Bursting Scientific Filter Bubbles: Boosting Innovation via Novel Author Discovery

Jason Portenoy University of Washington jporteno@uw.edu

Eric Horvitz
Microsoft
horvitz@microsoft.com

Marissa Radensky University of Washington radensky@cs.washington.edu

Daniel S. Weld Allen Institute for AI University of Washington danw@allenai.org Jevin West University of Washington jevinw@uw.edu

Tom Hope Allen Institute for AI University of Washington tomh@allenai.org

ABSTRACT

Isolated silos of scientific research and the growing challenge of information overload limit awareness across the literature and hinder innovation. Algorithmic curation and recommendation, which often prioritize relevance, can further reinforce these informational "filter bubbles." In response, we describe Bridger, a system for facilitating discovery of scholars and their work. We construct a faceted representation of authors with information gleaned from their papers and inferred author personas, and use it to develop an approach that locates commonalities and contrasts between scientists to balance relevance and novelty. In studies with computer science researchers, this approach helps users discover authors considered useful for generating novel research directions. We also demonstrate an approach for displaying information about authors, boosting the ability to understand the work of new, unfamiliar scholars. Our analysis reveals that Bridger connects authors who have different citation profiles and publish in different venues, raising the prospect of bridging diverse scientific communities.

CCS CONCEPTS

• Human-centered computing \rightarrow User studies; • Information systems \rightarrow Recommender systems; Document representation.

KEYWORDS

scholarly recommendation, filter bubbles, author discovery

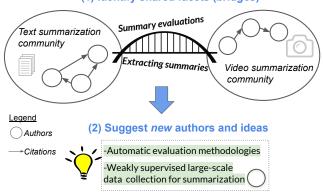
ACM Reference Format:

Jason Portenoy, Marissa Radensky, Jevin West, Eric Horvitz, Daniel S. Weld, and Tom Hope. 2022. Bursting Scientific Filter Bubbles: Boosting Innovation via Novel Author Discovery. In CHI Conference on Human Factors in Computing Systems (CHI '22), April 29-May 5, 2022, New Orleans, LA, USA. ACM, New York, NY, USA, 13 pages. https://doi.org/10.1145/3491102.3501905

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CHI '22, April 29-May 5, 2022, New Orleans, LA, USA © 2022 Association for Computing Machinery. ACM ISBN 978-1-4503-9157-3/22/04...\$15.00 https://doi.org/10.1145/3491102.3501905

(1) Identify shared facets (bridges)



erea otherwise aue to their disparities.

1 INTRODUCTION

"Opinion and behavior are more homogeneous within than between groups... Brokerage across structural holes provides a vision of options otherwise unseen." (Burt, 2004)

The volume of papers in computer science continues to skyrocket, with the DBLP computer science bibliography listing hundreds of thousands of publications in the year 2020 alone. In particular, the field of AI has seen a meteoric growth in recent years, with new authors entering the field every hour [55]. Research scientists rely largely on search and recommendation services like Google Scholar and Semantic Scholar to keep pace with the growing literature and the authors who contribute to it. The literature retrieval services algorithmically decide what information to serve to scientists [1, 14], using information such as citations and textual content as well as behavioral traces such as clickthrough data, to inform machine learning models that output lists of ranked papers or authors. By relying on user behavior and queries, these services adapt and reflect human input and, in turn, influence subsequent search behavior. This cycle of input, updating, engagement, and response can lead to an amplification of biases around searchers' prior awareness and knowledge [29]. Such biases include selective exposure [17],

homophily [35], and the aversion to information from novel domains that require more cognitive effort to consider [24, 30]. By reinforcing these tendencies, systems that filter and rank information run the risk of engendering so-called *filter bubbles* [45] that fail to show users novel content outside their narrower field of interest.

These bubbles and silos of information can be costly to individual researchers and for the evolution of science as a whole. They may lead scientists to concentrate on narrower niches [31], reinforcing citation inequality and bias [44], limiting cross-fertilization among different areas that could catalyze innovation [24, 26, 30], and preventing knowledge brokerage across groups that has been associated with creativity and success [5]. Addressing filter bubbles in general, in domains such as social media and e-commerce recommendations, is a hard and unsolved problem [12, 19, 67]. The problem is especially difficult in the scientific domain. The scientific literature consists of complex models and theories, specialized language, and an endless diversity of continuously emerging concepts. Connecting blindly across these cultural boundaries requires significant cognitive effort [58], translating to time and resources most researchers are unlikely to have at their disposal to enter unfamiliar research territory.1

Our vision in this paper is to develop an approach that **boosts scientific innovation and builds bridges across scientific communities**, by helping scientists **discover authors that spark new ideas** for research directions. Working toward this goal, we developed Bridger, illustrated in Figure 1. Our main contributions include:

- A multidimensional author representation for matching authors along specific facets. Our novel representation includes information extracted automatically from papers, specifically tasks, methods and resources, and automatically inferred personas that reflect the different focus areas on which each scientist works. Each of these aspects is embedded in a vector space based on its content, allowing the system to identify authors with commonalities along specific dimensions and not others, such as authors working on similar tasks but not using similar methods.
- Boosting discovery of useful authors and ideas from novel areas. We explore the utility of our author representation in experiments with computer science researchers interacting with Bridger. We find that this representation helps connect users with authors considered novel and relevant, assisting users in finding potentially useful research directions. Bridger outperforms a strong neural model currently employed by a public scholarly search engine for search and recommendation²— despite Bridger's focus on surfacing novel content and the built-in biases associated with this novelty. We conduct in-depth interviews with researchers, studying the tradeoffs between novelty and relevance in scientific content recommendations and discussing challenges and directions for author discovery systems.
- Exploring how to effectively depict recommended authors. In addition to assessing *what* authors to recommend to spark new ideas for research directions, we also consider *how* to display authors in a way that enables users to rapidly understand what

- new authors work on. We employ Bridger as an experimental platform to explore which facets should be displayed to users, investigating various design choices and tradeoffs. We obtain substantially better results in terms of user *understanding of profiles of unknown authors*, when displaying information taken from our author representation.
- Evidence of bridging across research communities. Finally, we conduct in-depth analyses revealing that Bridger surfaces novel and valuable authors and their work that are unlikely to be discovered in the absence of Bridger due to publishing in different venues, citing and being cited by non-overlapping communities, and having greater distances in the social co-authorship network.

Taken together, the ability to uncover novel and useful authors and ideas for research directions, and to serve this information to users in an effective and intuitive manner, suggests a future where automated systems are put to work to build bridges across communities, rather than blindly reinforcing existing filter bubbles.

2 RELATED WORK

Inspirational Stimuli. Our work is related to literature focused on computational tools for boosting creativity [9, 20, 24, 26, 30]. Experiments in this area typically involve giving participants a specific concrete problem, and examining methods for helping them come up with creative solutions [24, 26]. In our efforts reported in this paper, we do not assume to be given a single concrete problem. Rather, we are given *authors* and their papers, and automatically identify personalized inspirations in the form of other authors and their contributions. These computationally complex objects — authors can have many papers with different themes, each paper with many facets and authored by multiple co-authors — are very different to the short, single text snippets typically used in this line of work [24, 26], or even to paper abstracts [9]. A recurring theme in this area is the notion of a "sweet spot" for inspiration: not too similar to a given problem that a user aims to solve, and not too far afield [18]. Finding such a sweet spot remains an important challenge. Some work attempts to find this sweet spot by identifying analogies as inspirations — abstract structural relations between ideas [24, 26, 30]. We study a related notion, balancing commonalities and contrasts between researchers for discovering authors that spark new research directions, trading off relevance and novelty. In our work, commonalities represent shared facets between authors (e.g., similar tasks) intended to help surface relevant authors, while contrasts along other dimensions (e.g., dissimilar methods) help promote novelty.

Filter Bubbles and Recommendations. How to mitigate the filter bubble effect is a challenging open question for algorithmic recommendation systems [43], explored recently for movies [67] and in e-commerce [19] by surfacing content that is aimed at being both novel and relevant. One approach that has been explored for mitigating these biases is judging recommendations not only by accuracy, but with other metrics such as diversity (difference between recommendations) [12, 65], novelty (items assumed unknown to the user) [66], and serendipity (a measure of relevance and surprise associated with a positive emotional response) [62]. The notion of serendipity is notoriously hard to quantitatively define and measure [11, 28, 62, 64]; recently, user studies have explored human

 $^{^{1}}$ The challenge of limited time to explore novel directions is also discussed in our interviews with researchers; see \$6.

