

Scientific Evolutionary Pathways: Identifying and Visualizing Relationships for Scientific Topics

Yi Zhang

Decision Systems & e-Service Intelligence Research Lab, Centre for Artificial Intelligence, Faculty of Engineering and Information Technology, University of Technology Sydney, Australia

School of Management and Economics, Beijing Institute of Technology, Beijing, P. R. China. E-mail: yi.zhang@uts.edu.au

Guangquan Zhang

Decision Systems & e-Service Intelligence Research Lab, Centre for Artificial Intelligence, Faculty of Engineering and Information Technology, University of Technology Sydney, Australia. E-mail: guangquan.zhang@uts.edu.au

Donghua Zhu

School of Management and Economics, Beijing Institute of Technology, Beijing, P. R. China. E-mail: zhudh111@bit.edu.cn

Jie Lu

Decision Systems & e-Service Intelligence Research Lab, Centre for Artificial Intelligence, Faculty of Engineering and Information Technology, University of Technology Sydney, Australia. E-mail: jie.lu@uts.edu.au

Whereas traditional science maps emphasize citation statistics and static relationships, this paper presents a term-based method to identify and visualize the evolutionary pathways of scientific topics in a series of time slices. First, we create a data preprocessing model for accurate term cleaning, consolidating, and clustering. Then we construct a simulated data streaming function and introduce a learning process to train a relationship identification function to adapt to changing environments in real time, where relationships of topic evolution, fusion, death, and novelty are identified. The main result of the method is a map of scientific evolutionary pathways. The visual routines provide a way to indicate the interactions among scientific subjects and a version in a series of time slices helps further illustrate such evolutionary pathways in detail. The detailed outline offers sufficient statistical information to delve into scientific topics and routines and then helps address meaningful insights with the assistance of expert

knowledge. This empirical study focuses on scientific proposals granted by the United States National Science Foundation, and demonstrates the feasibility and reliability. Our method could be widely applied to a range of science, technology, and innovation policy research, and offer insight into the evolutionary pathways of scientific activities.

Introduction

Science maps are spatial representations of relationships among disciplines, fields, specialties, and individual papers or authors (Small, 1999), and are promising for visualizing the extent and structure of large-scale data to help understand scientific activities, innovative pathways, and interactive relationships (Börner, 2014). Similarity measures are the main analytic approach for science maps, and include many bibliometric indicators (Rafols, Porter, & Leydesdorff, 2010), such as bibliographic coupling (Kessler, 1963), citation analysis (Garfield, Sher, & Torpie, 1964), co-citation analysis (Small, 1973), co-word analysis (Callon, Courtial, Turner, & Bauin, 1983), and co-author analysis (Glänzel,

Received March 10, 2016; revised July 6, 2016; accepted September 7, 2016

© 2017 ASIS&T • Published online 0 Month 2017 in Wiley Online Library (wileyonlinelibrary.com). DOI: 10.1002/asi.23814

2001). Since the 1970s, science maps have achieved great success in both theoretical research and practical applications. The rapid growth of information technology (IT), especially information visualization techniques, helps further support the increasing need for these kinds of relational studies.

Citation- and term-based techniques are considered the two main parallel technical categories for mapping science. The former was first used in Garfield's historiographic map (Garfield et al., 1964) and Small strengthened the approach by exploring the relationships between the publications that were co-cited frequently (Small, 1973). However, critical comments exist, for example, the relationship between a reference and the originality, importance, or even the quality of that work (Okubo, 1997), the difference between the source documents within technical and applied fields (Rip, 1988), and the bias of databases (De Bellis, 2009). As natural language processing (NLP) techniques have developed, terms have become a fast way to uncover meaningful concepts from the large volume of textual records (Porter & Detampel, 1995), and the time information of the terms derived from multilevel entities would allow an analysis of the dynamics of scientific systems (Tijssen & Van Raan, 1994). Although citation-based approaches demonstrate great power for creating science maps, it is widely believed that the Pandora's box of term-based approaches, in the age of big data, will be opened by modern IT techniques such as large-scale quantum computing, machine learning, and convolutional networks.

Despite the ability of capturing scientific interactions, it is still challenging to identify evolutionary relationships using science maps. Clearly, the evolution of science, technology, and innovation (ST&I) never stops, and accumulative changes may induce evolution or perhaps even disruptive revolution (Kostoff, Boylan, & Simons, 2004). Addressing these concerns, the main objective of this paper is to introduce a learning technique to identify the evolutionary relationships (e.g., topic evolution, fusion, death, and novelty) between scientific topics and visualize such relationships in a map of scientific evolutionary pathways (SEP), where terms and their co-occurrence distributions are involved. In terms of technical bibliometric issues, we attempt to handle the challenges via the following efforts: (a) applying a term-clumping process (Zhang, Porter, Hu, Guo, & Newman, 2014) to effectively clean noisy terms; (b) introducing a K-means clustering approach (Zhang et al., 2016) for high-accurate topic gathering; and, in particular, (c) designing a learning process to identify the evolutionary relationships and train related functions to adapt to changing environments in real time.

The main contributions of this study are (a) a term-based science map to trace scientific evolutionary pathways with both visual routines and detailed outlines and (b) a learning process to introduce machine-learning techniques to investigate topic analysis in a changing environment over time. Considering the varied backgrounds of our audience, we have simplified the technical details and demonstrate our

method via a case study of the United States (US) National Science Foundation (NSF) Awards. The resulting analysis illustrates the landscape of the evolution of the US's scientific activities and could be a beneficial tool in a wide range of ST&I policy research, for example, locating hotspots, tracing emerging trends, and collecting insights for specific R&D needs.

Science Maps and Related Bibliometric Techniques

Science maps are rooted in Small's co-citation studies in the 1970s (Small & Griffith, 1974), which could be generally described as graphic references for scientific fields or portfolios that are used to chart strategic trajectories, locate emerging research frontiers, and profile insights into specific portfolios (Börner et al., 2012; Rafols et al., 2010). This article reviews science maps via related bibliometric techniques, where citation and co-citation analysis and co-word analysis are highlighted.

Citation and Co-Citation Analysis

The origin of citation analysis is Garfield's historiographic map (Garfield et al., 1964). Its informetric analytic software *HistCite* (Garfield & Pudovkin, 2004) further enhanced such efforts and Lucio-Arias and Leydesdorff (2008) applied main-path analysis to retrieve the structural backbone of a selected scientific field. *CiteSpace* is another significant contribution (Chen, 2006), which divides a time span into several slices and detects certain emerging trends (Chen, Ibekwe-SanJuan, & Hou, 2010).

Co-citation analysis continues the traditions of the Institute of Scientific Information (ISI) and includes some notable work. Small (1999) first presented a science map for entire scientific subjects using a large-scale publication dataset. Boyack, Klavans, and Börner (2005) extended this approach by addressing insights into both the intercitation and co-citation links between scientific journals. Leydesdorff and Rafols (2009) introduced ISI subject categories and proposed a citation-based overlay map to analyze interdisciplinary relationships. Their effort was further associated with multidimensional indicators (e.g., international patent classification [IPC] code; Kay, Newman, Youtie, Porter, & Rafols, 2014; Leydesdorff, Kushnir, & Rafols, 2014) and geographical maps (Leydesdorff & Bornmann, 2012).

Co-Word Analysis

Co-word information was first used to measure relationships between clusters of topics that represent a series of scientific subdomains and areas (Noyons & van Raan, 1998b). Significant pioneering endeavors are mostly credited to Van Raan and his colleagues at Leiden University (Noyons & van Raan, 1998a; Noyons, 2001; Peters & van Raan, 1993). However, compared to citations, the relationships derived from semantic structures are not as direct and clear as citation links. IT development provides significant opportunities

