# Asymmetry of social interactions and its role in link predictability: the case of coauthorship networks

Kamil P. Orzechowski, Maciej J. Mrowinski, Agata Fronczak and Piotr Fronczak

*Faculty of Physics, Warsaw University of Technology, Koszykowa 75, PL-00-662 Warsaw, Poland*

The paper provides important insights into understanding the factors that influence tie strength in social networks. Using local network measures that take into account asymmetry of social interactions we show that the observed tie strength is a kind of compromise, which depends on the relative strength of the tie as seen from its both ends. This statement is supported by the Granovetter-like, strongly positive weight-topology correlations, in the form of a power-law relationship between the asymmetric tie strength and asymmetric neighbourhood overlap, observed in three different real coauthorship networks and in a synthetic model of scientific collaboration. This observation is juxtaposed against the current misconception that coauthorship networks, being the proxy of scientific collaboration networks, contradict the Granovetter's strength of weak ties hypothesis, and the reasons for this misconception are explained. Finally, by testing various link similarity scores, it is shown that taking into account the asymmetry of social ties can remarkably increase the efficiency of link prediction methods. The perspective outlined also allows us to comment on the surprisingly high performance of the resource allocation index – one of the most recognizable and effective local similarity scores – which can be rationalized by the strong triadic closure property, assuming that the property takes into account the asymmetry of social ties.

## I. INTRODUCTION

Social networks representing patterns of human interactions have been the subject of both empirical and theoretical research since at least the middle of the last century [1, 2]. And although, over the past two decades, due to the rise of the Internet followed by the increased availability of large datasets on human interactions, methods of social network analysis have changed a lot, the basic challenge behind these analyses remained the same: To understand human behaviour [3–5]. In particular, questions that keep recurring in the literature of the field are: How, depending on the context studied, people establish social ties? To what extent can the evolution of a social network be modelled using features intrinsic to the network itself? Is it possible to predict existing but undisclosed or intentionally hidden connections based on those recorded? Finally, what influences the strength of social ties, and can these strengths be inferred from the binary link structure?

In what follows, building upon results of our previous paper [6], we refer to last two of the above questions. We show that in order to better understand weight-topology correlations in social networks, it is necessary to use measures that formally take into account asymmetry of social interactions, which may arise, for example, from differences in ego-networks of connected nodes. A simple argument in favour of this statement can be drawn from the theory of complex networks (more specifically, from the degree-based mean field approach [7–9]). In particular, in social networks with fat tailed node degree distributions the sizes of ego-networks of two connected nodes may differ considerably. This means that their common neighbours can be a significant part of the neighbourhood of one node and an insignificant part of the neighbourhood of the other, resulting in a completely different perception of the size of the common neighbourhood on both ends of the connection. This indicates that the observed *absolute* tie strength is a kind of compromise, which depends on the *relative* strength of the tie as seen from its both ends.

Recently, similar findings have been made in Ref. [10], where the concept of the *social bow tie* has been introduced. Bow tie consists of a focal tie and all nodes connected to either or both of the two focal nodes. In the mentioned study, a number of topological metrics quantifying properties of such a bow tie (including sum and absolute difference of clustering coefficients of connected nodes) have been investigated through machine learning and regression models in two different types of social networks (e.g. call network of mobile phone users). The main conclusion from this study was that in the considered networks tie strength depends not only on the properties of shared friends but also on those tied to only one person, hence introducing a fundamental asymmetry to social interaction. Despite interesting conclusions, the authors however failed to identify the most predictive, quantitative indicators of tie strength, basing their findings on a broad spectrum of different structural properties of bow ties. From this perspective, in the face of the growing interest in measuring and predicting social ties (see e.g. [11–14]), an important step towards finding such an informative metric has been made in our recent paper on Granovetter's theory in coauthorship networks [6].

Historically, the Granovetter's theory [15, 16] is of importance to weight-topology correlations in social networks, as Mark Granovetter was the first to distinguish between strong and weak social ties. He treated ties as if they were positive and symmetric, and suggested that, from a network structure perspective, tie strength between any two people should increase with the number

of their mutual friends. In line with this hypothesis, several intuitive network measures, such as the neighbourhood overlap [17], have been proposed to characterize the aforementioned correlations. Unfortunately, contrary to expectations, performance of these indicators turned out to be not very satisfactory: sometimes confirming [18–20], and sometimes saying nothing [14, 21], or even contradicting [22, 23] the Granovetter's hypothesis.

And while no systematic attempts have been made to explain the poor performance of these indicators to date, recent studies [6, 24, 25] may point to some reasons of their failure. For example, in Ref. [24], by analysing a large mobile-phone dataset, it has been shown that temporal features of social ties (such as the number of days with calls, number of bursty cascades, typical time of contacts, etc.) are related to both their strength and topological features of their nearest network neighbourhood. In Ref. [25], analysis of population-scale mobile-telephone and Twitter data has revealed that unembedded long-range connections (i.e. with no nearest neighbours and long second-nearest paths) can be as strong as embedded ones (with non-zero neighbourhood overlap). Finally, in Ref. [6], using a large scale real coauthorship network, we have provided evidence that the key to understand weight-topology correlations in social networks is to reject the assumption of the symmetry of social ties that is commonly used in scientific research.

It is no wonder then that such indicators as the number of common neighbours or neighbourhood overlap, when used in link prediction methods, gave results comparable (and often even worse) to the typical measures of nodes' similarity [26, 27], such as: the Adamic-Adar index [28] or the resource allocation index [29]. In fact, the above-mentioned problems are particularly evident in coauthorship networks, in which many independent studies [6, 14, 21–23] have confirmed non-monotonic (instead of strictly growing), U-shaped relation between tie strength and neighbourhood overlap of adjacent nodes that is contrary to the Granovetter's hypothesis. In our last paper [6], using DBLP computer science bibliography database, we identified the source of this problem, pointing to inappropriate (i.e. symmetric instead of asymmetric) quantities used to study the weight-topology correlations. We have introduced new measures: asymmetric neighbourhood overlap and asymmetric tie strength which allowed the successful verification of the Granovetter's theory, and which - we believe - may be helpful in developing new link prediction methods in social networks.

In this paper, to reinforce the message of our recent contribution [6], we investigate the weight-topology correlations in two more real coauthorship networks and in a synthetic model of scientific collaboration, which reproduces many of the properties of these networks. The motivation behind this study is twofold. First, the analyses with the use of different real data and synthetic networks are intended to validate our findings on the role of asymmetry in social ties, that were originally derived from analysis of just one dataset [6]; such validation is an important element of the research, as it shows that the results described in our previous paper are not an artefact resulting from the specificity of the only dataset used. Second, with this contribution, we would like to point out potential applications of the new network measures we introduced in [6] to the problem of link prediction in social networks; since most of the known link-prediction methods use symmetric network measures [11, 26, 27], contributions like this one are important because they increase the awareness of society that redefining traditional measures to account for link asymmetry can significantly improve their performance.

At this point, we would like to highlight the difference between our contribution and existing research on link prediction in directed networks [30–34]. We are dealing here with undirected networks. Howerer, despite the lack of link directions, we exploit a natural asymmetry in studied networks that can be used to predict links more effectively. This approach is completely new.

The reminder of this paper is organized as follows. In Section II, we study weight-topology correlations in three different real coauthorship networks and in a synthetic model of scientific collaboration. For this purpose, we use a new metric of local edge clustering - the asymmetric neighbourhood overlap, which extracts information about the asymmetry of social ties. In Section III, we provide an in-depth discussion of different similarity scores used in classical methods of link and weight prediction in complex networks. Understanding why some of these measures are successful allows us to design new, inherently asymmetric indices that outperform existing ones. Section IV draws conclusions of the paper.

## II.  ASYMMETRY-BASED WEIGHT-TOPOLOGY CORRELATIONS

### A.  Methods

#### 1.  Scientific collaboration networks

Coauthorship networks, with nodes representing all scientists in a particular discipline and edges joining pairs who have coauthored articles [35], are widely accepted as proxies of scientific collaboration networks [36, 37]. Accordingly, their properties are often compared to other proxies of social networks, such as mobile phone networks [14, 17, 18]. In this respect, when considered as binary networks - without any additional features assigned to nodes and connections - all these networks show numerous structural similarities (e.g. high clustering, small-world effect, and skewed degree distribution [38]). However, when the edges are assigned weights representing, depending on the network, the number of joint publications or the number of phone calls made then, although macroscopic features of these networks (such as distributions of node strengths and edge weights [39]) may still be similar, their weight-topology correlations arising

from the localization of strong and weak ties *seem to be* completely different [22].

Indeed, it is widely believed that coauthorship networks show atypical weight-topology correlations compared to other - let's say *typical* - social networks. Here, the term *typical* refers to networks that satisfy the Granovetter's hypothesis [15], according to which strong social ties are associated with densely connected groups of individuals, while weaker ties act as bridges between these groups. In what follows, we take a closer look at these issues. We show that the phrase we just used, namely: *seem to be* instead of *are*, is not accidental, because in fact coauthorship networks show weight-topological correlations *typical* of other social networks, provided that the measures used to analyse them are properly adapted to their structure.

