

# Visualization And Pattern Analysis on Cross-domain Scholarly Collaborations

Lingkun Kong, Bo Wang, Jiaqi Liu, Luoyi Fu, Xinbing Wang  
Shanghai Jiao Tong University, China

**Abstract**—Interdisciplinary collaborations, i.e., cross-domain scholarly collaborations play pivotal roles in scientific study, which have generated huge impact to society. However, what exactly are these cross-domain collaborations look like? Do classical social science features also exist in cross-domain collaborations? How to classify these collaborations? Do these collaborations impact research works’ quality in a certain pattern?

Due to theoretical and technical difficulties, there have been few studies that provide a systematic and practical understanding of these cross-domain collaborations. And no previous work properly answer these questions above. In this paper, we bridge this gap by analyzing cross-domain scholarly data, which well reflect and portray these cross-domain collaborations, by using real scholarly datasets – *Microsoft Academic Graph* [1] with 126 million papers collected from around 50 thousand domains.

By empirical exploration, we first visualize the cross-domain scholarly collaborations by using real-world scholarly data. And we find the power-law distribution, the classical distribution property in social science well fit our observation results. Then, in order to classify different collaborations, we design a quantification method to evaluate the distance, i.e., closeness relationship among different scientific domains and thus finding and evaluating the boundary in collaborations. And further with time data, we make case study to learn evolving closeness relationship, i.e., evolving boundary from different cross-domain collaborations. Moreover, we dig into the papers’ citation structure and compare papers’ citation distribution with papers’ domains distribution in cross-domain collaborations. We find these collaborations do impact the influence of papers, i.e., scientific works in a *peak* pattern. And therefore we can predict the influence of cross-domain collaborations.

Based on our empirical observations, we also make efforts in proposing novel models that can well reproduce the properties or patterns we discovered. To illustrate, we design a model of cross-domain power-law (*MCP*) to captures the power-law distributions in cross-domain data. And we reproduce the *peak* pattern by the help of gaussian distribution model. Moreover, through both theoretical analysis and empirical evaluations, we demonstrate that our models can accurately reproduce the features as well the patterns we probe in real-world dataset.

## I. INTRODUCTION

Interdisciplinary collaborations, i.e., cross-domain scholarly collaborations play pivotal roles in scientific study, which have generated huge impact to society.

For instance, collaborations between biology and computer science create the field of bioinformatics. Thanks to cross-domain collaborations, originally extremely expensive tasks such DNA sequencing have become scalable, which affords much broader usage. Another example goes to the field of medical informatics. Today medicine and data mining are working together in medical informatics, which is a big growth area that is expected to have huge impact on medicine [2] [3].

In literature, cross-domain scholarly collaboration is a theme that runs through large parts of scientific research [4] [5] [6]. However, previous works still have some limitations and leave four fundamental questions unanswered due to theoretical and practical difficulties.

The first fundamental question is **what are cross-domain scholarly collaborations look like?** Since cross-domain scholarly collaboration is a quite abstract concept, previous works only give definition without further illustration. Also, the paucity of real-world labelled interdisciplinary scholarly data adds difficulty to portray these collaborations comprehensively with fancy examples. However, to explore the features of cross-domain scholarly collaborations as well the hidden patterns, concise explanation is far away from what we need. In this case, visualizing the cross-domain scholarly collaborations is essential for further analysis and can help to offer researchers a better understanding of cross-domain collaborations.

The second fundamental question is **do classical social science features also exist in cross-domain collaborations?** In fact, papers [5] which study cross-domain collaborations make modeling with the intuitive assumption that the classical social science features should also exist in cross-domain collaborations’ situation. However, few of them make experiments to validate these properties by using real-world data with a large scale, which reduces models’ rigorousness. Therefore, using experimental results to verify the existence of classical properties is necessarily essential as it can convincingly demonstrate that many classical theorems can also be implemented in studying these collaborations, and can ensure the model could use refer previous studies in classical cases.

The third question goes to **how to classify cross-domain scholarly collaborations?** Researchers seldomly consider the difference of cross-domain scholarly collaboration itself. In other words, previous studies neglect the fact that the collaboration between different domains should be different as domains where collaboration establishes might have totally different relationship. For instance, the collaboration between mathematical and computer science should be different with the collaboration between mathematical and history, since the later two domains seems to have really far-away relation. Therefore, in order to classify different kinds of cross-domain scholarly collaboration, we need to find a proper way to reveal and evaluate this difference.

The fourth question is **do these collaborations impact research works’ quality in a certain pattern?** In order

to promote the communication of knowledge from different scientific fields and to solve interdisciplinary problems, research communities advocate interdisciplinary collaborations for many years. [6]–[9] And people believe that cross-domain scholarly collaborations could generate work of high quality and broad influence. However, no study probes this by specific data analysis.

In this paper, we bridge this gap – answering the questions above by implementing elaborate data-mining methods on big scholarly data. And we also build a model to fit properties we discover in this paper, which could give answers to all above questions.

For the first question, we study a real-world database – *Microsoft Academic Graph* (MAG) of large-scale which thoroughly explores the data of cross-domain scholarly collaboration. Based on the massive data, we visualize these collaborations in section 3, where we take domains in computer science as an example to illustrate the collaborations between different domains.

For the second question, based on the data from real-world database, we prove the existence of classical power-law distribution in cross-domain scholarly collaboration by experiment results, which are presented in section 3. And therefore, we can found our model with classical social science theorems and refer to previous social network modeling results.

For the third question, by using the domains’ hierarchy information in real-world dataset, in section 3, we visualize scientific domains’ hierarchy structure for understanding the closeness relationship, or boundary between different domains. Based on our observation, we propose the concept of *domains’ distance* to quantify the boundary between different domains and thus use quantified value to classify different kinds of cross-domain scholarly collaborations. Moreover, we add time slots to the relationship we try to explore, and find interesting evolving pattern in domains’ relationship by making case study in the field of data-mining.

For the fourth question, in section 3, we explore the papers’, i.e., works influential performance with the concern of cross-domain scholarly collaborations. That is based on paper’s membership relation with domains it belongs to, we study paper’s citation distribution. And surprisingly we get a *peak* pattern – the paper’s citation number is likely to get a maximum value when paper’s domains number comes to a certain amount and drop when domains number exceeds this amount. Further, we dig into the paper’s citation, dividing these citation into four parts according to the network structure of paper’s membership hierarchy, and thus know the decreasing citing possibility with the increasing cross-domain distance, which complies with our intuitive thinking.

Further, combining all these results together, we purpose a theoretical model to try to capture all properties we find in cross-domain scholarly collaborations. Also, we mathematically analyze the correctness of our model.

