

Research Topics Map: rtopmap

Md Iqbal Hossain and Stephen Kobourov

Department of Computer Science
University of Arizona

Abstract. In this paper we describe a system for visualizing and analyzing worldwide research topics, `rtopmap`. We gather data from google scholar academic research profiles, putting together a weighted topics graph, consisting of over 35,000 nodes and 646,000 edges. The nodes correspond to self-reported research topics, and edges correspond to co-occurring topics in google scholar profiles. The `rtopmap` system supports zooming/panning/searching and other google-maps-based interactive features. With the help of map overlays, we also visualize the strengths and weaknesses of different academic institutions in terms of human resources (e.g., number of researchers in different areas), as well as scholarly output (e.g., citation counts in different areas). Finally, we also visualize what parts of the map are associated with different academic departments, or with specific documents (such as research papers, or calls for proposals). The system itself is available at <http://rtopmap.arl.arizona.edu/>.

1 Introduction

Cataloguing and organizing science often involves taxonomies, ontologies, and knowledge graphs, but most often research topics are categorized in hierarchical trees [1, 23]; see Fig. 1. For example, “Hardware” and “computer systems organization” are subfields of “computer science.” Knowledge graphs make it possible to see more of the connections between topics than can be embedded in a tree, however, the ability to show clearly the underlying hierarchical structures is compromised.

Maps have guided human exploration for many centuries and recently there have been several efforts to visualize scholarly knowledge and research expertise using topic maps. Basemaps of science can be generated, for example, by analyzing citation links between publications and placing similar records next to each other. Such maps can be used to compare expertise profiles, to understand career trajectories, and to communicate emerging areas, as illustrated in the special *PNAS* issue on “Mapping Knowledge Domains” [49], and Börner’s “Atlas of Science” [13] and “Atlas of Knowledge” [14].

Most maps of science are really node-link diagrams with one level of detail; a few support two or more levels (e.g., the UCSD map of science which is the current standard has two levels of detail) [15]. However, people have difficulty reading large-scale networks [58] and few can derive knowledge from multi-level representations of networks. Given encouraging results about the effectiveness of map-like visualization of large graphs [45, 46], we adopt the Graph-to-Map

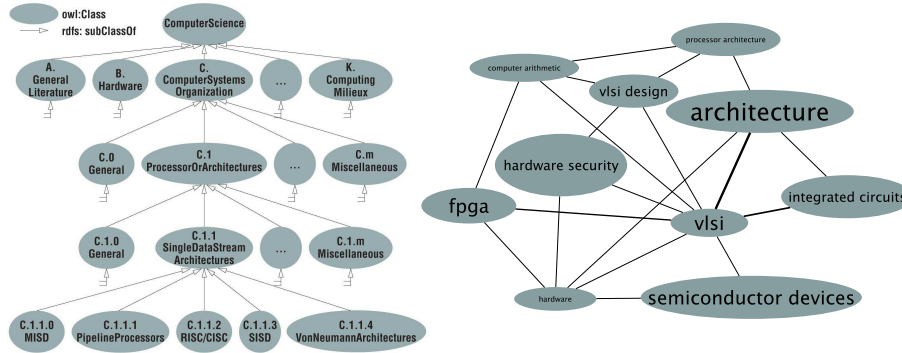


Fig. 1: Part of the ACM classification (left) and a more realistic network showing different types of connections between topics (right).

(GMap) framework [27] to visualize and explore our research topics map. The GMap visualization of relational data was introduced in the context of visualizing recommendations, where the underlying data is TV shows and the similarity between them [28] and has already been used to visualize research topics in computer science publications [26].

Our research topics map system, `rtopmap`, covers all research topics indexed by google scholar and provides the ability to show human resource investments (e.g., number of researchers in different areas) and scholarly output (e.g., citation counts in different areas) of different universities. The system supports zooming/panning/searching and other google-maps-based interactive features, including several map overlays showing what parts of the map are associated with different academic departments, or with specific documents; see Fig. 2 for an overview of the system.

We gather data from google scholar and then clean, split, merge, and correct the research topics (which become the nodes in the graph). We next compute a similarity matrix based on co-occurrence of topics in scholar profiles, which is used to place edges between topics that are frequently listed together. This gives us the topic network. We reduce the size of the graph by removing rarely occurring topics and weak connections. We then use a multi-level force-directed placement, node overlap removal, and clustering algorithms to represent the graph as a map. Nodes, node labels, polygon colors, and edges are transformed into google map objects, which are then drawn in the browser using the google maps API. Eight different level-of-detail (zoom levels) are precomputed, determining which nodes are present on a given level, computing label font sizes, and ensuring no label overlaps. Different overlays are added on demand.

2 Related work

Today, the most comprehensive map of science and classification uses ten years of paper-level data from Thomson Reuters' Web of Science and Elsevier's Scopus to