## 2. Datasets used

We analyse coauthorship networks built from three scientific databases containing publication records in the field of computer science and physics: the DBLP Computer Science Bibliography, the American Physical Society (APS) journal articles, and the Condensed Matter (CondMat) section of the preprint server ArXiv. In detail:

- DBLP is a digital library of article records published in computer science [40]. In this study, we use the 12th version of the dataset (DBLP-Citation-network V12; released in April 2020 [41]), which contains information on approximately 4.9 M articles published mostly during the last 20 years. We ourselves processed the raw DBLP data into the form of coauthorship network and, following previous studies of similar data, we focused on the largest component of this network, consisting of 2.9 M nodes (which is 65% of all nodes) and 12.5 M weighted links.

- APS dataset comprises of over 450 k articles published in all journals of the American Physical Society since 1893 [42]. In this study, we use the preprocessed APS data [14] covering the period between January 1970 and December 2006 and containing 315 k documents with up to 11 co-authors, from which we built coauthorship network with the largest component consisting of 184 k nodes (which is 96% of all APS authors recorded in this period) and 1 M weighted edges.

- CondMat is a weighted coauthorship network between scientists who published preprints on the Condensed Matter e-print archive between January 1995 and December 1999 [36]. To built the network we used a preprocessed bipartite dataset from [43]. The largest component of this network, which is taken into account for in-depth analysis, covers 14 k

authors (which is 83% of all authors) interacting via 45 k weighted links.

In this contribution, as in the previous one [6], our main dataset is DBLP, from which we derive the key findings and to which we relate analysis made in the other two databases and in the synthetic model of scientific collaboration. The leading role of DBLP in our research is due not only to its largest size compared to the other two datasets. Rather, it results from the care of the authors of this database to disambiguate the names of the authors of publications [44]. In DBLP, author names are disambiguated by the combination of algorithms and human curation [45], and not, as in many other bibliographic data - including APS and CondMat - represented by a string of characters corresponding to the surname(s) and initials of all forenames (or only the first one), which can lead to ambiguity of authors through merging or splitting their output [12]. However, the two additional datasets (APS and CondMat), although smaller and less accurate than DBLP, allow us to expand the scope of performed analyses by testing noise-resilience (e.g. due to incomplete data and problems with disambiguation of authors' names) of the recently observed weight-topology patterns [6].

## 3. Coauthorship network model

From various different models of scientific collaboration proposed so far (see e.g. [46–49]), for the analysis presented in this paper, we have chosen the model introduced in Ref. [23]. The choice of this particular model was dictated by several reasons. First, the model reproduces the weight-topology correlations observed in real networks, which we wanted to address and comment on in this paper. Second, despite its simplicity, the model takes into account many important features of the evolution of real scientific collaboration networks that can be easily verified by examining readily available coauthorship networks. These features include: i. growth over time by adding new nodes - students, ii. emergence of new research groups, in which junior scientists (former students) become group leaders (the evolution of career stages [50]), iii. creation of new publications based on intra- and inter-group relations, and finally iv. high probability that the young scientist will give up a further scientific career.

As indicated above, in the model studied, nodes are assigned to specific research groups in which they perform various functions. More specifically, each group consists of exactly one leader and a number of students, with the latter being "active" or "inactive" depending on the time elapsed since they were added to the network. It is assumed that the group leader is established at the time of group formation and remains in function until the end of the network evolution. The situation of students is a bit more complicated. After a node is added to the network,

it is assigned to one of the existing research groups as its active student, who can participate in scientific research and coauthor publications. However, just like in the real world, after some time such a student may cease to be active, giving up further research activity, or may pursue a scientific career as a leader of a new research group.

In the considered model, inter-node connections result from common (i.e. coauthored) publications, the number of which translates into the edge weight. The model has two mechanisms of producing new publications, through intra- or inter-group collaboration, with the number of coauthors taken from a certain distribution $P(l)$. This distribution and any other parameters of the model are determined on the basis of real data. We comment on them later in the text, in the part devoted to simulation results.

To be more specific, the evolution of the network model under study proceeds as follows:

(0) *Beginning of the evolution:* The network starts to grow with a single research group consisting of a leader and one student.

Then, in successive time steps, the following actions are performed (cf. description given in Ref. [23]):

(1) *Intra-group publications:* With probability $c$, each group publishes one paper by itself. The paper is written by the group's leader and $l-1$ active students preferentially chosen from the same group based on student's scientific expertise, which corresponds to the length of student's activity period.

(2) *Inter-group publications:* Each group may publish up to $\alpha$ papers with another group. External collaboration always takes place with the same group, which is randomly chosen right after a new group appears. Same as for intra-group publications, each of these $\alpha$ papers is realized with probability $c$ and is coauthored by $l$ individuals (two leaders and $l-2$ students, which are preferentially chosen from the pool of all active students of the two groups).

(3) *Groups' resource update:* Active students whose activity period has exceeded the threshold value $G$, with probability $f$ become leaders of new research groups, and with probability $1-f$ become inactive and no longer participate in network dynamics. A new active student is added to each group.

Although the model under study has several free parameters, the values of most of them can be approximated from various real data. For example, since in this paper we primarily use the coauthorship network extracted from DBLP Computer Science Bibliography to compare with the model results, we assume the distribution $P(l)$ of the number of coauthors in publications consistent with the corresponding distribution in the mentioned database [1] (alternatively, one could use real data to train a model describing the distribution of coauthors

in a similar way to [51]), see Fig. 1(a). Furthermore, although we examined a wide range of model parameter values in our studies, we ultimately decided to keep the values provided in Ref. [23] (where the model was originally introduced), as they result in the best agreement between simulations and real data. In Ref. [23] the following values have been taken as a reference: $c = 0.4$ for the probability to publish a paper; $\alpha = 3$ for the number of inter-group publications; $f = 0.2$ for the probability of a student becoming a group leader; and finally $G = 7$ for the length of the students' activity period, which also determines the maximum number of active students in a research group. The rationale for these parameters is described in more detail in Ref. [23], to which we refer interested readers (for a broader perspective, see also the recent studies: [49, 52, 53]).

### B. Results

#### 1. Real data vs. simulation results

In Fig. 1(b)-(f), we show basic structural characteristics of real coauthorship networks (DBLP, APS, and CondMat) and the model network with $N \simeq 10^4$ nodes (averaged over 100 realizations). Two obvious conclusions arise after analysing this figure. First: The examined features of real networks are very similar to each other. It is reasonable to claim that the slight differences in the range of the data shown are primarily related to the size of the analysed networks, which varies from millions of nodes (in DBLP), through hundreds (in APS) to tens of thousands (in CondMat). Second: The synthetic model of scientific collaboration reflects very well the basic features of the reference coauthorship networks, including their skewed distributions of node degrees, $P(k_i)$, and strengths, $P(s_i)$ (where the node strength is given as the sum of the weights of its edges: $s_i = \sum_j w_{ij}$), as well as the fat-tailed distributions of edge weights (tie strengths) $P(w_{ij})$ and $P(v_{ij})$ (where $w_{ij}$ represents the number of joint papers, and $v_{ij} = w_{ij}/p_i \neq v_{ji}$ (4) stands for an asymmetric tie strength, which is discussed afterwards). The good agreement between the model and real data, as can be seen in this figure, is all the more convincing as we checked that the differences between them decreased as the size of the model networks increased. The above results make the considered model a promising test-bed to study weight-topology correlations in scientific collaboration networks, which enables the formulation of well-established conclusions, based not only on

---

[1] By drawing the number $l$ of coauthors, we limit the range from which we randomly select to the size of the group (i.e. $G+1$ for