²https://twitter.com/SemanticScholar/status/1267867735318355968.

perceptions of serendipity [11, 62], yet this problem remains very much open. A distinct, novel feature of our work is the focus on the scientific domain, and that unlike the standard recommendation system setting we measure our system's utility in terms of boosting users' ability to discover authors that spur new ideas for research directions. In experiments with computer science researchers, we explore interventions that could potentially help provide bridges to authors working in diverse areas, with an approach based on finding faceted commonalities and contrasts between researchers. Our approach is broadly related to literature-based discovery [53], where the goal is to surface scientific hypotheses by identifying potential links between concepts (e.g., drugs and diseases) that are not apparent by reading individual papers. Work in this area typically does not evaluate in the context of inspiring human users but rather in the ability to predict future links between biomedical entities (e.g., new links between drugs and diseases) [39]. Furthermore, in our work we focus explicitly on surfacing potential links between authors, complex "objects" with many papers and lines of work.

Scientific Recommendations. Work in this area typically focuses on recommending papers, using proxies such as citations or coauthorship links in place of ground truth [2, 14, 48, 54]. In addition to being noisy proxies in terms of relevance, these signals reinforce existing patterns of citation or collaboration, and are not indicative of papers or authors that would help users generate novel research directions — the focus of Bridger. Furthermore, we perform controlled experiments with researchers to be able to better evaluate our approach without the biases involved in learning from observational data on citations or co-authorship. One related recent direction considers the problem of diversifying social connections made between academic conference attendees [56, 57, 63], by definition a relatively narrow group working in closely-related areas, using attendee metadata or publication similarity.

Visualization-aided Exploration of the Scientific Literature. There is a large body of work on the topic of mapping and visualizing networks of scholarly publications [6, 10, 13]. Recently, attempts have been made in the information visualization research community to build tools for exploring connected aspects of the literature using interactive visualizations [15, 22, 25, 27, 47]. One recent paper [41] designed a system to promote serendipitous discovery of new papers by finding semantically similar papers in a word embedding space; however, relying on such embeddings tuned for document similarity can reinforce filter bubbles, as we argue and demonstrate in this paper. In particular, we show that a state-of-art neural embedding model used by a popular scientific search engine for representing papers, underperforms our approach when it comes to discovering authors and ideas for research directions that are not only relevant, but also novel and more diverse. Finally, our work also studies novel design choices for displaying information on recommended authors, in a manner that increases users' ability to understand the work of scholars in unfamiliar areas.

3 BRIDGER: APPROACH OVERVIEW

In this section we present our novel faceted representation of authors, and methods for using this representation for author discovery by matching researchers along specific dimensions (Figure 2).

We also present methods for depicting the recommended authors when showing them to users. Bridger is designed to enable the study of different design choices for connecting authors and ideas across scientific filter bubbles and promoting discovery. We present the general framework, and the specific instantiations that we explore. We start by describing our representation for papers, and how Bridger represents authors by aggregating paper-level information and decomposing authors into *personas*.³

3.1 Paper Representations

Each paper P contains rich, potentially useful information. This includes raw text such as in a paper's abstract, incoming and outgoing citations, publication date, venues, and more. One key representation we derive from each paper P is a vector representation \tilde{P} , using a state-of-art scientific paper embedding model. This neural model captures overall coarse-grained topical information on papers, shown to be powerful in clustering and retrieving papers [14].

Another key representation is based on fine-grained facets obtained from papers. Let $\mathcal{T}_{P_i} = \{t_1, t_2, \ldots\}$ be a set of *terms* appearing in paper i. Each term is associated with a specific *facet* (category). Each term t is located in a "cell" in the matrix illustrated in Figure 2, with facets corresponding to the columns and papers to rows. Each term $t \in \mathcal{T}_{P_i}$ is also embedded in a vector space using a neural language model (see §3.5), yielding a \tilde{t} vector for each term.

We consider several categories of terms in this paper: coarse-grained paper topics inferred from the text [61], and fine-grained spans of text referring to *methods, tasks and resources* automatically extracted from paper *i* with a scientific named entity recognition model [59]. These three fine-grained categories are core aspects of computer science papers; in other words, they are key semantic concepts or "building blocks" with which computer scientists reason about research (developing new methods for a given task, developing new tasks, developing new datasets to support certain tasks, etc.) [8, 34]. These facets can help users find authors who spark ideas for new methods they can apply to their tasks, new tasks where their methods may be relevant, or new resources to explore. This relates to the fundamental role "functional aspects" play in science [23, 26] and in linking between distant ideas and areas [9, 24, 26].

3.2 Author Representations

We represent an author, \mathcal{A} , as a set of *personas* in which each persona is encoded with facet-wide aggregations of term embeddings across a set of papers. Figure 2 illustrates this with outlines of "slices" in bold — subsets of rows and columns in the illustrated matrix, corresponding to personas (subsets of rows) and facets (columns).

Author Personas. Each author $\mathcal A$ can work in multiple areas. In our setting, this can be important for understanding the different interests of authors, enabling more control on author suggestions. We experiment with a clustering-based approach for constructing personas, $P_{\mathcal A}$, based on inferring for each set of author papers $P_{\mathcal A}$ a segmentation into K subsets reflecting a common theme — illustrated as subsets of rows in the matrix in Figure 2. We also experiment with a clustering based on the network of co-authorship

³The source code for data processing, author representation and ranking, and the user-facing application displaying this data, can be found in the supplementary materials.

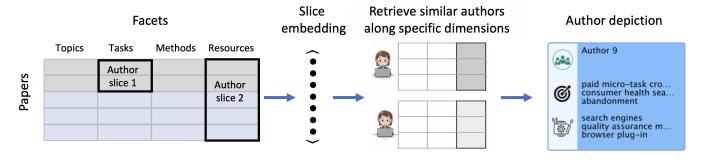


Figure 2: Overview of Bridger's author representation, retrieval, and depiction. Users are represented in terms of a matrix with rows corresponding to papers, and columns corresponding to facets. Bridger finds suggested authors who match along certain "slices" of the user's data – certain facets, subsets of papers, or both.

collaborations in which \mathcal{A} takes part. See §3.5 for details on clustering. As discussed later (§4), we find that the former approach in which authors are represented with clusters of papers elicits considerably better feedback from scholars participating in our experiments.

Co-authorship Information. Each paper P is in practice authored by multiple people, i.e., it can belong to multiple authors \mathcal{A} . Each author assumes a position k for a given paper, potentially reflecting the strength of affinity to the paper. As discussed below (§3.5), we make use of this affinity in determining what weight to assign terms \mathcal{T}_{P_i} for a given paper and given author.

Author-level Facets. Finally, using the above information on authors and their papers, we construct multiple author-level facets that capture different aggregate aspects of \mathcal{A} . More formally, in this paper we focus our experiments on author facets $\mathcal{V}_{\mathcal{A}} = \{\mathbf{m}, \mathbf{t}, \mathbf{r}\}$, where \mathbf{m} is an aggregate embedding of \mathcal{A} 's method facets, \mathbf{t} is an embedding capturing \mathcal{A} 's tasks, and \mathbf{r} represents \mathcal{A} 's resources. In addition, we also construct these facets separately for each one of the author's personas $P_{\mathcal{A}}$ — corresponding to "slice embeddings" over subsets of rows and columns in the matrix illustrated in Figure 2. In analyses of our experimental results (§5), we also study other types of information such as citations and venues; we omit them from the formal notations to simplify presentation.

3.3 Approaches for Recommending Authors

For a given author $\mathcal A$ using Bridger, we are interested in automatically suggesting new authors working on areas that are relevant to $\mathcal A$ but also likely to be interesting and spark new ideas for research directions. We are given a user $\mathcal A$, their set of personas $P_{\mathcal A}$, and for each persona its faceted representation $\mathcal V_{\mathcal A} = \{m,t,r\}$. We are also given a large pool of authors across computer science, $\{\mathcal A_1,\mathcal A_2,\ldots\}$, from which we aim to retrieve author suggestions to show $\mathcal A$.

Baseline Model. We employ Specter, a strong neural model to which we compare, trained to capture overall topical similarity between papers based on text and citation signals (see Cohan et al. [14] for details) and used for serving recommendations as part of a large public academic search system. For each of author \mathcal{A} 's papers P, we use this neural model to obtain an embedding \tilde{P} . We then

derive an aggregate author-level representation \tilde{p} (e.g., by weighted averaging that takes author-term affinity into account, see §3.5). Similar authors are computed using a simple distance measure over the dense embedding space. As discussed in the introduction and §2, this approach focuses on retrieving authors with the most overall similar papers to \mathcal{A} . Intuitively, the baseline can be thought of as "summing over" both the rows and columns of the author matrix in Figure 2. By aggregating across all of \mathcal{A} 's papers, information on finer-grained sub-interests may be lost. In addition, by being trained on citation signals, it may be further biased and prone to favor highly-cited papers or authors.

To address these issues, we explore a formulation of the author discovery problem in terms of matching authors along specific dimensions that allow more fine-grained control – such as by using only a subset of views in $\mathcal{V}_{\mathcal{A}}$, or only a subset of \mathcal{A} 's papers, or both — as in the row and column *slices* seen in Figure 2. This decomposition of authors also enables us to explore *contrasts* along specific author dimensions, e.g., finding authors who use similar tasks to \mathcal{A} but use very different methods or resources.