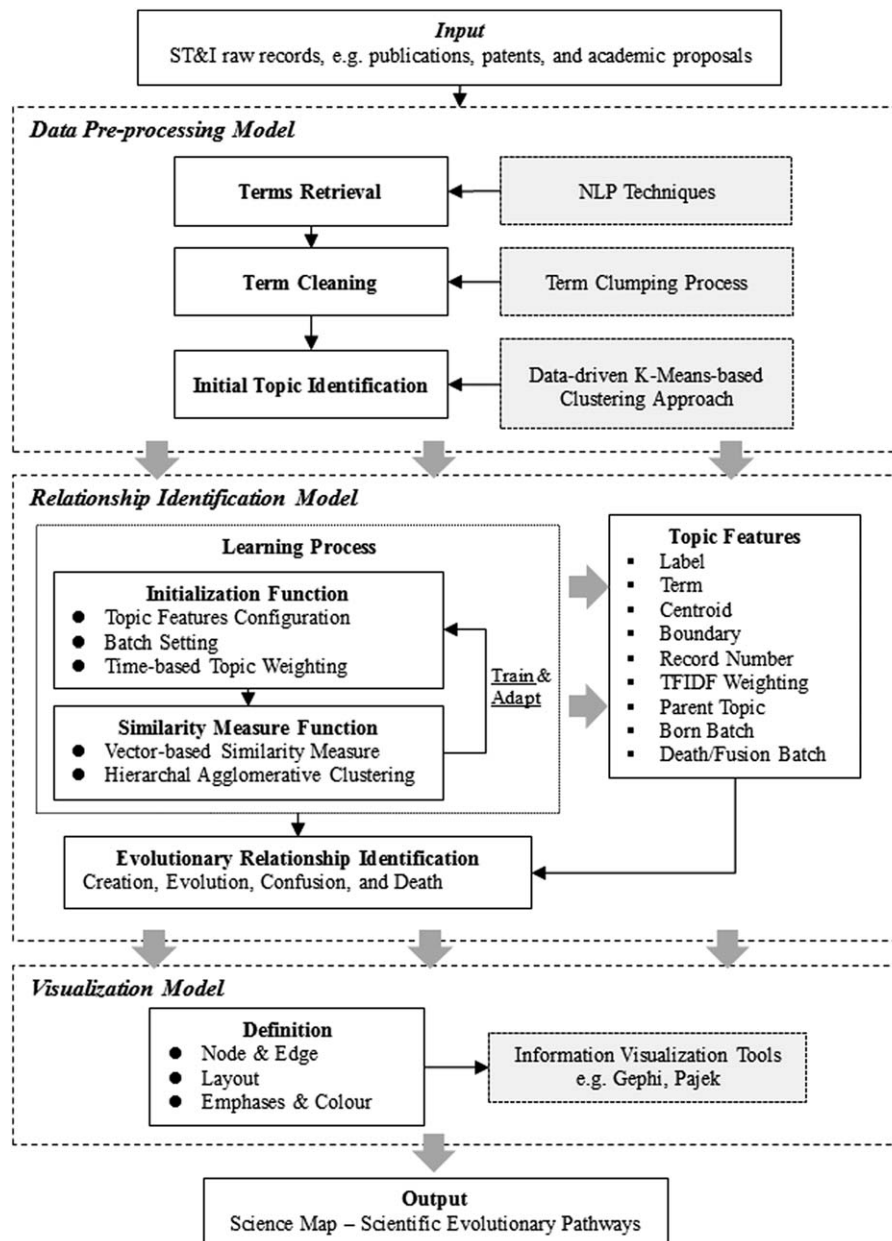


FIG. 1. The framework of the scientific evolutionary pathways.

to improve the computation of term identification and clustering (Noyons & van Raan, 1998b; van Eck, Waltman, Noyons, & Buter, 2010; Zhu & Porter, 2002). In recent years, the latent Dirichlet allocation approach (Blei, Ng, & Jordan, 2003) has played an active role in science maps (Suominen & Toivanen, 2015; Yau, Porter, Newman, & Suominen, 2014), as well as ontology and semantic webs (Wang, Lu, & Zhang, 2007). One important piece of software is *VantagePoint*, which provides systematic functions for term-based analysis. The attempts to combine co-word analysis with citation and co-citation analysis began around 2000 (Noyons, Moed, & Luwel, 1999), and were then associated with main-path analysis (Calero-Medina & Noyons, 2008). The software *VOSviewer* (Waltman, van Eck, & Noyons, 2010), which defines association links with co-

word, co-citation, and bibliographic couplings and visualizes grouped nodes as networks, is an important representation of such efforts.

The Validation of Science Maps

Boyack and Klavans contributed ground-breaking work in the validation of science maps (Boyack et al., 2005, 2011; Boyack & Klavans, 2010; Klavans & Boyack, 2006). They compared the accuracy of similarity measurements for topic identification, and extended these comparisons of science maps to represent research fronts via diverse indicator coupling. Their most recent study indicates that direct citation analysis is better for direct communication and detecting disciplines than co-citation analysis (Klavans & Boyack, 2016).

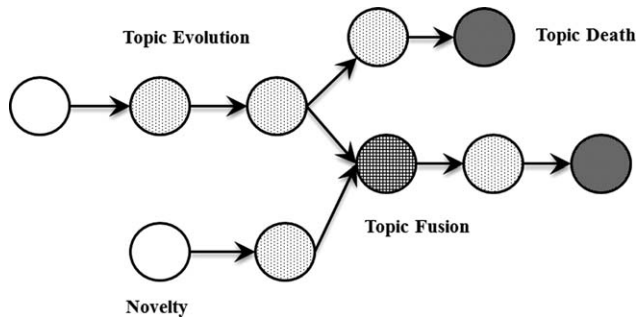


FIG. 2. The types of evolutionary relationships.

The Methodology of Scientific Evolutionary Pathways

This article proposes a method for composing a science map that identifies and visualizes the evolutionary relationships of scientific activities. Our method, called scientific evolutionary pathways (SEP), includes a data preprocessing model, a relationship identification model, and a visualization model, as shown in Figure 1.

Data Preprocessing Model

Our method is designed for ST&I textual data, for example, publications, patents, and academic proposals, but it is still necessary to note that diverse emphases and styles of term composition might exist in diverse ST&I textual data, and when applying our method to different databases, sufficient investigations are required. Generally, the data preprocessing model formats the raw data into our favored structures and draws upon several of our previous endeavors to achieve this goal.

An NLP technique is used to retrieve terms, after which a term-clumping process (Zhang, Porter, et al., 2014) removes meaningless noise and consolidates synonyms to reduce the number of terms from millions (or billions) to thousands. The main steps of the term-clumping process include the following: thesaurus-based term removal, knowledge-based term removal and consolidation (e.g., to consolidate certain technological synonyms), and association rule-based term consolidation (e.g., to consolidate terms that share certain individual words).

A K-means-based clustering approach (Zhang et al., 2016) follows, which includes a validation measure function and a feature selection and weighting function. We design methods to blend the general features (e.g., title and abstract terms) and special features (e.g., specific classification codes of a database), and highlight special features by weights. Then, based on a labeled sample dataset, we exhaustively run all possible combinations of a feature assembled set and decide a K value in a given interval with a preference for clustering accuracy; hence, the validation measurement is used to identify the best combination. This clustering approach groups records of the initial batch into several topics.

Relationship Identification Model

This model identifies the following evolutionary relationships, with a sample shown in Figure 2:

- **Novelty (white node):** a new topic is generated without any predecessors. Our model identifies topics that are generated suddenly and which share low or no similarities with existing topics as novelty. Thus, a novelty could be one which is related to something new that has never appeared before, or it also could be noise.
- **Topic evolution (node filled with spots):** a new topic is generated from existing topics. When the similarity between existing topics and a new record cannot be maintained above a threshold, but it still stays at a controllable interval, we set the topics grouped by these records as evolutions and the relationships between a predecessor and its evolved topics are identified as parent-child pairs.
- **Topic Fusion (node filled with small grid):** an existing topic fuses with another existing topic. However, identifying knowledge fusion is another difficult, but promising, research topic for current ST&I research, but we leave this task for further study.
- **Topic Death (gray node):** topic death occurs if no new records are added for several sequential batches.

Topics are the main focus, and we list the features of a topic in Figure 3. Our computation is mainly based on these features and detailed algorithms are presented in related functions.

The learning process runs through the whole model, and the algorithm is described as follows:

(1) Initialization function.

A. Batch setting. In the initialization function, a batch setting determines how to divide the whole dataset into a number of small batches. Despite a strategy of setting window size,¹ we generally divide data by time, e.g., all records published in the same year are gathered in one batch. Given that the whole dataset Θ is divided into n batches, each batch B contains a number of records, where $B(i, j)$ denotes the j -th record of the i -th batch. We denote that Num_t is the total number of records in the t -th batch and $1 \leq t \leq n$, and the input data stream of the learning process can be described as:

$$B(1, 1), \dots, B(1, Num_1), B(2, 1), \dots, B(2, Num_2), \dots, \\ B(t, 1), \dots, B(t, Num_t), \dots, B(n, 1), \dots, B(n, Num_n)$$

A record $B(t, x)$, $1 \leq x \leq Num_t$ can be a term frequency-based vector $\vec{T}_{B(t, x)} = \{v_{t_1}, v_{t_2}, \dots, v_{t_k}, \dots, v_{t_{m-1}}, v_{t_m}\}$, where t_k is the k -th term of the record $B(t, x)$ and v_{t_k} is its value, for example, raw term frequency or TFIDF value.

B. Topic initialization. Once a batch has been sequenced, the learning process iterates for each batch and accesses the records of the batch one-by-one. In each

¹The batch setting is a general window-size problem. While reading a data stream, we need to determine the window size to decide how many data items we should read at one time, and deciding an ideal window-size requires experiments.

No	Feature	Description
1	Label	The name of a topic and usually it is the highest frequency term/terms in the topic.
2	Centroid	The mathematical representation of a topic, and usually it is a term-frequency-based vector.
3	Boundary	Our method hypothesizes a topic is a circle, and the radius is the largest Euclidean distance between the centroid and the records.
4	TFIDF Weighting	A TFIDF-based algorithm is used to weight topics, and the TFIDF weighting, in some sense, is considered as the importance of a topic.
5	Parent Topic	The predecessor which evolved the current one.
6	Born Batch	The batch where the topic is generated.
7	Dead/Fusion Batch	The batch where the topic is fused or set as dead.