intragroup publications and $2(G+1)$ for intergroup publications). For the value of $G = 7$ adopted in this study, this limitation makes sense because, in the database under consideration, less than one per mille of publications has more than $2(G+1) = 16$ authors.
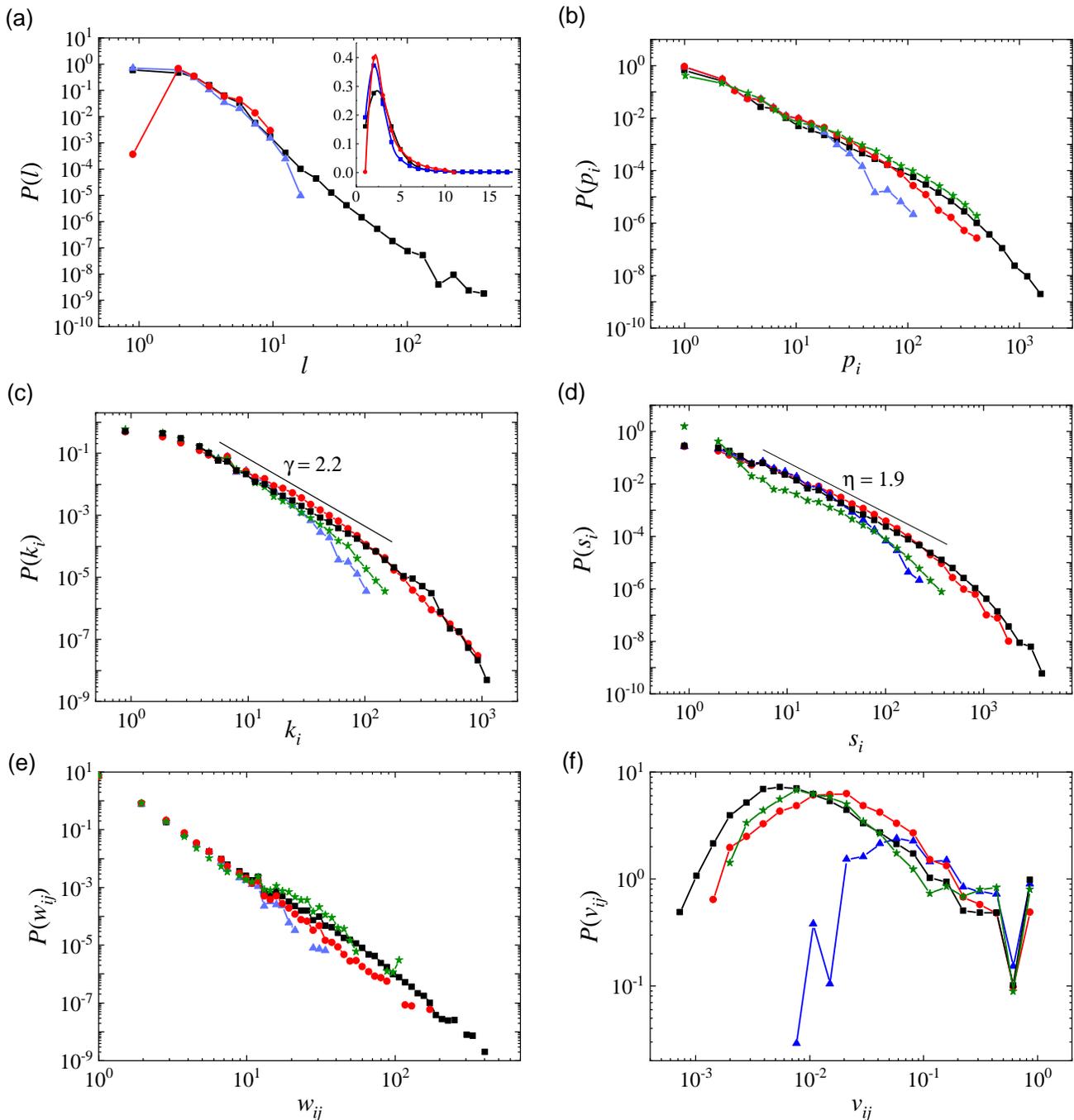
FIG. 1. **Structural properties of coauthorship networks.** The following graphs show distributions of: (a) the number of coauthors per paper, $P(l)$; (b) the number of publications per author, $P(p_i)$; (c) the node degrees, $P(k_i)$; (d) the node strengths, $P(s_i)$; (e) the edge weights (or symmetric tie strengths), $P(w_{ij})$; (f) the asymmetric tie strengths, $P(v_{ij})$. The symbols used are: black squares for DBLP, red circles for APS, blue triangles for CondMat, and green stars for results of numerical simulations obtained from the model network.

real datasets but also on results of repeatable numerical simulations.

## 2. *Granovetter's hypothesis in coauthorship networks*

As mentioned in the introduction, according to the Granovetter's hypothesis, strong social ties are expected to be associated with densely connected groups of individuals, while weaker ties act as bridges between these
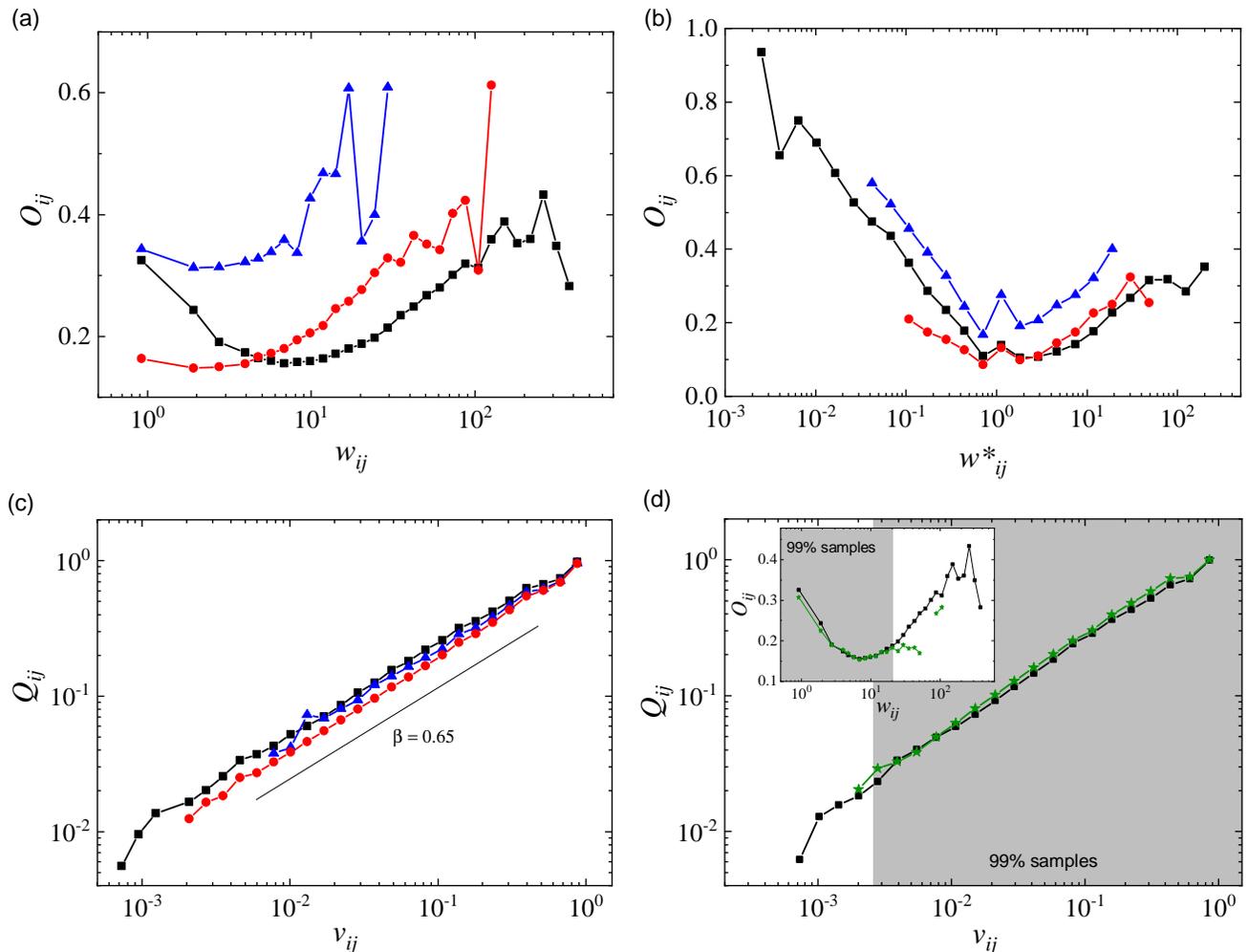
FIG. 2. **Weight-topology correlations in coauthorship networks** as observed on the basis of relationship between variously defined tie strengths ($w_{ij}$, $w_{ij}^*$, and $v_{ij}$) and neighbourhood overlaps ($O_{ij}$ and $Q_{ij}$). Detailed description of the notation used is given in the main test. Graphical symbols used are the same as in Fig. 1, i.e. black squares for BDLP, red circles of APS, etc. Note that the panels (a), (b), and (c) show only real coauthorship networks. Panel (d) presents results of numerical simulations of the synthetic network model and DBLP data for comparison.

groups. To quantitatively characterize such weight-topology correlations, in Ref. [17], the relationship between tie strength, $w_{ij}$, connecting two nodes ($i$ and $j$) and their neighbourhoods' overlap, $O_{ij}$, has been used, with the overlap defined as the ratio of the number of common neighbours, $n_{ij}$, of this node pair to the number of all their neighbours:

$$O_{ij} = \frac{n_{ij}}{(k_i - 1) + (k_j - 1) - n_{ij}}. \quad (1)$$

Correspondingly, clear empirical support for the Grannovetter's hypothesis, manifested a monotonically increasing dependence between $w_{ij}$ and $O_{ij}$, has indeed been observed in many social networks, but not in coauthorship networks (see Fig. 2(a) and (d - inset graph)), making the latter a flagship example of systems in which the hypothesis fails.