The paper is organized as follows. In Section 2, we discuss relevant literatures. In Section 3, We briefly introduce dataset we use and list all visualization and experimental results.

We give our model in Section 4 and analyze our model mathematically in Section 5. Finally we draw conclusion in Section 6.

## II. RELATED WORK

In literature, cross-domain scholarly collaboration is a theme that runs through large parts of scientific research [4]–[6]. However, previous works still have some limitations and leave four fundamental questions unanswered due to theoretical and practical difficulties.

The first fundamental question is what are cross-domain scholarly collaborations look like. In fact, though some works shows the authors’ collaboration by co-author network [10], papers’ collaboration by citation network [11] or even keywords network [12] no previous work which studies the cross-domain collaborations in research area directly visualize the domains’ collaborations. In order to answer this question, we observe the real-world database – *Microsoft Academic Graph* (MAG) [1] and use the labeled domains information to visualize the interdisciplinary collaborations.

The second question is do classical social science features also exist in cross-domain collaborations. Based on the data from real-world database, we prove the existence of classical power-law distribution [13]–[15] in cross-domain scholarly collaboration by experiment results, which provide solid foundation for further modeling.

The third question is how to classify cross-domain scholarly collaborations. Researchers [5] seldomly consider the difference of cross-domain scholarly collaboration itself, that is previous studies neglect that the collaboration between different domains is different as domains where collaboration establishes might have really different relationship. We bridge this gap by visualizing the domains’ hierarchy structure and then quantifying the distance, or boundaries between domains, which helps us to classify those collaborations.

The fourth question is do these collaborations impact research works’ quality in a certain pattern. Literatures [16] [17] studies research works’ impact by many different approaches.

However, few of them combine the interdisciplinary collaborations’ information with scientific works’ quality. Therefore, we explore works’ influential performance with the concern of cross-domain scholarly collaborations. To illustrate, we visualize the citations’ structure after adding domains’ information, and further observing the interesting peak pattern.

Moreover, we propose a novel model called *model of cross-domain power-law* – MCP to capture properties we find in cross-domain scholarly collaborations. And our model absorbs many previous studies knowledge and design, including random graph model built by Chakrabarti and Faloutsos [18], preferential attachment model [19] proposed by Barabasi et al., edge-copy model [20] created by R. Kumar et al., and affiliation network model [21] advanced by Silvio Lattanzi et al.

## III. VISUALIZATION AND PATTERN STUDY

In this part, by the help of data-mining methods, we properly explore the cross-domain scholarly data in real world

big scholarly data, and discover several stimulating results. Further, we implement visualization to our discoveries to make our result easy to view.

We make our observation both on domain-oriented aspect and paper-oriented or literature-oriented aspect. In the domain's perspective, we observe firstly the domain's number power-law distribution with its papers number as well its subdomain number. And we also study the closeness among different domains, i.e. using co-paper's ratio to describe the relationship between different domain. Moreover, we add time information to the relationship we try to explore, by mining paper's publication date, and thus find interesting evolving pattern in domains' relationship.

From the paper's perspective, we explore the paper's cross-domain performance, which includes its membership relation based on domains paper belong to and its cross-domain citation distribution. We firstly get the power-law distributed paper number with the number of domains paper belongs to. And then, based on paper's membership relation with domains it belongs to, we study paper's citation distribution. And surprisingly we get a "peak" distribution – the paper's citation number is likely to get a maximum value when paper's domain number comes to a certain amount. Further, we dig into the paper's citation, dividing these citation into four parts according to the network structure of paper's membership hierarchy, and thus know the decreasing citing possibility with the increasing cross-domain distance, which complies with our intuitive thinking. And we also studies successor phenomenon ...

#### A. Brief Introduction For MAG

We give our experiments based on *Microsoft Academic Graph* (MAG) which is an official and authoritative scholarly dataset containing massive scholarly information of publications such as titles, authors, conferences, fields of study and citations. Around 126 million papers in 50 thousand domains are included in this database and the published years of them vary from 1800 to 2016.

As we study the cross-domain collaborations of scholarly data, we mainly launch research on the scholarly data related with fields of study, i.e. domains in MAG. And the fields in MAG can be divided into four layers and we call them from L0, L1, L2, L3, where layer with lower number, which represents bigger scientific domain contains the layer with higher number, i.e. smaller scientific domain of study. E.g., we get a domain labeled with L0 layer in MAG called "Computer Science", which contains several domains labeled by L1 layers according to the MAG's hierarchy table, including "Artificial Intelligence", "Database", "Data-mining", "Computer Hardware", and etc. And "Data-mining" can contains several domains with lower layer label such as "Big data" in L2 Layer and "K-optimal pattern discovery" in L3 layer. Moreover, one thing supposed to be noticed is that the hierarchy in MAG is heterogeneous, which means the domain of L1 layer can directly contains or relates domains of L2 and L3 layers.

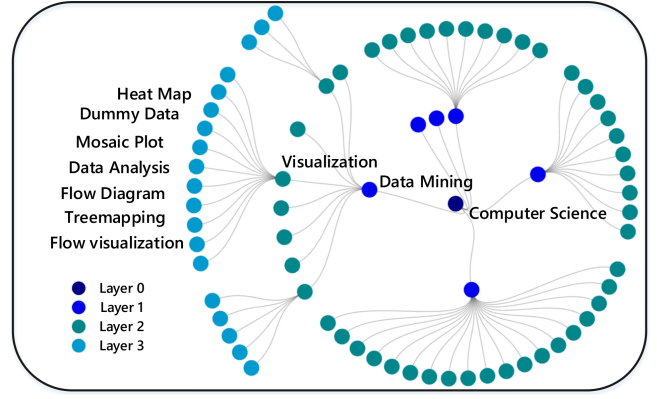


Fig. 1. Domain Hierarchy Structure in MAG

In figure 1, hierarchy example of MAG dataset can be viewed. This figure illustrates a part of the hierarchy structure of the domain "Computer Science". We pick up several nodes and mark their names besides the nodes. As can be seen from the figure, the lower the layer is, the more specifically the domain is.

#### B. Domain-oriented Exploration

Power-law distributed degree is a common feature of social networks, which is also well studied by many existing literatures. And XXX's work also find the power-law distribution also exists in scholarly network – using network structure to present the scholarly data. Stimulated by this result, when we studying domains' information in scholarly data, we also get power-law distribution. In figure 2, it is clearly can be viewed that there exists two kind of power-law distribution.