FIG. 3. The situations of evolution.

iteration (e.g., for the batch B_{t+1}), the initial topics $\Psi_t = \{\phi_1, \phi_2, \dots, \phi_p, \dots, \phi_{l-1}, \phi_l\}$, $1 \leq p \leq l$ are derived from the number of existing topics (i.e., not dead or fused) in previous batches, except for the first iteration, where l is the local optimum K defined in the K-means clustering approach. A series of topic features is then initialized. Let topic ϕ_p contain z records $B_i(t, x)$, $1 \leq i \leq z$, and, thus, the selected topic features are defined and calculated as:

- The label: a prevalence value $P(t)$ (Zhang, Zhou, Porter, & Gomila, 2014) is used to select terms to represent a topic. When a topic is born, we label it with the highest prevalence-value term (or the text of certain special features) at that current time.
- The centroid $\bar{C}(\phi_p)$: the mean of the term frequency-based vectors of all records in the topic, and is the mathematical representation of the topic,

$$\bar{C}(\phi_p) = \text{Avg}(B_i(t, x)) = \frac{1}{z} \times \sum_{i=1}^z \bar{T}_{B_i(t, x)}$$

- The boundary $R(\phi_p)$: the largest Euclidean distance between the centroid and the records,

$$R(\phi_p) = \text{Max}(Ed(\phi_p, B_i(t, x))) = \text{Max}(|\bar{C}(\phi_p) - \bar{T}_{B_i(t, x)}|)$$

- The TFIDF weighting $w_{\text{tf-idf}}(\phi_p)$: we apply a classical TFIDF formula (Salton & Buckley, 1988) to weight a topic,

$$w_{\text{tf-idf}}(\phi_p) = \frac{\text{Term amount in } \phi_p}{\text{Term amount in } \Psi_t} \times \log \frac{\text{Record amount in } \Psi_t}{\text{Record amount in } \phi_p}$$

The initialization function is called at the beginning of each batch, and the features (e.g., centroid and boundary) of all existing topics are recalculated to adapt to changing environments, since these newly assigned records can be accumulated to change the main content of a topic. Obviously, such efforts are within the scope of “training.” It is promising that such efforts directly improve the performance of the following analyses.

(2) Similarity measure function.

A. Similarity-based record assignment. For a new record $B(t+1, x)$, we first apply Salton’s cosine similarity measurement (Salton & McGill, 1986) to calculate the similarity between $B(t+1, x)$ and all the existing topics. In particular, a time-based topic weighting approach is used to take the term’s timeliness into consideration, that is, new terms are preferable to old ones. The similarity $S(\phi_p, B(t+1, x))$ between the record $B(t+1, x)$ and the topic ϕ_p is calculated as:

$$\begin{aligned} S(\phi_p, B(t+1, x)) &= w_{\delta}(\phi_p) \times \cos(\bar{C}(\phi_p), \bar{T}_{B(t+1, x)}) \\ &= \exp(-\delta \times T(\phi_p)) \times \frac{\bar{C}(\phi_p) \times \bar{T}_{B(t+1, x)}}{|\bar{C}(\phi_p)| |\bar{T}_{B(t+1, x)}|} \end{aligned}$$

where $T(\phi_p)$ is the time-gap (e.g., year) between the born batch of the topic ϕ_p and batch B_{t+1} , and δ is used as a sensitive parameter, that is, the larger δ is, the sooner old batches become unimportant.

The record $B(t+1, x)$ is then assigned to the topic that shares the highest similarity value, and we assume the topic is ϕ_p and denote the Euclidean distance between them as $Ed(\phi_p, B_i(t+1, x))$. We hypothesize a topic is a circle and the boundary is its radius, and then we set a threshold τ — evolutionary range — to draw a lower range and an upper range for the boundary. As shown in Figure 4, three situations exist:

- *Situation 1* - $Ed(\phi_p, B_i(t+1, x)) \in (0, R(\phi_p) \times (1-\tau)]$: the distance is less than the lower range of the boundary, so we set the record as a normal newcomer to the topic ϕ_p ;
- *Situation 2* - $Ed(\phi_p, B_i(t+1, x)) \in (R(\phi_p) \times (1-\tau), R(\phi_p) \times (1+\tau)]$: the distance is between the lower and upper range, so we set the record as evolved;
- *Situation 3* - $Ed(\phi_p, B_i(t+1, x)) \in (R(\phi_p) \times (1+\tau), +\infty)$: the distance is larger than the upper range, so we set the record as novel.

B. Hierarchical agglomerative clustering. There is no further operation for records in situation 1, but a hierarchical agglomerative clustering (HAC) approach is used to group

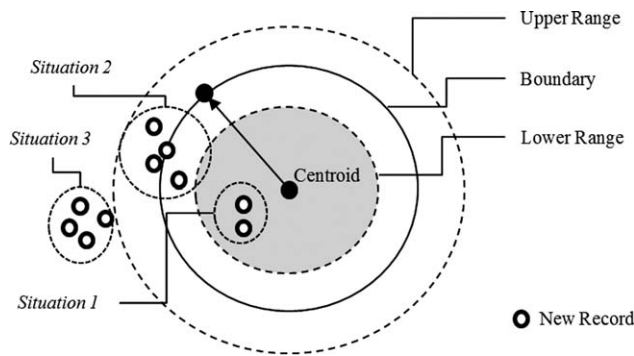


FIG. 4. The statistical information of the topics.

the records in Situations 2 and 3 separately. A basic HAC approach is applied, which we briefly describe as follows: (a) to set each record as a cluster; (b) to calculate the similarities among clusters and group the two clusters with the highest similarity value; and (c) to iterate step 2 until the terminal condition. Generally, the threshold has a positive correlation with the topic number, that is, the higher it is, the more evolved or novel the topic results, but we use a relatively objective way to eliminate such influence—the minimum similarity value among existing topics of the last is used, and once the highest similarity value between clusters is less than the threshold, HAC stops.

The concerns about choosing an HAC approach include the following: (a) the K-means approach needs labeled sample data to decide the best K value in a given interval, but this condition does not exist for the small datasets; and (b) it is not reasonable to set the topic number to a fixed value, for example, in some batches, very few records are assigned to Situations 2 or 3, but in some other batches, the record number is a little larger. Therefore, a dynamic topic number varying with actual conditions is more promising.

C. Evolutionary relationship identification. Generally, the topics grouped by the records in Situation 2 are identified as “evolved topics” and their predecessors are set as parent topics, whereas the topics in Situation 3 are identified as “novel topics” and have no parent topic. Specifically, referring to an interesting assumption that a “sleeping beauty” is an article that goes unnoticed (sleeps) for a long time and then, almost suddenly, attracts a lot of attention (van Raan, 2004, 2016), we propose a design of “death and resurgence” to capture this phenomenon. After reading all the records in a batch and finishing the tasks of evolution and novelty identification, the topics that have not been assigned any new records in certain sequential batches are set as “dead topics.” Then we measure the similarities² between the new topics generated in the current batch and all the existing/dead topics. Once the previous topic that shares the

highest similarity value with a new topic is not its parent, we consolidate the new topic with the previous one. In this circumstance, if the previous topic is a dead one, it is then resurged; if it is an existing topic, we set it as a confused topic. Note that in our current design, such consolidations miss all the information of a consolidated topic; thus, there is no topic with multiple parents.

The iteration stops when the learning process has read all the records in the data stream, and the output of this model is a list of topics (with their features) and their relationships.

Visualization Model

Using the basic format of a network, we identify the topics as nodes and their evolutionary relationships as arcs, and weigh a node via its related topic’s TFIDF weighting. Then, Gephi (Bastian, Heymann, & Jacomy, 2009) is used to visualize the results and the SEP is generated. The SEP is considered as an objective exhibition of quantitative data and analysis, and we then suggest involving experts for further understanding and implementation, which would increase the potential of SEP for addressing insights into specific problems and phenomenon.

Empirical Study: Scientific Evolutionary Pathways of the US NSF Awards

This paper proposes an SEP map to provide a graph of the visual routines of the US’s scientific activities and a detailed outline of the interactions of a series of scientific topics. The US NSF receives ~40,000 proposals per year, of which ~11,000 are granted. Considering the diverse purposes of different award types, for example, education, travel funding, or academic conference organization, this paper concentrates on the standard grant, the largest part of the NSF Awards. The NSF categorizes awards given prior to 1976 as historical data and reports on the possible features missing from these “old” awards. This article downloaded 388,909 awards from 1976 to 2014, and after some irrelevant awards (e.g., travel supports and sponsorship of academic activities) were removed, a total of 243,606 awards remained.

Data Preprocessing

The data preprocessing model was applied to the title and abstract. We first ran the NLP function of VantagePoint³ to retrieve 2,898,868 abstract terms and 312,639 title terms. Then, Zhang, Porter, et al.’s (2014) term-clumping process is used for term cleaning, which derived 178,262 distinct terms including 132,191 abstract terms and 62,089 title terms. Similar to IPCs, the US NSF Awards also have a systematic classification, for example, program element code

²Here we skip the aging function and directly apply the Cosine measure to the centroids of two topics, and focusing on topics with the same label, generally we add a weight to highlight our preference for consolidating such topics (but they are not always consolidated). Obviously, the weight will have a negative correlation with the final number of topics.

³Because the term-clumping process is to remove noise, the possible impacts resulting from diverse NLP functions can be ignored—a smart NLP function can lighten the stress, but a normal NLP function is also fine. The only requirement is that the output needs to be terms rather than individual words.