Interestingly, the characteristic U-shape non-

monotonic relation between tie strength and symmetric overlap $O_{ij}$, which is interpreted as the evidence of atypical weight-topology correlations in scientific collaboration networks becomes even more apparent, when the Newman's definition [22, 23, 39] of tie strength, $w_{ij}^*$, is taken into account (see Fig. 2(a),(b)). To grasp the difference between $w_{ij}$ and $w_{ij}^*$, recall that the tie strength $w_{ij}$, which is the standard used throughout this paper, stands for the number of joint publications, that is also the number of times a collaboration between two scientists has been repeated. Correspondingly, the Newman's tie strength is defined as:

$$w_{ij}^* = \sum_{p_{ij}} \frac{1}{l_{ij} - 1} \neq \sum_{p_{ij}} 1 = w_{ij}, \quad (2)$$

where the sum runs over the set of papers $p_{ij}$ co-authored by $l_{ij}$ scientists, including $i$ and $j$. The motivation behind

the Newman's formula for $w_{ij}^*$ is that an author divides his/her time and other resources between $l_{ij} - 1$ collaborators, and thus the tie strength of such a collaboration should vary inversely with $l_{ij} - 1$.

The possible cause of the failure of the Granovetter's hypothesis in scientific collaboration networks has only recently been clarified in Ref. [6], where it was suggested that the non-monotonic $O_{ij}(w_{ij})$ and/or $O_{ij}(w_{ij}^*)$ relations characterizing these networks are due to the definition of the neighbourhood overlap, Eq. (1) (hereafter called *symmetric overlap*) which is not properly suited to be a local network measure in networks with scale-free node degree-distributions.

The above mentioned problem with the symmetric overlap is particularly acute in the case of links connecting nodes with significantly different degrees. In such cases, for $k_i \ll k_j$, Eq. (1) can be simplified to $O_{ij} \simeq n_{ij}/k_j$, which shows that it is strongly biased towards nodes with high degrees, distorting the image of the common neighbourhood as seen from the perspective of nodes with small degrees. This drawback of symmetric overlap gains importance in networks with highly skewed, fat-tailed node degree distributions $P(k_i)$. In such networks, as brilliantly exploited by the degree-based mean-field theory of complex networks [7–9], node degree distributions for nearest neighbours are even more fat-tailed than the original distributions $P(k_i)$. As a result, the number of edges in such networks connecting nodes with high and low degrees can be very high, leading to an unintended overrepresentation of strongly connected nodes by Eq. (1).

To overcome the aforementioned problems with the symmetric overlap $O_{ij}$, the concept of *asymmetric overlap* has been introduced in Ref. [6]:

$$Q_{ij} = \frac{n_{ij}}{k_i - 1} \neq Q_{ji}, \qquad (3)$$

and it was used to describe the overlap between the neighbourhoods of two connected nodes from the perspective of each node separately [2]. In the context of coauthorship networks, this new definition is free from the shortcomings of the previous one. In particular, it copes well with collaborating scientists whose degrees (ego-networks) differ significantly in size - that is, when their common neighbours (if any) are a significant part of the neighbourhood of one node and an insignificant part of the neighbourhood of the other. The relevant situation is illustrated in Fig. 3(a).

The concept of asymmetric overlap naturally leads to the idea of directed networks and justifies the introduc-

---

[2] Note that the definition of the asymmetric neighbourhood overlap, Eq. (3), that we use in our manuscript is similar to the so-called *edge clustering coefficient*: $C_{ij} = n_{ij}/\min[k_i - 1, k_j - 1]$ [54]. However, the difference between the two measures is essential, because $Q_{ij} \neq Q_{ji}$ while $C_{ij} = C_{ji}$.
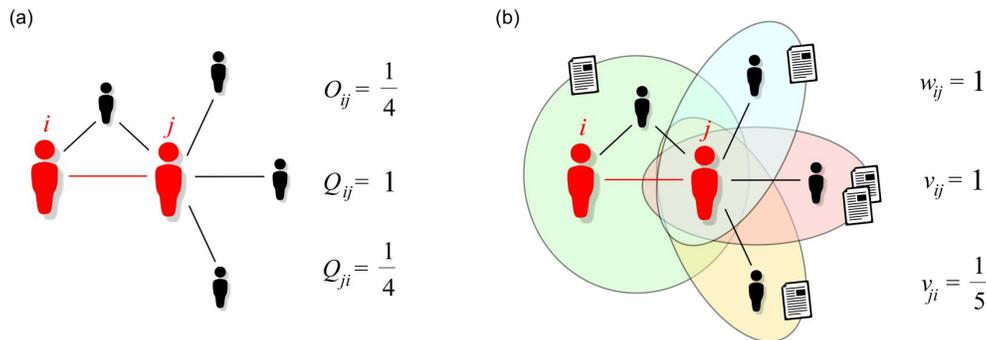
tion of the *asymmetric tie strength*:

$$v_{ij} = \frac{w_{ij}}{p_i} \neq v_{ji}, \qquad (4)$$

where $p_i$ stands for the number of all publications of the $i$-th author (note that the number of publications does not have to be equal to the strength of the node: $p_i \neq s_i = \sum_j w_{ij}$). The intuitive rationale behind Eq. (4) is illustrated in Fig. 3(b) and it proceeds as follows: For a young scientist, with a small number of publications, each publication makes a significant contribution to his or her publication output, just as each coauthor is an important part of his or her research environment (cf. Eqs. (3) and (4)). However, the importance of each publication and collaboration from the perspective of an established scientist with a large number of publications and an extensive network of collaborators is completely different. Depending on the circumstances, a given number of joint publications (e.g., $w_{ij} = 1$) may have a completely different meaning.

In Fig. 2(c) and (d-main panel), relationship between asymmetric neighbourhood overlap, $Q_{ij}$, is shown against the corresponding asymmetric tie strength, $v_{ij}$. Remarkably, although the relationships with the use of symmetric network measures (see Fig. 2(a),(b) and (d - inset graph)) are cumbersome to interpret (e.g. due to U-shape relation observed for $O_{ij}(w_{ij}^*)$ and significant differences between the real data observed for $O_{ij}(w_{ij})$), the relationship between asymmetric measures $v_{ij}$ and $Q_{ij}$ seems to be universal for all studied networks. The reasonable explanation for this observation is that, the relationship between symmetric tie strengths, $w_{ij}$ and $w_{ij}^*$, and the symmetric overlap $O_{ij}$ is not an informative measure for weight-topology correlations in coauthorship networks. On the other hand, the result for asymmetric measures is of particular importance as it aspires to be a universal scaling law that would require verification in other social networks, not only collaborative.

In fact, Fig. 2(c) confirms that the Granovetter's hypothesis holds in coauthorship networks. In other words, from the point of view of an individual scientist, strong ties do really correspond to dense local neighbourhoods, contrary to what has been suggested in other studies on these networks. The perspective of a single node, being one of the two ends of an edge, is important here. The new measures introduced ($Q_{ij}$ and $v_{ij}$) quantitatively capture the so far elusive concept of relativity in social relations. Following this line of reasoning, it may be tempting to say that the measured absolute tie strength (i.e. $w_{ij}$) is a kind of compromise and depends on relative strengths of the tie as seen from its both ends (i.e. $v_{ij}$ and $v_{ji}$). Moreover, it seems reasonable that similar thinking should also apply to the connection probability, not just to its weight. In this context, it is surprising that none of the so far proposed network measures that are used in link-prediction methods take into account, at least not explicitly, the asymmetry of these links. In the rest of the publication, we refer to these issues.

FIG. 3. **a) Difference between symmetric and asymmetric neighbourhood overlap**. In the figure, a pair of scientists ($i$ and $j$) is shown, with different numbers of collaborators (respectively, $k_i = 2$ and $k_j = 5$) and one neighbour in common (i.e. $n_{ij} = 1$). The figure shows that even in such a simple situation the neighbourhood overlap, as seen from the perspective of each of the two nodes, is significantly different: $Q_{ij} < Q_{ji} = O_{ij}$. **b) Difference between symmetric and asymmetric tie strength**. The figure presented in part a) is supplemented with additional data, which allow to determine appropriate edge weights. In particular, it can be seen from this figure that the node $i$ stands for an individual who coauthored (with two other individuals, including $j$) only one paper. The corresponding tie strengths are: $v_{ij} = w_{ij} > v_{ji}$.

## III. ASYMMETRY-BASED LINK PREDICTION

### A. Methods

Link prediction refers to the problem of finding missing or hidden links that are likely to exist in networks or will appear there in the near future [26, 27]. Predicting new friendships in social media or new collaborations in coauthorship networks [55–57], discovering previously unknown interactions in biological networks [58], predicting scientific research trends [59], or providing bibliography recommendations [60] are a few examples showing the importance and the diversity of the applications that can benefit from link prediction. It is also worth noting that link prediction is not limited to single-layer networks and prediction methods can utilise data from multiple layers [61] representing various types of interactions.

The simplest predicting methods are based on nodes' neighbourhood-related structural information that is used to compute the so-called similarity score, $s_{ij}$, of each pair of nodes in the network. Then, by ranking the pairs based on this score, an inference is made as to the existence or absence of edges. In the literature on link prediction, one can find dozens of such scores (or indicators) that perform better or worse depending on the network under study. In particular, the following examples of symmetric indicators have been widely employed by the research community due to their simplicity, computational efficiency and performance: the common neighbours index [62], the Jaccard's index [63], the Adamic-Adar (AA) index [28], and the resource allocation (RA) index [29]. Further in this subsection, using the Jaccard's index as an example, we will show that by redefining this index to take into account the asymmetry of network connections, one can significantly increase its prediction efficiency. The perspective will also allow us to comment on the surprisingly high performance of

the RA index, pointing to the asymmetry as a promising direction for further research on effective prediction methods.