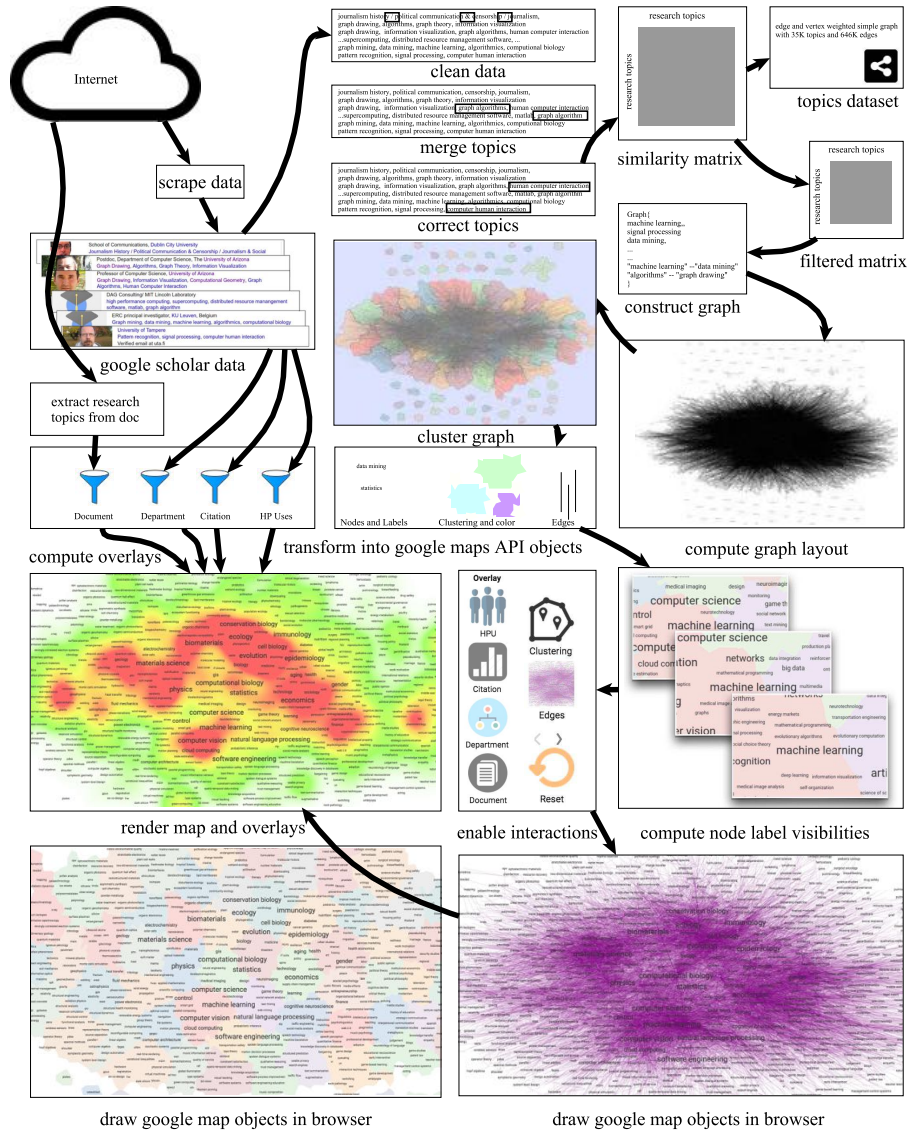


Fig. 2: Overview of the rtopmap system.

group about 25,000 journals into 554 subdisciplines that are further aggregated into 13 disciplines; see data and detailed procedure in [16]. However, the two-level map of 13 disciplines and 554 subdisciplines is too coarse for organizing, navigating, managing, and making sense of millions of publications.

Microsoft’s Academic Graph database has 50,000 fields of study (FOS) [50]. Three levels of relationships are present among the fields, although field importance is not measured or quantified. A FOS score based on researcher and citation counts has been proposed for computer science [23]. Hug *et al.* analyze

FOS and report that they tend to be dynamic and too specific, while field hierarchies are incoherent [33]. Liu *et al.* [36] use a hierarchical latent tree model (HLTM) to extract a set of hierarchical topics to summarize a corpus at different levels of abstraction. In HLTM, a topic is determined by identifying words that appear with high frequency in the topic and with low frequency in other topics. Yang *et al.* [59] use a HLTM in their visual analytics system, VISTopic. Mane and Börner [37] visualize 50 frequent and bursty words in their analysis of publication of the Proceeding of the National Academy of Sciences.

Words from paper titles have also been used as indicators for the content of a research topic, and visualizations based on this approach have been studied [12, 26, 61]. Many earlier approaches focus on analyzing specific journals, conferences, or research areas, e.g., analyzing computer science conferences and journals [26], trends in computer science research [23], the International Conference on Data Mining (ICDM) [38], publications in data visualization [31]. Domenico *et al.* [20] quantify attractive topics (i.e., topics that attract researchers from different areas). Sun *et al.* [53] build a network, with computer science conferences as nodes and edges between two conferences with common authors. Map-based visualization has been used for document visualization [51, 57].

Citations are considered an important contribution measurement [60] and are used in visualizations of scholar profiles [42] and paper recommendation systems [55]. Citation analysis with data from the Web of Science [41] and from Microsoft’s Academic Graph [33] have been considered. CiteRivers [31] and CiteVIS [52] analyze and visualize IEEE VIS conference citations, as do Ke *et al.* [35].

Also related to our work are many of the graph visualization techniques and tools. Graph layout algorithms are also provided in several libraries, such as GraphViz [4], OGDF [19], MSAGL [40], and VTK [47], which however, do not support interaction, navigation, and data manipulation. Visualization toolkits such as Prefuse [30], Tulip [9], Gephi [11], and yEd [56] support visual graph manipulation, and while they can handle large graphs, their rendering does not: even for graphs with a few thousand vertices, the amount of information rendered statically on the screen makes the visualization difficult to use.

There are research papers that describe interactive multi-level interfaces for exploring large graphs such as ASK-GraphView [8], topological fisheye views [29], and Grokker [44]. Software applications such as Pajek [21] for social networks, and Cytoscape [48] for biological data provide limited support for multi-level network visualization. These approaches rely on meta-graphs, made out of meta-vertices and meta-edges, which make interactions such as semantic zooming, searching, and navigation counter-intuitive. Not many of the tools and systems above provide browser-level navigation and interaction for large graphs.

Our work differs from earlier related work in several important ways: (1) we collect the underlying data using a bottom-up approach, based on self-reported data from the actual researchers, rather than using top-down taxonomies and ontologies, (2) our visualization provides map functionality (multiple zoom levels, searching, overlays), and (3) the ability to customize both the underlying base map and the overlays.