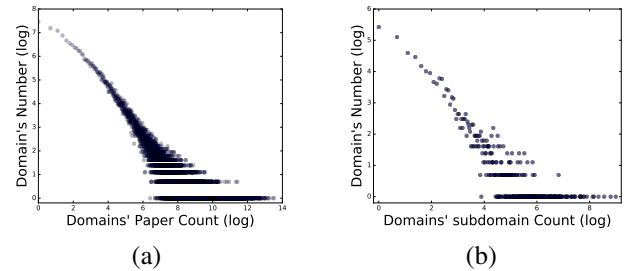


Fig. 2. Domain-oriented power-law distribution

First is power-law distributed domains number with papers' count in these domains, as shown in the figure 2.a. And second is power-law distribution between domains number with their sub-domains' count, which is drawn in figure 2.b.

We intuitively use co-paper's number among different domains to describe the correlation of domains, i.e. the closeness of relationships. Though behind this papers, a more complex networking topology, i.e. the paper's reference and citation network might include useful information for judging closeness between different domains, we originally use the number of co-papers to verify the correlation of these domains as co-paper is the bridge which levels up the gap among different

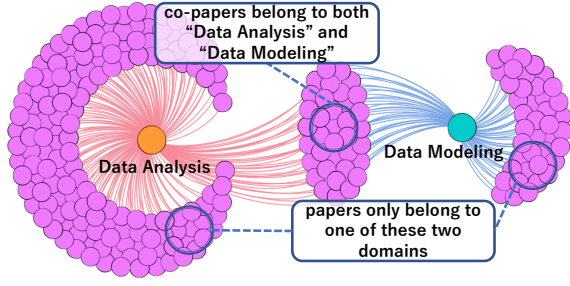


Fig. 3. The paper number distribution among target domains. Where are two domains in L3 layer, called "Data Analysis" and "Data Modeling".

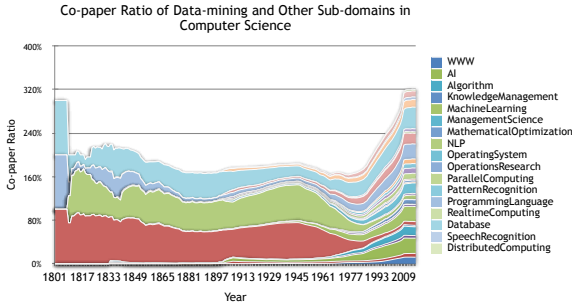


Fig. 4. Domain relationship in "Computer Science"

domains and the number of co-paper, the basic and essential feature of co-papers can linearly evaluate the closeness among domains.

First of all, we simply visualize this relationship in figure 3. And we can clearly find that different domains in fact have many co-papers. And therefore we can use the co-paper information to quantify the relationship between different domains.

Besides, we can add the time information by the data in MAG's paper publication time table to further explore the evolving pattern of the relationship among different domains.

What we do is to label every paper of two domains by their publication date, and thus co-papers are also labeled by time info. After that, we calculate the evolving co-paper ratio of one specific domain vs other different domains. And by that, we can find the changing relationship among one domain and other different domains. For instance, in figure 4, we study in the "computer science" domain, all L1 domains' relationship with "data-mining" domain.

In this figure, we can view the evolving relationship – the fluctuating correlation in different years among Data Mining and other different computer science related domains, such as AI, Algorithm, NLP and etc. which are listed in the figure's legend. By this figure, the affiliation of "Data Mining" domain is revealed, i.e. we can see the switching closeness of data mining to other domains. For instance, in the early stage, "Data Mining" is highly related with "Bio-information" – "Data

Mining" has almost 90% co-papers with "Bio-information" which refers high correlation while in current years, the co-papers' ratio of "Bio-information" has been decreased to a very low level, which indicates the degeneration of relationship between "Data Mining" and "Bio-information". Moreover, it can be viewed that recently the "Data Mining" domain has always preferred to combine knowledge in publication from "Artificial Intelligence" and "Machine Learning" domain as their co-paper ratio is rapidly growing.

### C. Paper-oriented Exploration

We also explore several features of paper's cross-domain performance, including papers' domain distribution, and paper's cross-domain citation distribution. First of all, we find that the paper's number is power-law distributed depending on paper's domain number.

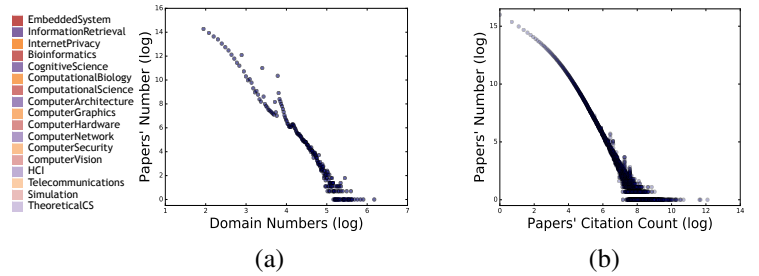


Fig. 5. Paper-oriented power-law distribution

In figure 5, the left graph presents relation between Art's paper domains count and paper number. Though the domain count is not large, it seems to be still power-law distribution. Right graph is computer science's. This graph more smoothly simulates the power-law distribution since there are much more papers in computer science and much more cross-domain collaboration.

### Paper citation distributions

We believe cross-domain paper citation performance is an essential indicator for paper's quality, i.e. whether the paper is good or not, since a good paper might cross several different domains as nowadays research emphasizes on the study's width. More specifically, a scientific study combining several fields of studies knowledge together, might be more likely to catch others' attention and generate stimulating results. For instance, paper in computer science domain with solid mathematical foundation or theory, i.e. also in mathematical domains are always more likely to be a good paper. Moreover, recently, scientific study launched in biology domain make many refreshing breakthroughs by combining computer science, especially data mining and machine learning's knowledge.

We extract papers' citation information and use their citations domain information to draw paper citation distribution map, which intuitively study the structure of paper citation distribution with domain information.

In figure 6, the green node in the center of the graph indicates the paper we want to study, and the pink node

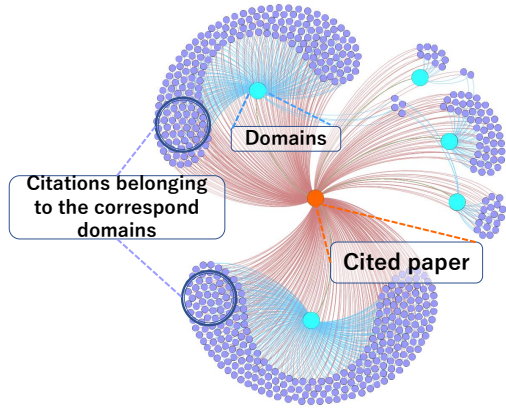


Fig. 6. Paper's citation domain structure

represents paper which cite the paper we study. The light green node refers to L1 layer domain, red node is L2 layer domain while blue is the L3 layer domain. And the line between green node to pink node means the paper's citation, while other lines among papers to domains describes membership relation between papers and domains. [More specific information need to be updated.](#)

What we can see is that the paper's citation can be divided into several cliques according to the paper's domain info. And there exists clique overlapping as citation papers might have more than single domains which are similar with the paper being cited. And we consider these papers to be successors of the original paper, which we will discuss later as these papers are more likely to be papers which inherit studying paper's idea or methods. [successor phenomenon](#)

#### Paper's average citation distribution over domain number.

We believe that paper's citation is related to its domain count, as the paucity of domains might constrain a paper's impact in a small domain, while too many domains might decrease paper's quality as too many domains might distract or diffuse author's attention and harms the paper's depth. Therefore, we plot the paper's average citation number over paper's domain count in below figures.