The Jaccard's index is widely used in information retrieval systems to compare the similarity and diversity of sample sets. In the context of link prediction methods, the index measures the proportion of common neighbours of two nodes ($i$ and $j$) in the total number of their neighbours. Correspondingly, it is given by the expression:

$$\text{JC}_{ij} = \frac{|\Gamma(i) \cap \Gamma(j)|}{|\Gamma(i) \cup \Gamma(j)|}, \qquad (5)$$

where $\Gamma(i)$ is the set of nearest neighbours of $i$, and $|\Gamma(i)|$ is the cardinality of this set. Since $|\Gamma(i)| = k_i$ and $|\Gamma(i) \cap \Gamma(j)| = n_{ij}$, the above formula can be rewritten as

$$\text{JC}_{ij} = \frac{n_{ij}}{k_i + k_j - n_{ij}} \simeq O_{ij}, \qquad (6)$$

which shows that the definition of the Jaccard's index is almost identical to the definition of the symmetric neighbourhood overlap, cf. Eq. (1).

Having the similarity score $s_{ij}$ defined (e.g. $s_{ij} = \text{JC}_{ij}$), the link prediction proceeds as follows: One ranks the non-connected pairs of nodes in descending order according to the values of $s_{ij}$ and then the pairs for which the score exceed some established threshold $T$ are considered to be connected. At this point, at least two problems arise, that we comment on below.

First, to validate the ranking method used and to evaluate the similarity score chosen, one has to know *a priori* which identified links are indeed present in the network. Thus, for the testing purposes, one has to construct a set of node pairs which include those pairs that are connected in the original network (labelled as positive links – P) and those that are not (labelled as negative links – N). Construction of such a set is not a trivial task per se, because the studied dataset is highly imbalanced - the

| | predicted links | predicted non-links |
|---|---|---|
| positive links | True Positives (TP) | False Negatives (FN) |
| negative links | False Positives (FP) | True Negatives (TN) |

TABLE I. Confusion matrix for link prediction.

number of negative (i.e non-existing) links is much larger than the number of positive ones. Moreover most of node pairs have no common neighbours, which may result in their similarity scores equal to zero (e.g. $JC_{ij} = 0$ for $n_{ij} = 0$). To overcome this problem, in this study, we construct our testing set by selecting $d$ existing links and $d$ non-existing links, both from those node-pairs which share at least one common neighbour.

Given such a correctly balanced testing set of existing and non-existing links, one can perform ranking on this set and create the confusion matrix (see Tab. I), whose elements - representing the numbers of connections labelled as: true positive (TP), true negative (TN), false positive (FP), and false negative (FN) - allow to derive several useful metrics to asses performance of a link-prediction method. In what follows, we use three, perhaps the most commonly used, metrics (see axis description in Fig. 4):

(1) *sensitivity*, also referred to as recall or the true positive rate (TPR), which measures the proportion of existing links that are correctly identified as true positive matches to all existing links: TPR=TP/(TP+FN),

(2) *specificity*, also known as the true negative rate (TNR), measures the proportion of non-links that are correctly identified as negative links: TNR=TN/(TN+FP), and finally

(3) *precision*, also called the link accuracy, that is the proportion of all correctly identified links to all node-pairs in the set that are classified (correctly or not) as existing links: PR=TP/(TP+FP).

Although the above measures seem similar, they relate to different aspects of link prediction and thus are complementary to each other. In particular, precision evaluates the correctness, while recall evaluates the completeness of both the similarity score and the method used. Generally, there is a trade-off between precision and recall, whereby a larger threshold $T$ increases precision and decreases recall.

The second problem with the method described is that the above measures use a fixed threshold to rank node pairs, but the value of this threshold may not be necessarily available or be the most optimal one. For example, it may be domain dependent. To deal with such cases threshold curve based metrics are used. This is where one of the most important goals of link prediction research emerges, which is optimizing these curves to find the most effective similarity score, $s_{ij}$. It is also where our research on the role of asymmetry in social ties contributes to the already huge research area of link prediction.

There are generally two such curves in use. The first one, the ROC curve (from: receiver operating character-

istic) represents the performance trade-off between true positives TPR and false positives FPR=1−TNR at different decision boundary thresholds. It can be interpreted as the probability that a randomly chosen true positive link will be ranked higher than a randomly chosen true negative link [64]. The area under the ROC curve, AUC, is always between 0 and 1, and, generally, the performance of any realistic classifier at AUC measure should be higher than 0.5 (which corresponds to completely random classifier) [65]. The second curve, the PR curve (from: precision recall) considers only prediction of the positive links, which, in some situations, can be more useful and informative, e.g. when negative links are not interesting [66]. Similarly to the AUC value characterizing the ROC curve, to get one number that describes performance of the method, one can also calculate the PRAUC, being the area under the precision-recall curve. Such a single value can be understood as the average of precision scores calculated for each recall threshold. The higher values of both threshold curve metrics, AUC and PRAUC, correspond to those similarity scores, $s_{ij}$, that rank better the pairs of nodes towards the discovery of existing edges between them.

### B. Results

#### 1. Similarity scores accounting for link asymmetry

The Jaccard's index, Eq. (5), which is equivalent to the symmetric overlap, Eq. (1), was one of the first similarity measures used in link prediction methods. Nevertheless, it was quickly realized that the link prediction based on this measure is not much better than a random classifier. Correspondingly, in the case of coauthorship networks analysed in this study, the values of AUC and PRAUC are close to 0.7, which is a rather poor result compared to other similarity scores reported in Tables II and III, which have values up to 0.97. Reasons of such a poor performance of the Jaccard's index, however, have not been thoroughly investigated. In this context, the results on weight-topology correlations in social networks presented in this paper shed some light on the problem. The relevant reasoning is as follows: If the tie strength of the missing link is high, the probability of its existence, based on the corresponding similarity score, should also be high. Likewise, low-strength ties should be less likely to be realized. Consequently, the prediction methods based on the Jaccard's index will have poor performance in networks where the correlation between tie strengths, $w_{ij}$, and the symmetric neighbourhood overlap, $O_{ij}$, are weakly positive or absent at all. This is the situation we deal with in social networks with fat-tailed degree distributions, of which coauthorship networks are a particularly vivid example.

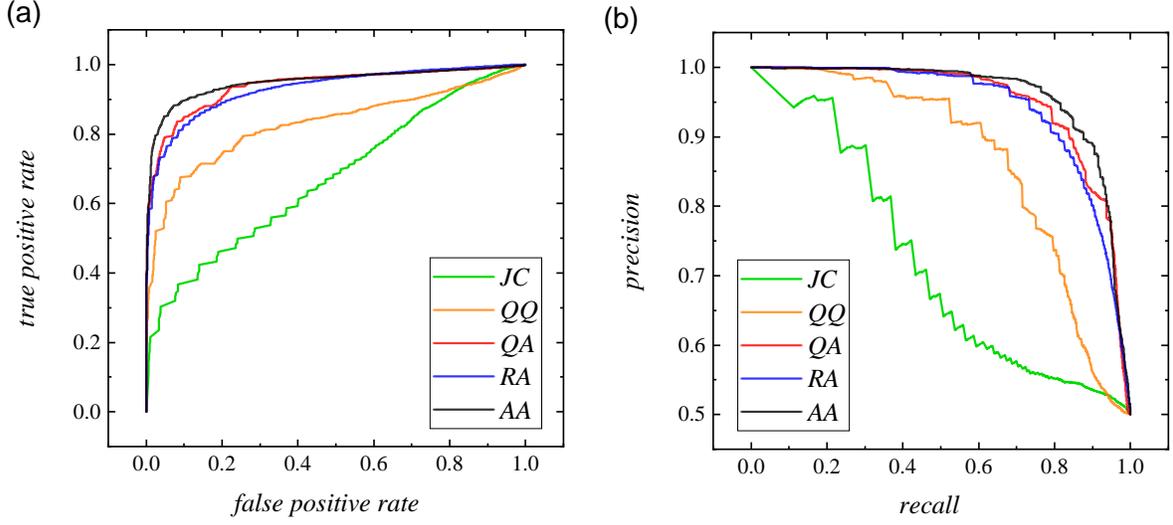To support the above reasoning, we have introduced

FIG. 4. **Performance of similarity-based link prediction methods in coauthorship networks**. ROC curves (a) and precision-recall curves (b) for different score measures in the DBLP dataset.