3 Network Generation

The set of research topics is not fixed or even well defined, as new topics are continuously created while old ones fade away. Automatically extracting keyword based topics from the research literature is a popular approach [26, 59], but has many limitations, such as identifying general topics (e.g., mathematics, physics) and specific sub-topics (e.g., graph drawing, network visualization). We use self-reported research topics from google scholar. Before we can build the topic graph, we scrape the data and then clean, split, merge, and correct the research topics. Next we build a similarity matrix M with topics as rows and columns. The value $M(i, j)$ represents the similarity between the pair of topics (i, j) , based on co-occurrence of the two topics in scholar profiles. The complete network is quite large, containing about 35,000 topics and 646,000 edges.

We reduce the size of the network by removing nodes and edges with low weights. Node weight is directly proportional to the number of scholar profiles that contain that topic, and edge weights are directly proportional to the number of scholar profiles that contain both topics. We remove a large number of infrequent topics, topics that contain typos, and topics listed in languages other than English.

Data Scraping: While some analysis of google scholar data exists [10, 25, 34], there is not much work based on data extracted from google scholar. Data retrieval is laborious due to the lack of API and metadata scarcity [17]. We started with a list of 1,000 universities [3] and then requested google scholar IDs for each university (e.g., MIT's ID is 16345133980181568013). We next collect research profiles from each university, by scraping the URL associated with that university (e.g., https://scholar.google.com/citations?view_op=view_org&org=16345133980181568013&hl=en&oi=io). Finally, we extract the name, affiliation, citations, and research topics of each individual researcher at that university, using a regular expression to match the relevant fields from the html file.

Data Cleaning: In the early days, researchers creating google scholar profiles manually created their own research topics. This might account for the large number of typos and acronyms in the dataset. These days google auto-suggests relevant topics and allows up to five research topics per profile. We use comma as the primary topics separator and a regular expression to replace other separators (e.g., ... / ; . #) with a comma. For html tags we use beautifulsoup [2], a python package for cleaning up html tags.

Topic Splitting and Merging: Once the data above is collected and analyzed, it is easy to see that many topics should be split; see Table 1. We split topics by pattern-matching conjunctions (i.e., or, and).

Merging is appropriate for topics that are similar but listed slightly differently. For example, out of out of half a million researchers in our dataset, 100 list *algorithm*, 20 list *algorithmics*, and 1,087 list *algorithms*. To handle this problem we need to determine the main topic with which the other topics should be merged. We use snowball [43] to find the root word by applying stemming

... methods for longitudinal or clustered data, statistics for neuroscience, ...
... data and model management, data mining, bioinformatics, algorithms ...
... new energy materials, Supercapacitor, photo and electro-catalysis of water ...
... Thyroid, Nuclear Cardiology and Neurology, Gluconeogenesis, ...

Table 1: Examples of records listing multiple topics that should be split.

(which removes endings such as $-s$, $-ed$, $-ing$). In the example above, snowball converted *algorithm*, *algorithmics*, and *algorithms* to the stemmed word *algorithm*, however, applying snowball may result in stems that are difficult to understand (e.g., *applied* and *applications* are converted to *appli*). With this in mind, we set the main topic to be the one with highest frequency among all the topics with the same stem.

Topic Correction: Topic splitting and merging does not resolve all topic issues, as further modifications might be needed due to leading and trailing spaces, lower and upper case letters, punctuation and control characters, and duplicate words. Other issues of this type include “Human Computer Interaction” and “Computer Human Interaction,” which are really the same topic. We try to address such issues with Google’s Openrefine [6] fingerprint key collision method, which attempts to find alternate representations of the same topic [18, 24, 32].

Network Statistics: After the steps above, our topic graph contains 34,774 topics and 646,582 edges. There are 17 components, but just one giant connected component (34,741 nodes and 646,565 edges). The average shortest path length is 3.141 which shows that the topic network is highly connected. The graph has a low global clustering coefficient of 0.09 (defined as the ratio of the number of triangles over the total number of node triples); see degree distribution in Fig. 11. The node “machine learning” has the highest degree and more researchers are reporting working on this topic than any other. Figure 3 shows the top ten topics by degree, by number of researchers, by citations per person.

Figure 9 shows which institutions contributed the most scholar profiles. Interestingly, some universities seem to have more profiles than academic staff, (likely due to doctoral and postdoctoral students with university affiliations), although the majority of the universities are associated with fewer profiles than the size of their academic staff.

Topics	Degree	Topics	Researchers	Topics	Cite/Person
machine learning	3314	machine learning	10726	particle physics	15906
artificial intelligence	2404	artificial intelligence	5766	high energy physics	15768
neuroscience	2033	neuroscience	5655	cosmology	7037
modeling	1902	computer vision	5372	clinical trials	6348
bioinformatics	1878	bioinformatics	4943	meta-analysis	6068
climate change	1846	robotics	3398	astronomy	6031
optimization	1827	data mining	3334	astrophysics	5501
education	1808	ecology	3281	genetic epidemiology	5389
nanotechnology	1788	materials science	3193	psychiatry	5205
statistics	1659	genetics	2951	nuclear physics	5130

Fig. 3: Top ten topics: highest degree, number of researchers, number of citations.

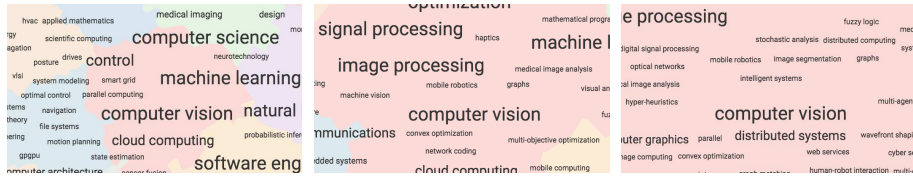


Fig. 4: Three zoom-level views near the “computer vision” topic.