	#T	%Proportion	#R	%Coverage	Survival Length	Mean TFIDF	S.D. TFIDF
Total topics	562	100%	18,3086	100%	9.3274	0.0419	0.0492
Always alive	177	31.49%	64,059	34.99%	10.0226	0.0541	0.0492
Alive with resurgence	60	10.68%	36,862	20.13%	17.5167	0.0796	0.0580
Dead without resurgence	214	38.08%	63,432	34.65%	9.1822	0.0323	0.0460
Dead with resurgence	30	5.34%	16,343	8.93%	12.3667	0.055	0.0340
Dead when born (with child)	28	4.98%	1694	0.93%	1	0.0148	0.0381
Dead when born (no child)	53	9.43%	696	0.38%	1	0.0039	0.0070

Note: 1) #T – the number of topics; %Proportion – the proportion of #T; #R – the number of the records within the topics; %Coverage – the coverage of #R; Survival Length – the total batches in which a topic lived; S.D. TFIDF – the standard deviation of the TFIDF value of the topics. 2) The deadline by which we decided a topic was alive or dead was 2014, and if records after 2014 are engaged, the results will vary.

FIG. 5. The list of topic features.

(PEC) and program reference code (PRC). We are aware of the value of the PEC and PRC for topic clustering and labeling, which provide a clearer taxonomy on naming classifications. Thus, we decide to label identified topics in all our models by the text of the highest prevalence-value PEC or PRC within the topic.

In this case, we imported 3,418 distinct PECs and 1,803 distinct PRCs. We assembled the abstract terms, title terms, PECs, and PRCs as the features of a record, and a record-feature matrix was generated, involving 4,414,646 cells. Since 1976 is the starting point of this data, we chose the 3,670 awards granted in 1976 to identify the initial topics. Following the general process of the clustering approach, we accumulated a 500-award labeled sample dataset and set the parameters as “K = 9,” “TFIDF-weighted term,” and “abstract terms + inverse-ratio-weighted (title terms + PEC + PRC).” The initial nine topics are: statistics, molecular biophysics, computer application research, renewable resources, polymers, metal and metallic nanostructure, geophysics, evolutionary processes cluster, and algebra and number theory.

Relationship Identification

We set the batch by year, resulting in 39 batches. The number of records in each batch varies and is $\sim 10,000$. Using the algorithm detailed in the methodology section: (a) for the exponential aging function $\exp(-T(\phi_p))$, we choose a conservative option to let δ be 1—we prefer the new terms but do not immediately ignore the old ones; (b) the threshold τ to determine the upper/lower range of the boundary for evolution/novelty identification was set as 10%; (c) the initial threshold for the terminal condition of the HAC approach, as we defined, is the minimum similarity among the initial topics. We also designed a function to detect changes in the minimum similarity in a series of batches, and once the amplitude is larger than 10%, we set the new value as the threshold, but the bottom line of the threshold cannot be less than the half of the initial one; (d) once a topic maintains its record number in two sequential batches (including the born batch), we set the topic as

dead; and (e) a strategy of labeling topics is used where we use the text of the highest prevalence-value PEC or PRC to label a topic but we use the highest prevalence-value term if the PEC and PRC are not ranked in the top-3 list.

As the results show, we generated 553 topics after the initial nine ones, which included 30 novel and 523 evolved topics. Each topic, except the initial ones and the novel ones, had unique parent topics, and so were identified as evolutionary pathways. The batch when the topic was set as dead was also recorded. The statistical information of the 562 topics is given in Figure 5.

We divided the 562 topics into six categories: (a) *always alive*: these topics have been alive since being born and are considered the backbone of the US NSF’s scientific activities; (b) *alive with resurgence*: using the idea of “sleeping beauties” (van Raan, 2004, 2016), it is interesting to assume that these topics might contain great potential for innovation; (c) *dead with resurgence*: an extension of the category “alive with resurgence,” but these topics had been dead and had not been resurged by 2014; (d) *dead without resurgence*: these topics could have been important at a previous time but their importance has decreased and their significance has not currently been recognized; (e) *dead when born with child*: these topics could be meaningless but we keep them because they have children and could act as certain conjunctions; and (f) *dead when born with no child*: we treat these topics as noise and delete them before visualization.

Based on the six categories and the statistical information shown in Figure 5, we observed certain interesting findings: (a) the topics *dead when born with no child* have the smallest mean of TFIDF values, and the smallest standard deviation, which further indicates a relatively stable dynamic of these values. Thus, this could support our hypothesis that the topics belonging to this category could be noise; (b) the topics *dead with resurgence* and *alive with resurgence* up until 2014 had the highest means of TFIDF values and the longest survival length, and the topics *dead with resurgence* were ranked higher than the topics *always alive* and were even more stable. Such observations could be a good endorsement for the idea of “sleeping beauties”; (c) despite occupying nearly 70% of both topic number and record

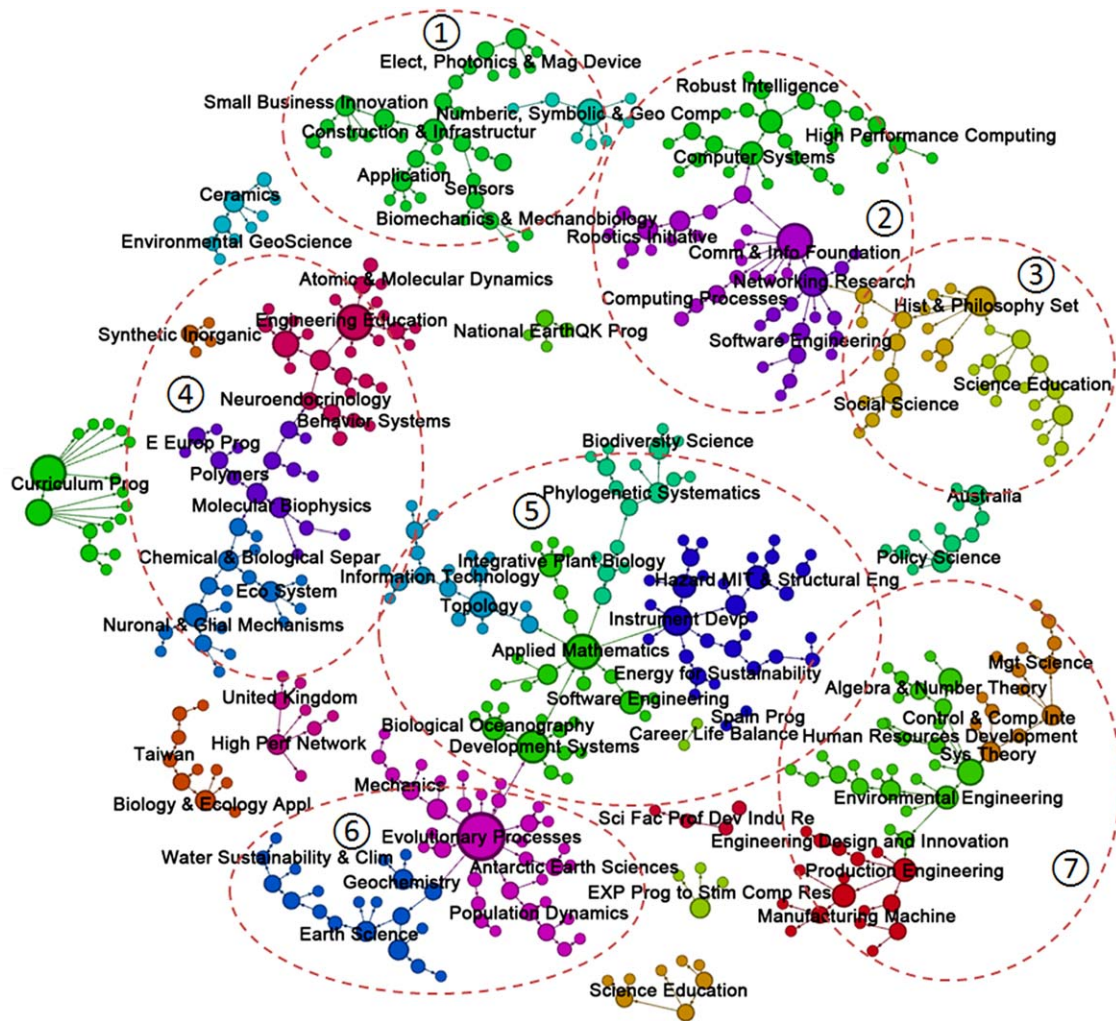


FIG. 6. The scientific evolutionary pathways for the US NSF award data. [Color figure can be viewed at wileyonlinelibrary.com]

number, the topics *always alive* might be not as important as we imagined and the same may apply to the topics *dead without resurgence*, which had an even lower TFIDF value than the mean of the total topics; and (d) generally a new topic survives 9 to 10 years, and if resurgence occurs, the survival length could extend to 12 to 17 years. When considering the life cycle of scientific activities, we might consider the length of the life cycle of a scholarly innovation granted by the US NSF to be around 9 years, that is, an innovation and its follow-up enhancement could last about 9 years, after which it might die until further significant innovation in this area appears again.

