
Joint Content-Context Analysis of Scientific Publications: Identifying Opportunities for Collaboration in Cognitive Science

Lu Cheng*

University of California, Los Angeles
Los Angeles, CA 90095
lucheng@g.ucla.edu

Girish Ganesan

Rutgers University
New Brunswick, NJ 08901
gg655@rutgers.edu

William He

Northwestern University
Evanston, IL 60208
WilliamHe@u.northwestern.edu

Daniel Silverston

Brown University
Providence, RI 02912
daniel_silverston@brown.edu

Harlin Lee

Department of Mathematics
University of California, Los Angeles
Los Angeles, CA 90095
harlin@math.ucla.edu

Jacob G. Foster

Department of Sociology
University of California, Los Angeles
Los Angeles, CA 90095
foster@soc.ucla.edu

Abstract

This work studies publications in the field of cognitive science and utilizes mathematical techniques to connect the analysis of the papers' content (abstracts) to the context (citation, journals). We apply hierarchical topic modeling on the abstracts and community detection algorithms on the citation network, and measure content-context discrepancy to find academic fields that study similar topics but do not cite each other or publish in the same venues. These results show a promising, systemic framework to identify opportunities for scientific collaboration in highly interdisciplinary fields such as cognitive science and machine learning.

1 Introduction

As scientific fields have grown larger and more specialized, researchers may be missing potentially-lucrative avenues of collaboration. For example, researchers may be pursuing similar paths in parallel while lacking a common language and literary academic foundation to connect their works. Uncovering such situations will enable more productive, coordinated research efforts, which is one of the principal goals of science of science [2, 6, 7, 8].

Science of science, or metascience, is the branch of science that uses quantitative measurements and scientific techniques to understand the interactions between scientific agents with the aim to refine and improve scientific practices and progress [4]. Yet currently, most metascience studies have focused on investigating either the content or context of research in relation to other publications without bridging the gap between them [3]. In this paper, we investigate the field of cognitive science through the twin lenses of content and context; information is extracted from both 1) paper abstracts through natural language processing (NLP) and 2) the citation network via graph community detection techniques.

*The first four undergraduate student authors are listed alphabetically.

We then propose a simple but effective criteria to determine which subdivisions within cognitive science are similar in content but not in context, and suggest what barriers may lie between them.

We focus on cognitive science, in part because it has been claimed that cognitive science has failed to achieve its intention of integrating the six disciplines of which it was to be comprised (psychology, linguistics, artificial intelligence, anthropology, philosophy and neuroscience) [10]. Hence, it will be revealing to discover which interdisciplinary connections are missing in the field and investigate how this gap could be filled. Beyond cognitive science, our approach and methods can provide a framework for the joint study of content and context in other interdisciplinary fields such as applied mathematics and machine learning.

2 Methods

We introduce NLP and graph methods that were used to analyze the publications dataset, as well as metrics used to quantify cluster similarities. Please see Appendix A for data acquisition and preprocessing details. Python code is available at <https://github.com/HarlinLee/cogsci-missed-connections>.

Dataset We used 59,384 papers in the field of “cognitive science” from the Microsoft Academic Graph [13], where the field tags of a paper are identified from its text and sometimes citations [12]. The papers are assigned a unique ID and include metadata such as title, author(s), journal of publication, year of publication, abstract text, and references.

Content Analysis We construct the word-by-abstract matrix \mathbf{X} using the bag-of-words model and tf-idf weighting, and apply Hierarchical Non-Negative Matrix Factorization (Hierarchical NMF) [5] to detect topics. NMF [9] approximates $\mathbf{X} \approx \mathbf{WH}$, where the dictionary matrix \mathbf{W} and the coding matrix \mathbf{H} are two rank- r non-negative matrices. The i th column of \mathbf{W} gives the weights of the words in the i th topic, while the j th column of \mathbf{H} gives the weights of the topics in the j th abstract. This allows us to represent a topic as a combination of words, and an abstract as a combination of topics. We describe each topic using its top three weighted words, and assign each paper to its most weighted topic. Next, we column-wise split \mathbf{X} into r sub-matrices, $\mathbf{X}_1^{(1)}, \dots, \mathbf{X}_1^{(r)}$, such that columns of $\mathbf{X}_1^{(i)}$ correspond to abstracts assigned to the i th topic. Then we perform NMF on each sub-matrix to obtain *subtopics*.