TABLE II. Performance results of analysed score measures tested on the DBLP, APS, CondMat datasets and the discussed model of scientific collaboration network.

| similarity score | DBLP | | APS | | Cond-Mat | | model | |
|---|---|---|---|---|---|---|---|---|
| | AUC | PRAUC | AUC | PRAUC | AUC | PRAUC | AUC | PRAUC |
| JC | 0.684 | 0.714 | 0.698 | 0.684 | 0.764 | 0.766 | 0.634 | 0.659 |
| QQ | 0.829 | 0.865 | 0.758 | 0.795 | 0.847 | 0.870 | 0.798 | 0.823 |
| QA | 0.939 | 0.949 | 0.925 | 0.941 | 0.901 | 0.920 | 0.860 | 0.885 |
| RA | 0.934 | 0.945 | 0.922 | 0.927 | 0.910 | 0.924 | 0.869 | 0.902 |
| AA | 0.948 | 0.952 | 0.937 | 0.944 | 0.918 | 0.929 | 0.888 | 0.920 |

and analysed a simple similarity score:

$$QQ_{ij} = \frac{|\Gamma(i) \cap \Gamma(j)|}{|\Gamma(i)|} + \frac{|\Gamma(i) \cap \Gamma(j)|}{|\Gamma(j)|}, \qquad (7)$$

the form of which refers to the sum of the asymmetric neighbourhood overlaps (3) of the considered pair of nodes:

$$QQ_{ij} = \frac{n_{ij}}{k_i} + \frac{n_{ij}}{k_j} \simeq Q_{ij} + Q_{ji}. \qquad (8)$$

As one can see in Figs. 4 and 5 and Tab. II, for all studied coauthorship networks (real and synthetic), the link prediction results obtained with this similarity score are much better than with the Jaccard's index, Eq. (6). Moreover, the results can be further improved by making adjustments learned from the Adamic-Adar index (see Eq. (11)), in which the importance of nodes is expressed not by their degree, but by the logarithm of the degree:

$$QA_{ij} = \frac{n_{ij}}{\log k_i} + \frac{n_{ij}}{\log k_j}. \qquad (9)$$

As compared to other similarity scores listed in Tab. II, the results obtained using the QA index are significant for at least two reasons. Firstly, the accuracy of predictions

according to QA is similar to the accuracy obtained using measures such as the resource allocation index [26]:

$$RA_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} \frac{1}{k_z}, \qquad (10)$$

and the Adamic-Adar index:

$$AA_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} \frac{1}{\log k_z}, \qquad (11)$$

both of which require much more detailed information about the neighborhood of the possible edge between $i$ and $j$ than just the number of their common neighbors, $n_{ij} = |\Gamma(i) \cap \Gamma(j)|$, as QA does. Secondly, it shows that taking the asymmetry of social ties into account can significantly improve the effectiveness of link prediction methods. In the rest of the paper, we refer to the last remark, showing that the high efficiency of RA and AA indicators in social networks is obviously related to the asymmetry of social ties. In showing this, in order to avoid distraction, we limit ourselves to the RA index only.
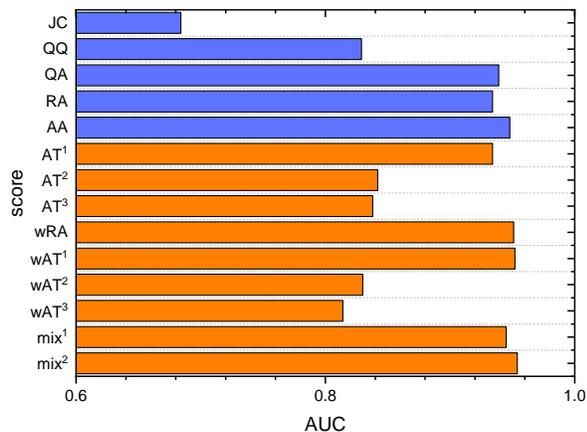
FIG. 5. **AUC measures for different scores** calculated for DBLP dataset. Blue and orange colors correspond to the data presented in the Table II and III respectively.

### 2. Asymmetry-based perspective on RA-like scores

All the local similarity indices we have discussed so far have one thing in common: They are designed based on the assumption that two nodes are more likely to have a link if they have many common neighbors. In particular, the RA index (10) is the greater the more common neighbors the nodes $i$ and $j$ have, but it also reduces the contribution of high degree common neighbors by assigning more importance to those less-connected [3]. In what follows, we show that the punishment of high degree nodes in Eq. (10) can be justified by the strong triadic closure property, assuming that the property takes into account the asymmetry of social ties.

Triadic closure property states that: If two people $i$ and $j$ in a social network have a friend $z$ in common, then there is an increased *likelihood* that they will become friends themselves at some point in the future (see [2], p. 44). Accordingly, the strong triadic closure property completes the previous statement saying that: If a node $z$ has edges to $i$ and $j$, then the connection between $i$ and $j$ is *especially likely* to form if $z$'s edges to $i$ and $j$ are both *strong ties* (see [2], p. 49).

Clearly, both triadic closure properties are intuitively very natural. Furthermore, although their original wordings apply to single triads, i.e. those in which nodes $i$ and $j$ have only one common neighbor $z$, it is reasonable to assume that the more common neighbors, the higher the *likelihoods* in question should be. In the context of link predictability - the primary concern of which is to

————

[3] Note that the AA index also punishes the high-degree common neighbors but to a lesser extent than RA. It depends on the network under study which of the approaches to punish nodes with high degrees is better (RA or AA). In scientific collaboration networks we study in this paper, AA performs better, but this is not a rule, because in other networks (not only the social ones) it may be the other way around [26, 67].
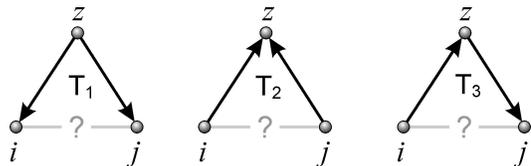


FIG. 6. **Three types of directed triads** that are analysed in the text.

correctly estimate the mentioned likelihoods - the above reasoning leads to the simplest similarity measures defined as a bare number of common neighbors [26, 67]:

$$\mathrm{CN}_{ij} = n_{ij}, \tag{12}$$

and its weighted version:

$$\mathrm{wCN}_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} (w_{zi} + w_{zj}), \tag{13}$$

where $w_{ij}$ stand for the symmetric tie strength between $i$ and $j$. However, as numerous studies show, these indices are not as efficient as the previously introduced RA (10) and AA (11) scores. Below we show, where this inefficiency comes from and how to improve it without going beyond the concept of triadic closure.

The idea - which is consistent with our findings about asymmetry of social ties - is to replace the symmetric weights $w_{ij}$ in Eq. (13) with the asymmetric tie strengths $v_{ij}$ (4) or with the asymmetric neighbourhood overlaps $Q_{ij}$ (3) that show a high correlation with asymmetric tie strengths. In the case of binary networks, the corresponding similarity score - referring to the newly addressed concept of the directed triad closure [68, 69] - can be defined as:

$$\mathrm{AT}^1_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} (Q_{zi} + Q_{zj}) \tag{14}$$

$$= \sum_{z \in \Gamma(i) \cap \Gamma(j)} \frac{n_{zi} + n_{zj}}{k_z - 1}, \tag{15}$$

and for the weighted networks, as:

$$\mathrm{wAT}^1_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} (v_{zi} + v_{zj}), \tag{16}$$

where the acronyms AT and wAT account for 'asymmetric triad' and 'weighted asymmetric triad', and the superscript - here '1' - refers to the triad ordering $T_1$ as shown in Fig. 6.

Now, when we compare Eq. (15) with the definition of the RA index (10), we see that they are very similar. In particular, in both indices, RA and $\mathrm{AT}^1$, each neighbor of the pair $(i, j)$ is weighted inversely proportional to its degree. The difference between the two measures is that our $\mathrm{AT}^1$ index additionally takes into account the entire triad's neighborhood, which is represented by nodes adjacent to the edges $(z, i)$ and $(z, j)$. Comparing the

TABLE III. Performance results of analysed RA-like similarity measures tested on the DBLP, APS, CondMat datasets and the discussed model of scientific collaboration network.

| similarity score | DBLP | | APS | | Cond-Mat | | model | |
|---|---|---|---|---|---|---|---|---|
| | AUC | PRAUC | AUC | PRAUC | AUC | PRAUC | AUC | PRAUC |
| RA | 0.934 | 0.945 | 0.922 | 0.927 | 0.910 | 0.924 | 0.869 | 0.902 |
| $AT^1$ | 0.934 | 0.949 | 0.924 | 0.941 | 0.883 | 0.915 | 0.874 | 0.911 |
| $AT^2$ | 0.842 | 0.885 | 0.813 | 0.862 | 0.805 | 0.853 | 0.747 | 0.803 |
| $AT^3$ | 0.838 | 0.888 | 0.814 | 0.865 | 0.802 | 0.857 | 0.803 | 0.862 |
| wRA | 0.951 | 0.959 | 0.940 | 0.944 | 0.929 | 0.943 | 0.906 | 0.923 |
| $wAT^1$ | 0.952 | 0.961 | 0.945 | 0.954 | 0.922 | 0.941 | 0.909 | 0.931 |
| $wAT^2$ | 0.830 | 0.881 | 0.802 | 0.849 | 0.814 | 0.859 | 0.814 | 0.868 |
| $wAT^3$ | 0.814 | 0.862 | 0.779 | 0.828 | 0.797 | 0.836 | 0.708 | 0.765 |
| $mix^1$ | 0.945 | 0.956 | 0.928 | 0.943 | 0.903 | 0.925 | 0.916 | 0.937 |
| $mix^2$ | 0.954 | 0.964 | 0.940 | 0.952 | 0.933 | 0.948 | 0.912 | 0.934 |

effectiveness of both measures (see Fig. 5 and Tab. III), there is a slight argument in favour of the $AT^1$ index, which proves the legitimacy of the presented reasoning. In the same vein, and for the sake of completeness, it is worth noting that our weighed $wAT^1$ index is almost identical to the RA weighted index [67, 70]:

$$\text{wRA}_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} \frac{w_{zi} + w_{zj}}{s_z}, \qquad (17)$$

although our $wAT^1$ score performs slightly better than wRA (see Tab. III).