- Single-facet matches: For each author \mathcal{A}_i in the pool of authors $\{\mathcal{A}_1,\mathcal{A}_2,\ldots\}$, we obtain their respective aggregate representations $\mathcal{V}_{\mathcal{A}_i} = \{\mathbf{m},\mathbf{t},\mathbf{r}\}$. We then retrieve authors with similar embeddings to \mathcal{A} along one dimension (or matrix column in Figure 2; e.g., \mathbf{r} for resources), ignoring the others. Unlike the baseline model, which aggregates all information appearing in \mathcal{A} 's papers tasks, methods, resources, general topics, and any other textual information this approach is able to disentangle specific aspects of an author, potentially enabling discovery of more novel, remote connections that can expose users to more diverse ideas and cross-fertilization opportunities.
- Contrasts: Finding matches along *one* dimension does not guarantee retrieving authors who are *distant* along the others. As an example, finding authors working on *tasks* related to *scientific knowledge discovery* and *information extraction from texts*, could be authors who use a diverse range of *resources*, such as *scientific papers*, *clinical notes*, etc. While the immense diversity in scientific literature makes it likely that focusing on similarity along one dimension only will still surface diverse results in terms of the other (see results in §5), we seek to further ensure this. To do so, we apply a simple approach inspired by recent work on

retrieving inspirations [26]: We first retrieve the top K authors $\{\mathcal{A}_1, \mathcal{A}_2, \ldots, \mathcal{A}_K\}$ that are most similar to \mathcal{A} along one dimension (e.g., t), for some relatively $large\ K$ (e.g., K=1000). We then rank this narrower list inversely by another dimension (e.g., r), and show user \mathcal{A} authors from the top of this list. Intuitively, this approach helps balance relevance and novelty by finding authors who are similar enough along one dimension, and within that subset find authors who are relatively distant along another.

Persona-based matching: Finally, to account for the different focus areas authors may have, instead of aggregating over all of an author's papers, we perform the same single-view and contrast-based retrieval using the author's personas P_A — or, in other words, row-and-column slices of the matrix in Figure 2.

3.4 Depicting Recommended Authors

Our representation allows us to explore multiple design choices not only for *which* authors we show users, but also *how* we show them. In our experiments (§4, §5), we evaluate authors' facets and personas in terms of their utility for helping researchers learn about new authors, and for controlling how authors are filtered.

Term Ranking Algorithms to Explain What Authors Work On. Researchers, flooded with constant streams of papers, typically have a very limited attention span to consider whether some new author or piece of information is relevant to them. It is thus important that the information we display for each author (such as their main methods, tasks, resources, and also papers) is ranked, such that the most important or relevant terms appear first. We explore different approaches to rank the displayed terms, balancing between relevance (or centrality) of each term for a given author, and coverage over the various topics the author works on. We compare several approaches, including a customized relevance metric we design, in a user study with researchers (§4). We discuss in more detail the ranking approaches we try in §3.5.

Retrieval Explanations. In addition to term ranking approaches aimed at explaining to users of Bridger what a new suggested author works on, we also provide users with two rankings that are geared for explaining how the retrieved authors relate to them. First, we allow users to rank author terms $\mathcal T$ by how similar they are to their own list of terms (for each facet, separately). Second, users can also rank each author's papers by how similar they are to their own — showing the most similar papers first. These explanations can be regarded as a kind of "anchor" for increasing trust, which could be especially important when suggesting novel, unfamiliar content.

3.5 Implementation Details

3.5.1 Data. We use data from the Microsoft Academic Graph (MAG) [51]. We use a snapshot of this dataset from March 1, 2021. We also link the papers in the dataset to those in a Semantic Scholar, a large public academic search engine. We limit the papers and associated entities to those designated as Computer Science papers. We focus on authors' recent work, limiting the papers to those published between 2015 and 2021, resulting in 4,650,474 papers from 6,433,064 authors. Despite using disambiguated MAG author data, we observe the challenge of author ambiguity still persists

[52]. In our experiments, we thus exclude participants with very few papers (see §5), since disambiguation errors in their papers stand out prominently.

3.5.2 Term Extraction. We extract terms (spans of text) referring to tasks, methods, and resources mentioned in paper abstracts and titles, using the state-of-art DyGIE++ IE model [59] trained on SciERC [34]. We extracted 10,445,233 tasks, 20,705,854 methods, and 4,978,748 resources from 3,594,975 papers. We also use MAG topics, higher-level coarse-grained topics available for each paper in MAG. We expand abbreviations in the extracted terms using the algorithm in [50] implemented in ScispaCy [42].

3.5.3 Scoring Papers by Relevance to an Author. The papers published by an author have varying levels of importance with regard to that author's overall body of publications. To capture this, we use a simple heuristic that takes into account two factors: the author's position in a paper as a measure of affinity (see §3.2), and the paper's overall impact in terms of citations. More formally, for each author \mathcal{A} , we assign a weight $w_{\mathcal{A},P}$ to each paper P in $P_{\mathcal{A}}$, $w_{\mathcal{A},P} = \mathsf{pos}_{\mathcal{A},P} \times \mathsf{Rank}_P$, where $\mathsf{pos}_{\mathcal{A},P}$ is 1.0 if \mathcal{A} is first or last author on P and 0.75 otherwise, and Rank_P is MAG's assigned paper Rank (a citation-based measure of importance, see [61] for details), normalized by min-max scaling to a value between .5 and 1.

3.5.4 Author Similarity. We explore several approaches for author similarity and retrieval, all based on paper-level aggregation as discussed in §3.3. For the document-level Specter baseline model discussed in §3.3, we obtain 768-dimensional embeddings for all of the papers. To determine similarity between authors, we take the average embedding of each author's papers, weighted by the paper relevance score described above. We then compute the cosine similarity between this author and the average embedding of every other author. For our faceted approach, we compute similarities along each author's facets, using embeddings we create for each term in each facet. The model used to create embeddings was CS-RoBERTa [21], which we fine-tuned for the task of semantic similarity using the Sentence-BERT framework [49]. For each author or persona, we calculate an aggregate representation along each facet by taking the average embedding of the terms in all of the papers, weighted by the relevance score of each associated paper.

3.5.5 Identification of Personas. We infer author personas using two different approaches. For the first approach we cluster the coauthorship network using the ego-splitting framework in [16]. In a second approach, we cluster each authors' papers by their Specter embeddings using agglomerative clustering with Ward linkage [37] on the Euclidean distances between embedding vectors. In our user studies, we show participants their personas and the details of each one (papers, facets, etc.). To make this manageable, we sort the clusters (personas) based on each cluster's most highly ranked paper according to MAG's assigned rank, and show participants only their top two personas.

⁴https://www.semanticscholar.org/

⁵This specific implementation reflects the norms around author position in computer science research. While many fields share these same norms, they are not universal, and so these methods can be adjusted when this system is applied to fields with different conventions.

⁶Implemented in the scikit-learn Python library [46]. Distance threshold of 85.

⁷Some authors do not have detected personas; we observe this to often be the case with early-career researchers.

- 3.5.6 Term Ranking for Author Depiction. We evaluate several different strategies to rank terms (methods, tasks, resources) shown to users in Experiment I (§4):
- **TextRank:** For each term t in an author's set of papers, we create a graph $G_F = (V, E)$ with vertices V the terms and weighted edges E, where weight w_{ij} is the euclidean distance between the embedding vectors $\tilde{t_i}$ and $\tilde{t_j}$. We score each term t_i according to its PageRank value in G_F [36].
- **TF-IDF:** For each t, we compute TF-IDF across all authors, considering each author as a "document" (bag of terms) in the IDF (inverse document frequency) term, counting each term once per paper. We calculate the TF-IDF score for each term for each author, and use this as the term's score.
- Author relevance score: For each *t*, we calculate the sum of the term's relevance scores (§ 3.5.3) derived from their associated papers. If a term is used in multiple papers, the associated paper's score is used for each summand.
- **Random:** Each term *t* is assigned a random rank.

4 EXPERIMENT I: AUTHOR DEPICTION

In systems that help people find authors, such as Microsoft Academic Graph, Google Scholar, and AMiner [60], authors are often described in terms of a few high-level topics. In advance of exploring how we might leverage facets to engage researchers with a diverse set of authors, we performed a user study to gain a better understanding of what information might prove useful when depicting authors. We started from a base of Microsoft Academic Graph (MAG) topics, and then added their extracted facets (tasks, methods, resources). We investigated the following research questions:

- **RQ1**: Do tasks, methods, and/or resources complement MAG topics in depicting an author's research?
- RQ2: Which term ranking best reflects an author's interests?
- **RQ3**: Do tasks, methods, and/or resources complement MAG topics in helping users gain a better picture of the research interests of *unknown* authors?
- **RQ4**: Do personas well-reflect authors' different focus areas?