Context Analysis After assigning papers to nodes and citations between those papers to edges, our citation data yields a graph with 59,384 nodes and 191,871 directed edges. We then isolate the largest weakly-connected component, which leaves us with 41,465 nodes (69.8% of original papers) and 190,997 edges (99.5% of original citations). We employ Degree-Discounted Symmetrization [11] to get the degree-discounted and symmetric adjacency matrix, and use Louvain’s Algorithm [1] to find a community scheme that maximizes the modularity of the final graph.

Content-Context Discrepancy Let c_i be the i th largest community of publications in the citation network. We measure topic similarity $T(c_i, c_j)$ and journal similarity $J(c_i, c_j)$ as proxies for content and context similarity, respectively. Then, we calculate the discrepancy $\rho(c_i, c_j)$ and use these metrics to identify communities that are more similar in content than they are in context.

Recall that every paper in c_i is assigned to an NMF topic, and has its journal of publication known. Let \mathbf{t}_i be the frequency distribution of the topics of the papers in c_i . Similarly, \mathbf{p}_i is the frequency distribution of journals that the papers in c_i were published in. Normalize them by $\hat{\mathbf{t}}_i = \mathbf{t}_i / \|\mathbf{t}_i\|_2$, $\hat{\mathbf{p}}_i = \mathbf{p}_i / \|\mathbf{p}_i\|_2$, and define the similarity metrics as their dot product. Our proposed discrepancy index combines these two metrics such that topic similarity is considered more heavily:

$$T(c_i, c_j) = \langle \hat{\mathbf{t}}_i, \hat{\mathbf{t}}_j \rangle, J(c_i, c_j) = \langle \hat{\mathbf{p}}_i, \hat{\mathbf{p}}_j \rangle, \rho(c_i, c_j) = T(c_i, c_j) - J(c_i, c_j)/2. \quad (1)$$

3 Results and Discussion

We display topic modeling and community detection results on the publications dataset, and discuss how it may relate to missed opportunities for scientific collaboration in cognitive science.

3.1 Hierarchical Topics in Cognitive Science



Figure 1: Hierarchical topics of cognitive science according to paper abstracts. Labels are the topics’ keywords, and wedge size is proportional to number of papers in the topic.

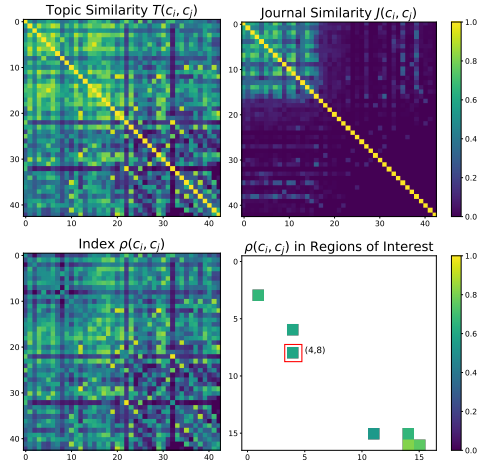


Figure 2: Heatmaps of metrics T , J and ρ . Axes are indices i, j .

Figure 1 presents the hierarchical topics extracted from abstracts. The inner circle contains 15 NMF topics, and each topic is further split into 8 or 10 subtopics in the outer circle. Some keywords suggest connections to known fields of cognitive science. For example, “language, linguistic, communication” \sim linguistics, “human, social, behavior” \sim anthropology, and “consciousness, conscious, mind” \sim philosophy. It is notable that neither “computer science” nor “psychology” seem to exist as keywords to a main topic even though they are claimed to dominate the field of cognitive science in [10]. A hypothesis is that as those fields have become so broad and popular, researchers avoid those terms and instead use specific subtopics or methods under the field to describe their work. Alternatively, these fields could be so prevalent and diffused within cognitive science that they would not appear as a distinct topic.