To complete the discussion on the role of the asymmetry of social interactions in link predictability, we also examined other similarity scores accounting for the concept of 'asymmetric triad' (see Fig. 6). More precisely, we tested the following 'structural' measures:

$$\text{AT}^2_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} (Q_{iz} + Q_{jz}), \qquad (18)$$

$$\text{AT}^3_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} (Q_{iz} + Q_{zj}), \qquad (19)$$

and their 'weighted' counterparts:

$$\text{wAT}^2_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} (v_{iz} + v_{jz}), \qquad (20)$$

$$\text{wAT}^3_{ij} = \sum_{z \in \Gamma(i) \cap \Gamma(j)} (v_{iz} + v_{zj}), \qquad (21)$$

and found that they give significantly worse results (see Fig. 5 and Tab. III). This observation clearly confirms the general message of our paper that the asymmetry of social interactions is important, even in networks where the edges are not formally assigned a direction.

At this point, we would like to refer to the RA score once again. The original rationale for this index is related to the efficiency of *resource transmission between the nodes i and j, through their nearest neighbours* [29]. There is, however, a subtle inaccuracy between the above reasoning and the RA definition (10). Namely, the reasoning seems to apply to the directed triad $T_3$ shown in

Fig. 6 and quantified by our $AT^3$ index rather than by the RA-like, our $AT^1$ index that describes the $T_1$ triad. Correspondingly, the success of the measure $AT^1$ (and also $wAT^1$) can be justified as follows: The asymmetric overlap $Q_{zi}$ acts as a proxy for the amount of resources the node $z$ invests in collaborating with the node $i$. If two such investments, e.g. $Q_{zi}$ and $Q_{zj}$, are resource-intensive, they probably cannot be independent, simply because personal resources of $z$ are limited (e.g. there are only 24 hours in a day). This observation stands in line with the study by Dunbar [71], who concluded that humans, even very active, have a limited capacity to maintain significant interpersonal relationships. In the case of the collaboration network such a dependency between both $z$'s investments simply means a joint project or publication co-authored by $z$, $i$ and $j$, so the link between $i$ and $j$ has to be expected. Please note that such a dependency cannot be simply derived from investments $Q_{iz}$ and $Q_{jz}$, or $Q_{iz}$ and $Q_{zj}$ standing behind the indices $AT^2$ or $AT^3$, respectively. Please also note that such a perspective of resource allocation is completely different from the original one, which we have already questioned.

### 3. Towards intrinsically asymmetric similarity scores

Finally, while it was not our intention in this paper to search for better than already existing similarity scores for link prediction, nor to argue that the local network measures we introduced, such as asymmetric neighborhood overlap, are generally better to design such indices, we cannot pass over the fact that through trial and error, we have managed to find indices whose effectiveness exceeds (at least in the analysed networks of scientific collaboration) the effectiveness of all others that were examined in this paper. These measures are defined as linear combinations of the previous ones (see Fig. 5 and Tab. III):

$$\text{mix}^1_{ij} = \text{wAT}^1_{ij} + \text{QQ}_{ij}, \qquad (22)$$

and

$$\text{mix}^2_{ij} = \text{wAT}^1_{ij} + \text{QA}_{ij}, \qquad (23)$$

where QQ, QA and wAT[1] are given by Eqs. (8), (9) and (16), respectively.

## IV. CONCLUDING REMARKS

The leitmotif of this paper is the problem of "relativity" (or the lack of symmetry) in social relations. To draw attention to this problem, we focused on coauthorship networks, and used the known controversy regarding their atypical weight-topology correlations to show that taking the asymmetry into account can change the understanding of even well-established findings, such as that scientific collaboration networks do not satisfy the Granovetter's strength of weak ties hypothesis.

More precisely, in this paper, by analysing three different real coauthorship networks (DBLP, APS, and Cond-Mat) and their reliable synthetic model, we show that the networks show strong positive correlations between tie strength, $v$, and neighbourhood overlap, $Q$, of the connected nodes only when both measures take into account the lack of symmetry of the relationship. The observed correlations satisfy the power law scaling: $Q \propto v^\beta$, with the same characteristic exponent $\beta \simeq 0.65$ for all studied networks.

In light of the noticed strong correlations, research on link prediction methods that would take advantage of link asymmetry seems particularly interesting. By testing various link scores used in similarity-based unsupervised link and weight prediction methods [67, 72–74], we argue that taking into account the asymmetry of social ties can remarkably increase efficiency of these methods. We are also convinced that taking into account the asymmetry of social ties can also improve more advanced prediction methods, especially those supervised [75]. Finally, since in many ways, scientific collaboration networks are very specific, a natural continuation of the research presented here would be to check whether similar results can be obtained by analysing other (more typical) social networks, or even other complex networks, not necessarily social ones.

## SUPPLEMENTARY MATHERIALS

All the data used in this paper as well as the software developed for calculation of overlaps and for prediction can be found at [76].

[1] S. Wasserman and K. Faust, *Social Network Analysis – Methods and Applications* (Cambridge Univ. Press, 1994).

[2] D. Easley and J. Kleinberg, *Networks, Crowds, and Markets. Reasoning about a Highly Connected World* (Cambridge Univ. Press, 2010).

[3] D. M. J. Lazer, A. Pentland, D. J. Watts, and et al., Computational social science: Obstacles and opportunities, Science **369**, 1060 (2020).

[4] R. Conte, N. Gilbert, G. Bonelli, and et al., Manifesto of computational social science, Eur. Phys. J. Spec. Top. **214**, 325 (2012).

[5] J. Giles, Computational social science: Making the links, Nature **488**, 448 (2012).

[6] A. Fronczak, M. J. Mrowinski, and P. Fronczak, Scientific success from the perspective of the strength of weak ties, Sci. Rep. **12**, 5074 (2022).

[7] S. N. Dorogovtsev, A. V. Goltsev, and J. F. F. Mendes, Critical phenomena in complex networks, Rev. Mod. Phys **80**, 1275 (2008).

[8] A. Barrat, M. Barthélemy, and A. Vespignani, *Dynamical Processes on Complex Networks* (Cambridge Univ. Press, 2008).

[9] R. Pastor-Satorras, C. Castellano, P. Van Mieghem, and A. Vespignani, Epidemic processes in complex networks, Rev. Mod. Phys. **87**, 925 (2015).

[10] H. Mattie, K. Engo-Monsen, R. Ling, and J.-P. Onnela, Understanding tie strength in social networks using a lo-

cal "bow tie" framework, Scie. Rep. **8**, 9349 (2018).

[11] L. Lü, L. Pan, T. Zhou, and et al., Toward link predictability of complex networs, Proc. Natl. Acad. Sci. USA **112**, 2325 (2015).

[12] J. Kim, Scale-free collaboration networks: An author name disambiguation perspective, J. Assoc. Inf. Sci. Technol. **70**, 685 (2019).

[13] Z. Zuo and K. Zhao, Understanding and predicting future research impact at different career stages — a social network perspective, J. Assoc. Inf. Sci. Technol. **72**, 454 (2021).

[14] E. Ubaldi, R. Burioni, V. Loreto, and F. Tria, Emergence and evolution of social networks through exploration of the adjacent possible space, Commun. Phys. **4**, 28 (2021).

[15] M. Granovetter, The strength of weak ties, Am. J. Sociol. **78**, 1360 (1973).

[16] M. Granovetter, *Getting a Job: A Study of Contacts and Careers* (Univ. Chicago Press, 1995).

[17] J.-P. Onnela, J. Saramäki, J. Hyvönen, G. Szabó, and et al., Structure and tie strengths in mobile communication networks, Proc. Natl. Acad. Sci. USA **104**, 7332 (2007).

[18] J.-P. Onnela, J. Saramäki, J. Hyvönen, G. Szabø, and et al., Analysis of a large-scale weighted network of one-to-one human communication, New J. Phys. **9**, 179 (2007).

[19] M. Szell and S. Thurner, Measuring social dynamics in a massive multiplayer online game, Soc. Netw. **32**, 313

(2010).

[20] Šuvakov M., M. Mitrović, V. Gligorijević, and B. Tadić, How the online social networks are used: dialogues-based structure of myspace, J. R. Soc. Interface **10**, 20120819 (2013).

[21] S. Pajevic and D. Plenz, The organization of strong links in complex networks, Nature Phys. **8**, 429 (2012).

[22] R. K. Pan and J. Saramäki, The strength of strong ties in scientific collaboration networks, EPL **97**, 18007 (2012).