Visualization

Based on the topics (we removed the 53 topics in the category *dead when born with no child* and used the remaining 509 topics) and their evolutionary relationships, we generated a node list and an arc list as the inputs for Gephi (Bastian et al., 2009), where the TFIDF value of topics was used as the weights for the size of the nodes and the Gephi's modularity function was applied to determine the color. The SEP

is shown in Figure 6 and is one of the outputs of our method.

As the coloring strategy might influence readers, we discuss our selection as follows: the modularity function is based on a modularity optimization-based heuristic method (Blondel, Guillaume, Lambiotte, & Lefebvre, 2008), which makes good sense when objectively seeking and grouping similar communities via the similarity values among topics (identified as the weights of arcs in Gephi). Other options for assigning colors could be based on the results of the HAC approach. When considering the granularity of these coloring strategies, only the “brother” topics that are generated by the same parent topic are painted the same color and such granularity could be too trivial.

The Insights of the Scientific Evolutionary Pathways

Visual Routines

The graph of the visual routines is the main output of the SEP, which continues the traditions of science maps and

address concerns on scientific activities and their evolutionary relationships in a landscape. We discuss the insights gained from the SEP in two parts: the distribution of disciplines and the version in time slices.

The distribution of disciplines. As shown in Figure 6, we grouped the routines into seven clusters and briefly discuss each.

- Cluster 1: Infrastructure: Described as one of the main responsibilities of the US NSF, this topic covers the construction of infrastructure for high school, universities, and academic institutions. It ran through the whole period, and, as time went on, this infrastructure diversified into detailed fields, for example, electronic devices, sensors, and other applications.
- Cluster 2: Computer Science: It is interesting that one of the branches in the topic “history & philosophy” is “networking research,” and the computer science subject stems almost entirely from this. Several pathways were identified: software engineering, computing processes, computer and information foundation, and computer systems.
- Cluster 3: Social Science: Evolving from the general topic “history & philosophy,” this topic, also associated with science education, comprised many related topics, for example, science, technology and society.
- Cluster 4: Molecular Biophysics: it is difficult to completely distinguish between biology and chemistry, where strong interactions and interdisciplinary integrations tend to occur. Molecular biophysics was an initial topic in our settings, and its evolved topics included “biological and chemical separations,” “ecosystems,” “neuroendocrinology,” and “behavior systems,” where neuroscience is a significant direction for future biology studies.
- Cluster 5: Applied Mathematics: Mathematics is widely used in multiple disciplines, for example, software engineering and information technology, biology, and instrument development. This cluster includes a large range of scientific applications that deal with real-world needs through mathematical approaches.
- Cluster 6: Evolutionary Processes: This cluster also originated from an initial topic, which shared similarities with Cluster 4, but focused on Darwin’s evolution theory and on earth science, for example, geochemistry, population dynamics, and biological oceanography.
- Cluster 7: Industrial Engineering: A combination of management principles, system theory, mathematics, and engineering applications are grouped in this cluster, which are fundamental topics for modern industrial engineering research. Starting with the topic “thermal transport processes,” these topics are divided into two groups: theoretical research and applications.

As shown in Figure 6, education-related topics appear three times: “science education” in Cluster 3 and between Clusters 6 and 7, and “engineering education” located at the top of Cluster 4. This is definitely not due to confusion resulting from our topic-labeling approach. One important responsibility of the US NSF in progressing the US’s academic research and advanced technique development is to

fund education. Despite an independent Directorate of Education & Human Resource, a large number of funding programs for fundamental education are allocated by other specific-subject-oriented directorates. As an example, a division of the Engineering Education and Centers is under the Directorate of Engineering, and obviously, the biological, chemical, and neuroscientific topics in Cluster 4 are important branches of engineering. Therefore, an explanation for such a phenomenon is that education-related programs were widely proposed in different directorates and different fields, which were evolved by or then evolved to related specific scientific topics.

The SEP in time slices. Aiming to better demonstrate the advantages of SEP in identifying evolutionary pathways, we set four time slices, shown in Figure 7, to trace the evolution of scientific topics. The slice before 1980 is the beginning of the US NSF programs, and shows some early-stage topics, or, as we say, the predecessors of some routines. The slices for the 1990s and 2000s follow, and the newly generated topics indicate how the routines evolved. The fourth slice covering the period from 2001 to 2014 indicates the new ideas and innovation in recent years, and it is interesting to examine the in-depth implications.

- Slice 1 (1976 to 1980): This period does not show any impressive innovations, when the US NSF focused on fundamental research, infrastructure construction, and education.
- Slice 2 (1981 to 1990): Commercial innovation in infrastructure-related programs is one significant topic in this period, computer science started to appear, attempts at blending system theory with mathematics appeared, and earth science received increasing support.
- Slice 3 (1991 to 2000): Biology, especially neuroscience, was the most rapidly developing discipline in the last decade of the 20th century, since a large proportion of newly granted proposals related to this discipline. In the same period, engineering was also important and was widely applied to related subjects.
- Slice 4 (2001 to 2014): Information technology and computer science dominated research at the beginning decade of the 21st century. Another significant research was sustainability, as energy concerns had become an emerging need for the modern world, and engaging new techniques, materials, and products for environmental sustainability were on the US NSF’s agenda.

We present the SEP in time slices to further demonstrate the feasibility of our method for exploring the evolutionary pathways for entire scientific fields and specific subjects. Similar to Figure 6, the SEP provides an effective method to identify the relationships among scientific topics and visualize such topics and relationships in a landscape-type manner.

A Detailed Outline of Specific Topics

SEP’s major advantage is its ability to trace the evolutionary pathways of scientific topics; however, the graph of the visual routines of the SEP is only able to reveal changes in topics and these changes mostly relate to labels. We selected the topic “evolutionary processes” and its related

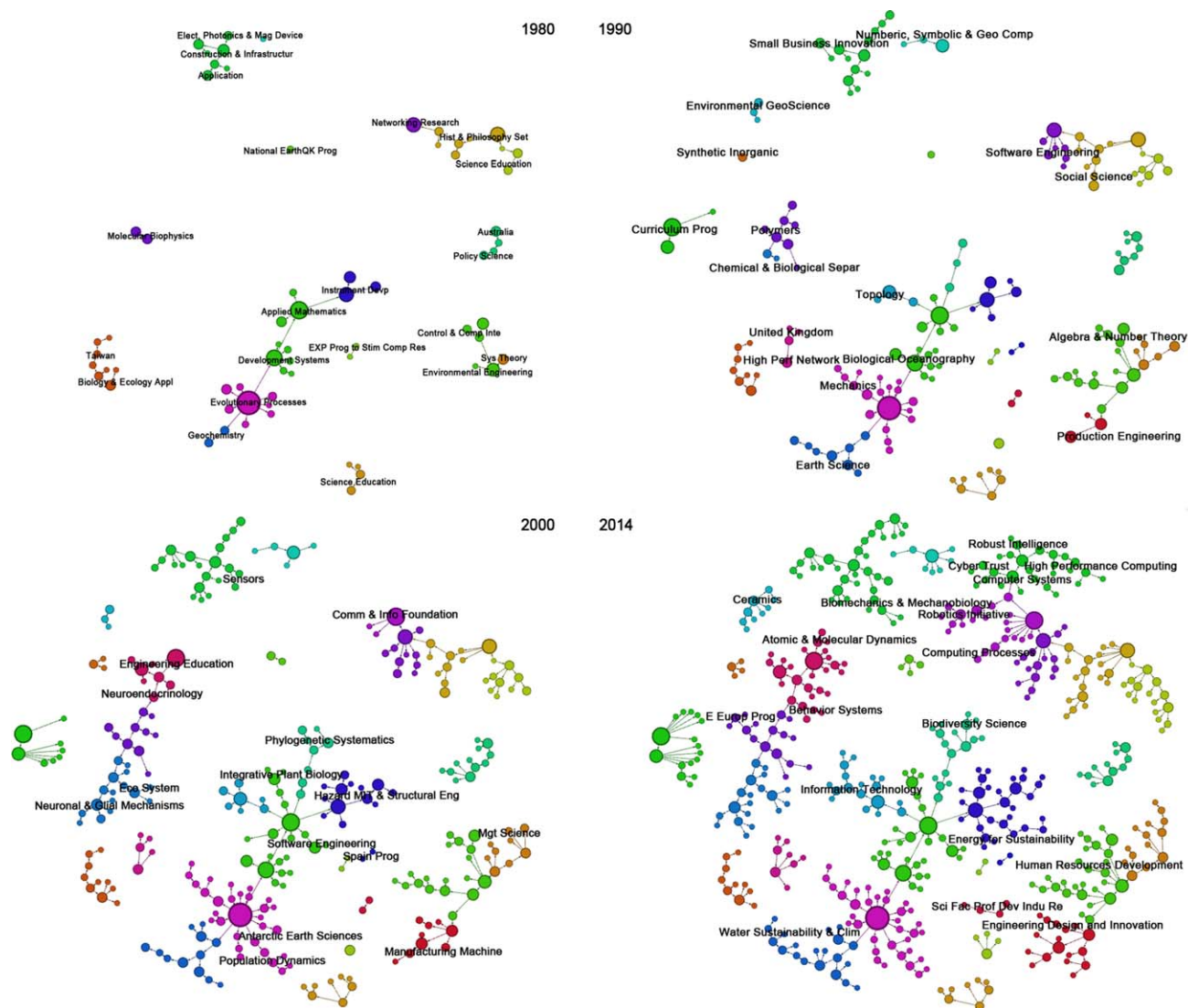


FIG. 7. The scientific evolutionary pathways for the US NSF award data in time slices. [Color figure can be viewed at wileyonlinelibrary.com]

generations as an example to demonstrate changes in feature space (e.g., terms) and data distribution (e.g., the number of top terms) to indicate such evolutions.

