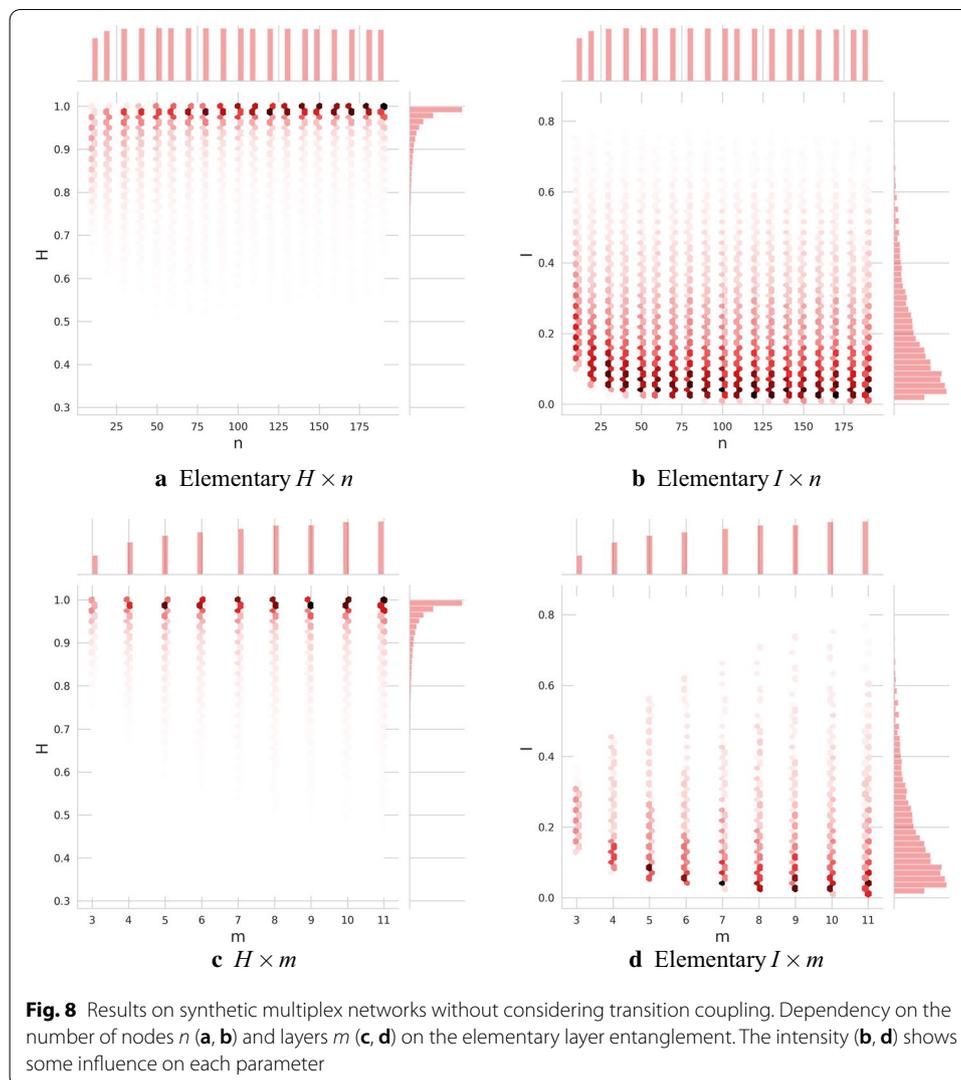
Multiplex networks with transition coupling

A second experiment is focusing on multiplex graphs with transition coupling, i.e. considering only the coupling edges in our 1,329,696 generated networks (illustrated in Fig. 10). This experiment reproduces the previous one, but focusing on the transition coupling entanglement. Results are shown in Figs. 11 and 12, dependency on the number of nodes and layers is illustrated in “Appendix”. From Fig. 11, the shape is globally the same, with the difference in a skewed density of high-homogeneity without a dense production of very low intensity generated networks (from the density lines).

The profile is sensibly the same than that of the previous experiment, except that the layer assignment probability o appears to have a more diffuse impact, and the direct dependency is this time observed on the coupling edge probability q . Comparison with parameter p obviously does not influence entanglement, but can be found in “Appendix” for additional inspection.

Overall, the networks with transition coupling are more *saturated* when compared to the ones without transition. The reason may be that we only consider here transition coupling edges that only connect the *same node* across layers.

For the interested reader, we also illustrate in the “Appendix” the independence of parameters q over the elementary layer entanglement and p over the transition coupling entanglement. We also report there the computation of entanglement over the combined elementary layers and transition coupling, which displays a dependency on both p and q parameters. Finally, we have computed the layer correlation coefficient, as



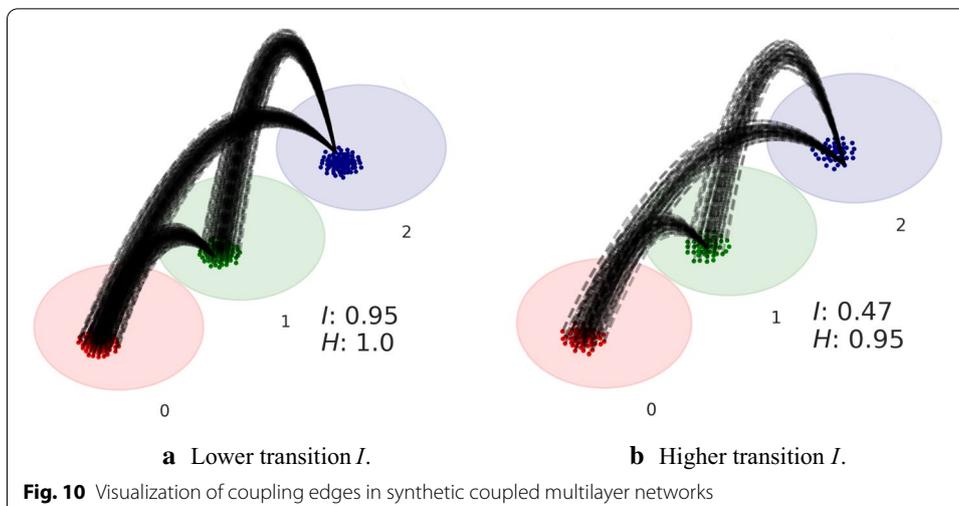
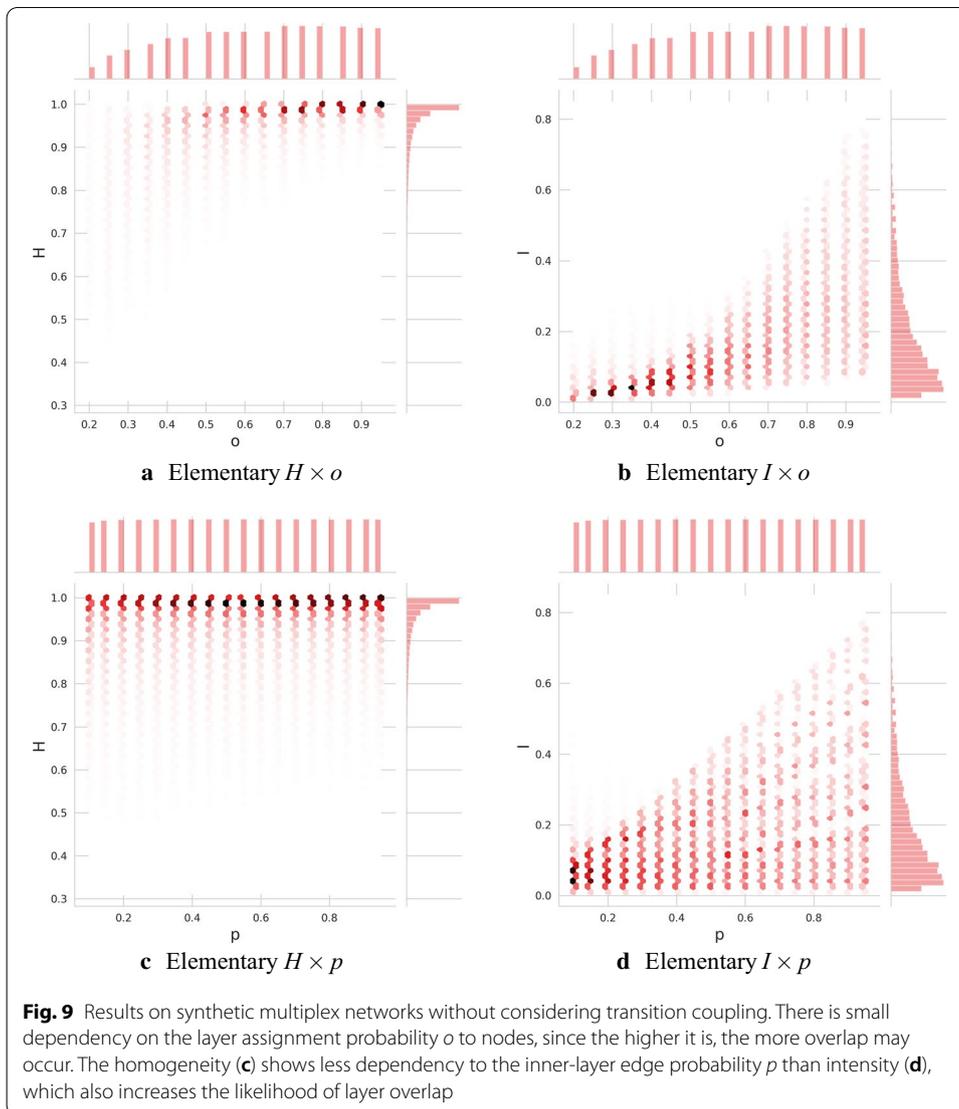
suggested in Nicosia and Latora (2015), confirming the role of the different parameters of our generator.