4 Map Generation

We use the GMap framework [27] to generate map layouts of the research topic graph and extend it to support semantic zooming. There are three high-level steps: (1) embedding the topic graph in the plane, (2) grouping vertices into clusters, and (3) creating the geographic map representation. We embed the graph using a scalable force-directed algorithm (`sfdp` from `graphviz`) and then group the vertices using k -means clustering. To create the geographic map look, we use a modified Voronoi diagram based on the obtained embedding and clustering. The geographic regions are colored such that no two adjacent countries have colors that are too similar, using the spectral vertex labeling method [27].

We use the GraphViz implementation of node-overlap removal provided by `prism`, but that provides non-overlapping labels only for the complete basemap, and not for the other 7 level-of-detail views, needed for semantic zooming. Semantic zoom requires modifications to nodes, edges, clusters, and heatmaps. The google maps API handles all of these issues except for node-overlap (and hence node-label overlap), which is a natural side effect of zooming-in. To ensure that neither nodes nor labels overlap on any zoom-level, we compute different *node visibilities* for different zoom-levels. For each level, we sort the nodes by their weight (recall, that node weight is proportional to the number of researchers working on the topic associated with the node). We make i -th node visible on the j -th level if the bounding box of the i -th node does not overlap with the bounding boxes of nodes $1, 2, \dots, (i - 1)$. This algorithm takes $O(n^2)$ time but it can be improved by using [22]. Figure 4 shows how the local neighborhood of “computer vision” is changing in different zoom levels.

The size of the font label for topic t is directly proportional to the number of researchers working on that topic, denoted by the weight: $w(t)$. We assign font size from the range 80% to 200% of the default browser font size, as follows:

$$\mathcal{F}_t = \begin{cases} 80 & \text{if } w_t/10 \leq 80 \\ 200 & \text{if } w_t/10 \geq 200 \\ w_t/10 & \text{otherwise} \end{cases}$$

Web Interface and User Interaction: GMap produces a “basemap” from the given graph which is a static image that is not ideal for user interaction, such as zooming, panning, and searching. We enable interactions with the basemap with the help of the google maps API [54]. Specifically, we take the output from GMap and convert it into google map objects, i.e., `google.maps.SymbolPath`,

google.maps.Polygon, *google.maps.Polyline*, etc. For the web interface we provide 8 levels of details, showing different subgraphs, depending on the zoom level. We provide basic search functionality, which finds topics containing the query words. Clicking on a node shows the number of people who work on that topic in the underlying dataset and highlights edges to adjacent nodes (other topics that are frequently co-listed with that topic); see Fig. 10.

Basemaps: We compute several different basemaps, covering a different set of universities. The full map is determined by researchers in 1,000 universities around the world [3], but we also provide basemaps for universities in the United States and universities in Europe. Changing the basemap results in different node-weights and hence different label font-sizes. This is a useful feature when comparing a specific US university with universities around the world, with universities in the US, or with universities in Europe.

Overlays: After creating the basemap and all level-of-detail maps, we use overlays to show additional information, such as human resource investments of and citations associated with a specific university. Overlays can also be used to highlight topics associated with different departments (e.g., Computer Science, History) and even individual text documents (e.g., research paper, call for proposals). The overlay requests are collected from the browser (client) and the request is processed on the server, which then returns the necessary data to produce the overlay in the client. This is discussed in more detail in the next section.

5 Knowledge Strengths and Weaknesses

Quantifying knowledge strengths and weaknesses is a non-trivial challenge. We provide a university-level search feature which allows a specific university to be selected. We then attempt to visualize the strength and weakness of that university by computing the number of people working on different research topics and the number of citations associated with different topics for researchers from that university. Figure 5 illustrates this using the University of Arizona (UofA), Arizona State University (ASU), and the California Institute of Technology (Cal-Tech). It is easy to see that UofA has a significant human resource investment in ecology/evolution (large green circles around these topics) and this translates to many citations in these topics. Such visualizations also make it possible to see that UofA has not invested human resources in computer science (purple circles around CS topics) but CS is still associated with a large number of citations.

Citation Overlays: Let r be a researcher and let $topics(r)$ be the set of topics associated with researcher r . We denote number of citation received by r as $cite(r)$. Then we can define T as the set of topics that associated with university X as follows: $T = \bigcup_{\forall r \in X} topics(r)$. Then the citations associated with university X for each each topic $t \in T$ is determined by the sum of citations of researchers at university X who work on topic t : $c_X(t) = \sum_{r \in X \& t \in r} cite(r)$.

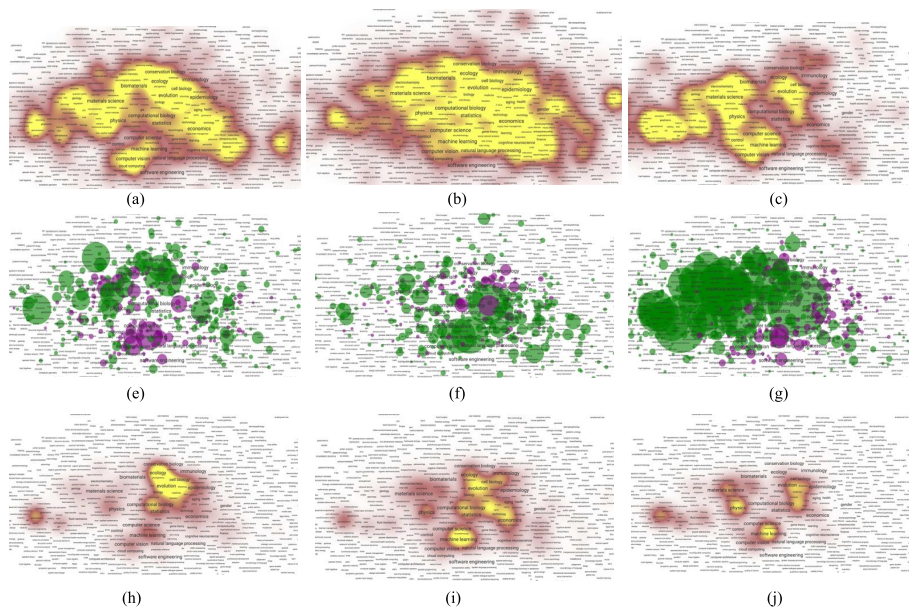


Fig. 5: University of Arizona (left column), Arizona State University (middle column), and California Institute of Technology (right column). (a-c) Heatmap overlays, based on citations. (e-g) Showing the number of people associated with different topics: green (purple) circles represent higher (lower) than average number of people working in this area. (h-j) Normalized citation heatmaps. Higher resolution images can be found in the image gallery at rtopmap.arl.arizona.edu.