4.1 Experiment Design

After the experiment was approved for IRB exemption, thirteen computer-science researchers were recruited for the experiment through Slack channels and mailing lists. Participants were compensated \$20 over PayPal for their time. Study sessions were one-hour, semi-structured interviews recorded over Zoom. The participants engaged in think-aloud throughout the study. They evaluated a depiction of a known author (e.g., research mentor) for accuracy in depicting their research, as well as depictions of five *unknown* authors for usefulness in learning about new authors.

Throughout all parts of the experiment, the interviewer asked follow-up questions regarding the participant's think-aloud and reactions. To address **RQ1** and **RQ2**, the participants first evaluated the accuracy of a known author's depiction.

Step I. To begin, we presented the participant with only the top 10 MAG topics for the known author. We asked them to mark any topic that was unclear, too generic, or did not reflect the author's

research well. Next, we provided five more potential lists of terms. One of these lists consisted of the next 10 top topics. The other four presented 10 tasks, each selected as the top-10 ranked terms using the strategies described in §3.5. We asked participants to rank the five lists (as a whole) in terms of how well they complemented the first list (with an option to select none).

Step II. The process then repeated for five more potential lists to complement the original topics and the highest-ranked second list selected in Step I — this time, with methods instead of tasks. If the participant ranked a methods list highest, we then presented the participant with a resources list that used the same ranking strategy preferred by the participant for methods, and asked whether or not this list complemented those shown so far.

Step III. To address **RQ3**, participants next evaluated the utility of author depictions for five unknown authors. To describe each unknown author, we provided topics, tasks, methods, and resources lists with 10 terms each. The non-topics lists were ranked using TF-IDF as a default. The participant noted whether or not each additional non-topics list complemented the preceding lists in helping them understand what kind of research the unknown author does.

Step IV. Finally, for **RQ4**, we asked participants to evaluate the known author's distinct personas presented in terms of tasks, which were ranked using TF-IDF. On a Likert-type scale of 1-5, participants rated their agreement with the statement, "The personas reflect the author's different research interests (since the year 2015) well."

4.2 Results

The think-aloud results were evaluated using thematic analysis [4]. The transcripts of each participant were reviewed and paraphrased with quotes and contextual notes. Codes and themes were generated from the review. The notes for each participant (and full transcripts as necessary) were then reviewed again, and relevant codes were connected to each participant.

4.2.1 Results for RQ1. The majority of participants found that tasks, methods, and resources complemented topics to describe a known author's research. For both tasks and methods, 11 of 13 participants felt that seeing information about that facet, more so than additional top MAG topics or no additional information, complemented the original top ten MAG topics. The prevailing grievance with the additional MAG topics was that they were too general. For example, looking at the topics column that was an alternative to the potential tasks columns, P7 said, "Some of it like 'healthcare' and 'applied psychology' are too high-level." In the same situation, P2 commented, "AI is very general and implied by [the author's other topics] machine learning or NLP or IR." Furthermore, 7 of 9 participants who evaluated a resources list thought that it complemented the preceding lists.

4.2.2 Results for RQ2. Participants overall preferred the relevance score ranking strategy for tasks and methods. We compared the four ranking strategies and MAG topics baseline strategy for both tasks and methods. For each participant, we awarded points to each strategy based on its position in the participant's ranking of the five strategies (Figure 3a, b). We awarded the least favorite strategy one point and the most favorite strategy five points. Since there were 13 participants, a strategy could accumulate up to 65

⁸The script for Experiment I can be found in our supplementary materials.

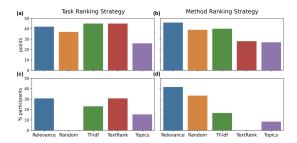


Figure 3: Points awarded to each ranking strategy for tasks (a) and methods (b), and percentage of participants who favored each strategy most for tasks (c) and methods (d).

points. Separately, we counted how many times each strategy was a participant's favorite strategy (Figure 3c, d). With regards to tasks, TextRank and TF-IDF accrued the most points from participants, with the relevance score trailing close behind (Figure 3a). Meanwhile, the MAG topics baseline accrued the least points, even fewer than the random task ranking strategy. In addition, relevance score and TextRank were chosen most often as the favorite task ranking strategy (Figure 3c). With regards to methods, the relevance score ranking strategy performed best in terms of both total points (Figure 3b) and favorite strategy (Figure 3d).

4.2.3 Results for RQ3. Participants generally found tasks, methods, and resources helpful to better understand what kind of research an unknown author does. To calculate how many participants were in favor of including tasks, methods, and resources to help them better understand an author, we determined the average of each participant's binary response per facet. Adding up the 13 responses for each facet, we saw that the majority of participants thought each additional facet helped them understand the unknown author better. All 13 participants found the tasks helpful, eight found the methods helpful, and 12 found the resources helpful. As an example, P12 connected an unknown author's topics, tasks, and methods to better understand them: "I wouldn't have known they were an information retrieval person from the [topics] at all.... The previous things [in topics and tasks] that mentioned translation and information retrieval and kind of separately... This [methods section] connects the dots for me, which is nice." Interestingly, methods were not viewed to be as useful as tasks or resources. The majority of participants cited unfamiliar terms as a key issue. For instance, after looking at the first four methods such as "Experience Sampling Method" and "TapSense" in an unknown author's methods column, P9 looked at the fifth one and noted, "'Body posture calculations'- I think that's the first phrase that I can pick out maybe what it means"

4.2.4 Results for RQ4. Participants indicate preference for personas selected based on papers rather than co-authorship. After the experiment, six participants were informally asked to compare the experiment's personas selected based on co-authorship with the personas based on paper-based clustering (see §3.5). Four of them preferred the updated version. Furthermore, one of the users who preferred the old version still thought the updated version had better personas themselves and merely did not like the

updated personas' ordering. In addition, all six participants liked seeing the personas in terms of papers. In our experiment in §5, we observed much higher satisfaction with the updated personas in comparison to the original personas of this experiment.

5 EXPERIMENT II: AUTHOR DISCOVERY

We now turn to our main experiment, exploring whether facets can be employed in Bridger to spur users to discover valuable and novel authors and their work. We use our two author-ranking strategies (§3.3), one based on similar tasks alone (sT) and the other on similar tasks with contrasting (distant) methods (sTdM). We compare these strategies to the Specter (ss) baseline. More specifically, we investigated the following research questions:

- RQ5: Do sT and sTdM, in comparison to Specter, surface suggestions of authors that are considered novel and valuable, coming from research communities more distant to the user?
- RQ6: Does sorting based on personas help users find more novel and valuable author suggestions?

5.1 Experiment Design

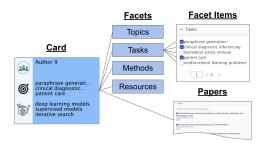


Figure 4: Illustration of information shown to users in Experiment II, §5. When the user clicks on an author card, an expanded view is displayed with 5 sections: papers, topics, and our extracted facets — tasks, methods, and resources.

Twenty computer-science researchers participated in the experiment approved for IRB exemption, after recruitment through Slack channels and mailing lists. Participants were compensated \$50 over PayPal for their time.

All participants were shown results based on their overall papers (without personas) consisting of 12 author cards they evaluated one by one. Four cards were included for each of sT, sTdM, and ss. We only show cards for authors who are at least 2 hops away in the co-authorship graph from the user, filtering authors with whom they had previously worked.

For participants who had at least two associated personas, we also presented them with authors suggested based on each separate persona: four author cards for each of their top two personas (two under sT and two under sTdM). Whether the participants saw the personas before or after the non-persona part was randomized.

Each author card provides a detailed depiction of that author (see Figure 2). The author's name and affiliation is hidden in this experiment to mitigate bias. As shown in Figure 4, cards showcase five sections of the author's research: their papers, MAG topics, and our extracted facet terms (i.e., tasks, methods, and resources). We

also let users view the tasks and methods ranked by *similarity* to them, which could be helpful to explain why an author was selected and better understand commonalities.

The cards showed up to five items for each section, with some sections having a second page, depending upon data availability, for a maximum of ten items per section. Papers could be sorted based on recency or similarity to a participant / persona. To avoid biasing participants, the only information provided for each paper was its title, date, and the suggested author's position on each paper (e.g., first, last).

Each of these items (papers and terms) had a checkbox, which the user was instructed to check if it fulfilled two criteria: 1) potentially interesting and valuable for them to learn about or consider in terms of utility, and 2) not too similar to things they had worked on or used previously. Following a short tutorial, 9 participants evaluated each author shown by checking the aforementioned checkboxes (see Figure 4, right). While evaluating the first and last author (randomized), the participant engaged in a protocol analysis methodology (sharing their thinking as they worked). Participants with personas were also asked, based on each persona's top five associated papers, whether they each reflected a coherent focus area, and whether they seemed useful for filtering author suggestions. 10

5.2 Quantitative Results

For each author card evaluated by a user, we calculate the ratio of checked boxes to total boxes (typically 5-10 for papers, 10 for topics and facets; see §5.1) in that card. Then, for each user, we calculate the average of these ratios per condition (sT, sTdM, ss), and calculate a user-level preference *S* specifying which of the three conditions received the highest average ratio. Using this score, we find the proportion of users who preferred each of the sT and sTdM conditions in comparison to ss. This metric indicates the user's preference between Bridger- and Specter-recommended authors in terms of novelty and value (**RQ5**).