As one of the initial nine topics, “evolutionary processes” was set to be dead in 1996, 2001, and 2011, respectively, but resurged twice before 2014. This topic is a good representative example to illustrate an entire evolutionary pathway with the dynamics of both feature space and data distribution. The topic contained 5,922 records from 24 batches, from which we retrieved 10,475 title terms, 10,353 abstract terms, 6,567 PECs, and 805 PRCs. We combined the title and abstract terms to arrive at 3,643 distinct terms that constitute the topic’s feature space. Because only 498 terms appeared more than 10 times, we consolidated the related terms into a group where a simple association rule was applied, that is, terms sharing the same words or the same stem were grouped. As an example, terms such as “gene engineering” was first consolidated with the term “gene” since both terms shared the word “gene,” and then,

based on the stem “*gene*,” all related terms were consolidated with the term “genetics” (the decision to choose either “gene” or “genetics” as the representation requires human intervention). Finally, we selected seven groups to represent their principal features—ecological factors, molecular analyses, geography, population, genetics, evolution, and species.

We recorded the term frequency of the principal features in each batch, and created a 100% stacked area chart, shown in Figure 8. Given that the records in 1976 were grouped as the initial topics, we skipped this batch and started the chart from 1977. Although the seven principal features dominated the top of the feature space, it is obvious that their distributions kept changing over time. The label of the topic was “evolutionary processes,” and we set it to stable, but the main focus of this topic was different in diverse time intervals. As shown in Figure 8, the proportion of species reduced dramatically, while some emerging subjects, for example, genetics and molecular analyses, played increasingly important roles. Geography-related studies grew well

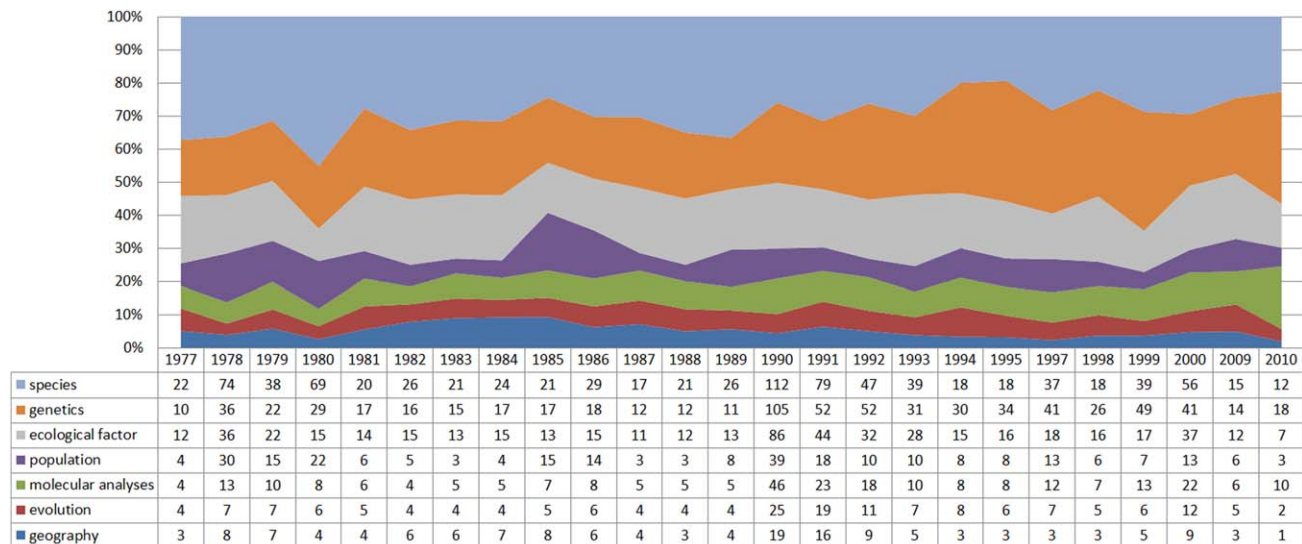


FIG. 8. The dynamics of the principal features' distribution in the topic "evolutionary processes." [Color figure can be viewed at wileyonlinelibrary.com]

at the beginning, but this tendency weakened rapidly after 1984, and a similar situation also occurred with population.

Aiming to better illustrate the change in the distributions, we magnified Figure 6 and focused only on our selected topics, as shown in Figure 9. Thirteen child topics were generated in diverse batches, most of which could be easily matched with the seven principal features, for example, population dynamics, geochemistry, and animal behavior. This phenomenon might explain how a new topic generates; that is, while new records are being assigned to a topic, the term composition of the topic also keeps changing (shown as Figure 8), and such accumulative change finally leads to the generation of a child topic.

Such discoveries act as a good way to: address more details and concerns for specified subjects and fields; provide sufficient statistical information to trace detailed

evolutionary pathways; and aid the understanding of the visual routines of the SEP.

Case Study on the Routines of Big Data Research From 2009 to 2014

Considering the breadth of our empirical study on the US NSF Awards from 1976 to 2014, we found it was not realistic to validate our results by manually reviewing the topics. Moreover, although Klavans and Boyack (2009) generated a consensus map to provide certain indicators for validating science maps, this approach was specifically designed for multidisciplinary studies and its main criteria can be summarized as the following: How many disciplines can be indicated from a science map (16 areas of science were presummarized), and can the real linkages between these disciplines be displayed correctly? However, these criteria do not exactly match the SEP, where we focus on the evolutionary relationships between scientific topics. Obviously, our definition of a scientific topic might cover certain disciplines and the evolutionary relationships between scientific topics could be both multidisciplinary and interdisciplinary interactions. Moreover, while comparing their 16 areas of science, the US NSF Awards exclude certain disciplines, for example, infectious disease, medical specialties, and health services. At this stage, instead of quantitatively measuring the validation of our results, we instead apply a case study on the routines of big data research from 2010 to 2014. The main objectives of this case study are (a) to demonstrate the reliability and robustness of the SEP, (b) to examine the role that the US NSF plays in implementing the Obama Administration's *Big Data Research and Development Initiative*, and (c) to specifically focus on the evolutionary pathways of big data-related topics over the recent 5 years, which can be a practical application of the SEP and might be meaningful

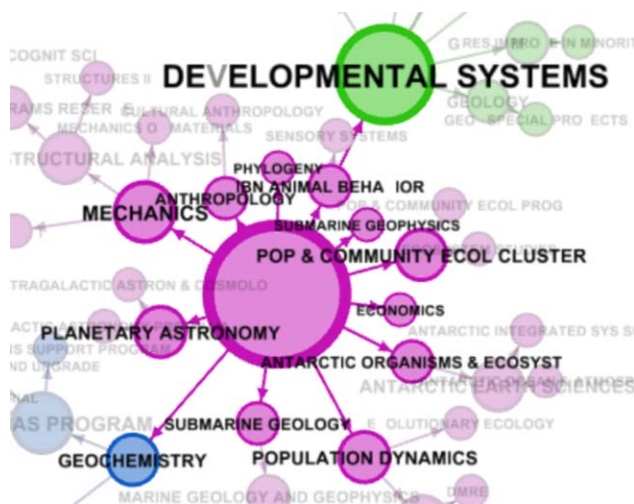


FIG. 9. The next generation of the topic "evolutionary processes." [Color figure can be viewed at wileyonlinelibrary.com]

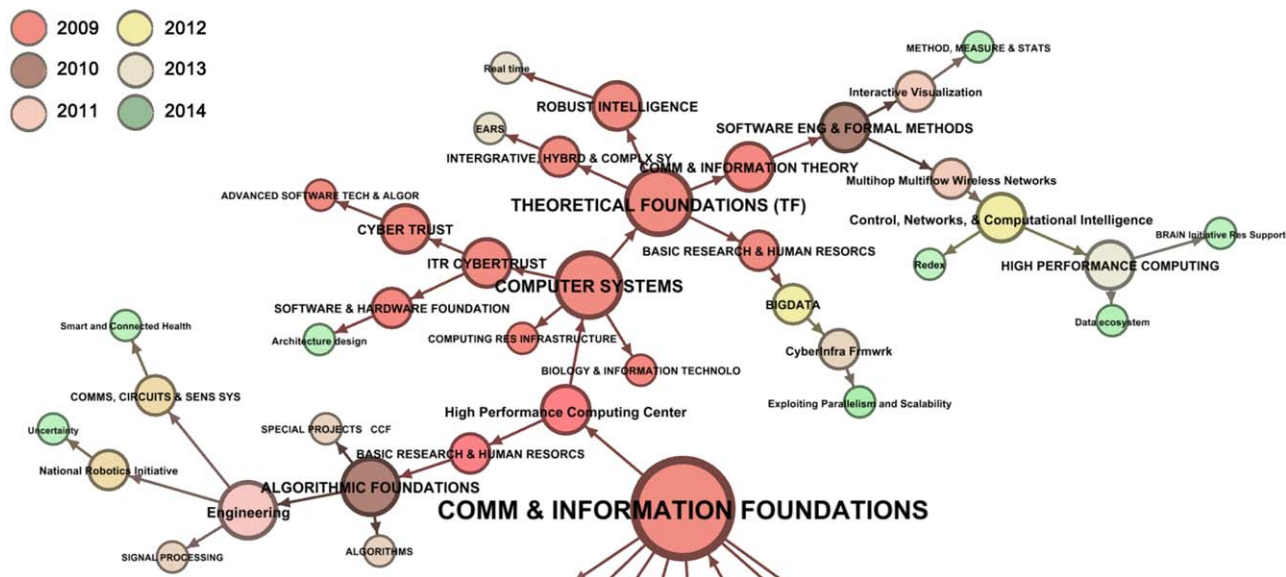


FIG. 10. The routines of big data-related topics (from 2010 to 2014). [Color figure can be viewed at wileyonlinelibrary.com]

for science policy makers, program officers, and academic researchers in related areas.