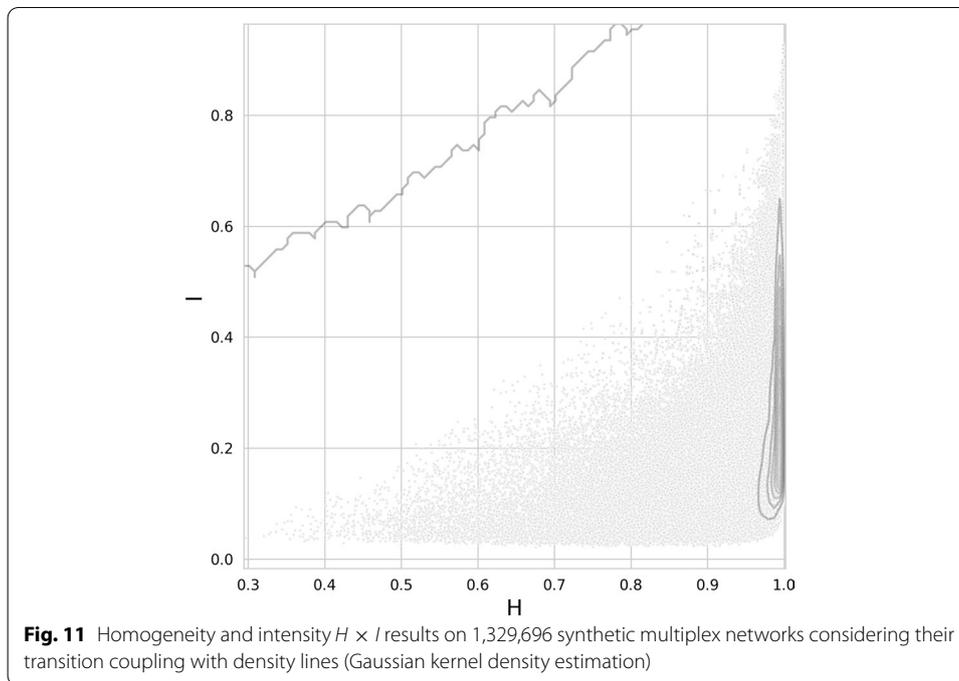
Multiplex network comparison across disciplines

We now consider real world static networks. All considered networks are summarised with their main characteristics in Table 1¹. Unfortunately, we have not found a real case with a large number of transition coupling edges, so we limit this evaluation to elementary layer entanglement. For each network, we computed elementary layer homogeneity and intensity for all connected components.

We first investigate individual results through the distributions of each metric across network types, Fig. 13. We then compare individual networks across entanglement intensity and homogeneity Fig. 14.

¹ The networks are hosted at <https://comunelab.fbk.eu/data.php>





Two main observations are apparent when studying the results on real networks. First, the difference between social and genetic (biological) multiplex networks becomes obvious when both entanglement intensity and homogeneity are considered (Fig. 14). To confirm these differences, we further compare their distributions, i.e., the intensity and homogeneity of social versus genetic networks, in Fig. 15.

In addition, from Fig. 14, we may observe that many genetic networks sit in relatively low intensity/homogeneity places, whereas social networks sit in the top right corner: the high entanglement homogeneity of social networks is quite noticeable. This suggests a few interpretations:

- genetic networks show in general very little layer overlap;
- some genetic networks could be matched to synthetic networks of low inner-layer edge probability, especially when homogeneity is low, being very sparse, potentially pointing at low layer assignment probability too;
- layers in social networks tend to overlap a lot;
- social networks tend to be quite dense and may be simulated by synthetic networks with a high inner-layer edge probability;

The results on social networks indicate a high level of layer overlap and it may be due to the overall behaviour of people, which is rather similar across different networks, whatever their means of interaction. *Simmelian ties*, *triadic closure*, and *homophily* (which are well studied in social sciences) are probably strong drivers of this layer overlap.

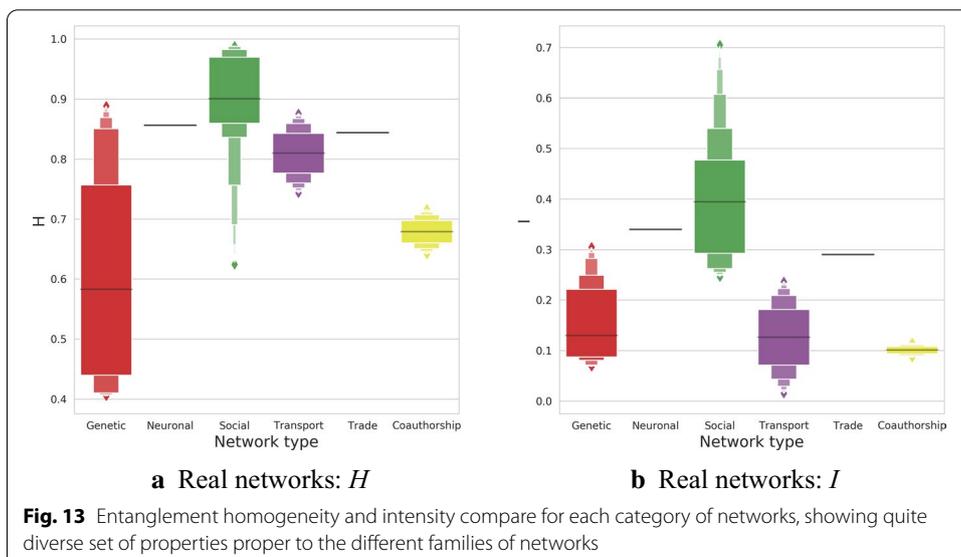
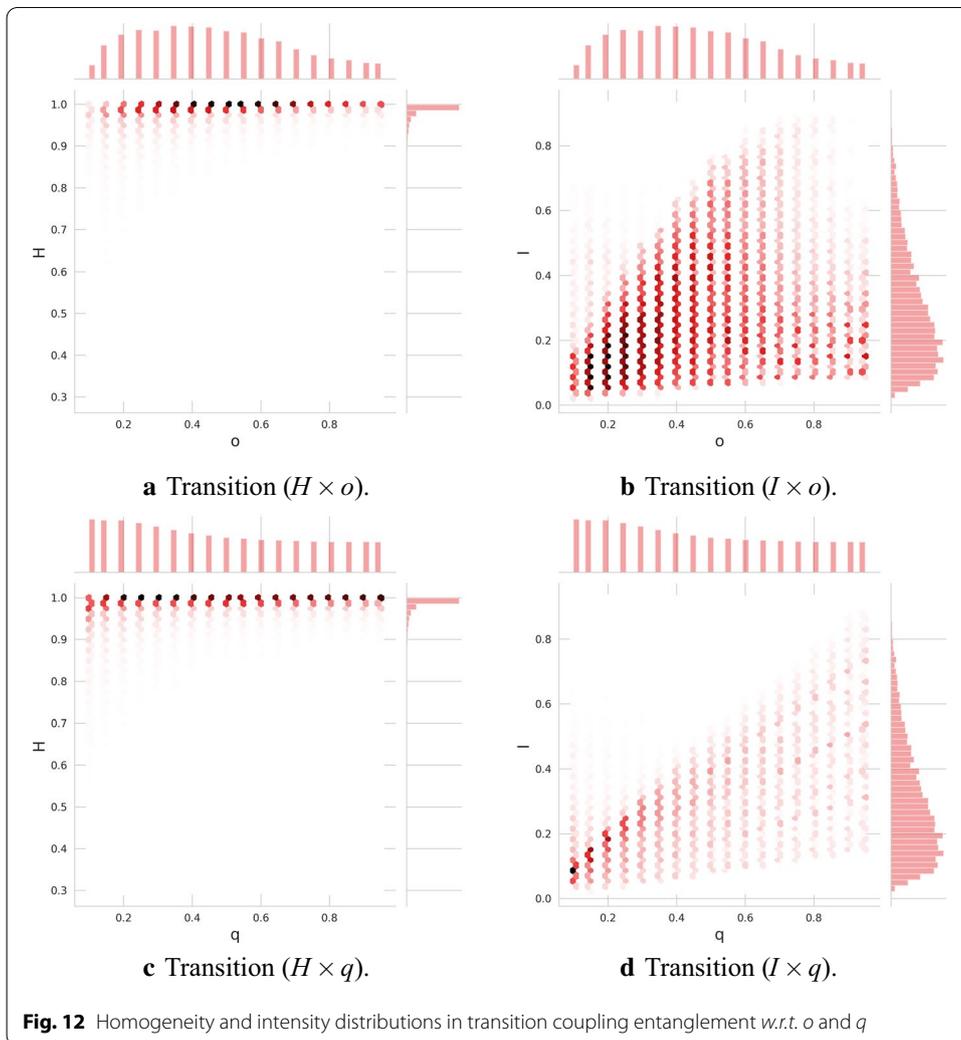
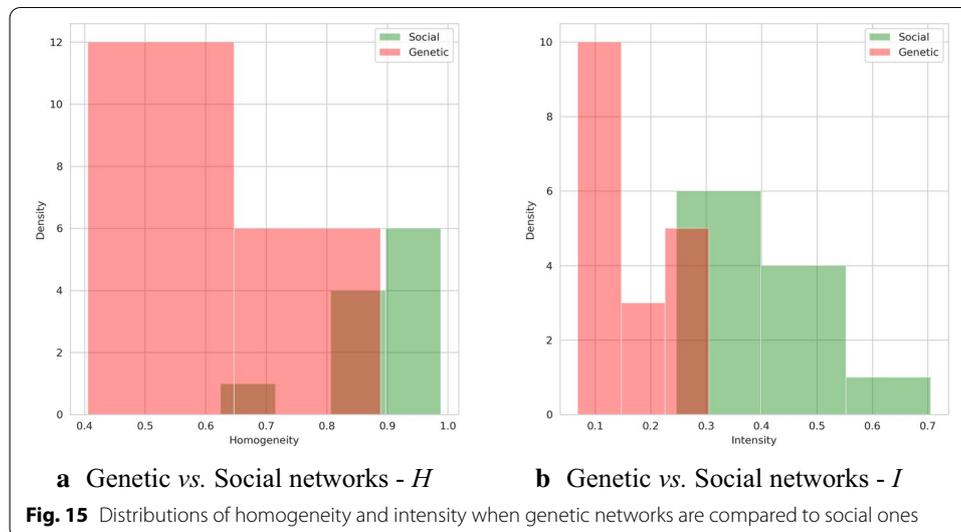