Figure 5(a), shows results by this metric. The facet-based approaches lead to a boost over the non-faceted ss approach, with users overall preferring suggestions coming from the facet-based conditions. This is despite comparing against an advanced baseline geared at relevance, to which users are naturally primed.

We break down the results further by slightly modifying the metric to account for the different types of items users could check off. In particular, we distinguish between the task/method/resource/topic checkboxes, and the paper checkboxes. For each of these two groups, we compute *S* in the same way, ignoring all checkboxes that are not of that type (e.g., counting only papers). This breakdown reveals a more nuanced picture. For the task, method, resource and topic facets, the gap in favor of sT grows considerably (Figure 5b). In terms of papers only, ss, which was trained on aggregate paper-level information, achieves a marginally better outcome compared to sT, with a slightly larger gap in comparison to sTdM (Figure 5c). Aside from being trained on paper-level information, Specter also benefits from the fact that biases towards filter bubbles can be particularly strong with regard to papers. Unlike with facets, users

Item type	sT	sTdM
Overall	58%	75%
Paper	83%	67%
Topic	58%	75%
Task	42%	50%
Method	67%	58%
Resource	50%	67%

Table 1: Percentage of users with personas (N=12), for which the average proportion of checked items was higher for the persona-matched authors than for the overall-matched authors. Users saw suggested authors based on two of their personas. The suggestions came from either the sT or sTdM conditions. Reported here are counts of users who showed preference for one or both personas. "Overall" shows results for all checkboxes considered in aggregate; this is followed by results for individual item categories (papers and facets).

must tease apart aspects of papers that are new and interesting to them versus aspects that are relevant but familiar. See §5.4 for more discussion and concrete examples.

Importantly, despite obtaining better results overall with the faceted approach, we stress that our goal in this paper is not to "outperform" Specter, but mostly to use it as a reference point — a non-faceted approach used in a real-world academic search and recommendation setting. We also note that Specter and other existing alternative baselines we could use are not tuned to our task of building bridges for authors across filter bubbles.

Personas. We also compare the results from sT and sTdM conditions based on personas P for user \mathcal{A} , versus the user's non-personabased results presented above (**RQ6**). We compare the set of authors found using personas with authors retrieved without splitting into personas (equivalent to the union of all personas). Table 1 shows the number of users for which the average proportion of checked items was higher for the persona-matched authors than for the overall-matched authors (for at least one of the personas). For most participants, users signalled preference for persona-matched authors when looking at one or both of their personas. Interestingly, for papers we see a substantial boost in preference for both conditions, indicating that by focusing on more refined *slices* of the user's papers, we are able to gain better quality along this dimension too.

5.3 Evidence of Bursting Bubbles

The matched authors displayed to users were identified based either on sT and sTdM or the baseline Specter-based approach (ss). These two groups differed from each other substantially according to several empirical measures of similarity. We explore the following measures, based on author dimensions in our data that we do not use as part of the experiment: (1) Citation distance: A measure of distance in terms of citations that the user has in common with the matched author (Jaccard distance: 1 minus intersection-overunion). This is calculated both for incoming and outgoing citations. (2) Venue distance: The Jaccard distance between user and matched author for publication venues. (3) Coauthor shortest path: The shortest path length between the user and the matched author in

 $^{^9{\}rm The}$ tutorial slides are available in our supplementary materials.

¹⁰ See supplementary materials for the source code used for generating the data for Experiment II, as well as the code for the interactive application used in the evaluation, and the script used to direct the participants.

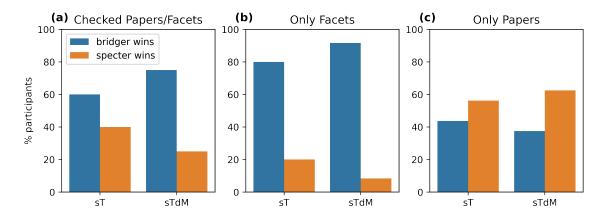


Figure 5: More users prefer Bridger for suggesting novel, interesting authors. Percent of the participants who preferred author suggestions surfaced by faceted conditions (sT and sTdM, blue bars) compared to a baseline non-faceted paper embedding (ss, orange bars). On average, users prefer the former suggestions, leading to more discovery of novel and valuable authors and their work (a). When broken down further, we find users substantially preferred the facet items shown for authors in our condition (b), and preferred the paper embedding baseline when evaluating papers (c). See §5 for discussion.

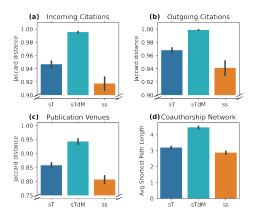


Figure 6: Bridger suggests authors that are more likely to bridge gaps between communities. In comparison to the baseline, facet-based (Bridger) author suggestions link users to broader areas. Clockwise: (a, b) Jaccard distance between suggested authors' papers and the user's papers for incoming citations (a) and outgoing citations (b); greater distance means that suggested authors are less likely to be cited by or cite the same work. (c) Jaccard distance for publication venues. (d) Shortest path length in the coauthorship graph between author and user (higher is more distant). Bridger conditions (sT and, especially, sTdM) show higher measures of distance.

the coauthorship graph. Findings of this analysis, shown in Figure 6, suggest that Bridger surfaces novel authors from more diverse, distant fields and research communities than Specter (RQ5).

5.4 Qualitative Findings: User Interviews

Experiment II's qualitative results were evaluated similarly to Experiment I (§4), using thematic analysis [4]. The interviews support

our quantitative results, affirming that Bridger authors encourage more diverse connections (RQ5). Under the Bridger conditions, participants noted diverse potentially useful research directions that connected their work to other authors not only within their own subareas, but also other areas. This was especially true under the sTdM condition. For instance, P9, who works on gradient descent for convex problems, saw a sTdM author's paper discussing gradient descent but for deep linear neural networks, which imply non-convex problems. They remarked, "This is a new setup. It's very different, and it's super important ... definitely something I would like to read ..." Considering a paper under a sTdM author, P6 observed an interesting contrast with their work: "I think my work has been bottom-up, so top-down would be an interesting approach to look at." As another example, P2 drew a connection between the mathematical area of graph theory and their area of human-AI decision-making under the sTdM condition: "This could be interesting mostly because ... they're using graph theory for decision making ... something I have not considered in the past." P19 remarked of an sTdM author's paper, "This one actually seems quite interesting because it seems like explicitly about trying to bridge the gap between computational neuroscience models, understanding the neocortex, and computing. So that seems like it's... going to actually chart the path for me between my work and the stuff I think about like artificial neural networks and machines."

In reacting to sTdM authors, many participants were able to go further than simply stating their interest in a connection and also describe *how* they would utilize the connection. Looking at a sTdM author, P6 explained how the author's neuroscience work could motivate work in their area of natural language processing: "I might learn from that [paper] how people compose words, and that might be inspiring for work on learning compositional representation ..." P18 checked off a paper titled "Multidisciplinary Collaboration to Facilitate Hypothesis Generation in Huntington's Disease" under a

sTdM author "because new ways to think about generating hypotheses could be interesting." Seeing the topic 'spike-timing-dependent plasticity' under a sTdM author, P19 mused, "I would like to understand how spike-timing-dependent plasticity works and whether that could lead to a better learning rule for other types of neural nets, like the ones I work with on language, so that seems fun." P12 described a sTdM author's paper about knowledge-driven search applications as useful to them because "One of my primary research areas is knowledge base completion. However, that's not an end application. An end application would be a search application which kind of uses my method to complete the knowledge base, and gives the user the end result. …" Though the sTdM condition presented more of a risk in terms of surfacing authors with which the user could draw connections, it also surfaced the more far-reaching connections.

The sT condition also helped participants ponder new connections, though perhaps not as distant. P8 said of a sT author's work, "I've worked a bit on summarization, so I want to know whether the approaches that I've worked on are applicable to real-time event summarization, which is a task I don't know about." Also reflecting on a sT author, P1 explained, "I've done a lot of work with micro tasks and these seem more maybe larger scale, like physical tasks or like planning travel.... There are so many problems ... that I could apply my techniques to." Other times, participants would connect one facet of their work to a different facet of the suggested author's work. In discussing a question-answering paper from a sT author, P8 explained, "I don't have experience with [the method] adversarial neural networks [used in this paper], but question answering is a task that I've worked on, so I would want to check this." Conversely, if participants found new connections with Specter, they tended to be more immediate connections to authors in their area. As an example, when checking off the paper "Efficient Symmetric Norm Regression via Linear Sketching" from a Specter-suggested author, P9 observed, "I have used sketching techniques and I have [also] used norm regression, but [on] this specific problem I have not." P9 also identified some of the papers from the suggested author as co-authored by their advisor.