As shown in Figure 7, the routines evolved by the topic “Communication & Information Foundations (CIF)” started after 2000 but grew rapidly during the next 15 years. Professionally, big data and related analytic techniques can only be a subcategory of CIF, but considering the extensive and in-depth influence of big data, we have named the entire routine big data research. These routines are shown in Figure 10, where the red nodes indicate all topics generated before 2010, and the topics identified from 2010 to 2014 are shown in diverse colors.

We first reviewed the CIF studies before 2010, which mainly related to topics on computer systems and theoretical foundations, and the most significant pilot studies during this period were the topics “cyber trust” and “robust intelligence.” On one hand, launched by the US NSF in late 2003, the cyber trust program involved certain divisions of the Directorate for Computer & Information Science & Engineering (Landwehr, 2003). The archived programs in the US NSF indicated that cyber trust programs are open for applicants each year and they are available until 2016. Despite no significant new topics being generated after 2010, cyber trust programs integrated well with the scope of big data, and the state-of-the-art terms include: trustworthy computing, secure and trustworthy cyberspace, cyber security, digital privacy, and so on.⁴ On the other hand, the robust intelligence program⁵ belongs to the division of Information & Intelligent Systems (IIS), and includes research areas such as artificial intelligence, robotics, human language research, and computational neuroscience. One new topic “real time”

was generated from robust intelligence in 2013, and is one of the main attributes of big data analytic techniques. In addition, a nearby topic, “Enhancing Access to the Radio Spectrum (EARS),” was a collaborative program involving three US NSF directorates,⁶ and such collaboration was also motivated by the emerging requirements of big data analysis: exploring insights from unstructured data (Manyika et al., 2011).

Interesting topics generated in 2010 or after include (a) high performance computing—despite appearing in the early 2000s as “high-performance computing center,” evolving from topics related to wireless networks, control, and computational intelligence, this topic was consistent with the US NSF’s plan entitled “*Advanced Computing Infrastructure: Vision and Strategic Plan*” in February 2012,⁷ and specifically focused on expanding the use of high-end resources to larger and more diverse communities; (b) exploiting parallelism and scalability—parallel computing is a big data-related technique that concerns not only new hardware architectures but also new concepts, models, and algorithms; (c) robotics—a “national robotics initiative” program was operated by the IIS but it collaborated with not only other directorates of the US NSF but also with a wide range of US agencies, for example, National Institutes of Health (NIH) and NASA⁸; and (d) smart and connected health—we specifically picked this topic since it involved direct collaboration⁹ between the US NSF and NIH, which was advanced by the Obama Administration’s Big Data Program in 2012.

⁶Resources can be seen at the website https://www.nsf.gov/funding/pgm_summ.jsp?pims_id=503480.

⁷Resources can be seen at the website <http://www.nsf.gov/pubs/2012/nsf12051/nsf12051.pdf>.

⁸Resources can be seen at the website <http://www.nsf.gov/pubs/2015/nsf15505/nsf15505.htm>.

⁹Resources at the website https://www.nsf.gov/funding/pgm_summ.jsp?pims_id=504739.

⁴More details can be found at the website https://www.whitehouse.gov/sites/default/files/docs/big_data_privacy_report_may_1_2014.pdf.

⁵A general history of the NSF’s robust intelligence program was recalled by (54).

As discussed earlier, the SEP can effectively identify emerging scientific topics and the case study provides meaningful evidence to endorse the reliability of our results. Meanwhile, the results address in-depth insights into the latest programs of the US NSF in big data research and offer efficient guidance for stakeholders in these areas.

Discussion and Conclusions

This paper proposed an SEP method to identify and visualize the relationships among scientific topics. We focused on the potential of terms derived from ST&I textual data via NLP techniques, and a learning process was introduced to trace the evolutionary relationships. We demonstrated our method using a case study of all proposals granted by the US NSF over the time period between 1976 and 2014. The results include a graph of visual routines to visualize the scientific evolutionary pathways of several grouped routines and a detailed outline of scientific topics with statistical information and insights.

Potential for Term-Based Science Maps

Terms, as shown in the SEP, first provide an easy and meaningful way to label topics. A set of terms can constitute comprehensive semantic meanings and help to better understand the related topics, and, compared to citation linkages, the relationships between terms are easy to recognize manually. In particular, as one of the basic hypotheses of our method, a term is defined as the feature of SEP, and the dynamics of terms and their frequency result in the evolution of related scientific topics. In addition, as semantic elements, it also makes good sense to blend terms with other ST&I entities on science maps to link active agents with objects, for example, who was the key player in this scientific arena and what is the relationship between the players, competitive or complementary? The development of NLP techniques greatly assists the retrieval of accurate terms, and clustering approaches in multiple dimensions effectively reduce negative influences from single or limited scopes.

The Benefits of Learning Process-Based Bibliometrics

Our endeavor to introduce machine-learning techniques to deal with bibliometric problems is another exciting effort. Traditional bibliometric approaches analyze data in a stable environment, and the time and topic label is the only factor to identify the relationships between time interval and its forward or afterward topics. This design has been widely adopted in many bibliometric studies. However, it is not always correct to link topics with the same label yet diverse time intervals.

Although our empirical data was offline, we simulated a data stream. Once we identified the initial topics in batch 0, our method followed the sequence of the batch and read its records one-by-one. Obviously, once new records were classified into an existing topic, the features of the topic changed. Our learning process focused on changing

environments and adjusted topics in real time. Thus, the evolutionary relationship identified by our method not only pays attention to time but also takes these dynamic interactions into prior consideration.

Implementation

It is our belief that implementation of SEP is a beneficial development. The synergy between the visual routines and the detailed outline works well for tracing evolutionary pathways for both entire scientific subjects and selected ones. The SEP could be an effective tool by which to investigate a wide range of ST&I policy research, for example, multidisciplinary interactions, scientific outputs evaluation for selected entities, and competitive technical intelligence studies. We list two possible implementations as follows:

- Multidisciplinary interactions: SEP has a great capability for exploring multidisciplinary interactions, as most science maps do. The empirical study and Figure 6 demonstrated the feasibility of our SEP for such needs. It is easy to explore the interactions between diverse subjects and, in particular, detect interdisciplinary activities, e.g., topic evolution, confusion, and death. There are many substitutes to replace the US NSF Awards, and these data options provide different scopes and insights for multidisciplinary needs.
- Competitive and technical intelligence: SEP is helpful for locating the core scientific nodes in evolutionary pathways for selected paths, which would be significant materials, techniques, or algorithms. If we narrow the focus to specific entities, for example, individuals, organizations, countries, and regions, the comparison study between different SEPs generated from diverse entities' records, it would be promising to analyze potential collaborators and competitors. With the assistance of strategic analysis, such a combination would be of great interest to stakeholders.

The Limitations of the SEP

There are several limitations to our current methods: (a) the level of term cleaning heavily influences the results of the similarity measure. This is a general problem with term-based analysis. Although the term-clumping process effectively helps reduce the term amount, there is still a gap between the ideal situation and the current one; (b) we apply a clustering algorithm to identify initial topics, but the parameter configuration of the clustering algorithm highly depends on the labeled sample dataset; and (c) because most ST&I data is unlabeled and the interaction between semantic structures is complex, we could not find an efficient approach to validate the accuracy of our method or compare it with other approaches, except by expert knowledge.

Acknowledgments

This work was partially supported by the Australian Research Council under Discovery Grant DP150101645 and the National High Technology Research and Development Program of China under Grant 2014AA015105. We

thank Dr. Alan Porter at the Georgia Institute of Technology for valuable suggestions and insights, and to the three anonymous referees for comments, criticisms, and suggestions.