Table 1 Real multiplex networks and their properties

Dataset	ID	Type	Nodes	Edges	Number of layers	Mean degree	CC	Intensity	Homogeneity
arXiv-Netscience (De Domenico et al. 2015a)	34	Coauthorship	26,796	59,026	13	4.41	3660	0.114786	0.641670
PierreAuger (De Domenico et al. 2015a)	26	Coauthorship	965	7153	16	14.82	131	0.086551	0.716156
Arabidopsis (Stark et al. 2006)	0	Genetic	8765	18,655	7	4.26	387	0.111636	0.408940
Bos (Stark et al. 2006)	1	Genetic	369	322	4	1.75	82	0.160341	0.582015
Candida (Stark et al. 2006)	5	Genetic	418	398	7	1.90	50	0.284783	0.888476
Celegans (Stark et al. 2006)	7	Genetic	4557	8182	6	3.59	193	0.115718	0.420231
DanioRerio (Stark et al. 2006)	8	Genetic	180	188	5	2.09	45	0.068219	0.870304
Drosophila (Stark et al. 2006)	9	Genetic	11,970	43,367	7	7.25	346	0.082283	0.405509
Gallus (Stark et al. 2006)	12	Genetic	367	389	6	2.12	54	0.151845	0.433374
HepatitisCVirus (Stark et al. 2006)	13	Genetic	129	137	3	2.12	4	0.304679	0.777382
Homo Sapiens (Stark et al. 2006)	14	Genetic	36,194	170,899	7	9.44	785	0.101047	0.519648
HumanHerpes4 (Stark et al. 2006)	16	Genetic	261	259	4	1.98	21	0.245979	0.595037
HumanHIV1 (Stark et al. 2006)	15	Genetic	1195	1355	5	2.27	13	0.158347	0.583648
Oryctolagus (Stark et al. 2006)	24	Genetic	151	144	3	1.91	21	0.241322	0.635943
Plasmodium (Stark et al. 2006)	27	Genetic	1206	2522	3	4.18	27	0.249623	0.853694
Rattus (Stark et al. 2006)	28	Genetic	3263	4268	6	2.62	296	0.126889	0.457888
SacchCere (Stark et al. 2006)	29	Genetic	27,994	282,755	7	20.20	432	0.070428	0.695150
SacchPomb (Stark et al. 2006)	30	Genetic	10,178	63,677	7	12.51	286	0.079756	0.407135
Xenopus (Stark et al. 2006)	32	Genetic	582	620	5	2.13	109	0.082539	0.829466
YeastLandscape (Costanzo et al. 2010)	33	Genetic	17,770	8,473,997	4	953.74	4	0.132035	0.534030
CElegans (Chen et al. 2006)	7	Neuronal	791	5863	3	14.82	6	0.339461	0.856373
Cannes2013 Omodei et al. 2015	6	Social	659,951	991,854	3	3.01	48,375	0.269159	0.900587
CKM-Physicians-Innovation (Coleman et al. 1957)	3	Social	674	1551	3	4.60	12	0.394666	0.988309
CS-Aarhus (Magrani et al. 2013)	4	Social	224	620	5	5.54	13	0.341388	0.894766
Kapferer-Tailor-Shop (Kapferer 1972)	17	Social	150	1018	4	13.57	5	0.438509	0.910168



Entanglement in temporal multiplex networks

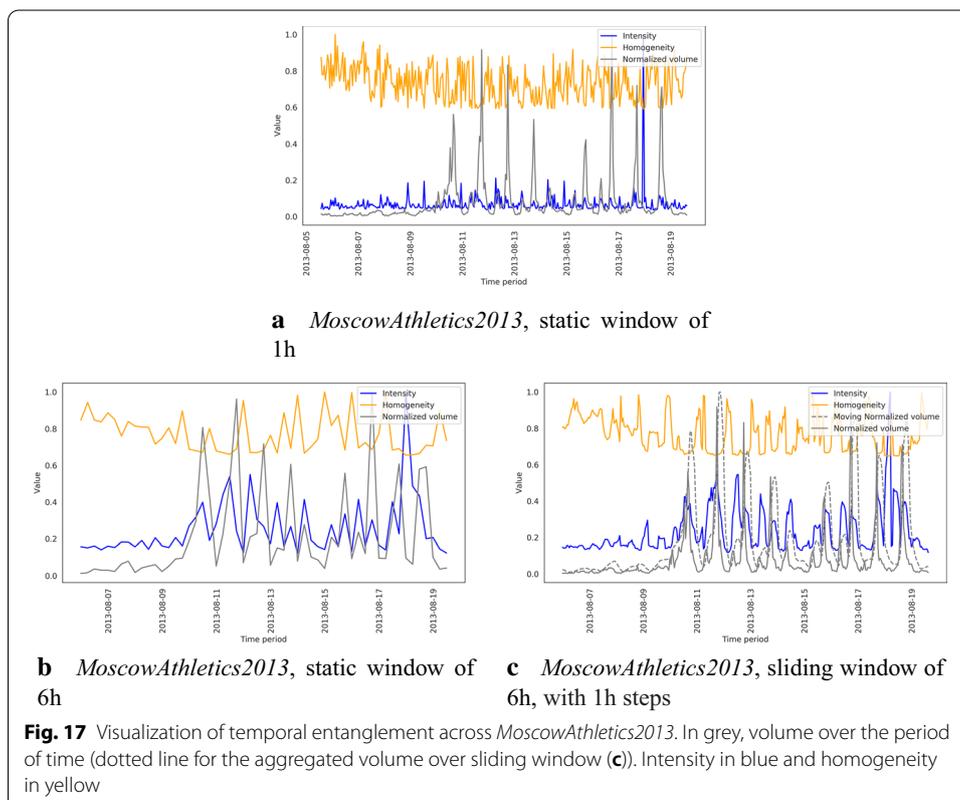
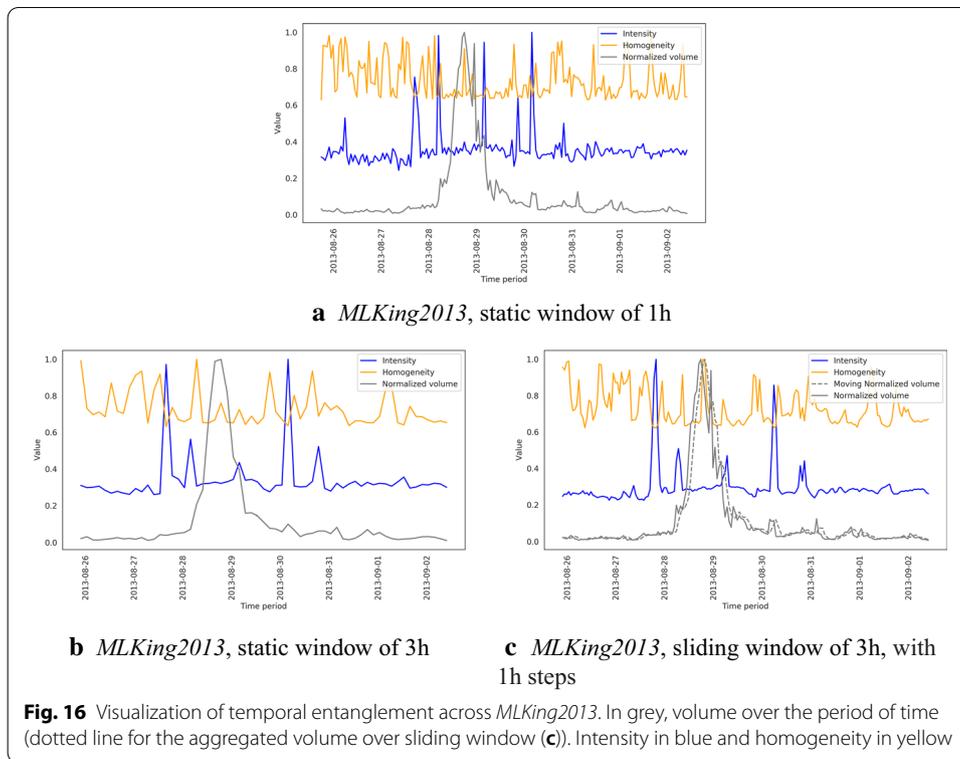
In our last experiment, we investigate entanglement across time slices of three real-life temporal multiplex networks: *MLKing2013*, *MoscowAthletics2013*, and *Cannes2013* [as found in Omodei et al. (2015)]. Each network consists in a collection of Twitter activity related to some event. The networks are comprised of three layers of connection, namely *retweets*, *replies* and *comments*. They can be summarised as follows. The *MLKing2013* data set consists of 421,083 events covering a week of celebration of M.L. King’s speech “*I have a dream*” in 2013, forming 396,671 edges between 327,708 nodes. The *MoscowAthletics2013* data set consists of 303,330 events covering two weeks of the World Championships of Athletics held in Moscow in 2013, forming 210,250 edges between 88,805 nodes. The *Cannes2013* network consists of 1,297,545 events (temporal edges) covering a month of the 2013 Cannes Film Festival, together forming a network of 930,419 edges and 438,538 nodes. Note that the networks are not trivially small, offering additional evidence of the stability of the entanglement computation.

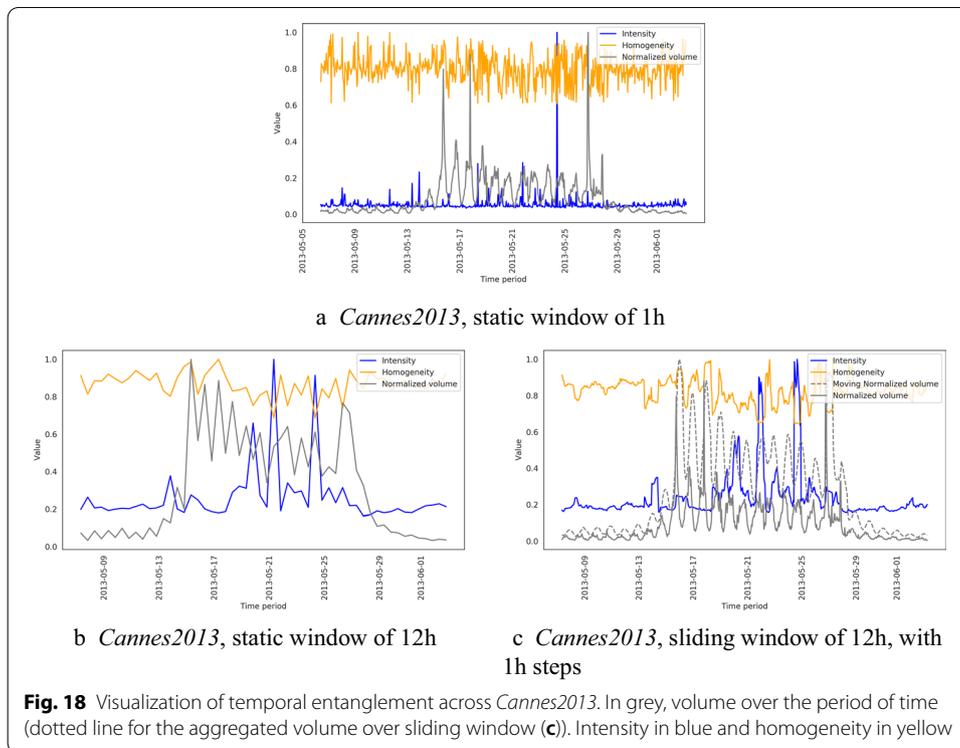
The networks were analysed following the methodology introduced in “[Layer entanglement in temporal multiplex networks](#)” section. We propose two experiments with regard to time segmentation.