The above formulas produce raw citation counts, although not all research fields cite at the same rate, e.g., “particle physics” is associated with more citations than average (due to high number of co-authors and citations per paper); see also the citations-per-person table in Fig. 3. With this in mind we provide the option to normalize citation counts by total number of citations associated with a specific field t : *normalized citation of t* = $c_X(t) \frac{c(t)}{C}$, where $c(t) = \sum_{r \in X \& t \in r} cite(r)$ and C is the total number of citations for all topics.

When a researcher lists multiple topics, it is not easy to determine which citation contributes to which topic. In this case each citation contributes equally to each topic. A more careful analysis of the citation meta-data might allow us to distribute the contribution of each citation to different topics.

Human Resources Overlays: We calculate human resource investment of a particular university by simply counting researchers and comparing the results to the averages. That is, to determine the human resource investment in topic t at university X , given a base set (top 1,000 universities, US universities, European universities), we calculate the difference in the percentage of researchers at

university X who work on topic t and the percentage of researchers in the base set who work on topic t . If this difference is positive (negative) then we consider this a human resource strength (weakness) of university X . This is illustrated with circles of different color: green for strength and purple for weakness. The size of the circles is proportional to the magnitude of the difference.

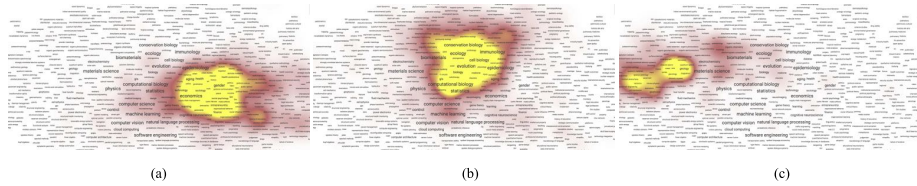


Fig. 6: Department visualization on the topic map. (a) Department of Economics, (b) Department of Biology and (c) Department of Geology.

Department Overlays: We can also visualize what part of the map is associated with what academic department using map overlays. Figure 6 shows three such department heatmap overlays. Different universities give different names to the same department, or group together academic units in colleges and schools. We match a given a keyword, such as economics, biology, geology, with the affiliation-information in our database. Google scholar affiliations are a combination of designation, department name, university name, as shown in Table 2. We use a regular expression to match researchers based on department name, which gives us the topics associated with the department, and eventually the weights, which are used to build the heatmap.

Affiliation
Professor of Computer Science and Artificial Intelligence, Granada University
Professor of Chemistry and Chemical Biology
Computer Science, Virginia Commonwealth University
Professor of Biochemistry and Molecular Biophysics, Washington University

Table 2: Examples of affiliations in google scholar.

Using such visualizations, we can find research fields that are shared by different departments, e.g., “machine learning,” which is studied in mathematics, management and information systems, and statistics; see Fig. 7.

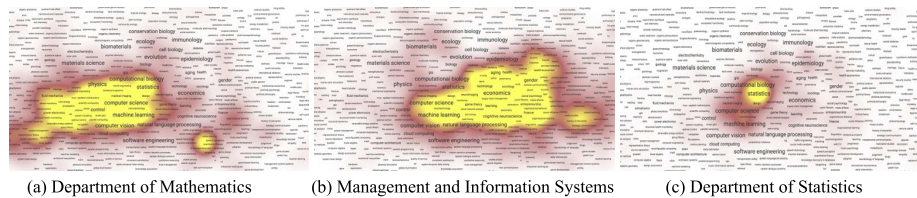


Fig. 7: Topics overlapped with different departments and institutes.

Document Overlays: We can also visualize documents as map overlays. We extract research terms and their frequency from a given input document, usually a URL. This requires a cleanup of the text, tokenization, and stemming. We next compute the term frequency of unigrams, bigrams, and trigrams, which become our *candidate research topics*. Then we match the candidate research topics with the research topics that are already in our database. The result is a collection of research topics associated with the document along with weights, which are used to create the heatmap overlay; see Fig. 8.

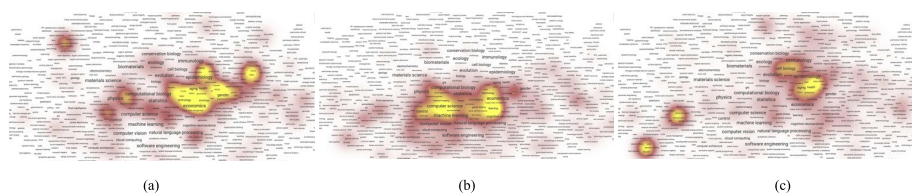


Fig. 8: Documents: (a) research paper on economics [39], (b) call for paper for CS conference NIPS [5], and (c) call for proposals for National Institute of Health [7].

6 Implementation

We use a variety of tools to process, clean, store, and process our data: mongodb scripts, sqlite, python, R, Java-Lucene, openrefine. Google maps API and jquery are used for map drawing and to handle user interaction in the web application. We run python-django for the webserver and mongodb for database storage and query.

For each researcher, our database stores name, google scholar id, university id, total citations, domain name of email address, affiliation, raw research areas, research phrases, and stemmed phrases.