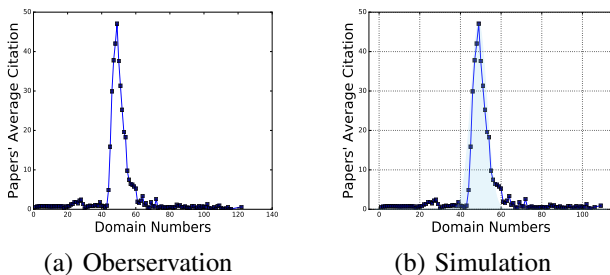


Fig. 7. Papers' average citation count with certain domain numbers of overall paper in the database

Figure III-C is the overall paper's citation count over their domain count. And we can clearly find a peak when paper's

domain count goes to about 50. And we check these paper's number in case paper's number is too small to represent a pattern when paper's domain count exceeds a certain number, e.g. 50. And we find when paper's domain count is less than 100, the paper's number is large enough to calculate the average citation count. In other words, this "peak" pattern does exist. And we use gaussian distribution to simulate this peak pattern, which is also shown in the figure III-C. And the modeling details and simulating results will be explained later.

Moreover, when we look into smaller domains, computer science's sub-domains, for instance, we find that this rule is almost generally valid. And in fact though in some sub-domains, Information Retrieval for instance, though paper's domain count is small, the peak still exist.

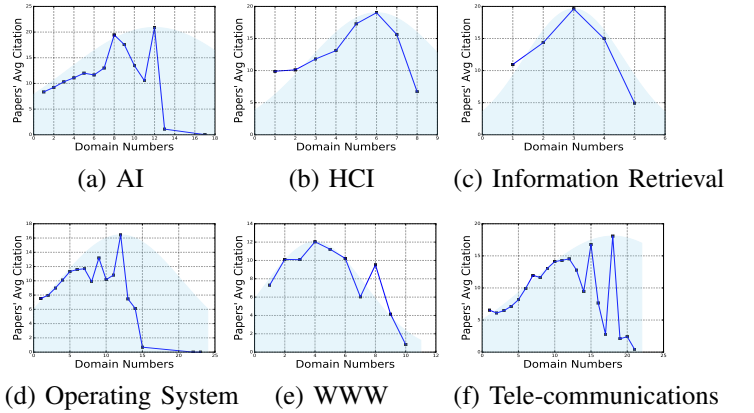


Fig. 8. Papers' average citation number influenced by their domain numbers in Computer Science Domain

In figure 8, we can easily view these properties. Also, we use gaussian distribution with trained parameters to simulate this peak pattern.

**Paper's citation performance over more specific domains**  
Further, we look into the paper's citation distribution. And we divide one paper's citation into several different types according to the cross-domain step. Intuitively speaking, the cross-domain step is a parameter we use to evaluate the paper's cross-domain distance. For instance, a paper in "Astrology" domain cites paper in "Data-mining" has longer cross-domain distance than paper in "Data Base" cite paper in "Data-mining". But how to quantify this difference? We take advantage of our dataset's level partition. And we set cross-domain step into four class, i.e. merge at L3 layer, L2 layer, L1 layer and L0 layer. Noticed that a paper may belong to several domains, the domains of two papers can merge at different layers. We choose the lowest merge layer as the cross-domain distance since the domains at lower layer can represent the paper more specifically. figure 9 presents our dividing rules:

In figure 9, the domains at layer 3 such as "Graphical Tools" and "Data Analyze" merge at layer 2, while "Graphical Tools" and "Network Simulation" merge at layer 1. If the cited paper and its citation belongs to "Graphical Tools" and "Network Simulation" respectively, their domains merge at layer 1 and layer 0, and we chose the lowest merge layer, i.e. layer 1, as the



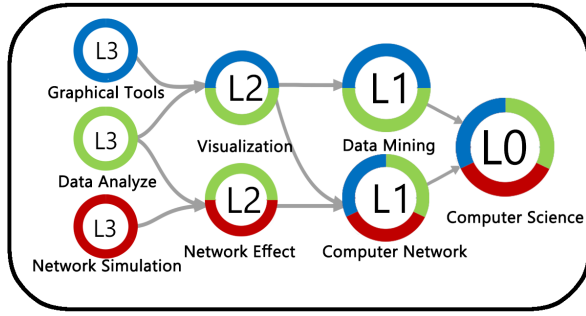


Fig. 9. The example of domains merging at different layers

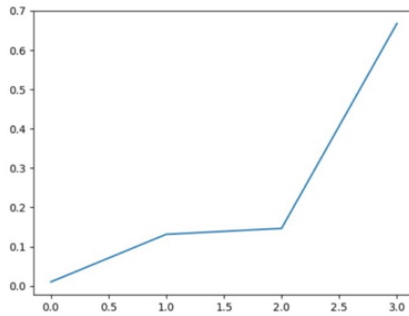


Fig. 10. 1000 papers cross L0, L1, L2, L3 domains ratio distribution

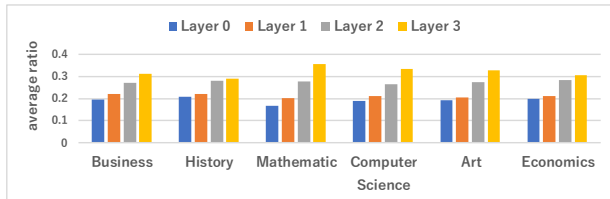


Fig. 11. The average cross domains ratio distribution over L0 domains

merge layer of this citation. And according to our rules, it is easily can be found that a citation type is supposed to only belong to one class. And the possibility of a citation belongs to these 4 class should be summed into exactly 1. **Drawbacks:** we currently cannot figure the distance between math to cs and art to cs which is longer.

As we can see in figure 10, we pick up some papers with high citations to draw this distribution and the citation in cross-L3 type's possibility is much higher than cross L1 and cross L2 type's. And in fact cross L2 type's possibility is slightly higher than cross L1 types. And the cross L0's possibility is very low, almost 0.