The first experiment considers fixed time windows of sizes 1h, 3h, 6h, and 12h. We compare with the activity volume in form of H of a total number of tweets—as found in Omodei et al. (2015), Figure 1 for a 1h window size, here reported in Figs. 16a, 17a, and 18a. We normalise here this volume so values are in $[0, 1]$.

We selected the coarse windows at their best readability for each dataset (3h for *MLKing2013* in Fig. 16b, 6h for *MoscowAthletics2013* in Fig. 17b, and 12h for *Cannes2013* in Fig. 18b)—each coarsening is further illustrated in “[Appendix](#)”. A second experiment considers a moving window of the size corresponding to these best windows, sliding by the hours (Figs. 16c, 17c, and 18c).

In the *MLKing2013* data set (Fig. 16), we can observe that spikes of intensity surround the main spike of volume activity. A smaller spike of intensity consistently coincides with a smaller spike of volume at the end of the main spike.





In the *MoscowAthletics2013* data set (Fig. 17), the 1h-time window does not show a consistent behaviour. However, we can see that spikes in coarser time windows coincide with the spikes in volume. A larger spike in intensity appears before the final spike in volume.

In the *Cannes2013* data set (Fig. 18), the 1h-time window shows some spikes in intensity, especially a major by the end of the period of activity in terms of volume. In coarser time windows, we can notice four main spikes: one before the beginning of volume of activity; the next two ones appear just before a slight increase in the daily volume; the last one appears the day before the last day of the volume activity. This last peak appears even more prominent from the sliding window example.

The volume captures Twitter activity, governed by the human activity following the day/night rhythm. Although entanglement intensity is also submitted to it, we see emerging patterns that seem proper to each type of event. The activity of entanglement shows definitely some relationship with volume while telling a different story. The sports event that is *MoscowAthletics2013* may be much more subject to the day-by-day routine in which different disciplines are at play. On the other hand, the speech celebration in *MLKing2013* has some very specific activity before (could it be anticipation?) and after (could it be ripples?) the event. The movie festival in *Cannes2013* may be governed by sub-events of different importance in terms of networking activity.

In accordance with the position of social networks in our evaluation of real-world networks in “[Multiplex network comparison across disciplines](#)” section, we see a decrease

in homogeneity whenever we see spiking of intensity. This may indicate that a lot of the network activity suddenly focuses on one specific modality of exchange (such as *replies*). Entanglement study may help in targeting when this is driven by a particular modality.

Further studies on the nature of the events, and the specific topologies of the *LIN* networks that gave rise to these entanglement values is necessary for a more in-depth analysis of each case. Since we see some spiking activity of entanglement before actual events took place, we may suspect that, beyond monitoring, there is a *predictive power* of modelling time series from entanglement in past data (sliding windows).

Discussion and conclusions

In this work, we have revisited the notion of layer entanglement and extended it to coupled multilayer networks and temporal networks. To investigate entanglement, we have proposed a random generator for coupled multilayer networks, and generated a large set of synthetic ones. We have evaluated entanglement intensity and homogeneity in all cases, and compared to static and temporal real world networks.

Our analysis of the synthetic networks outlined that entanglement intensity is directly correlated with edge probability parameter—the sparser the network, the lower the intensity. This result indicates the proposed generator indeed emits networks which adhere to this property. We have also observed that large parts of the generated networks are subject to high homogeneity with various degrees of entanglement intensity.

Entanglement in the synthetic networks appears very sensitive to the different probabilities characterising the model (o , p , and eventually q for the coupled multilayer networks). The influence of each parameter should also be investigated theoretically in future work.

The high homogeneity observed may be a byproduct of our computations. First, our random generation induces a lot of small connected components of the coupled multilayer networks, and small components tend to show higher homogeneity since there are not so many degrees of freedom for edges to overlap. Because we are averaging the entanglement intensity and homogeneity over all components, this may go in favour of high homogeneity. Understanding this effect deserves more investigation. Second, entanglement homogeneity is a cosine measure, and the observed values may suffer from the skewness of cosine values when distributed in a linear space, amplifying the effect of having large values. Furthermore, it might also suffer from the curse of dimensionality in the case of a high number of layers. It would be worth considering normalizing this homogeneity with respect to the number of layers involved and the number of edges they cover. Instead of cosine, a Shannon's entropy measure may overcome some of these limitations.

One of the aspects that was not extensively evaluated as a part of this work is the processing of the repeated links in a given time slice. The current implementation considers, for each time slice, the collection of *unique* links, which are not weighted by their possible multiple occurrences. This way, the diversity of connections is emphasized, instead of the link frequency. A more detailed study of how the links can be re-weighted will be considered in future work.

We further demonstrated that the two measures offer interesting insights when computed across a wide array of real-world networks. The observed relationship between the intensity and homogeneity of layer entanglement with the family of data-set was previously reported for clusters of documents [in Renoust et al. (2013, Figure 5)]. In this previous experiments, clusters of documents were mostly located at the left frontier of high intensity for a varying homogeneity. Our current experiments showed that real networks cluster based on their type (e.g. biological vs. social), also close to this frontier. We have observed (from Fig. 14) that the set of genetic networks tend to sit in areas with low entanglement intensity, which could correspond to lower edge probability p , but they also tend to show a wider span of entanglement homogeneity including our lowest values measured (from Fig. 13), which could correspond to lower layer assignment probability o . Further work should be invested on finding the reason why genetic networks tend to show lower homogeneity. This is opposed to social networks which tend to find their way in the higher probability area. This should be further investigated, but this may be related to *homophily* (McPherson et al. 2001; Borgatti et al. 2009). Homophily is the implied similarity of two entities in a social network, and the property of entities to agglomerate when *being similar*. If the reason of *'being similar'* could be modelled as a layer of interaction, the result of a group of entities in *'being similar'* would lead to the formation of a clique in this layer, hence locating social networks in high probability areas.

The proposed work offers at least two prospects of multiplex network study which are in our belief worth exploring further. The difference between the genetic and social networks is possibly subject to very distinct topologies which emerge in individual layers. This claim may further be investigated via other measurements, such as graphlets, communities or other structures. Next, genetic networks are less homogeneous. Future work includes exploration of this fact, as it can be merely a property of the networks considered, empirical methodology used to obtain the networks or some other effect.

We believe that theoretical properties of the proposed network generator can also be further studied, offering potential insights into how multiplex networks behave and whether the human-made aspects are indeed representative of a given system's state. The model that we are currently exploring only takes into account a probability of linkage through (or within) layers without guarantee of connectivity. We made this choice to be able to compare between different fields, without prior assumption which could, for example, rule in favour of similarity to social network. Our future work will investigate other generation models including Erdős-Rényi-based (Caimo and Gollini 2020) or other with preferential attachment (Nicosia et al. 2014).

The analysis of real-life temporal networks offers cues on evolution in layer entanglement which can happen prior to some other events. We have tested multiple time scales. Too small time windows mostly result in noisy time series carrying low amounts of useful information, while higher coarsening shows activity related to volume, but with a different light on the events that are captured. Future work will

dive deeper into these events, and consider testing entanglement as a predictor using approaches such as of Prophet (Taylor and Letham 2018).

When considering entanglement as a either a monitoring or a predictive variable, its utility largely depends on the time scale at which a given edge stream needs to be considered. We leave extensive, possibly automatic determination of a setting where entanglement would be of practical relevance for future work. To study the parameters driving the dynamics of entanglement in temporal networks, we will consider comparing entanglement measures with synthetic temporal networks in our future investigations.

Abbreviations

LIN: Layer interaction Network.

Acknowledgements

We acknowledge Dagstuhl seminar 19061 where many ideas implemented in this paper emerged.

Authors' contributions

Both authors have contributed equally to the theoretical background, design of the experiments, and the writing of the manuscript. BR contributed to experiments, but the most of the experiments was handled by BS. Both authors read and approved the final manuscript.

Funding

The work of the first author was funded by the Slovenian Research Agency through a young researcher grant. The work of other authors was supported by the Slovenian Research Agency (ARRS) core research programme *Knowledge Technologies* (P2-0103) and ARRS funded research project *Semantic Data Mining for Linked Open Data* (financed under the ERC Complementary Scheme, N2-0078).

Availability of data and materials

The code for reproduction of experiments is freely available at <https://gitlab.com/skblaz/entanglement-multiplex>.

Competing interests

The authors declare that they have no competing interests.

Author details

¹ Jožef Stefan International Postgraduate School and Jožef Stefan Institute, Ljubljana, Slovenia. ² Osaka University Institute for Data Science, Osaka, Japan.