Although participants were asked to only check off interesting papers that suggested something new for them to explore, biases towards filter bubbles can be particularly strong with regard to papers because users must tease apart papers' new and interesting aspects from their relevant but familiar aspects. Even if a paper is directly connected to a user's research, they may be tempted to check off a paper because they have not seen that exact paper or because it has minute differences from their work. For instance, P10 commented, "The problem is if it's so general a title, I assume there's something interesting happening, but I'm not completely sure." In contrast, when judging a particular facet item, participants need only contemplate the novelty of the term itself, without distraction or fixation on other terms [24, 26, 30]. As an example, P17 swiftly separated a task's general relevance from its lack of novelty to know not to check it. They explained, "'Scientific article summarization'-It is relevant, [but] I'm already familiar with it." This bias helps explain the overall preference for Specter when considering only papers (Figure 5(c)).

5.4.1 Personas. For the 12 participants who had personas, seven described their two personas as distinct, coherent identities that

would be useful for filtering author suggestions. As an example, P2 characterized their personas as related to "human-AI collaboration or decision-making" and "error analysis and machine learning debugging" respectively. The other 5 participants described one persona as coherent and seemingly useful for filtering authors. Concerns about their other personas were related to coherence, granularity, overlap with the other persona, and preference for the non-persona results after already looking through them and their first persona. Though the persona author suggestions performed relatively well in generating novel connections (Table 1), a few participants commented that they did not see the connection between suggested authors and their persona. For example, under a persona associated with lexical semantics, P6 commented on a sTdM paper, "'Causality' is not one of the topics that I would work on in lexical semantics."

6 DISCUSSION, LIMITATIONS & DESIGN IMPLICATIONS

We provide further analysis and discussion of our findings, including limitations of our proposed system and key challenges surfaced in the course of user interviews. We also discuss potential design implications for future author discovery systems.

6.1 Author Representations

6.1.1 Faceted Representation. Our experimental results point toward the advantages of a facet-based approach in the context of author discovery. We find that short, digestible items in the form of an author's tasks, methods, and resources can help participants consider interesting new research directions that they did not consider based on the author's papers alone. For instance, one participant (P14) expressed that a Bridger-suggested author's paper associated with medical image diagnosis would not be useful for them to consider because "breaking into that space for me would require a lot of work." However, when they later saw 'medical image diagnosis' as a task, they commented, "As a task, I could see some usefulness there. There could be other approaches that might more quickly catch my interest." Committing to interest in the overall task required much less effort. Moreover, participants were able to peruse more of an author's interesting tasks and methods that they did not necessarily find in their top papers. Reacting to one Bridger-suggested author, P3 did not see any papers related to 'biomedical question answering,' but they did see 'biomedical question answering system' as a method. They then noted, "I'm going to click 'biomedical question answering' because that's not what I have worked on before, but I'm interested in learning about it."

Our work in the computer science domain made use of *tasks*, *methods*, and *resources* — important functional aspects in this area [34]. Our results for Experiment I indicated that scientists find these more granular terms helpful both in describing their own work or the work of a known researcher, as well as for learning about unfamiliar researchers (§4.2.1, 4.2.3). This suggests that to extend our approach more broadly to other areas of science, categories of terms that have important semantic meaning in specific domains are required, as opposed to using generic keywords or considering text in aggregate. In biomedical research, for example, salient facets might include the drugs that researchers study, or the diseases

that their work addresses [25]. Alternatively, a future system could provide users with the ability to select or define which "author dimensions" matter to them, as opposed to assuming one universal pre-defined set.

6.1.2 Author Personas. Our results in experiments I and II suggest that directly accounting for authors' different lines of work, or personas, can help boost user satisfaction in discovery systems. In our work we focused on a specific notion of personas based on clustering authors' papers, but this can be extended and generalized. For example, we could allow users to more directly select themselves subsets of papers for which they want to find interests, allowing users to define their own lines of work. This could also help with challenges in automatically clustering papers (§ 5.4.1) by providing interactive feedback and supervision from users. Another important point for consideration is the temporal element, capturing author evolution over the years. Our current method simply looks 5 years back, but this could and should in principle be made adaptive per author, to account for different timelines for different authors. This too could potentially be made an interactive choice by the user, allowing them to segment their work temporally.

6.2 Challenges with Novel Information & Ideas

6.2.1 Novel Terminologies. An interesting design consideration that emerges in our experiments is that highly fine-grained terms (such as names of specific methods or datasets) can also introduce challenges in the context of discovery of novel, unfamiliar authors and their work. In particular, as discussed in Experiment I, while overall our faceted representation of authors was considered more useful for understanding the work of new authors, unfamiliar names of methods also hindered understanding.

This tradeoff was also reflected in interviews in our user study evaluating author recommendations (Experiment II). Participants commonly identified tasks, methods, and resources as interesting, even when they did not fully understand their meaning. When P4 saw the method 'least-general generalization of editing examples' from a Bridger-suggested author, they stated, "Don't know what this means exactly, but it sounds interesting." P13 marked their interest in the task "folksonomy-based recommender systems" under a Bridger-suggested author after commenting, "I'm curious [about folksonomy] simply because I'm ignorant." In seeing the resource 'synaptic resources' under a Bridger-suggested author, P19 simply said, "I'd like to know what that is." Nonetheless, many participants also struggled with indiscernible terms. For example, P20 said of the resource 'NAIST text corpus' under a Bridger-suggested author, "I'm not sure what this is, and I can't guess from the name. And it wasn't mentioned in the title of the papers." P2 explained that a paper did not "seem that interesting, but mostly because I don't understand all of these words." Thus, providing term definitions may be helpful. For additional context, multiple participants expressed interest in having abstracts available, and P15 suggested including automated summaries [7].

This problem of "unknown terms" encountered by our participants increases the effort required from users, potentially deterring users from considering certain authors/directions. An important line of future work will be addressing this problem by providing justin-time definitions of terms using extractive summarization [40] or

generative approaches [33]. In particular, an especially appealing idea is to develop methods for *personalized* explanations of new concepts, anchored in concepts with which the end user is already familiar (e.g., explaining a new neural network model by relating it to an older known one) [38].

6.2.2 Biases Toward Scientific Filter Bubbles. An important challenge reflected in our results is that of time constraints in the fast-moving world of research, inhibiting exploration beyond the filter bubble. Despite clear interest in an author's distant research, a couple of participants in Experiment II were hesitant to make connections. For example, in reacting to a Bridger-suggested author, P11 recognized, "There's just a bunch of really interesting kind of theory application papers in this list that I'm not familiar with. ... I would maybe scan a little bit of these, but it's so far off that it's harder to make room to read someone that far away, but still cool."

Unknown background knowledge can make it difficult to consider new areas. Engaging with distant authors' work requires a large cognitive load that can impede uncovering connections. As P18 in Experiment II noted: "Maybe there's some theoretical computer science algorithm that if I knew to apply it to my problem would speed things up or something like that, but I wouldn't know enough to recognize it as interesting." This further suggests that unfamiliar terms can especially hinder making interesting connections, and that a personalized system design that highlights the most useful aspects of a distant author's research may facilitate building far-reaching connections. This also further compounds aversion to novel ideas, and fixation on familiar frames of problems and solutions [9, 18, 24, 26, 30]. Because Bridger's authors are selected to be more distant from the user than Specter's authors, they sometimes met with hard-line resistance, without full consideration of potential links. Looking at a Bridger-suggested author, the natural language processing (NLP) researcher P20 said, "This is not really an NLP paper, so I would pass." Similarly, P17 rejected a paper from a Bridger suggestion, saying "I don't know anything about neuroscience, and I'm not going to start now probably."

6.3 Data & System

One limitation of our work is that many early-career researchers are excluded from our system because they do not have enough papers. This user group, however, could potentially especially benefit from this type of discovery system. This problem relates more broadly to the much-discussed "cold start" challenge in recommendation systems [3, 32]. One potential implication for future systems is to provide such users with alternative options, such as viewing recommendations to authors they consider relevant (e.g., advisors, mentors, etc.).

Another limitation lies in the accuracy of the information extraction methods we employ. This is an issue that impacts all work that relies on such methods; however, it is somewhat mitigated in our setting by aggregating over many spans extracted from authors' papers, which averages out some noise. Additionally, other work has found that even with moderate extraction accuracy, strong ideation utility can be obtained [26]. The problem of identifying connections between mentions of tasks, methods or resources across scientific papers is very much an open one [8]; as future models in this

area become more precise, our approach for matching authors is expected to become more accurate, too.

Finally, our evaluation focused primarily on surfacing authors who spark new ideas outside users' familiar areas. We design controlled studies measuring our approach's ability to surface inspirations in comparison to a real-world baseline scientific search model. An interesting and challenging direction for further evaluation is to measure Bridger's longer-term ability to provide useful inspirations that yield viable research directions and projects, and to measure the system's ability to help users in less controlled settings. This type of evaluation would require longer-term interaction with the user, and longitudinal observations.