In figure 11 we calculate the average cross domains ratio in several domains. The result shows that the citations of a paper

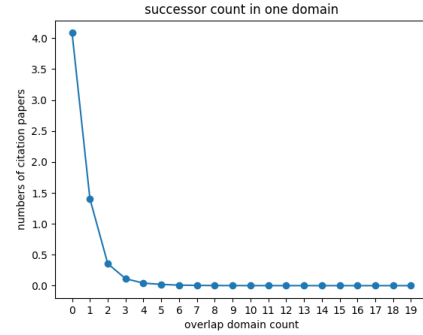


Fig. 12. overlap domain count and the number of cited papers in one domain

are tends to be in a relatively close distance of the cited paper. In some domains such as "Literature", the ratios of cross L1 level and L2 level are much lower than other domains, which indicates that in those domains, the interdisciplinary trend is not so strong.

**successors** In figure 6, we find some citation papers which might have more than one domains that overlap with the fields of cited paper. This phenomenon exposes that the cited paper introduces more papers to cross the domains as the cited paper does. It illustrates the real ability of the cited paper — the ability to lead more papers to be interdisciplinary. So we call this phenomenon as successor phenomenon.

In figure 12, we count the number of citation papers averagely for one cited paper and the number of the corresponding overlapping domains with the cited paper. We find that successors which have more overlapping domains with the cited paper are relatively less. It means that even the cited paper crosses many domains, only a few citation papers cross the same domains as the cited paper does. [Waiting for clearer explanation and more graphs](#)

#### IV. MODELING AND ANALYSIS ON OBSERVATION

In this section, based on the previous observations in real world dataset, we use three models to help to simulate or evaluate the scholarly data's cross-domain performance.

##### A. Power-law Distribution Modeling

In our observation, it can be viewed clearly that there exists two kinds of power-law distributions in our study. The first is the power-law distributed domain's number with paper's number in this domain. And the second is the power-law distribution between the paper's number and the this paper's relative domains' number. And here we construct a evolving model which can properly reproduce these power-law distribution in the cross-domain scholarly data. And we call this model as *model of cross-domain power-law* – MCP.

##### MCP Construction:

We use network structure to present the scholarly data. And in MCP, the graph is denoted as  $G(P, D)$ . Then, we use bipartite graph to present the inter-correlation between elements. Besides, we also focus on the intra-features of

every element. For an intuitive understanding, we illustrate the framework of our evolving scholarly model in Figure IV-A. It contains:

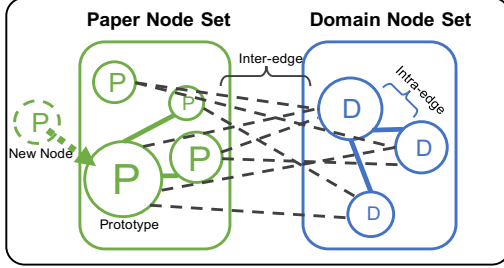


Fig. 13. Structure of evolving scholarly model

(1) **Two node sets:** Paper node set  $N_p$ , and domain node set  $N_d$ . The node in each node set is marked as  $n_p$  and  $n_d$ .

(2) **Inter-edge sets:** We denote the edges between every two node sets as inter-edge sets. And, we refer all edges between paper node and domain node as  $E_{pd}(E_{dp})$ , then an edge  $e_{n_p n_d}$  which belongs to  $E_{pd}$  means paper  $n_p$  belongs to the domain  $n_d$  or equivalently domain  $n_d$  has paper  $n_p$ .

(3) **Two intra-edge sets:** Intra-edge is the edge in the same node set, and our graph has two intra-edge sets, which we refer as  $E_{pp}$  and  $E_{dd}$ . If an edge  $e_{n_p^i n_p^j} \in E_{pp}$ , then we know that paper  $n_p^i$  and  $n_p^j$  have reference or citation relationship. While if an edge  $e_{n_d^i n_d^j} \in E_{dd}$ , then we know that domain  $n_d^i$  and  $n_d^j$  are directly connected in the domain hierarchy dataset.

In Figure IV-A, nodes are illustrated as colorful circles in each node set while intra edge and inter edge are labeled. And a new paper node is trying to preferentially attach himself with some heavily linked papers nodes (distinguished by their sizes) that are already in the paper set. With these nodes and edges of the model, we can well extract the structure of papers' cross-domain relationship in scholarly networks. We present notations in Table I for later convenience and describe the evolving process of the proposed model in the following subsection.

TABLE I  
NOTATIONS AND DEFINITIONS

Notations	Definitions
$N_p, N_d$	Node set of Paper and Domain
$E_{pd}$	Inter-edge set between nodes in $N_p$ and $N_d$
$E_{pp}, E_{dd}$	Intra-edge set of $N_p$ and $N_d$
$\alpha_p, \alpha_d$	Probability that a new node arrives in $N_p$ and $N_d$
$\beta_p, \beta_d$	Probability that an edge added in set $E_{pp}$ and $E_{dd}$
$c_{pd}, c_{dp}$	Number of edges added to set $E_{pd}$ at one time slot
$G(P, D)$	Graph of our cross-domain scholarly model
$B(N_p, N_d)$	Bipartite graph with sets $N_p, N_d$ , and $E_{pd}$

#### Evolving process:

While we defer the detailed evolving process of the proposed model to Algorithm 1, we would also like to provide a corresponding brief summary of the process. We first fix parameters including  $\alpha_i$ ,  $\beta_{ij}$  and  $c_{ij}$  where  $i \neq j \in \{p, d\}$ , and then assume that the evolution starts from an initial case

that can be modeled as an initial graph, showing that each node in the graph is linked to a number of nodes in other node sets. After initialization, for every time slot, we classified the process into five main steps: **1)** A new node, which can be randomly designated as a paper or a domain, is added to the graph. For clarity, here we only take the arrival of a new paper as example for explanation of the subsequent steps. And the symmetry also holds for the domain. **2)** With probability proportional to degree in  $B(N_p, N_d)$ , paper node  $n_p^d$  is chosen as prototype for the new node  $n_p$ . **3)**  $c_{pd}$  neighbors ( $n_d^1, \dots, n_d^{c_{pd}}$ ) of  $n_p^d$  in  $N_d$  are randomly chosen to have connections with node  $n_p$ . **4)**  $c_{pd}$  edges are added between  $n_d^1, \dots, n_d^{c_{pd}}$ . **5)** Edges between every two paper nodes are added with probability  $\beta_{pd}$  if they have a common domain.