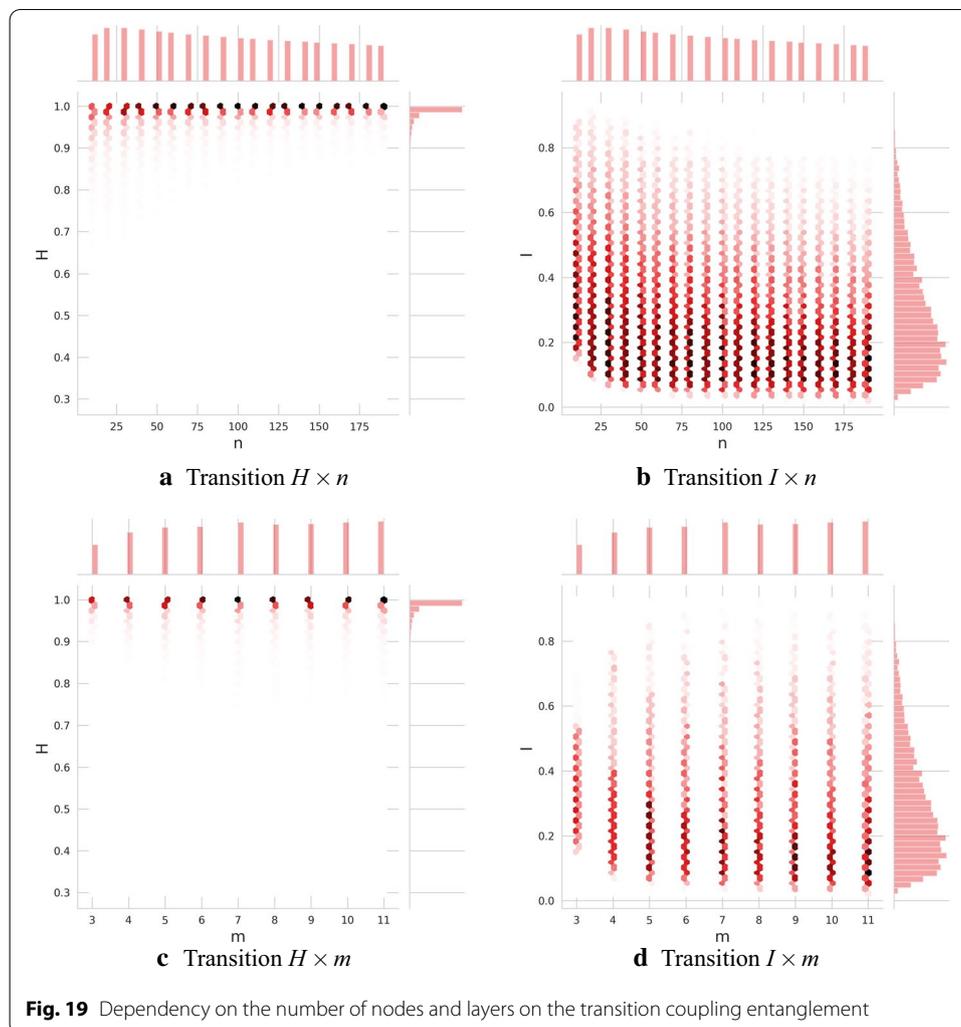
Appendix

Dependency in synthetic networks over nodes and layers

One can predict of course a level of dependency over the number of nodes n and layers m for the transition coupling case too. The dependency tends towards lower entanglement values since when increasing the number of nodes and layers, we increase the degree of freedom for layers to overlap. This trend, first illustrated in Fig. 8, is confirmed in Fig. 19.

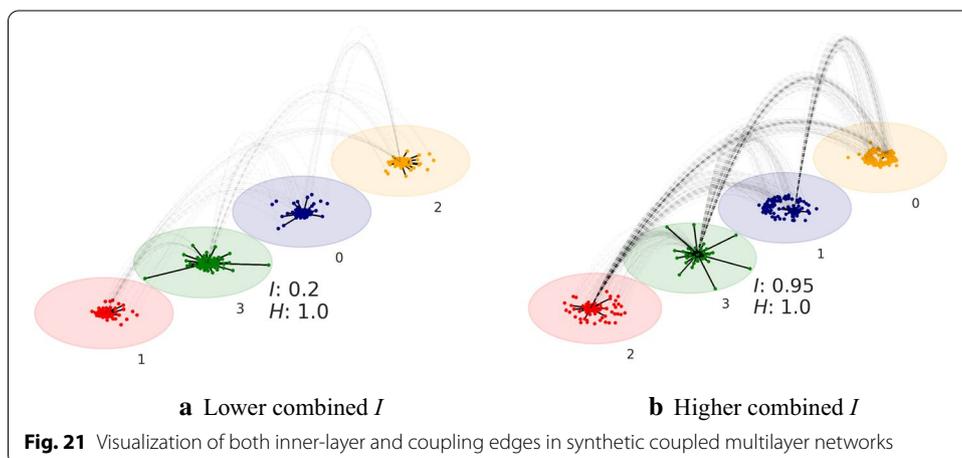
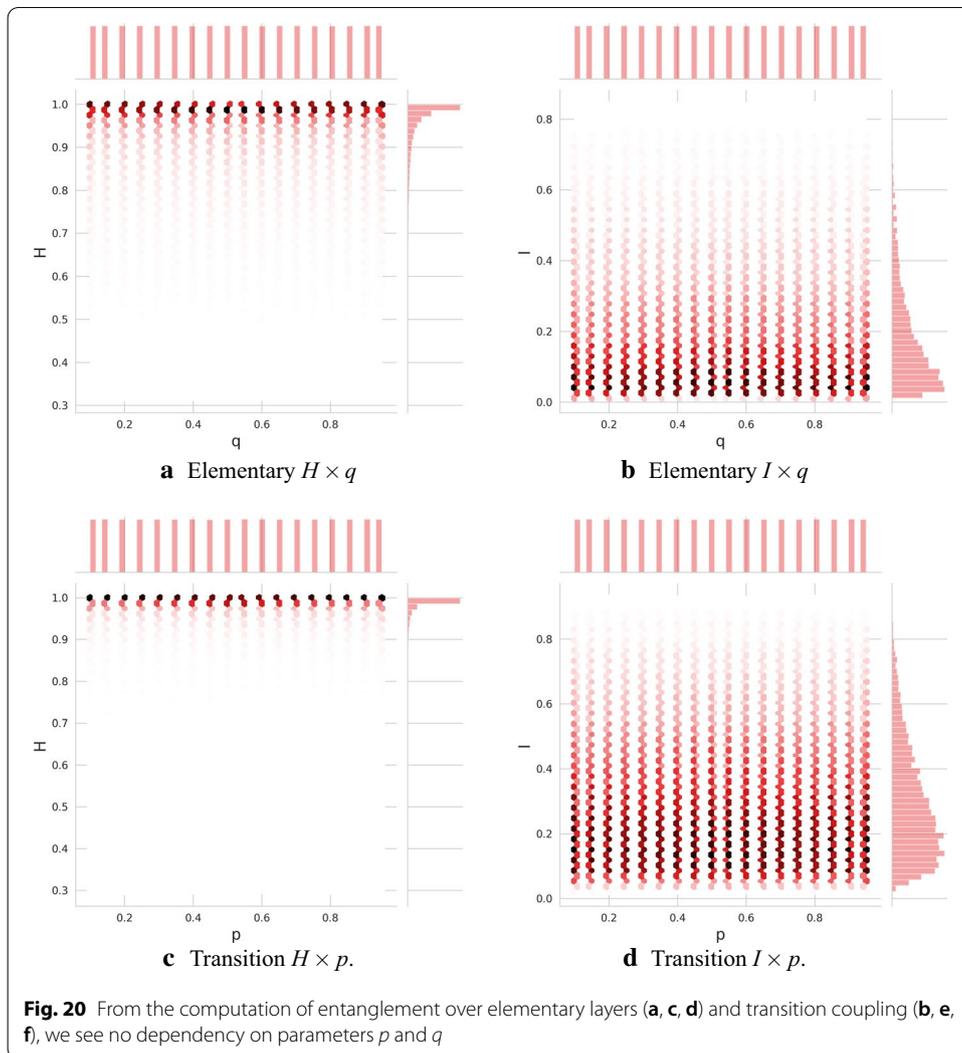
Independence of parameters

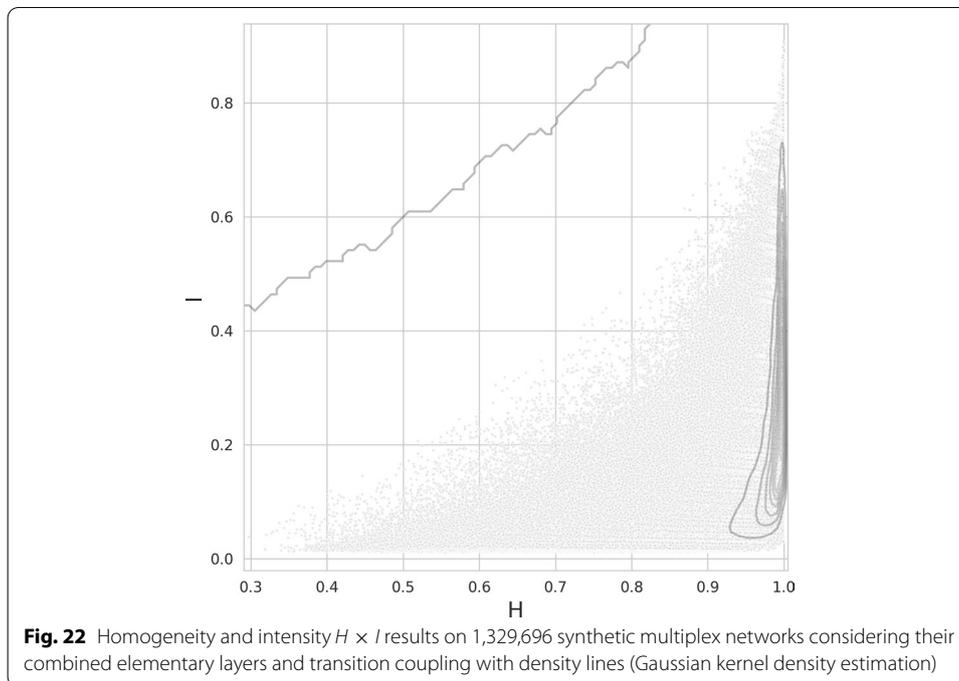
The distribution of parameters of transition coupling intensity and homogeneity over parameter p , and elementary layer intensity and homogeneity over parameter q , show no dependency as illustrated in Fig. 20.