7 CONCLUSION

We presented Bridger, a framework for facilitating discovery of novel and valuable scholars and their work. Bridger consists of a faceted author representation, allowing users to see authors who match them along certain dimensions (e.g., tasks) but not others. Bridger also provides "slices" of a user's papers, enabling them to find authors who match the user only on a subset of their papers, and only on certain facets within those papers. Our experiments with computer science researchers show that the facet-based approach was able to help users discover authors with work that is considered more interesting and novel, substantially more than a relevance-focused baseline representing state-of-art retrieval of scientific papers. Importantly, we show that authors surfaced by Bridger are indeed from more distant communities in terms of publication venues, citation links and co-authorship social ties. While our work only considers the domain of computer science research, we believe the techniques could generalize outside of computer science, potentially connecting people with ideas from even more disparate fields as we make steps toward bridging gaps across all of science. These results suggest a new and potentially promising avenue for mitigating the problem of isolated silos in science.

ACKNOWLEDGMENTS

We thank Kishore Vasan, Jonathan Borchardt, and Jonathan Bragg for their contributions to this project. We also thank our study participants, and the four anonymous reviewers for their helpful suggestions. This work is partially supported by NSF Grant OIA-2033558, NSF RAPID grant 2040196, and ONR grant N00014- 18-1-2193.

REFERENCES

- Jöran Beel and Bela Gipp. 2009. Google Scholar's ranking algorithm: an introductory overview. In Proceedings of the 12th international conference on scientometrics and informetrics (ISSI'09), Vol. 1. Rio de Janeiro (Brazil), 230–241.
- [2] Joeran Beel, Bela Gipp, Stefan Langer, and Corinna Breitinger. 2016. Paper recommender systems: a literature survey. *International Journal on Digital Libraries* 17, 4 (2016), 305–338.
- [3] Jesús Bobadilla, Fernando Ortega, Antonio Hernando, and Jesús Bernal. 2012. A collaborative filtering approach to mitigate the new user cold start problem. Knowledge-Based Systems 26 (Feb. 2012), 225–238. https://doi.org/10.1016/j. knosvs.2011.07.021
- [4] Virginia Braun and Victoria Clarke. 2012. Thematic analysis. (2012).
- [5] Ronald S. Burt. [n.d.]. Structural Holes and Good Ideas. 110, 2 ([n.d.]), 349–399. https://doi.org/10.1086/421787
- [6] Katy Börner, Chaomei Chen, and Kevin W. Boyack. 2005. Visualizing knowledge domains. Annual Review of Information Science and Technology 37, 1 (Jan. 2005), 179–255. https://doi.org/10.1002/aris.1440370106
- [7] Isabel Cachola, Kyle Lo, Arman Cohan, and Daniel S Weld. 2020. TLDR: Extreme Summarization of Scientific Documents. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings. 4766–4777.
- [8] Arie Cattan, Sophie Johnson, Daniel Weld, Ido Dagan, Iz Beltagy, Doug Downey, and Tom Hope. 2021. SciCo: Hierarchical Cross-Document Coreference for Scientific Concepts. arXiv preprint arXiv:2104.08809 (2021).
- [9] Joel Chan, Joseph Chee Chang, Tom Hope, Dafna Shahaf, and Aniket Kittur. 2018. Solvent: A mixed initiative system for finding analogies between research papers. Proceedings of the ACM on Human-Computer Interaction 2, CSCW (2018), 1–21.
- [10] Chaomei Chen. 2017. Expert Review. Science Mapping: A Systematic Review of the Literature. Journal of Data and Information Science 2, 2 (2017), 1–40. https://doi.org/10.1515/jdis-2017-0006 00001.
- [11] Li Chen, Yonghua Yang, Ningxia Wang, Keping Yang, and Quan Yuan. 2019. How serendipity improves user satisfaction with recommendations? a large-scale user evaluation. In The World Wide Web Conference. 240–250.
- [12] Wanyu Chen, Pengjie Ren, Fei Cai, Fei Sun, and Maarten de Rijke. 2020. Improving end-to-end sequential recommendations with intent-aware diversification. In Proceedings of the 29th ACM International Conference on Information & Knowledge Management. 175–184.
- [13] M.j. Cobo, A.g. López-Herrera, E. Herrera-Viedma, and F. Herrera. 2011. Science mapping software tools: Review, analysis, and cooperative study among tools. *Journal of the American Society for Information Science and Technology* 62, 7 (July 2011), 1382–1402. https://doi.org/10.1002/asi.21525
- [14] Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel Weld. 2020. SPECTER: Document-level Representation Learning using Citation-informed Transformers. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, Online, 2270–2282. https://doi.org/10.18653/v1/2020.acl-main.207
- [15] Marian Dörk, Nathalie Henry Riche, Gonzalo Ramos, and Susan Dumais. 2012. PivotPaths: Strolling through faceted information spaces. Visualization and Computer Graphics, IEEE Transactions on 18, 12 (2012), 2709–2718. http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=6327277
- [16] Alessandro Epasto, Silvio Lattanzi, and Renato Paes Leme. 2017. Ego-Splitting Framework: from Non-Overlapping to Overlapping Clusters. In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '17. ACM Press, Halifax, NS, Canada, 145–154. https://doi.org/10. 1145/3097983.3098054
- [17] Dieter Frey. 1986. Recent research on selective exposure to information. Advances in experimental social psychology 19 (1986), 41–80.
- [18] Katherine Fu, Joel Chan, Jonathan Cagan, Kenneth Kotovsky, Christian Schunn, and Kristin Wood. 2013. The Meaning of Near and Far: The Impact of Structuring Design Databases and the Effect of Distance of Analogy on Design Output. JMD (2013)
- [19] Yingqiang Ge, Shuya Zhao, Honglu Zhou, Changhua Pei, Fei Sun, Wenwu Ou, and Yongfeng Zhang. 2020. Understanding echo chambers in e-commence recommender systems. In Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. 2261–2270.
- [20] Kosa Goucher-Lambert, Joshua T Gyory, Kenneth Kotovsky, and Jonathan Cagan. 2020. Adaptive Inspirational Design Stimuli: Using Design Output to Computationally Search for Stimuli That Impact Concept Generation. Journal of Mechanical Design 142, 9 (2020).
- [21] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. In ACL. https://doi.org/10.18653/v1/2020.acl-main 740
- [22] Florian Heimerl, Qi Han, Steffen Koch, and Thomas Ertl. 2016. CiteRivers: Visual Analytics of Citation Patterns. IEEE Transactions on Visualization and Computer Graphics 22, 1 (Jan. 2016), 190–199. https://doi.org/10.1109/TVCG.2015.2467621 Conference Name: IEEE Transactions on Visualization and Computer Graphics.