For a better intuitive understanding of this evolving process, let us, for instance, consider the arrival of a new paper. This paper is likely to learn from an influential paper, which is thus selected as a prototype and influences the new paper on choosing research domains. To illustrate, this new paper will have high possibility to generate studies in the same domains like the old, influential paper. Besides, the papers belongs to the same domain are often relevant, indicating that these papers belong to these domains with a higher possibility to be connected, i.e. cite each other than those belong to different domains.

Similarly, when a new topic emerges in the literature, it is usually inspired by some existing topics (prototypes) and these topics are more likely to be related by the same papers they have.

#### Algorithm 1 Evolving Process

**Parameters:** Simulated time steps:  $T$ , Fixed probability  $\alpha_i$  that a new node arrives in  $N_i$ , fixed  $\beta_{ij} \in (0, 1)$  and integers  $c_{ij} > 0$  where  $i \neq j \in \{p, d\}$ .

**Initialisation:** In initial graph, the node in paper set has a certain number of neighbors with domain set. For example, a paper node  $n_p$  connects to at least  $c_{pd}$  domain nodes. So the inter-edge set  $E_{pd}$  has at least  $c_{pd}$  edges in the beginning.

1: **for**  $1 \leq t \leq T$  **do**

2: **1) Node arrival:** According to  $\alpha_p, \alpha_d$ , we decide the type of node to join the graph. In later discussion, we take the arrival of a new paper node  $n_p$  as example, and the symmetry also holds for domain.

3: **2) Preferentially chosen ProtoType:** A node  $n_p^d \in N_p$  is chosen as prototype for the new node, with probability proportional to its degree in  $B(N_p, N_d)$ .

4: **3) Edge copying:**  $c_{pd}$  edges are copied from  $n_p^d$ , that is,  $c_{pd}$  neighbors of  $n_p^d$ , denoted by  $n_d^1, \dots, n_d^{c_{pd}}$  in  $N_d$  are chosen uniformly at random, and the edges  $(n_p, n_d^1), \dots, (n_p, n_d^{c_{pd}})$  are added to the graph.

5: **4) Evolution inside:** For every two nodes  $n_p^x$  and  $n_p^y$  ( $x \neq y$ ), if they have a common domain, then with probability  $\beta_{pd}$ , an edge  $(n_p^x, n_p^y)$  is added in  $E_{pp}$ .

6: **end for**

## B. Peak

## C. Evaluation of Paper's Influence's Broadness

### V. THEORETICAL ANALYSIS

In this section, we mathematically analyze our MCP model and the **peak model** to confirm that our models can well reproduce properties in the real-world database, i.e. the scholarly network in this paper.

#### A. MCP Analysis

Here we prove that our MCP can well reproduce two kinds of the power-law distributions in our observation.

According to our model, we divide the nodes' degree into two types – the first is the *inter-degree*, i.e., the node degree between node sets, related with the growth of  $E_{pd}$ , we call it  $d^{ir}$  –  $ir$  here means inter. And the second is the *intra-degree*, i.e. the node degree inside node set, related with the growth of  $E_{ii}$ , i.e.,  $E_{pp}$  or  $E_{dd}$ , we call it  $d^{ia}$  –  $ia$  here means intra.

**Growth of inter-degree:** Assuming node  $n$  arrives at node set  $N_i$  at time  $t_0$  with initial inter-degree  $d_i^{ir}(t_0)$ , the inter-degree of  $n$  at time  $t > t_0$  is

$$d_i^{ir}(t) = \left(\frac{t}{t_0}\right)^{\lambda_i} d_i^{ir}(t_0),$$

where  $\lambda_i \in (0, 1)$  is a constant, and  $i \in \{p, d\}$ .

In fact, an implications can be deduced by this result, that the inter-degree  $d_i^{ir}(t)$  grows with polynomial rate in time  $t$ , following the power  $\lambda_i \in (0, 1)$ .

This implication gives the growth rate of node's inter-degree. And the detailed proof is given in Theorem 1.

**Growth of intra-degree:** Again, we set beginning time as  $t_0$  and the intra-degree of node set  $N_i$  at time  $t > t_0$  is  $d_i^{ia}(t)$ , then

$$d_i^{ia}(t) = \Theta\left(t^{\frac{1}{\lambda_j}+1}\right),$$

where  $\lambda_j$  represents the constant  $\lambda$  in  $N_j$ . For instance when  $i$  is  $p$ , i.e., the paper, the  $j$  represents the  $d$ , i.e., the domain.

The equation reveals that, in our model, the intra-degree of a node set actually is related with the inter-degree's growing rate variable  $\lambda$ . As in equation the intra-degree is positively related with the growing with time slot  $t$ , we can say the intra-degree also grows with time. The detailed proof is given in Theorem 2.

Also, we analyze nodes' power-law distribution in two cases – inter and intra-degree respectively.

**Distribution of inter-degree:** For the node  $n \in N_i$  in  $G(P, D)$  with  $t \rightarrow \infty$ , the inter-degree distribution of it follows

$$\mathbb{P}\{d_i^{ir}(t) = x\} \propto x^{-\frac{1}{\lambda_i}-1}.$$

And we find that the inter-degree  $d^{ir}$  follows the power-law distribution with exponent  $-\frac{1}{\lambda_i} - 1$ .

Results show our model well capture the power-law distribution of nodes' inter-degree, which are proved in Theorem 3 and verified by experimental measurements.

**Distribution of intra-degree:** For the node  $n$  in  $G(P, D)$  with  $t \rightarrow \infty$ , the intra-degree distribution of  $n \in N_i$  follows

$$\mathbb{P}\{d_i(t) = x\} \propto x^{-\omega_i},$$

where  $\omega_i$  is a constant which describes the exponential factor in power-law distribution.

This means our model well simulates the power-law distribution of nodes' intra-degree. And results are proved in Theorem 4.

Combining above four results together, it can be easily viewed in our model – MCP that the degree of the node in graph  $G(P, D)$  grows with polynomial rate in time  $t$ , and the growth rate differs from inter-degree to intra-degree of the node. Moreover, the inter-degree as well the intra-degree are proven to be powerlaw-distributed in MCP. Therefore, our MCP model can well simulate and reproduce our observation in real-world database.

**Theorem 1:** For graph  $G(P, D)$  generated after  $t$  time slots ( $t \geq t_0$ ), with the initial condition that a certain node  $n \in N_p$  is added to node set  $N_p$  at time  $t_0$  with the degree  $d^{ir}(t_0)$  from  $N_p$  to  $N_d$ , the inter-degree of  $n$  at time  $t$  satisfies

$$d_p^{ir}(t) = \left(\frac{t}{t_0}\right)^{\lambda_p} d_p^{ir}(t_0).$$

This result also holds for  $n \in N_d$  with symmetrical expressions.