Combining both elementary and transition coupling

As we mentioned in “[Transition coupling entanglement](#)” section, one can compute entanglement over all the network, combining elementary layers and transition coupling (as illustrated in Fig. 21). Although we have not identified practical use cases for this entanglement (often both categories of layers tell a different story), we report here the results over our synthetic networks in Figs. 22, 23, and 24. As expected we may observe a strong dependency over both p and q parameters combined (Fig. 24). Note that the current generator does not forbid the creation of loops enabling overlap between elementary layer and transition coupling. A generation of transition coupling edges that would connect *different nodes* between layers would create even more overlap between elementary layers and transition coupling. Such a parameter is actually available in the proposed code, but beyond the scope of this paper.





Choosing the right size of time window

Choosing the right size of time-window fundamentally depends on the dataset we observe. We report all variations of fixed time window coarsening we have explored, among 1h, 3h, 6h, and 12h-long windows for each of the *MLKing2013* (Fig. 25), *Moscow-Athletics2013* (Fig. 26), and *Cannes2013* (Fig. 27) events. Too fine selection displays a lot of noise, too coarse eludes most of the content.

Correlation analysis

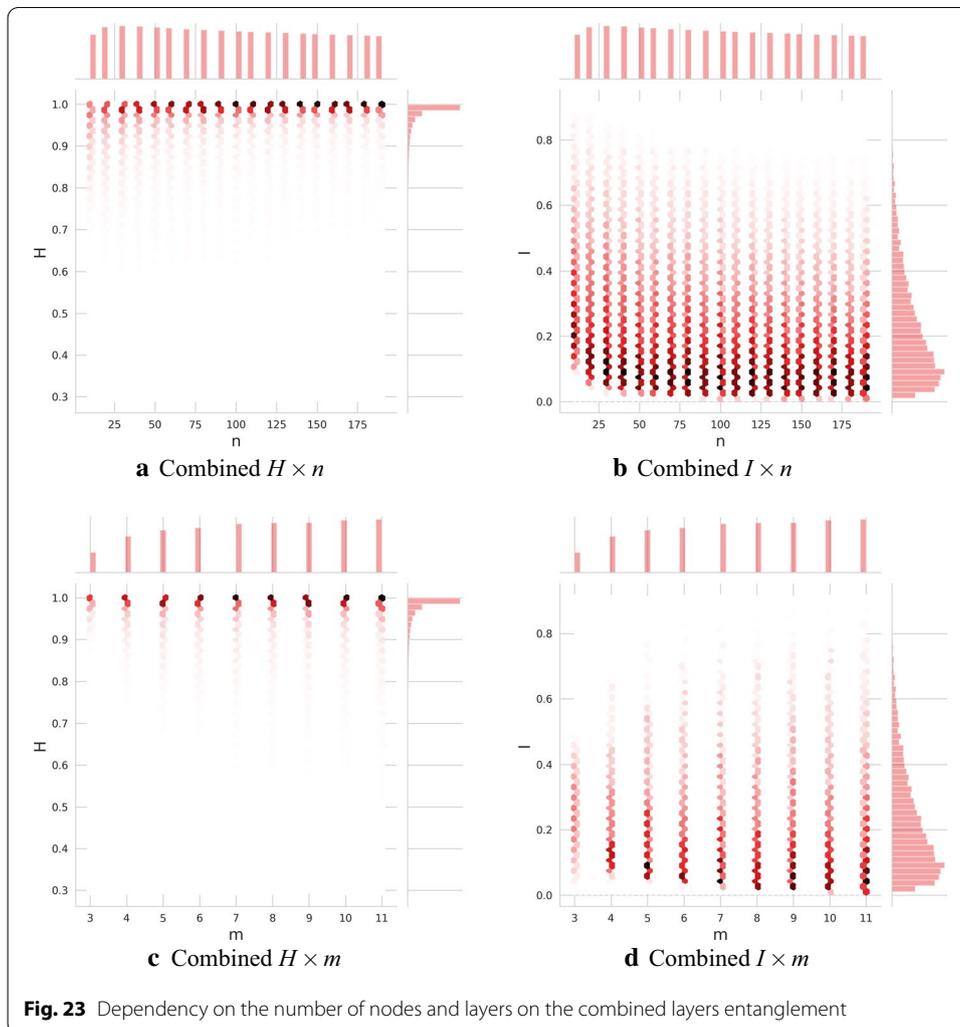
For completeness, we also computed the correlation between the occurrence of a given pair of nodes, as discussed in Nicosia and Latora (2015). As a measure of correlations of layer activity, we have computed the *pairwise multiplexity*. Following the original notation, we computed, for each synthetic network:

$$Q_{\alpha,\beta} = \frac{1}{|N|} \sum_i b_i^{[\alpha]} b_i^{[\beta]},$$

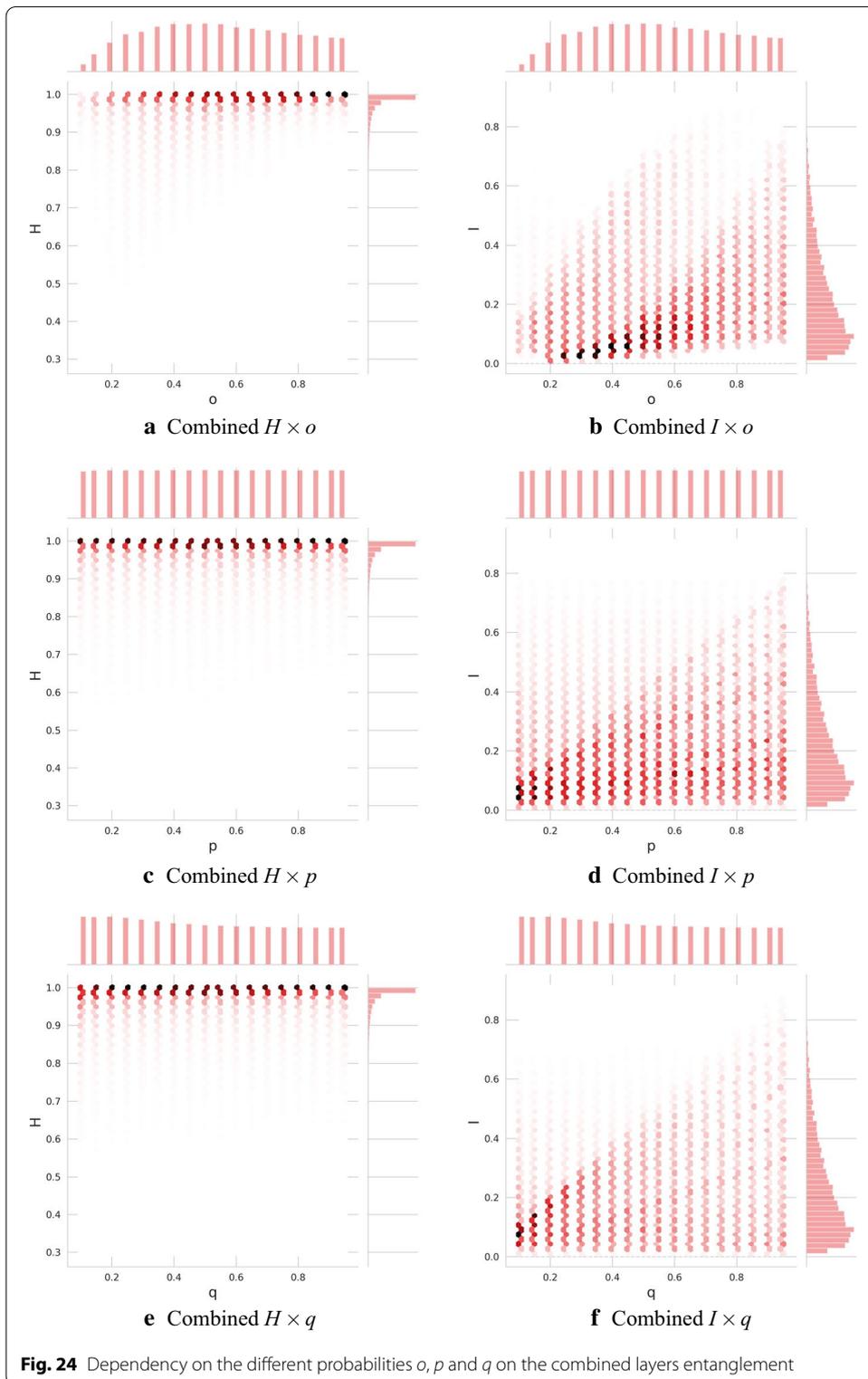
where $\alpha \in L$ and $\beta \in L$ are two distinct layers. Here, $b_i \in [0, 1]$ represents the presence of a given node i , hence, the product equals zero if a given node is not co-present on both considered layers. Given the large space of synthetic multiplex networks, it is not sensible to analyse individual ones, as performed in Nicosia et al. (2014), hence we further computed:

$$Q_d = \{Q_{\alpha,\beta}\}_{\alpha,\beta \in L^2 \wedge (\alpha \neq \beta)},$$

i.e., the distribution of all possible pairwise correlations. In Fig. 28, we present the statistical properties of this distribution further segmented according to the number of layers.



We further compare this distribution with each parameter of our generator in Fig. 29. We first could expect some level of relationship with the number of layers and the node layer probability. This is confirmed, and we can observe a direct dependency with the number of layers m , and the node-layer assignment probability o . The relationship with the number of layers slowly decreases with a minimum of 5 layers, since with more layers, there are more degrees of freedom for nodes to be assigned on layers. The relationship with o is almost linear. There is almost no dependency with the number of nodes (once enough nodes are assigned). The inner-layer edge probability p does not influence this measure, although we may observe a sudden increase first, it is an artefact of our algorithm because nodes with degree zero in a layer are discarded from the layer to speed up computations. The transition coupling edge probability q shows no correlation because this algorithmic artefact does not apply to it.