The default settings for filtering the network (removing nodes and edges with low weights) results in the more compact version of the network with 6,052 nodes and 26,162 edges. Generating the topics map in svg format (layout, clustering, node-overlap removal, etc.) takes 14 seconds. Loading the initial base map takes 12,638ms, including 7,836ms for scripts, 2,444ms for rendering, and 275ms for painting (in google chrome v. 58 browser). Hiding the edges results in much faster interaction and this is indeed the default setting. Interaction with the basemap, map navigation, zooming with edges takes 1,515ms. Computing human resource overlays takes 3,129ms, while citation heatmaps require 1,231ms.

The `rtopmap` system currently runs as a virtual machine on a Dell PowerEdge R430 server with 2 Intel(R) Xeon(R) CPU E5-2530 v4 @2.20GHz processors and 32GB of memory.

7 Discussion and Limitations

We use google scholar as the source for our data, with all of the advantages (e.g., a great deal of information) and disadvantages (e.g., the data is not curated)

that this choice is associated with. Further, different research areas differ in their representation on google scholar. For example, there seem to be many more computer science and physics profiles than history and psychology ones. Researchers from different universities also use google scholar profiles at different rates.

Once the data has been gathered, the set of universities used to create the basemap has a non-trivial impact. Focusing only on English language terms, also biases the results. Despite our attempts to clean, split and merge topics, many issues remain. For example, researchers sometimes use acronyms (NLP for natural language processing, or HCI for human computer interaction) which should be expanded and merged.

Our strengths and weaknesses calculation, based on human resource investments, is associated with other biases: the human resource counts (from Wikipedia) are not guaranteed to be accurate, we do not distinguish tenure-track faculty from other type of staff (e.g., doctoral and postdoctoral students). Our citation-based calculations are also biased and inaccurate, as google scholar often misattributes papers and we cannot match specific citations to specific topics associated with a researcher, or distribute the citation contribution among its co-authors.

8 Conclusions and Future Work

We presented `rtopmap`: a system for visualizing and analyzing research topics. Despite non-trivial limitations, we believe that the system is useful as it gathers this information in self-reported, bottom-up fashion, rather than the more traditional top-down hierarchical taxonomies and ontologies. With the help of map overlays, `rtopmap` makes it possible to visualize human resource investments and scholarly output for different academic institutions. Department profiles and documents can also be visualized via overlays. Finally, the system implements in-the-browser, map-based interactive navigation of the fairly large underlying network, supporting panning, zooming, and searching. The `rtopmap` system is available at <http://rtopmap.arl.arizona.edu/>.

A natural extension of this work would be to match calls for proposals with individual researchers or groups of researchers at a specific university, making it possible to quickly identify potential participants in multi-disciplinary research proposals. Adding more data (e.g., funding statistics from national funding agencies, patents, media coverage of research projects) can augment the picture of a specific university, or enable more detailed comparisons between different universities. The visualization system itself is in a prototype stage. Our goal is to make it more responsive, improve HCI aspects, and extend its functionality to smaller screens.

Acknowledgments We thank Nirav Merchant, Mihai Surdeanu, and the Data7 Institute at the University of Arizona.

References

1. 2010 Mathematics Subject Classification - MSC2010 database. www.ams.org/msc/msc2010.html, accessed: 05-04-2017
2. Beautiful soup: We called him tortoise because he taught us. <https://www.crummy.com/software/BeautifulSoup/>, accessed: 01-13-2016
3. Cwur — center for world university rankings. <http://cwur.org/>, accessed: 09-13-2016
4. Graphviz — Graphviz - Graph Visualization Software. <http://www.graphviz.org/>, accessed: 05-25-2017
5. NIPS 2017 Call for Papers. <https://nips.cc/Conferences/2017/CallForPapers>, accessed: 05-04-2017
6. Openrefine. <http://openrefine.org/>, accessed: 05-04-2017
7. PA-17-219: Mechanisms of Alcohol-associated Cancers (R21). <https://grants.nih.gov/grants/guide/pa-files/PA-17-219.html>, accessed: 05-23-2017
8. Abello, J., Van Ham, F., Krishnan, N.: Ask-GraphView: A large scale graph visualization system. *Visualization and Computer Graphics, IEEE Transactions on* 12(5), 669–676 (2006)
9. Auber, D., Archambault, D., Bourqui, R., Lambert, A., Mathiaut, M., Mary, P., Delest, M., Dubois, J., Mélançon, G.: The Tulip 3 framework: A scalable software library for information visualization applications based on relational data. *Tech. Rep. RR-7860, INRIA* (2012)
10. Bar-Ilan, J.: Which h-index? a comparison of wos, scopus and google scholar. *Scientometrics* 74(2), 257–271 (2007)
11. Bastian, M., Heymann, S., Jacomy, M.: Gephi: an open source software for exploring and manipulating networks. *ICWSM 8*, 361–362 (2009)
12. Van den Besselaar, P., Heimeriks, G.: Mapping research topics using word-reference co-occurrences: A method and an exploratory case study. *Scientometrics* 68(3), 377–393 (2006)
13. Börner, K.: *Atlas of Science: Visualizing What We Know*. The MIT Press, Cambridge, MA (2010)
14. Börner, K.: *Atlas of Knowledge: Anyone Can Map*. The MIT Press, Cambridge, MA (2015)
15. Börner, K., Klavans, R., Patek, M., Zoss, A.M., Biberstine, J.R., Light, R.P., Larivière, V., Boyack, K.W.: Design and update of a classification system: The UCSD map of science. *PLOS ONE* 7(7) (Jul 2012)
16. Börner, K., Klavans, R., Patek, M., Zoss, A.M., Biberstine, J.R., Light, R.P., Larivière, V., Boyack, K.W.: Design and update of a classification system: The UCSD map of science. *PLoS ONE* 7(7), e39464 (07 2012)
17. Bornmann, L., Thor, A., Marx, W., Schier, H.: The application of bibliometrics to research evaluation in the humanities and social sciences: An exploratory study using normalized google scholar data for the publications of a research institute. *Journal of the Association for Information Science and Technology* 67(11), 2778–2789 (2016)
18. Cavnar, W.B., Trenkle, J.M., et al.: N-gram-based text categorization. *Ann Arbor MI* 48113(2), 161–175 (1994)
19. Chimani, M., Gutwenger, C., Jünger, M., Klau, G.W., Klein, K., Mutzel, P.: The open graph drawing framework (OGDF). *Handbook of Graph Drawing and Visualization* pp. 543–569 (2011)