**Proof:** At each time slot  $t$ , the inter-degree of node  $n \in N_p$  in  $B(N_p, N_d)$ , i.e.  $d_p^{ir}(t)$ , can only increase in follow case: a new node arrives at  $N_d$  and is connected to  $n$ , which results in  $d_p^{ir}(t) = d_p^{ir}(t-1) + 1$ .

In edge copying, we choose the prototype node according to its inter-degree, while the endpoint of any edge is chosen with equal probability. Thus, the probability that a new added edge in  $B(N_p, N_d)$  points to a certain node  $n$  is  $\frac{d_p^{ir}(t-1)}{s_p(t-1)}$ , where  $s_p(t-1)$  denotes the sum number of edges in  $B(N_p, N_d)$  at time  $t-1$ , and we have

$$s_p(t-1) = (\alpha_p c_{pd} + \alpha_d c_{dp})(t-1).$$

Then, we get

$$d_p^{ir}(t) - d_p^{ir}(t-1) = \alpha_d c_{dp} \frac{d_p^{ir}(t-1)}{s_p(t-1)}$$

With the initial condition that

$$d_p^{ir}(t) = \left(\frac{t}{t_0}\right)^{\lambda_p} d_p^{ir}(t_0), \quad (1)$$

where  $\lambda_p = \frac{\alpha_d c_{dp}}{\alpha_p c_{pd} + \alpha_d c_{dp}}$ .

By same approach we can obtain the expression result of  $d_d^{ir}(t)$  for nodes in  $N_d$ . That is

$$d_d^{ir}(t) = \left(\frac{t}{t_0}\right)^{\lambda_d} d_d^{ir}(t_0),$$

where  $\lambda_d = \frac{\alpha_p c_{pd}}{\alpha_p c_{pd} + \alpha_d c_{dp}}$ . Thus we complete the proof. ■

**Theorem 2:** For graph  $G(P, D)$  generated after  $t$  time slots ( $t \geq t_0$ ), with the condition that inter-degree in node set  $N_d$



growing with the power  $\lambda_d$ , the intra-degree of  $n \in N_p$  at time  $t$  satisfies

$$d_p^{ia}(t) = \Theta\left(t^{\frac{1}{\lambda_d}+1}\right).$$

This result also holds for  $n \in N_d$  with symmetrical expressions.

*Proof:* The intra-degree in  $N_p$  is generated by common neighbors in  $N_d$ .

When a certain node  $d \in N_d$  has node degree  $x$  from  $N_d$  to  $N_p$ , it has exactly  $x$  neighbors in  $N_p$ . Thus, the expected intra-degree in  $N_p$  added by this node is  $2\beta_p\binom{x}{2}$ , where  $\beta_p$  is the linking probability when two nodes inside node set  $N_p$  have a common neighbor node in  $N_d$ . And the number of nodes in  $N_d$  who have  $x$  neighbors in  $N_p$  is expected as  $|N_d|\mathbb{P}\{d_{ap}^{ir}(t) = x\}$  where  $\mathbb{P}$  denotes the probability that node in  $N_a$  having  $x$  neighbors in  $N_p$  exists and  $|N_a|$  denotes the total nodes in  $N_a$ . Therefore, the intra-degree generated by nodes with  $x$  neighbors in  $N_a$  is

$$\text{Contribution}(x) = 2\beta_p\binom{x}{2}|N_d|\mathbb{P}\{d_d^{ir}(t) = x\}. \quad (2)$$

Considering we add certain number of nodes with a certain probability in the node set, we get  $|N_a| = \Theta(t)$ . Thus, combining the result of Theorem 3, we get the intra-degree  $d_p^{ia}(t)$  in node set  $N_p$  contributed is

$$\begin{aligned} d_p^{ia}(t) &= \sum_{x=1}^{\max} \text{Contribution}(x) \\ &= \sum_{x=1}^{\max} 2\beta_p\binom{x}{2}|N_d|\mathbb{P}\{d_d^{ir}(t) = x\} \\ &= \Theta\left(\sum_{x=1}^{\max} x^2 x^{-\frac{1}{\lambda_d}-1} t\right) \\ &= \Theta\left(\sum_{x=1}^t x^{-\frac{1}{\lambda_d}+1} t\right), \end{aligned}$$

where  $\max$  presents the maximum inter-degree in  $E_{pd}$  which satisfies  $\max = \Theta(t)$ . By using the sum of  $p$ -series, we get

$$\sum_{x=1}^t x^{-\frac{1}{\lambda_d}+1} = t^{1-(1-\frac{1}{\lambda_d})}.$$

Therefore, we have  $d_p^{ia}(t) = \Theta\left(t^{\frac{1}{\lambda_d}+1}\right)$ .

By same approaches, we can also obtain the expression result of  $d_d^{ia}$  for nodes in  $N_d$ , thus we complete the proof. ■

*Theorem 3:* For graph  $G(P, D)$  generated after  $t$  time slots, when  $t \rightarrow \infty$ , the inter-degree sequences of  $n \in N_p$  in  $B(N_p, N_d)$  follows power-law distribution that

$$\mathbb{P}\{d_p^{ir}(t) = x\} \propto x^{-\frac{1}{\lambda_p}-1},$$

where  $x$  is one node's total degree and  $\mathbb{P}$  presents the probability. This result also holds for node  $n \in N_d$  sharing symmetrical expressions.

*Proof:* First of all, we consider the distribution of  $d_p^{ir}(t)$  which denotes the degree of node  $n \in N_p$  in  $B(N_p, N_a)$ . According to Equation (1), the cumulative distribution function of  $d_p^{ir}(t)$  can be calculated as

$$\begin{aligned} \mathbb{P}\{d_p^{ir}(t) < x\} &= \mathbb{P}\left\{d_p^{ir}(t_0) \left(\frac{t}{t_0}\right)^{\lambda_p} < x\right\} \\ &= \mathbb{P}\left\{t_0 > \left(\frac{d_p^{ir}(t_0)}{x}\right)^{\frac{1}{\lambda_p}} t\right\} \\ &= 1 - d_p^{ir}(t_0)^{\frac{1}{\lambda_p}} x^{-\frac{1}{\lambda_p}}. \end{aligned}$$