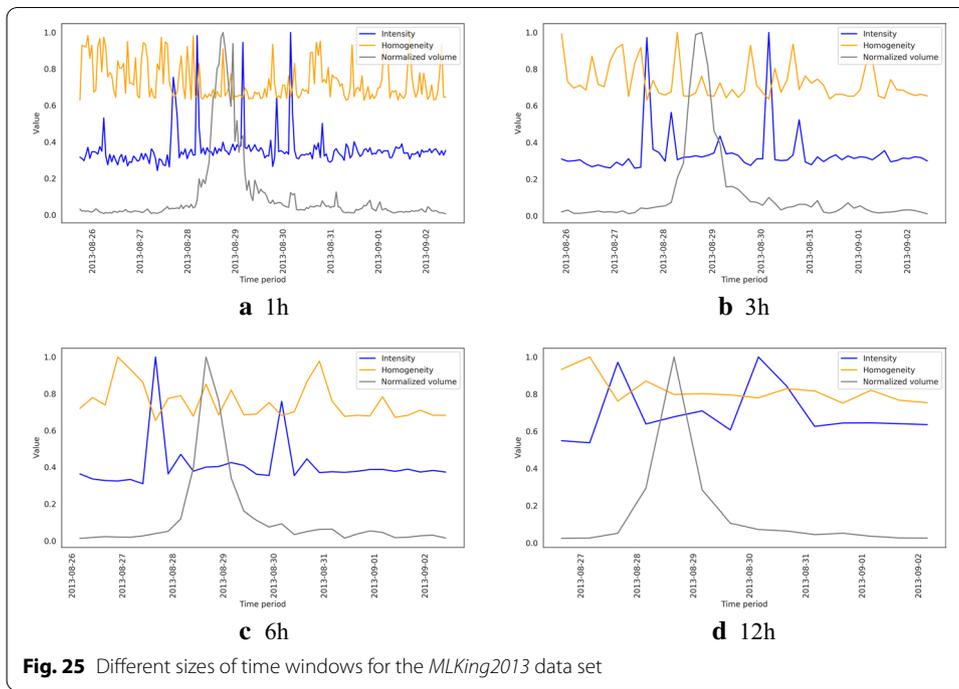


Fig. 25 Different sizes of time windows for the *MLKing2013* data set

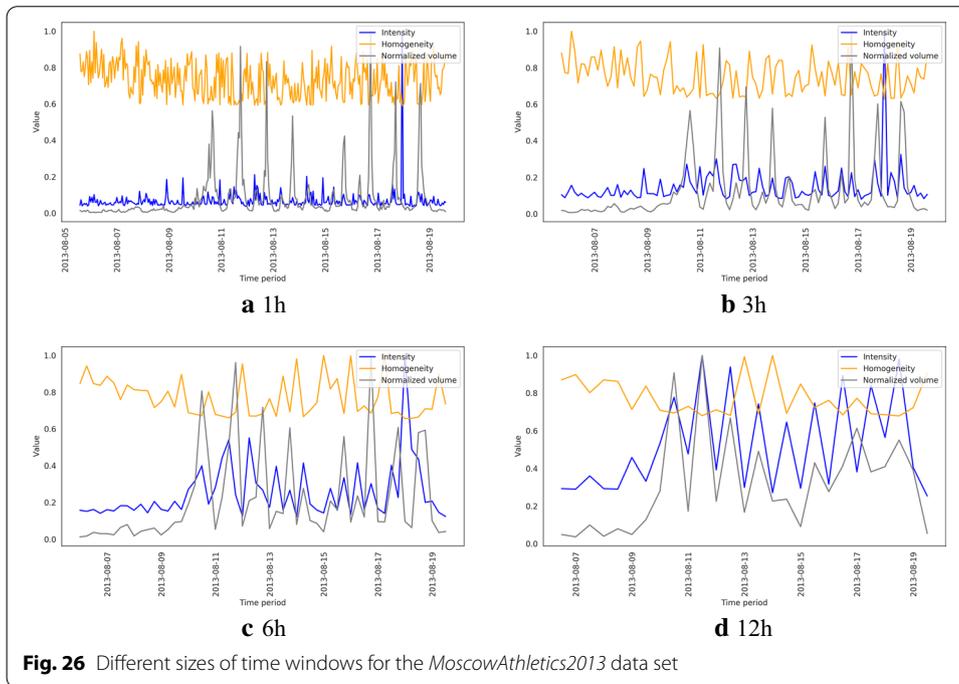


Fig. 26 Different sizes of time windows for the *MoscowAthletics2013* data set

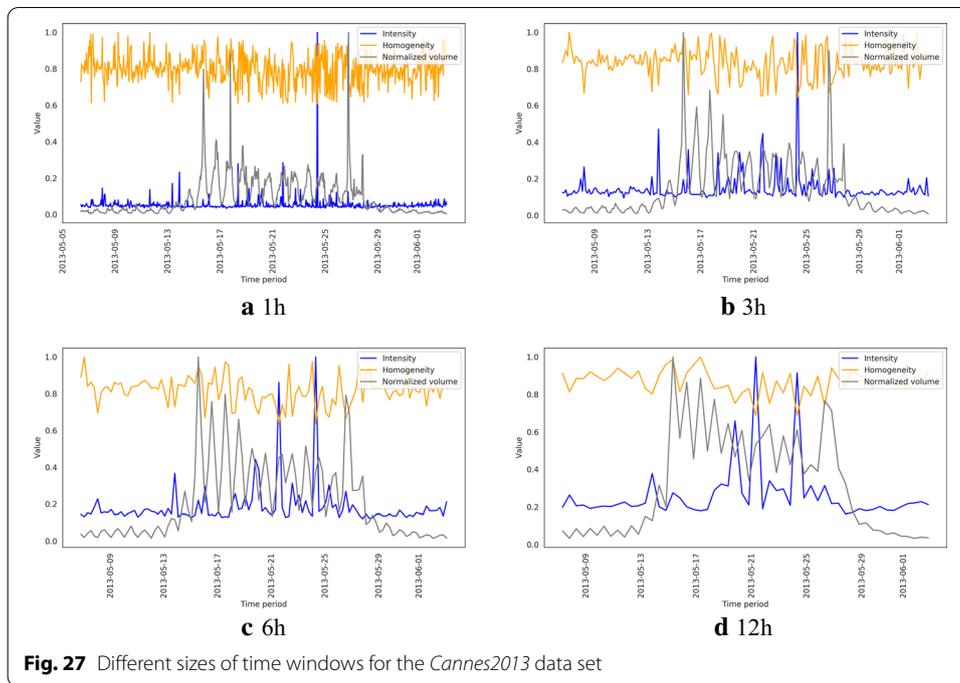


Fig. 27 Different sizes of time windows for the *Cannes2013* data set

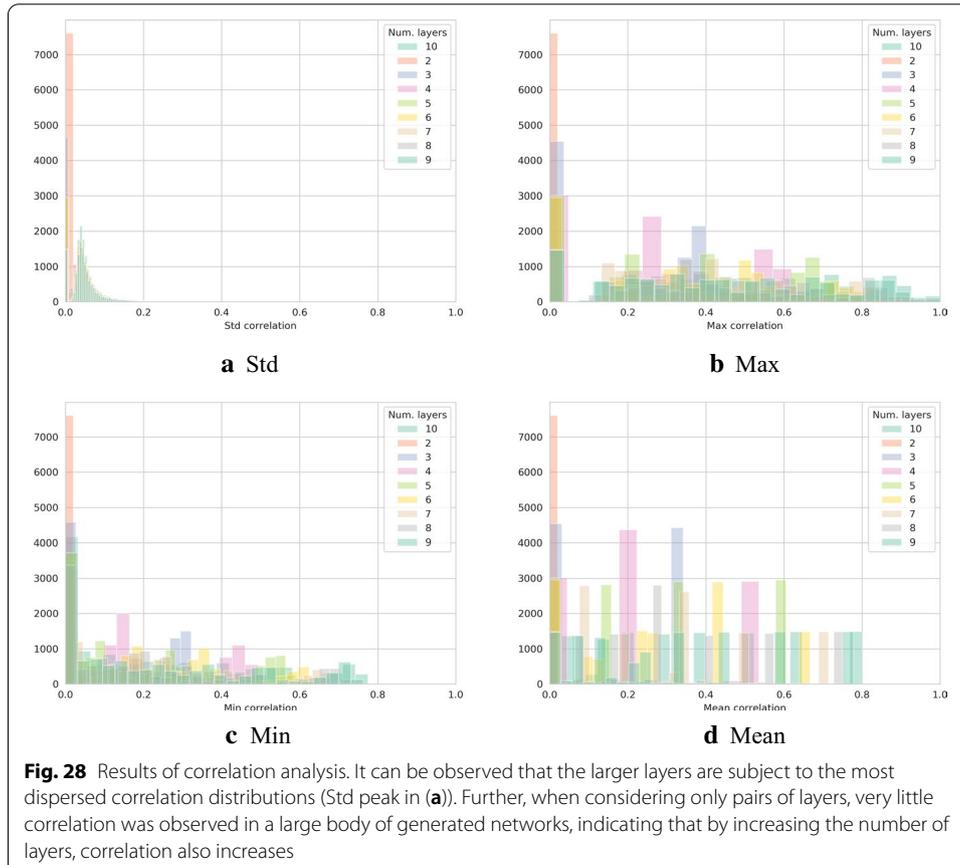
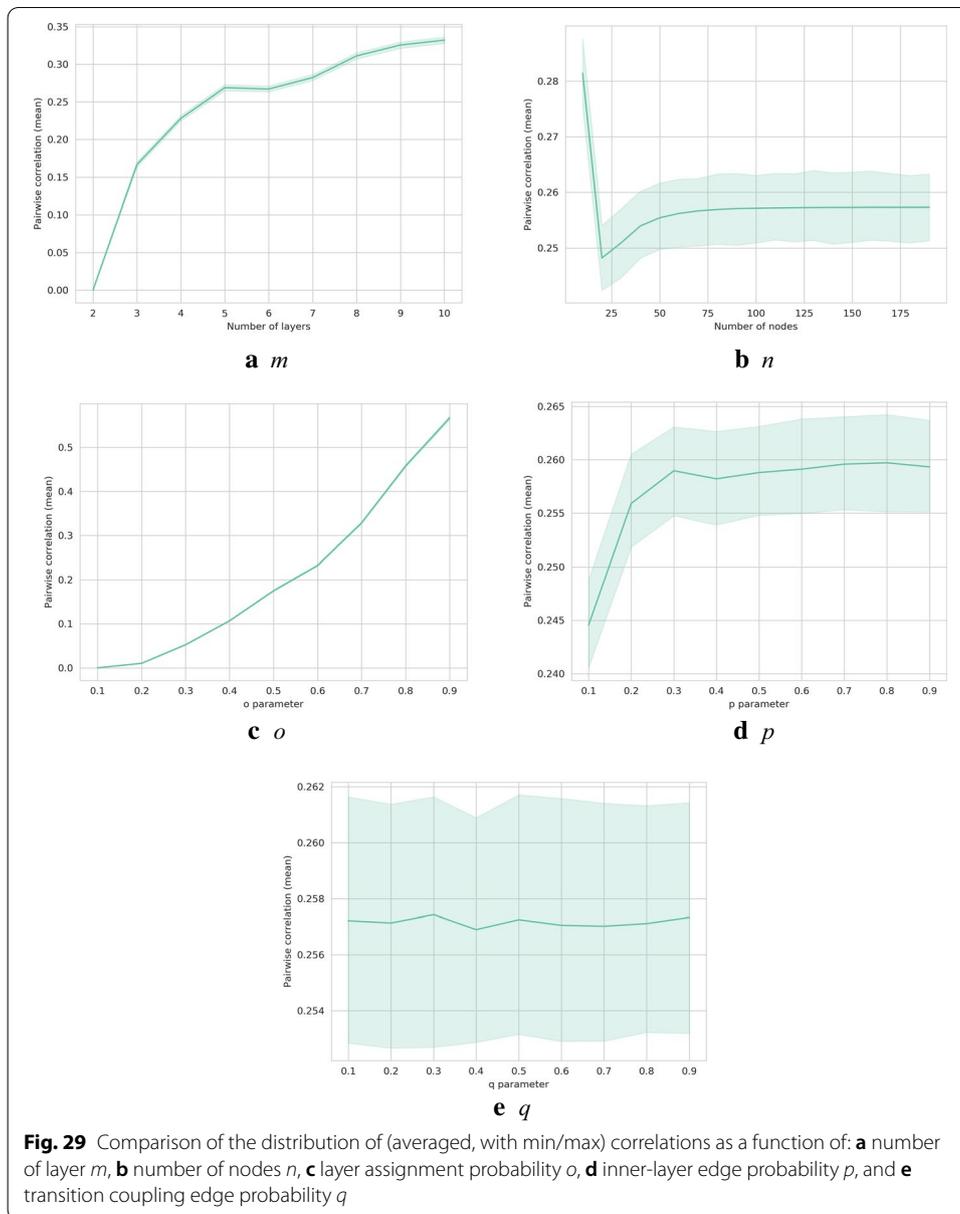