References

- Bastian, M., Heymann, S., & Jacomy, M. (2009). Gephi: An open source software for exploring and manipulating networks. *ICWSM*, 8, 361–362.
- Blei, D.M., Ng, A.Y., & Jordan, M.I. (2003). Latent dirichlet allocation. *The Journal of Machine Learning Research*, 3, 993–1022.
- Blondel, V.D., Guillaume, J.L., Lambiotte, R., & Lefebvre, E. (2008). Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008, P10008.
- Börner, K. (2014). *Atlas of knowledge*. Cambridge, MA: MIT Press.
- Börner, K., Klavans, R., Patek, M., Zoss, A.M., Biberstine, J.R., Light, R.P., . . . Boyack, K.W. (2012). Design and update of a classification system: The UCSD map of science. *PLoS One*, 7, e39464.
- Boyack, K.W., & Klavans, R. (2010). Co-citation analysis, bibliographic coupling, and direct citation: Which citation approach represents the research front most accurately?. *Journal of the American Society for Information Science and Technology*, 61, 2389–2404.
- Boyack, K.W., Klavans, R., & Börner, K. (2005). Mapping the backbone of science. *Scientometrics*, 64, 351–374.
- Boyack, K.W., Newman, D., Duhon, R.J., Klavans, R., Patek, M., Biberstine, J.R., . . . Börner, K. (2011). Clustering more than two million biomedical publications: Comparing the accuracies of nine text-based similarity approaches. *PLoS One*, 6, e18029.
- Calero-Medina, C., & Noyons, E.C. (2008). Combining mapping and citation network analysis for a better understanding of the scientific development: The case of the absorptive capacity field. *Journal of Informetrics*, 2, 272–279.
- Callon, M., Courtial, J.P., Turner, W.A., & Bauin, S. (1983). From translations to problematic networks: An introduction to co-word analysis. *Social Science Information*, 22, 191–235.
- Chen, C. (2006). CiteSpace II: Detecting and visualizing emerging trends and transient patterns in scientific literature. *Journal of the American Society for Information Science and Technology*, 57, 359–377.
- Chen, C., Ibekwe-SanJuan, F., & Hou, J. (2010). The structure and dynamics of cocitation clusters: A multiple-perspective cocitation analysis. *Journal of the American Society for Information Science and Technology*, 61, 1386–1409.
- De Bellis, N. (2009). *Bibliometrics and citation analysis: From the science citation index to cybermetrics*. Lanham, MD: Scarecrow.
- Fu, M.C., & Barton, R.R. (2012, June). O.R. & the National Science Foundation, Part 2. *ORMS-Today*.
- Garfield, E., & Pudovkin, A.I. (2004). The HistCite system for mapping and bibliometric analysis of the output of searches using the ISI Web of Knowledge. Paper presented at the Proceedings of the 67th Annual Meeting of the American Society for Information Science and Technology.
- Garfield, E., Sher, I.H., & Torpie, R.J. (1964). *The use of citation data in writing the history of science*. Philadelphia: Institute for Scientific Information Inc.
- Glänzel, W. (2001). National characteristics in international scientific co-authorship relations. *Scientometrics*, 51, 69–115.
- Kay, L., Newman, N., Youtie, J., Porter, A.L., & Rafols, I. (2014). Patent overlay mapping: Visualizing technological distance. *Journal of the Association for Information Science and Technology*, 65, 2432–2443.
- Kessler, M.M. (1963). Bibliographic coupling between scientific papers. *American Documentation*, 14, 10–25.
- Klavans, R., & Boyack, K.W. (2006). Identifying a better measure of relatedness for mapping science. *Journal of the American Society for Information Science and Technology*, 57, 251–263.
- Klavans, R., & Boyack, K.W. (2009). Toward a consensus map of science. *Journal of the American Society for Information Science and Technology*, 60, 455–476.
- Klavans, R., & Boyack, K.W. (2016). Which type of citation analysis generates the most accurate taxonomy of scientific and technical knowledge? Manuscript submitted for publication.
- Kostoff, R.N., Boylan, R., & Simons, G.R. (2004). Disruptive technology roadmaps. *Technological Forecasting and Social Change*, 71, 141–159.
- Landwehr, C.E. (2003). NSF activities in cyber trust. Paper presented at the Second IEEE International Symposium on Network Computing and Applications.
- Leydesdorff, L., & Bornmann, L. (2012). Mapping (USPTO) patent data using overlays to Google Maps. *Journal of the American Society for Information Science and Technology*, 63, 1442–1458.
- Leydesdorff, L., Kushnir, D., & Rafols, I. (2014). Interactive overlay maps for US patent (USPTO) data based on International Patent Classification (IPC). *Scientometrics*, 98, 1583–1599.
- Leydesdorff, L., & Rafols, I. (2009). A global map of science based on the ISI subject categories. *Journal of the American Society for Information Science and Technology*, 60, 348–362.
- Lucio-Arias, D., & Leydesdorff, L. (2008). Main-path analysis and path-dependent transitions in HistCiteTM-based historiograms. *Journal of the American Society for Information Science and Technology*, 59, 1948–1962.
- Manyika, J., Chui, M., Brown, B., Bughin, J., Dobbs, R., Roxburgh, C., & Byers, A.H. (2011). *Big data: The next frontier for innovation, competition, and productivity*. McKinsey Global Institute.
- Noyons, E. (2001). Bibliometric mapping of science in a policy context. *Scientometrics*, 50, 83–98.
- Noyons, E.C., Moed, H.F., & Luwel, M. (1999). Combining mapping and citation analysis for evaluative bibliometric purposes: A bibliometric study. *Journal of the Association for Information Science and Technology*, 50, 115.
- Noyons, E., & van Raan, A. (1998a). Advanced mapping of science and technology. *Scientometrics*, 41, 61–67.
- Noyons, E.C., & van Raan, A.F. (1998b). Monitoring scientific developments from a dynamic perspective: Self-organized structuring to map neural network research. *Journal of the American Society for Information Science*, 49, 68–81.
- Okubo, Y. (1997). Bibliometric indicators and analysis of research systems. *OECD Science, Technology and Industry Working papers*. doi: 10.1787/18151965
- Peters, H., & van Raan, A.F. (1993). Co-word-based science maps of chemical engineering. Part I: Representations by direct multidimensional scaling. *Research Policy*, 22, 23–45.
- Porter, A.L., & Detampel, M.J. (1995). Technology opportunities analysis. *Technological Forecasting and Social Change*, 49, 237–255.
- Rafols, I., Porter, A.L., & Leydesdorff, L. (2010). Science overlay maps: A new tool for research policy and library management. *Journal of the American Society for Information Science and Technology*, 61, 1871–1887.
- Rip, A. (1988). Mapping of science: Possibilities and limitations. In A.F.J. van Raan (Ed.), *Handbook of quantitative studies of science and technology* (pp. 253–273). North-Holland, Netherlands: Elsevier Science.
- Salton, G., & Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24, 513–523.
- Salton, G., & McGill, M.J. (1986). *Introduction to modern information retrieval*. New York: McGraw-Hill.
- Small, H. (1973). Co-citation in the scientific literature: A new measure of the relationship between two documents. *Journal of the American Society for Information Science*, 24, 265–269.
- Small, H. (1999). Visualizing science by citation mapping. *Journal of the American Society for Information Science*, 50, 799–813.
- Small, H., & Griffith, B.C. (1974). The structure of scientific literatures I: Identifying and graphing specialties. *Science Studies*, 4, 17–40.
- Suominen, A., & Toivanen, H. (2015). Map of science with topic modeling: Comparison of unsupervised learning and human-assigned subject classification. *Journal of the Association for Information Science and Technology*, 67, 2464–2476.

- Tijssen, R.J., & Van Raan, A.F. (1994). Mapping changes in science and technology bibliometric co-occurrence analysis of the R&D literature. *Evaluation Review*, 18, 98–115.
- van Eck, N., Waltman, L., Noyons, E., & Buter, R. (2010). Automatic term identification for bibliometric mapping. *Scientometrics*, 82, 581–596.
- van Raan, A.F. (2004). Sleeping beauties in science. *Scientometrics*, 59, 467–472.
- van Raan, A.F. (2016). Sleeping beauties cited in patents: Is there also a dormitory of inventions? Manuscript submitted for publication.
- Waltman, L., van Eck, N.J., & Noyons, E.C. (2010). A unified approach to mapping and clustering of bibliometric networks. *Journal of Informetrics*, 4, 629–635.
- Wang, C., Lu, J., & Zhang, G. (2007). Mining key information of web pages: A method and its application. *Expert Systems with Applications*, 33, 425–433.
- Yau, C.K., Porter, A., Newman, N., & Suominen, A. (2014). Clustering scientific documents with topic modeling. *Scientometrics*, 100, 767–786.
- Zhang, Y., Porter, A.L., Hu, Z., Guo, Y., & Newman, N.C. (2014). “Term clumping” for technical intelligence: A case study on dye-sensitized solar cells. *Technological Forecasting and Social Change*, 85, 26–39.
- Zhang, Y., Zhang, G., Chen, H., Porter, A.L., Zhu, D., & Lu, J. (2016). Topic analysis and forecasting for science, technology and innovation: Methodology and a case study focusing on big data research. *Technological Forecasting and Social Change*, 105, 179–191.
- Zhang, Y., Zhou, X., Porter, A.L., & Gomila, J.M.V. (2014). How to combine term clumping and technology roadmapping for newly emerging science & technology competitive intelligence: “Problem & solution” pattern based semantic TRIZ tool and case study. *Scientometrics*, 101, 1375–1389.
- Zhu, D., & Porter, A.L. (2002). Automated extraction and visualization of information for technological intelligence and forecasting. *Technological Forecasting and Social Change*, 69, 495–506.