Then, the probability density function of  $d_p^{ir}(t)$  can be calculated using  $\mathbb{P}\{d_p^{ir}(t) = x\} = \frac{\partial \mathbb{P}\{d_p^{ir}(t) < x\}}{\partial x}$ . Also, it can be expressed as

$$\mathbb{P}\{d_p^{ir}(t) = x\} = \frac{x^{-\frac{1}{\lambda_p}-1}}{\sum_{x=1}^n x^{-\frac{1}{\lambda_p}-1}},$$

where  $\sum_{x=1}^n x^{-\frac{1}{\lambda_p}-1}$  is a constant normalization coefficient. Therefore, we get

$$\mathbb{P}\{d_p^{ir}(t) = x\} \propto x^{-\frac{1}{\lambda_p}-1},$$

By same approaches, we can also calculate the distribution of  $d_d^{ir}(t)$ , and thus the proof is complete. ■

*Theorem 4:* For graph  $G(P, D)$  generated after  $t$  time slots, when  $t \rightarrow \infty$ , the nodes' intra-degree sequences of  $n \in N_p$  follow power-law distribution that

$$\mathbb{P}\{d_p^{ia}(t) = x\} \propto x^{-\omega_p},$$

where  $x$  is one node's total degree,  $\mathbb{P}$  presents the probability and  $\omega_p$  is a constant. This result also holds for node  $n \in N_d$  as they share symmetrical expressions.

*Proof:* The proof uses the result of *Silvio Lattanzi and D. Sivakumar's* research work.

[citelattanzi2009affiliation](#). In their work, the model's bipartite network's structure is similar to our model's bipartite networks' which are disconstructed from  $G(P, D)$ .

And by Theorem 4 and Theorem 8 in their paper, they fully prove the total degree distribution is similar to the inter-degree distribution when time slot  $t \rightarrow \infty$ . Which means the total degree is also power-law distributed.

Therefore, the total degree distribution in our model follows

$$\mathbb{P}\{d_p(t) = x\} \propto x^{-\omega_p},$$

where  $\omega_p$  is a constant.

Using same methods, we can obtain the distribution for node  $n \in N_d$  and thus complete the proof. ■

## VI. CONCLUSION

In general, this work helps to judge or explain the relation and the boundary between different domains. For instance, we can explore the similarity and distinction of Literature and Mathematic, two domains largely different from each other in common sense. Moreover, when adding publication

date as time slot, we can explore the evolving pattern of the domains relationship. Besides, according to the performance of literatures cross-domain studying, the literatures depth and breadth can be well measured by properties of domains which they belong to and their cross-domain citation distribution, which affords researchers more accurate browsing results when they want to wade into a new scientific field.

## REFERENCES

- [1] (2016) Microsoft academic graph. [Online]. Available: <https://www.microsoft.com/en-us/research/project/microsoft-academic-graph/>
- [2] M. F. Collen and M. J. Ball, *A history of medical informatics in the United States*. Springer, 2015.
- [3] A. Burnett-Hartman, P. A. Newcomb, C. X. Zeng, Y. Zheng, J. M. Inadomi, C. Fong, M. P. Upton, and W. M. Grady, "Abstract pr05: Using medical informatics to evaluate the risk of colorectal cancer in patients with clinically diagnosed sessile serrated polyps," 2017.
- [4] D. Hristovski, A. Kastrin, and T. C. Rindflesch, "Implementing semantics-based cross-domain collaboration recommendation in biomedicine with a graph database," *DBKDA 2016*, p. 104, 2016.
- [5] J. Tang, S. Wu, J. Sun, and H. Su, "Cross-domain collaboration recommendation," in *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2012, pp. 1285–1293.
- [6] S. J. Derry, C. D. Schunn, and M. A. Gernsbacher, *Interdisciplinary collaboration: An emerging cognitive science*. Psychology Press, 2014.
- [7] B. Taebe, A. Correlje, E. Cuppen, M. Dignum, and U. Pesch, "Responsible innovation as an endorsement of public values: The need for interdisciplinary research," *Journal of Responsible Innovation*, vol. 1, no. 1, pp. 118–124, 2014.
- [8] B. K. Sovacool, "Energy studies need social science," *Nature*, vol. 511, no. 7511, p. 529, 2014.
- [9] H. Rabbani, "Interdisciplinary researches in iran v: Toward interdisciplinary technologies," *Journal of medical signals and sensors*, vol. 6, no. 3, p. 129, 2016.
- [10] N. Aggrawal and A. Arora, "Visualization, analysis and structural pattern infusion of dblp co-authorship network using gephi," in *Next Generation Computing Technologies (NGCT), 2016 2nd International Conference on*. IEEE, 2016, pp. 494–500.
- [11] J. Portenoy and J. D. West, "Dynamic visualization of citation networks showing the influence of scholarly fields over time," in *International Workshop on Semantic Analytics, Visualization*. Springer, 2016, pp. 147–151.
- [12] P. Isenberg, T. Isenberg, M. Sedlmair, J. Chen, and T. Möller, "Visualization as seen through its research paper keywords," *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 1, pp. 771–780, 2017.
- [13] M. E. Newman, "Power laws, pareto distributions and zipf's law," *Contemporary physics*, vol. 46, no. 5, pp. 323–351, 2005.
- [14] C. Xie, L. Yan, W.-J. Li, and Z. Zhang, "Distributed power-law graph computing: Theoretical and empirical analysis," in *Advances in Neural Information Processing Systems*, 2014, pp. 1673–1681.
- [15] L. Muchnik, S. Pei, L. C. Parra, S. D. Reis, J. S. Andrade Jr, S. Havlin, and H. A. Makse, "Origins of power-law degree distribution in the heterogeneity of human activity in social networks," *arXiv preprint arXiv:1304.4523*, 2013.
- [16] T. Honicke and J. Broadbent, "The influence of academic self-efficacy on academic performance: A systematic review," *Educational Research Review*, vol. 17, pp. 63–84, 2016.
- [17] X. Zhu, P. Turney, D. Lemire, and A. Vellino, "Measuring academic influence: Not all citations are equal," *Journal of the Association for Information Science and Technology*, vol. 66, no. 2, pp. 408–427, 2015.
- [18] D. Chakrabarti and C. Faloutsos, "Graph mining: Laws, generators, and algorithms," *ACM computing surveys (CSUR)*, vol. 38, no. 1, p. 2, 2006.
- [19] A.-L. Barabási and R. Albert, "Emergence of scaling in random networks," *science*, vol. 286, no. 5439, pp. 509–512, 1999.
- [20] R. Kumar, P. Raghavan, S. Rajagopalan, D. Sivakumar, A. Tomkins, and E. Upfal, "Stochastic models for the web graph," in *Foundations of Computer Science, 2000. Proceedings. 41st Annual Symposium on*. IEEE, 2000, pp. 57–65.
- [21] S. Lattanzi and D. Sivakumar, "Affiliation networks," in *Proceedings of the forty-first annual ACM symposium on Theory of computing*. ACM, 2009, pp. 427–434.