Fig. 28 Results of correlation analysis. It can be observed that the larger layers are subject to the most dispersed correlation distributions (Std peak in (a)). Further, when considering only pairs of layers, very little correlation was observed in a large body of generated networks, indicating that by increasing the number of layers, correlation also increases



Received: 8 April 2020 Accepted: 23 October 2020

Published online: 11 November 2020

References

Battiston F, Nicosia V, Latora V (2014) Structural measures for multiplex networks. *Phys Rev E* 89(3):032804
 Beck F, Burch M, Diehl S, Weiskopf D (2014) The state of the art in visualizing dynamic graphs. In: *EuroVis (STARS)*. Citeseer
 Borgatti SP, Mehra A, Brass DJ, Labianca G (2009) Network analysis in the social sciences. *Science* 323(5916):892–895
 Burt RS, Schött T (1985) Relation contents in multiple networks. *Soc Sci Res* 14(4):287–308
 Caimo A, Gollini I (2020) A multilayer exponential random graph modelling approach for weighted networks. *Comput Stat Data Anal* 142:106825
 Cardillo A, Gómez-Gardenes J, Zanin M, Romance M, Papo D, Del Pozo F, Boccaletti S (2013) Emergence of network features from multiplexity. *Sci Rep* 3:1344

- Chen BL, Hall DH, Chklovskii DB (2006) Wiring optimization can relate neuronal structure and function. *Proc Natl Acad Sci* 103(12):4723–4728
- Chen X, Wang R, Tang M, Cai S, Stanley HE, Braunstein LA (2018) Suppressing epidemic spreading in multiplex networks with social-support. *New J Phys* 20(1):013007
- Coleman J, Katz E, Menzel H (1957) The diffusion of an innovation among physicians. *Sociometry* 20(4):253–270
- Costanzo M, Baryshnikov A, Bellay J, Kim Y, Spear ED, Sevier CS, Ding H, Koh JL, Toufighi K, Mostafavi S et al (2010) The genetic landscape of a cell. *Science* 327(5964):425–431
- Cozzo E, Moreno Y (2016) Characterization of multiple topological scales in multiplex networks through supra-Laplacian eigengaps. *Phys Rev E* 94(5):052318
- Cozzo E, Kivelä M, De Domenico M, Solé-Ribalta A, Arenas A, Gómez S, Porter MA, Moreno Y (2015) Structure of triadic relations in multiplex networks. *New J Phys* 17(7):073029
- De Domenico M, Solé-Ribalta A, Gómez S, Arenas A (2014) Navigability of interconnected networks under random failures. *Proc Natl Acad Sci* 111(23):8351–8356
- De Domenico M, Lancichinetti A, Arenas A, Rosvall M (2015a) Identifying modular flows on multilayer networks reveals highly overlapping organization in interconnected systems. *Phys Rev X* 5(1):011027
- De Domenico M, Nicosia V, Arenas A, Latora V (2015b) Structural reducibility of multilayer networks. *Nat Commun* 6:6864
- Gomez S, Diaz-Guilera A, Gomez-Gardenes J, Perez-Vicente CJ, Moreno Y, Arenas A (2013) Diffusion dynamics on multiplex networks. *Phys Rev Lett* 110(2):028701
- Kapferer B (1972) Strategy and transaction in an African factory: African workers and Indian management in a Zambian Town. University Press, Manchester
- Kivelä M, Arenas A, Barthelemy M, Gleeson JP, Moreno Y, Porter MA (2014) Multilayer networks. *J Complex Netw* 2(3):203–271
- Kivelä M, McGee F, Melançon G, Henry Riche N, von Landesberger T (2019) Visual analytics of multilayer networks across disciplines (seminar 19061). Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik
- Krackhardt D (1987) Cognitive social structures. *Soc Netw* 9(2):109–134
- Lazega E et al (2001) The collegial phenomenon: the social mechanisms of cooperation among peers in a corporate law partnership. Oxford University Press, Oxford
- Magnani M, Micenkova B, Rossi L (2013) Combinatorial analysis of multiple networks. arXiv preprint arXiv:1303.4986
- McPherson M, Smith-Lovin L, Cook JM (2001) Birds of a feather: homophily in social networks. *Annu Rev Sociol* 27(1):415–444
- Mittal R, Bhatia M (2019) Analysis of multiplex social networks using nature-inspired algorithms. In: Nature-inspired algorithms for big data frameworks. IGI Global, pp 290–318
- Nicosia V, Latora V (2015) Measuring and modeling correlations in multiplex networks. *Phys Rev E* 92(3):032805
- Nicosia V, Bianconi G, Latora V, Barthelemy M (2014) Nonlinear growth and condensation in multiplex networks. *Phys Rev E* 90(4):042807
- Omodei E, De Domenico MD, Arenas A (2015) Characterizing interactions in online social networks during exceptional events. *Front Phys* 3:59
- Padgett JF, Ansell CK (1993) Robust action and the rise of the medici, 1400–1434. *Am J Sociol* 98(6):1259–1319
- Renoust B, Melançon G, Viaud M.-L. (2013) Measuring group cohesion in document collections. In: 2013 IEEE/WIC/ACM international joint conferences on web intelligence (WI) and intelligent agent technologies (IAT), vol 1. IEEE, pp 373–380
- Renoust B, Melançon G, Viaud M.-L. (2014) Entanglement in multiplex networks: understanding group cohesion in homophily networks. In: Missaoui R, Sarr I (eds) Social network analysis. Springer, Berlin, pp 89–117
- Renoust B, Melançon G, Munzner T (2015) Detangler: visual analytics for multiplex networks. *Computer Graphics Forum* 34(3):321–330
- Renoust B, Kobayashi T, Ngo TD, Le D-D, Satoh S (2016a) When face-tracking meets social networks: a story of politics in news videos. *Appl Netw Sci* 1(1):4
- Renoust B, Le D-D, Satoh S (2016b) Visual analytics of political networks from face-tracking of news video. *IEEE Trans Multimedia* 18(11):2184–2195
- Sannino S, Stramaglia S, Lacasa L, Marinazzo D (2017) Visibility graphs for fMRI data: multiplex temporal graphs and their modulations across resting-state networks. *Netw Neurosci* 1(3):208–221
- Škrlić B, Kralj J, Lavrač N (2019) Cbssd: Community-based semantic subgroup discovery. *J Intell Inf Syst* 53(2):265–304
- Škrlić B, Renoust B (2019) Patterns of multiplex layer entanglement across real and synthetic networks. In: International conference on complex networks and their applications. Springer, pp 671–683
- Stark C, Breitkreutz BJ, Reguly T, Boucher L, Breitkreutz A, Tyers M (2006) Biogrid: a general repository for interaction datasets. *Nucleic Acids Res* 34(suppl-1):535–539
- Taylor SJ, Letham B (2018) Forecasting at scale. *Am Stat* 72(1):37–45
- Tejedor A, Longjas A, Fofoula-Georgiou E, Georgiou TT, Moreno Y (2018) Diffusion dynamics and optimal coupling in multiplex networks with directed layers. *Phys Rev X* 8(3):031071
- Valdeolivas A, Tichit L, Navarro C, Perrin S, Odelin G, Levy N, Cau P, Remy E, Baudot A (2018) Random walk with restart on multiplex and heterogeneous biological networks. *Bioinformatics* 35(3):497–505
- Vickers M, Chan S (1981) Representing classroom social structure. Victoria Institute of Secondary Education, Melbourne
- Wang W, Cai M, Zheng M (2018) Social contagions on correlated multiplex networks. *Physica A* 499:121–128
- Wasserman S, Faust K (1994) Social Network analysis, methods and applications. Structural analysis in the social sciences. Cambridge University Press, Cambridge. <https://doi.org/10.1017/CBO9780511815478>

Publisher's Note

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