20. De Domenico, M., Omodei, E., Arenas, A.: Quantifying the diaspora of knowledge in the last century. *Applied Network Science* 1(1), 15 (2016), <http://dx.doi.org/10.1007/s41109-016-0017-9>
21. De Nooy, W., Mrvar, A., Batagelj, V.: *Exploratory social network analysis with Pajek*, vol. 27. Cambridge University Press (2011)
22. Dwyer, T., Marriott, K., Stuckey, P.J.: Fast node overlap removal. In: *International Symposium on Graph Drawing*. pp. 153–164. Springer (2005)
23. Effendy, S., Yap, R.H.: Analysing trends in computer science research: A preliminary study using the microsoft academic graph. In: *Proceedings of the 26th International Conference on World Wide Web Companion*. pp. 1245–1250. WWW '17 Companion, International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, Switzerland (2017), <https://doi.org/10.1145/3041021.3053064>
24. Elmagarmid, A.K., Ipeirotis, P.G., Verykios, V.S.: Duplicate record detection: A survey. *IEEE Transactions on knowledge and data engineering* 19(1) (2007)
25. Falagas, M.E., Pitsouni, E.I., Malietzis, G.A., Pappas, G.: Comparison of pubmed, scopus, web of science, and google scholar: strengths and weaknesses. *The FASEB journal* 22(2), 338–342 (2008)
26. Fried, D., Kobourov, S.G.: Maps of computer science. 2014 IEEE Pacific Visualization Symposium (PacificVis) 00, 113–120 (2014)
27. Gansner, E.R., Hu, Y., Kobourov, S.: Gmap: Visualizing graphs and clusters as maps. In: *2010 IEEE Pacific Visualization Symposium (PacificVis)*. pp. 201–208 (March 2010)
28. Gansner, E., Hu, Y., Kobourov, S., Volinsky, C.: Putting recommendations on the map: Visualizing clusters and relations. In: *Proceedings of the Third ACM Conference on Recommender Systems*. pp. 345–348. RecSys '09, ACM, New York, NY, USA (2009), <http://doi.acm.org/10.1145/1639714.1639784>
29. Gansner, E., Koren, Y., North, S.: Topological fisheye views for visualizing large graphs. *TVCG* 11(4), 457–468 (July 2005)
30. Heer, J., Card, S.K., Landay, J.A.: Prefuse: a toolkit for interactive information visualization. In: *Proc. SIGCHI conference on Human factors in computing systems*. pp. 421–430. ACM (2005)
31. Heimerl, F., Han, Q., Koch, S., Ertl, T.: Citerivers: Visual analytics of citation patterns. *IEEE Transactions on Visualization and Computer Graphics* 22(1), 190–199 (Jan 2016)
32. Hjaltason, G.R., Samet, H.: Index-driven similarity search in metric spaces (survey article). *ACM Transactions on Database Systems (TODS)* 28(4), 517–580 (2003)
33. Hug, S.E., Ochsner, M., Brändle, M.P.: Citation analysis with microsoft academic. *CoRR abs/1609.05354* (2016), <http://arxiv.org/abs/1609.05354>
34. Jacsó, P.: Google scholar: the pros and the cons. *Online information review* 29(2), 208–214 (2005)
35. Ke, W., Borner, K., Viswanath, L.: Major information visualization authors, papers and topics in the acm library. In: *Information Visualization, 2004. INFOVIS 2004. IEEE Symposium on*. pp. r1–r1. IEEE (2004)
36. Liu, T., Zhang, N.L., Chen, P.: Hierarchical latent tree analysis for topic detection. In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. pp. 256–272. Springer (2014)
37. Mane, K.K., Börner, K.: Mapping topics and topic bursts in pnas. *Proceedings of the National Academy of Sciences* 101(suppl 1), 5287–5290 (2004), http://www.pnas.org/content/101/suppl_1/5287.abstract