[23] Q. Ke and Y.-Y. Ahn, Tie strength distribution in scientific collaboration networks, Phys. Rev. E **90**, 032804 (2014).

[24] J. Ureña Carrion, J. Saramäki, and M. Kavelä, Estimating tie strength in social networks using temporal communication data, EPJ Data Sci. **9**, 37 (2020).

[25] P. S. Park, J. E. Blumenstock, and M. W. Macy, The strength of long-range ties in population-scale social networks, Science **362**, 1410 (2018).

[26] L. Lü and T. Zhou, Link prediction in complex networks: A survey, Physica A: Statistical Mechanics and its Applications **390**, 1150 (2011).

[27] V. Martínez, F. Berzal, and J. C. Cubero, A survey of link prediction in complex networks, ACM Comput. Surv. **49**, 1 (2017).

[28] L. A. Adamic and E. F. Adar, Friends and neighbors on the web, Soc. Netw. **25**, 211 (2003).

[29] T. Zhou, L. Lü, and Y.-C. Zhang, Predicting missing links via local information, Eur. Phys. J. B **71**, 623 (2009).

[30] D. Schall, Link prediction in directed social networks, Soc. Netw. Anal. Min. **4**, 157 (2014).

[31] E. Bütün and M. Kaya, A pattern based supervised link prediction in directed complex networks, Physica A **525**, 1136 (2019).

[32] X. Wang, X. Zhang, C. Zhao, Z. Xie, S. Zhang, and D. Yi, Predicting link directions using local directed path, Physica A **419**, 260 (2015).

[33] Y. Liu, T. Li, and X. Xu, Link prediction by multiple motifs in directed networks, IEEE Access **8**, 174 (2020).

[34] Q.-M. Zhang, L. Lü, W.-Q. Wang, Yu-Xiao, and T. Zhou, Potential theory for directed networks, PLoS ONE **8**, e55437 (2013).

[35] Z. Xie, Z. Ouyang, and J. Li, A geometric graph model for coauthorship networks, J. Informetr. **10**, 299 (2016).

[36] M. E. J. Newman, The structure of scientific collaboration networks, Proc. Natl. Acad. Sci. USA **98**, 404 (2001).

[37] M. E. J. Newman, Coauthorship networks and patterns of scientific collaboration, Proc. Natl. Acad. Sci. USA **101**, 5200 (2004).

[38] M. E. J. Newman, Scientific collaboration networks. I. Network construction and fundamental results, Phys. Rev. E **64**, 016131 (2001).

[39] M. E. J. Newman, Scientific collaboration networks. II. Shortest paths, weighted networks, and centrality, Phys. Rev. E **64**, 016132 (2001).

[40] J. Tang, J. Zhang, L. Yao, J. Li, L. Zhang, and Z. Su, Arnetminer: Extraction and mining of academic social networks, in *Proceedings of the Fourteenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD'2008)* (2008) pp. 990–998.

[41] DBLP Citation Network Dataset, `https://www.aminer.org/citation`, accessed: 2022-08-30.

[42] APS data sets for research, `https://journals.aps.org/datasets`, accessed: 2022-08-30.

[43] T. K. project", `http://konect.cc/networks/opsahl-collaboration/`, accessed: 2022-08-30.

[44] J. Tang, A. C. M. Fong, B. Wang, and J. Zhang, A unified probabilistic framework for name disambiguation in digital library, IEEE Transactions on Knowledge and Data Engineering **24**, 975 (2012).

[45] M. Müller, F. Reitz, and N. Roy, Data sets for author name disambiguation: An empirical analysis and a new resource, Scientometrics **111**, 1467 (2017).

[46] A.-L. Barabási, H. Jeong, Z. Néda, and et al., Evolution of the social network of scientific collaborations, Physica A: Statistical Mechanics and its Applications **311**, 590 (2002).

[47] R. Guimera, B. Uzzi, J. Spiro, and L. A. N. Amaral, Team assembly mechanisms determine collaboration network structure and team performance, Science **308**, 697 (2005).

[48] X. Sun, J. Kaur, S. Milojević, and et al., Social dynamics of science, Sci. Rep. **3**, 1069 (2013).

[49] C. Zhang, Y. Bu, Y. Ding, and J. Xu, Understanding scientific collaboration: Homophily, transitivity, and preferential attachment, J. Assoc. Inf. Sci. Technol. **69**, 72 (2018).

[50] W. Lu, Y. Ren, Y. Huang, Y. Bu, and Y. Zhang, Scientic collaboration and career stages: An ego-centric perspective, J. Informetr. **15**, 101207 (2021).

[51] Z. Xie, Predicting the number of coauthors for researchers: A learning model, J. Informetr. **14**, 101036 (2020).

[52] Y. Bu, Y. Ding, X. Liang, and D. S. Murray, Understanding persistent scientific collaboration, J. Assoc. Inf. Sci. Technol. **69**, 438 (2018).

[53] J. Kim and J. Diesner, Coauthorship networks: A directed network approach considering the order and number of coauthors, J. Assoc. Inf. Sci. Technol. **66**, 2685 (2015).

[54] F. Radicchi, C. Castellano, F. Cecconi, and et al., Defining and identifying communities in networks, Proc. Natl. Acad. Sci. USA **101**, 2658 (2004).

[55] D. Liben-Nowell and J. Kleinberg, The link-prediction problem for social networks, J. Assoc. Inf. Sci. Technol. **58**, 1019 (2007).

[56] B. Schafer, D. Frankowski, J. Herlocker, and S. Sen, Collaborative filtering recommender systems, in *The Adaptive Web* (Springer, 2007) p. 291.

[57] R. Guns and R. Rousseau, Recommending research collaborations using link prediction and random forest classiers, Scientometrics **101**, 1461 (2014).

[58] A. Clauset, C. Moore, and M. E. J. Newman, Hierarchical structure and the prediction of missing links in networks, Nature **453**, 98 (2008).

[59] S. Behrouzi, Z. Shafaeipour Sarmoor, K. Hajsadeghi, and K. Kavousi, Predicting scientific research trends based on link prediction in keyword networks, J. Informetr. **14**, 101079 (2020).

[60] C. Pornprasit, X. Liu, P. Kiattipadungkul, and et al., Enhancing citation recommendation using citation network embedding, Scientometrics **127**, 233 (2022).

[61] F. Karimi, S. Lotfi, and H. Izadkhah, Community-guided link prediction in multiplex networks, J. Informetr. **15**, 101178 (2021).

[62] M. E. J. Newman, Clustering and preferential attachment in growing networks, Phys. Rev. E **64**, 025102 (2001).

[63] P. Jaccard, Étude comparative de la distribution florale dans une portion des alpes et des jura, Bull. Soc. Vaud. Sci. Nat. **37**, 547 (1901).

[64] J. A. Hanley and B. J. McNeil, The meaning and use of the area under a receiver operating characteristic (roc) curve, Radiology **143**, 29 (1982).

[65] S. J. Mason and N. E. Graham, Areas beneath the relative operating characteristics (roc) and relative operating levels (rol) curves: Statistical significance and interpretation, Q.J.R. Meteorol. Soc. **128**, 2145 (2002).

[66] J. Davis and M. Goadrich, The relationship between precision-recall and roc curves, in *PProceedings of the 23rd international conference on Machine learning (ICML '06)* (Association for Computing Machinery, 2006) p. 233.

[67] B. Zhu, Y. Xia, and X.-J. Zhang, Weight prediction in complex networks based on neighbor set, Sci. Rep. **6**, 38080 (2016).

[68] H. Yin, A. R. Benson, and J. Leskovec, The local closure coefficient: a new perspective on network clustering, in *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, WSDM '19 (2019) p. 303–311.

[69] H. Yin, A. Benson, and J. Ugander, Measuring directed triadic closure with closure coefficients, Network Science **8**, 551 (2020).

[70] L. Lü and T. Zhou, Link prediction in weighted networks: The role of weak ties, EPL **89**, 18001 (2010).

[71] R. I. Dunbar, Neocortex size as a constraint on group size in primates, J. Hum. Evol. **22**, 469 (1992).

[72] J. Zhao, L. Miao, J. Yang, and et al., Prediction of links and weights in networks by reliable routes, Sci. Rep. **5**, 12261 (2015).

[73] J. Li, J. Peng, S. Liu, X. Ji, X. Li, and X. Hu, Link prediction in directed networks utilizing the role of reciprocal links, IEEE Access **8**, 28668 (2020).

[74] J. Li, J. Peng, S. Liu, and Z. Li, Mining missing links in directed social networks based on significant motifs, in *2020 IEEE 10th International Conference on Electronics Information and Emergency Communication (ICEIEC)* (2020) pp. 31–38.

[75] C. Fu, M. Zhao, L. Fan, and et al., Link weight prediction using supervised learning methods and its application to yelp layered network, IEEE Transactions on Knowledge and Data Engineering **30**, 1507 (2018).

[76] P. Fronczak, A. Fronczak, M. J. Mrowinski, and K. P. Orzechowski, Supplementary matherials for "Asymmetry of social interactions and its role in link predictability:the case of coauthorship networks", Mendeley Data, V1, doi:10.17632/8x6y22jzfz.1 (2022).