- [23] Tom Hope, Aida Amini, David Wadden, Madeleine van Zuylen, Sravanthi Parasa, Eric Horvitz, Daniel Weld, Roy Schwartz, and Hannaneh Hajishirzi. 2021. Extracting a Knowledge Base of Mechanisms from COVID-19 Papers. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, Online, 4489–4503. https://doi.org/10.18653/v1/2021.naacl-main.355
- [24] Tom Hope, Joel Chan, Aniket Kittur, and Dafna Shahaf. [n.d.]. Accelerating Innovation Through Analogy Mining. In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (New York, NY, USA, 2017-08-04) (KDD '17). Association for Computing Machinery, 235–243. https://doi.org/10.1145/3097983.3098038
- [25] Tom Hope, Jason Portenoy, Kishore Vasan, Jonathan Borchardt, Eric Horvitz, Daniel S Weld, Marti A Hearst, and Jevin West. 2020. SciSight: Combining faceted navigation and research group detection for COVID-19 exploratory scientific search. In EMNLP.
- [26] Tom Hope, Ronen Tamari, Hyeonsu Kang, Daniel Hershcovich, Joel Chan, Aniket Kittur, and Dafna Shahaf. 2021. Scaling Creative Inspiration with Fine-Grained Functional Facets of Product Ideas. arXiv e-prints (2021), arXiv-2102.
- [27] Sanjay Kairam, Nathalie Henry Riche, Steven Drucker, Roland Fernandez, and Jeffrey Heer. 2015. Refinery: Visual Exploration of Large, Heterogeneous Networks through Associative Browsing. Computer Graphics Forum (Proc. EuroVis) 34, 3 (2015). http://idl.cs.washington.edu/papers/refinery
- [28] Marius Kaminskas and Derek Bridge. 2016. Diversity, serendipity, novelty, and coverage: a survey and empirical analysis of beyond-accuracy objectives in recommender systems. ACM Transactions on Interactive Intelligent Systems (TiiS) 7, 1 (2016), 1–42.
- [29] Lanu Kim, Jevin D West, and Katherine Stovel. 2017. Echo Chambers in Science?. In American Sociological Association.
- [30] Aniket Kittur, Lixiu Yu, Tom Hope, Joel Chan, Hila Lifshitz-Assaf, Karni Gilon, Felicia Ng, Robert E. Kraut, and Dafna Shahaf. [n.d.]. Scaling up analogical innovation with crowds and AI. 116, 6 ([n.d.]), 1870–1877. https://doi.org/10. 1073/pnas.1807185116 Publisher: National Academy of Sciences Section: Social Sciences.
- [31] Joel Klinger, Juan Mateos-Garcia, and Konstantinos Stathoulopoulos. 2020. A narrowing of AI research? arXiv preprint arXiv:2009.10385 (2020).
- [32] Xuan Nhat Lam, Thuc Vu, Trong Duc Le, and Anh Duc Duong. 2008. Addressing cold-start problem in recommendation systems. In Proceedings of the 2nd international conference on Ubiquitous information management and communication (ICUIMC '08). Association for Computing Machinery, New York, NY, USA, 208–211. https://doi.org/10.1145/1352793.1352837
- [33] Y. Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, M. Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. ArXiv abs/1907.11692 (2019).
- [34] Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-Task Identification of Entities, Relations, and Coreference for Scientific Knowledge Graph Construction. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, Brussels, Belgium. https://doi.org/10.18653/v1/D18-1360
- [35] Miller McPherson, Lynn Smith-Lovin, and James M Cook. 2001. Birds of a feather: Homophily in social networks. Annual review of sociology 27, 1 (2001), 415–444.
- [36] Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing Order into Text. In Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, Barcelona, Spain, 404–411. https://www.aclweb.org/anthology/W04-3252
- [37] Fionn Murtagh and Pierre Legendre. 2014. Ward's hierarchical agglomerative clustering method: which algorithms implement Ward's criterion? Journal of classification 31, 3 (2014), 274–295.
- [38] Sonia K. Murthy, Daniel King, Tom Hope, Daniel Weld, and Doug Downey. 2021. Towards personalized descriptions of scientific concepts. In The Fifth Widening Natural Language Processing Workshop at EMNLP.
- [39] Rahul Nadkarni, David Wadden, Iz Beltagy, Noah A Smith, Hannaneh Hajishirzi, and Tom Hope. 2021. Scientific Language Models for Biomedical Knowledge Base Completion: An Empirical Study. AKBC (2021).
- [40] Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Ranking Sentences for Extractive Summarization with Reinforcement Learning. In NAACL-HLT.
- [41] Arpit Narechania, Alireza Karduni, Ryan Wesslen, and Emily Wall. 2021. vitaLITy: Promoting Serendipitous Discovery of Academic Literature with Transformers & Visual Analytics. IEEE Transactions on Visualization and Computer Graphics (2021), 1–1. https://doi.org/10.1109/TVCG.2021.3114820 Conference Name: IEEE Transactions on Visualization and Computer Graphics.
- [42] Mark Neumann, Daniel King, Iz Beltagy, and Waleed Ammar. 2019. ScispaCy: Fast and Robust Models for Biomedical Natural Language Processing. undefined (2019). /paper/ScispaCy%3A-Fast-and-Robust-Models-for-Biomedical-Neumann-King/de28ec1d7bd38c8fc4e8ac59b6133800818b4e29
- [43] Tien T Nguyen, Pik-Mai Hui, F Maxwell Harper, Loren Terveen, and Joseph A Konstan. 2014. Exploring the filter bubble: the effect of using recommender systems on content diversity. In Proceedings of the 23rd international conference on World wide web. 677–686.

- [44] Mathias Wullum Nielsen and Jens Peter Andersen. 2021. Global citation inequality is on the rise. Proceedings of the National Academy of Sciences 118, 7 (2021).
- [45] Eli Pariser. 2011. The filter bubble: What the Internet is hiding from you. Penguin UK.
- [46] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in Python. the Journal of machine Learning research 12 (2011), 2825–2830.
- [47] Jason Portenoy, Jessica Hullman, and Jevin D. West. 2017. Leveraging Citation Networks to Visualize Scholarly Influence Over Time. Frontiers in Research Metrics and Analytics 2 (Nov. 2017), 8. https://doi.org/10.3389/frma.2017.00008
- [48] Jason Portenoy and Jevin D West. 2020. Constructing and evaluating automated literature review systems. Scientometrics 125 (2020), 3233–3251.
- [49] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In EMNLP/IJCNLP. https://doi.org/10.18653/v1/ D19-1410
- [50] Ariel S. Schwartz and Marti A. Hearst. 2002. A simple algorithm for identifying abbreviation definitions in biomedical text. In *Biocomputing 2003*. WORLD SCIENTIFIC, 451–462. https://doi.org/10.1142/9789812776303_0042
- [51] Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June (Paul) Hsu, and Kuansan Wang. [n.d.]. An Overview of Microsoft Academic Service (MAS) and Applications. ACM Press, 243–246. https://doi.org/10.1145/2740908.2742839
- [52] Shivashankar Subramanian, Daniel King, Doug Downey, and Sergey Feldman. 2021. S2AND: A Benchmark and Evaluation System for Author Name Disambiguation. arXiv preprint arXiv:2103.07534 (2021).
- [53] Don R Swanson and Neil R Smalheiser. 1996. Undiscovered Public Knowledge: A Ten-Year Update.. In KDD. 295–298.
- [54] Jie Tang, Sen Wu, Jimeng Sun, and Hang Su. 2012. Cross-domain collaboration recommendation. In Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining. 1285–1293.
- [55] Xuli Tang, Xin Li, Ying Ding, Min Song, and Yi Bu. 2020. The pace of artificial intelligence innovations: Speed, talent, and trial-and-error. *Journal of Informetrics* 14, 4 (2020), 101094.
- [56] Chun-Hua Tsai and Peter Brusilovsky. 2018. Beyond the ranked list: User-driven exploration and diversification of social recommendation. In 23rd international conference on intelligent user interfaces. 239–250.
- [57] Chun-Hua Tsai, Jukka Huhtamäki, Thomas Olsson, and Peter Brusilovsky. 2020. Diversity Exposure in Social Recommender Systems: A Social Capital Theory Perspective. work 5, 11 (2020), 22.
- [58] Daril Vilhena, Jacob Foster, Martin Rosvall, Jevin West, James Evans, and Carl Bergstrom. [n.d.]. Finding Cultural Holes: How Structure and Culture Diverge in Networks of Scholarly Communication. 1 ([n.d.]), 221–238. https://doi.org/10. 15195/v1.a15
- [59] David Wadden, Ulme Wennberg, Yi Luan, and Hannaneh Hajishirzi. 2019. Entity, Relation, and Event Extraction with Contextualized Span Representations. In EMNLP/IJCNLP. https://doi.org/10.18653/v1/D19-1585
- [60] Huaiyu Wan, Yutao Zhang, Jing Zhang, and Jie Tang. 2019. Aminer: Search and mining of academic social networks. Data Intelligence 1, 1 (2019), 58–76.
- [61] Kuansan Wang, Zhihong Shen, Chiyuan Huang, Chieh-Han Wu, Darrin Eide, Yuxiao Dong, Junjie Qian, Anshul Kanakia, Alvin Chen, and Richard Rogahn. 2019. A Review of Microsoft Academic Services for Science of Science Studies. Frontiers in Big Data 2 (2019). https://doi.org/10.3389/fdata.2019.00045 Publisher: Frontiers.
- [62] Ningxia Wang, Li Chen, and Yonghua Yang. 2020. The Impacts of Item Features and User Characteristics on Users' Perceived Serendipity of Recommendations. In Proceedings of the 28th ACM Conference on User Modeling, Adaptation and Personalization. 266–274.
- [63] Wei Wang, Jiaying Liu, Zhuo Yang, Xiangjie Kong, and Feng Xia. 2019. Sustainable collaborator recommendation based on conference closure. IEEE Transactions on Computational Social Systems 6, 2 (2019), 311–322.
- [64] Jevin D West, Ian Wesley-Smith, and Carl T Bergstrom. 2016. A recommendation system based on hierarchical clustering of an article-level citation network. *IEEE Transactions on Big Data* 2, 2 (June 2016), 113–123. https://doi.org/10.1109/ TBDATA.2016.2541167
- [65] Mark Wilhelm, Ajith Ramanathan, Alexander Bonomo, Sagar Jain, Ed H Chi, and Jennifer Gillenwater. 2018. Practical diversified recommendations on youtube with determinantal point processes. In Proceedings of the 27th ACM International Conference on Information and Knowledge Management. 2165–2173.
- [66] Pengfei Zhao and Dik Lun Lee. 2016. How much novelty is relevant? it depends on your curiosity. In Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval. 315–324.
- [67] Ziwei Zhu, Jianling Wang, and James Caverlee. 2020. Measuring and Mitigating Item Under-Recommendation Bias in Personalized Ranking Systems. In Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. 449–458.