38. Misue, K.: Visual exploration of a series of academic conferences. In: 2014 IEEE International Conference on Data Mining Workshop. pp. 314–320 (Dec 2014)
39. Motesharrei, S., Rivas, J., Kalnay, E.: Human and nature dynamics (handy): Modeling inequality and use of resources in the collapse or sustainability of societies. *Ecological Economics* 101, 90 – 102 (2014), <http://www.sciencedirect.com/science/article/pii/S0921800914000615>
40. Nachmanson, L., Robertson, G., Lee, B.: Drawing graphs with GLEE. In: *Graph Drawing*. pp. 389–394. Springer (2008)
41. Perianes-Rodriguez, A., Ruiz-Castillo, J.: University citation distributions. *Journal of the Association for Information Science and Technology* (2015)
42. Portenoy, J., West, J.D.: Visualizing scholarly publications and citations to enhance author profiles. In: *Proceedings of the 26th International Conference on World Wide Web Companion*. pp. 1279–1282. WWW '17 Companion, International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, Switzerland (2017), <https://doi.org/10.1145/3041021.3053058>
43. Porter, M.F.: Snowball: A language for stemming algorithms (2001)
44. Rivadeneira, W., Bederson, B.B.: A study of search result clustering interfaces: Comparing textual and zoomable user interfaces. *Studies* 21, 5 (2003)
45. Saket, B., Scheidegger, C., Kobourov, S.G., Börner, K.: Map-based visualizations increase recall accuracy of data. *COMPUTER GRAPHICS FORUM* 34(3), 441–450 (2015)
46. Saket, B., Simonetto, P., Kobourov, S., Börner, K.: Node, node-link, and node-link-group diagrams: An evaluation. *IEEE Transactions on Visualization & Computer Graphics* 20(12), 2231–2240 (2014)
47. Schroeder, W.J., Avila, L.S., Hoffman, W.: Visualizing with VTK: a tutorial. *Computer Graphics and Applications* 20(5), 20–27 (2000)
48. Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N., Schwikowski, B., Ideker, T.: Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome research* 13(11), 2498–2504 (2003)
49. Shiffrin, R.M., Börner, K.: Mapping knowledge domains. *Proceedings of the National Academy of Sciences of the United States of America* 101, 5183–5185 (2004), <http://www.pnas.org/cgi/lookup/doi/10.1073/pnas.0308030101>
50. Sinha, A., Shen, Z., Song, Y., Ma, H., Eide, D., Hsu, B.j.P., Wang, K.: An overview of microsoft academic service (mas) and applications. In: *Proceedings of the 24th international conference on world wide web*. pp. 243–246. ACM (2015)
51. Skupin, A.: A cartographic approach to visualizing conference abstracts. *IEEE Computer Graphics and Applications* 22(1), 50–58 (2002)
52. Stasko, J., Choo, J., Han, Y., Hu, M., Pileggi, H., Sadana, R., Stolper, C.D.: Citevis: Exploring conference paper citation data visually. *Posters of IEEE InfoVis* 2 (2013)
53. Sun, X., Ding, K., Lin, Y.: Mapping the evolution of scientific fields based on cross-field authors. *Journal of Informetrics* 10(3), 750 – 761 (2016), <http://www.sciencedirect.com/science/article/pii/S1751157715302352>
54. Svennerberg, G.: *Beginning Google Maps API 3*. Apress (2010)
55. West, J.D., Wesley-Smith, I., Bergstrom, C.T.: A recommendation system based on hierarchical clustering of an article-level citation network. *IEEE Transactions on Big Data* 2(2), 113–123 (June 2016)
56. Wiese, R., Eiglsperger, M., Kaufmann, M.: yFiles: Visualization and automatic layout of graphs. In: *GD*. pp. 453–454 (2001)

57. Wise, J.A., Thomas, J.J., Pennock, K., Lantrip, D., Pottier, M., Schur, A., Crow, V.: Visualizing the non-visual: spatial analysis and interaction with information from text documents. In: Proceedings of Visualization 1995 Conference. pp. 51–58 (Oct 1995)
58. Wojton, M.A., Heimlich, J.E., Burris, A., Tramby, Z.: Sense-making of big data spring break 2013 visualization recognition and meaning making. Report, Indiana University, Lifelong Learning Group (2014)
59. Yang, Y., Yao, Q., Qu, H.: Vistopic: A visual analytics system for making sense of large document collections using hierarchical topic modeling. Visual Informatics pp. – (2017), <http://www.sciencedirect.com/science/article/pii/S2468502X17300074>
60. Zhao, D., Strotmann, A.: Analysis and visualization of citation networks. Synthesis Lectures on Information Concepts, Retrieval, and Services 7(1), 1–207 (2015)
61. Zhou, N., Saltz, J., Mueller, K.: Maps of Human Disease: A Web-Based Framework for the Visualization of Human Disease Comorbidity and Clinical Profile Overlay, pp. 47–60. Springer International Publishing, Cham (2016), http://dx.doi.org/10.1007/978-3-319-41576-5_4

Appendix

We place several additional tables and figures in this section.

Name of University	Academic Staff	# Profiles
Stanford University	2118	8104
University of Washington	5803	5562
Harvard University	4671	5356
Massachusetts Institute of Technology	1021	3527
University of Michigan	6771	3413
University of Toronto	2547	3148
University of Cambridge	6645	2669
Texas A&M University	2700	2515
University of Minnesota	3804	2511
Pennsylvania State University	8864	2368

Fig. 9: Academic staff numbers according to Wikipedia entries.

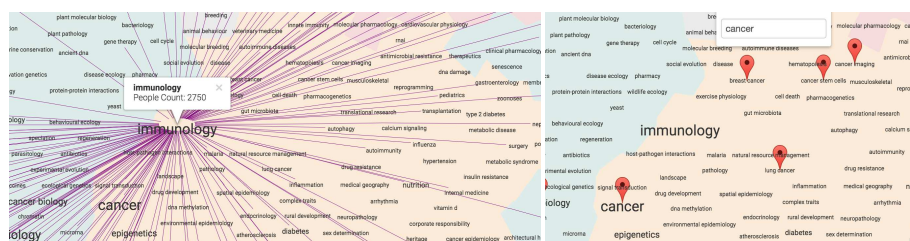


Fig. 10: Topic selection and search in rtopmap.

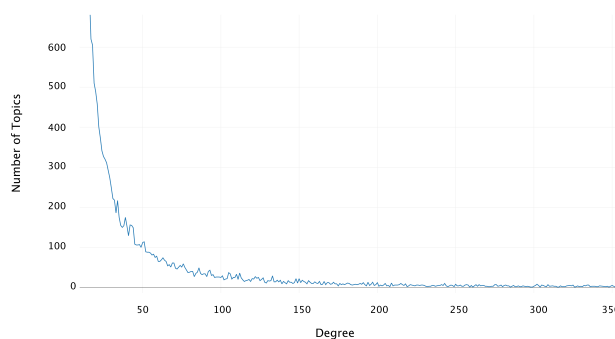


Fig. 11: Degree distribution of the topics network.