3.2 Content-Context Discrepancy Criteria

After uncovering 15 topics in the abstracts and 43 communities in the citation network, we examined and visualized in Figure 2 the metrics $T(c_i, c_j)$ (top left), $J(c_i, c_j)$ (top right), and $\rho(c_i, c_j)$ (bottom left). The color of each pixel represents the metric value for the pair of publication clusters. Note that $J(c_i, c_j)$ drops significantly at $i, j = 17$. The sample space of journal distribution in this dataset is large, but many communities are very small, often with merely tens of papers; see Figure 4. This means the journal distribution vectors are necessarily sparse, leading to a flawed comparison between smaller communities. Therefore, we limit our analysis to the 17 largest communities and compare only those close to each other in size to minimize other size effects.

We use the following criteria to identify regions of interest, i.e. communities in cognitive science that may discuss similar themes but do not cite each other or publish in the same venues:

- Similar topic distribution: $T(c_i, c_j) > 0.75$
- Dis-similar journal distribution: $J(c_i, c_j) < 0.5$
- High discrepancy: $\rho(c_i, c_j) > 0.5$
- Similar size: $|i - j| \leq 5$
- Large enough size: $i, j \leq 16$.

The bottom right of Figure 2 marks the 7 identified pairs, which we can then examine in detail.

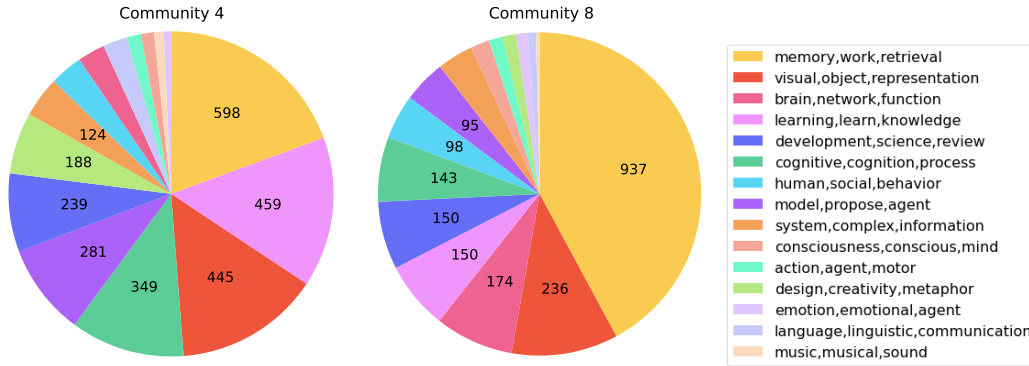


Figure 3: Topic distributions in communities 4 and 8. Wedge labels are numbers of papers in the topic. Legend shows keywords.

3.3 Case Study on Communities 4 and 8

Communities 4 and 8 (boxed in red in Figure 2 bottom right) yielded $T(c_4, c_8) = 0.826$, $J(c_4, c_8) = 0.479$, and $\rho(c_4, c_8) = 0.586$. According to the pie charts in Figure 3, the two communities have a very similar topic composition—both are a mix of “memory” + “visual” + “learning”. At the same time, the fact that they are split into two graph communities indicates that they are not very connected in the citation network. In fact, there are approximately 15,000 intra-community edges in these two communities, and only 800 inter-community edges. Furthermore, we find very little overlap in the top 10 published-in journal sets in these communities.

Community 4		Community 8	
Advances in Psychology	78	Trends in Cognitive Sciences	84
Memory & Cognition	66	Behavioral and Brain Sciences	50
Journal of Experimental Psychology	63	BiorXiv	47
Applied Cognitive Psychology	61	Frontiers in Human Neuroscience	35
Educational Psychologist	52	Neuropsychologia	34
Educational Psychology Review	44	Journal of Cognitive Neuroscience	33
Psychology of Learning and Motivation	43	Neuron	30
Journal of Educational Psychology	35	Current Biology	30
Psychonomic Bulletin & Review	34	Neuroscience & Biobehavioral Reviews	30
Memory	32	Memory	29

Community 4 is mostly published in (educational) psychology journals, whereas community 8 is associated with neuroscience journals. Clearly, there is a citational and academic disconnect between them, even though they share similar topic distributions. Initiating conversation between them could help further our understanding of complex subjects like memory, as it can provide a more holistic view of the theme, and even inspire fresh research questions and methods.

4 Conclusions and Future Work

We outlined a method that connects the analysis of the content and context of scientific papers in cognitive science. We extracted topics from paper abstracts using hierarchical NMF, detected communities in the citation network, and analyzed their journal publication distributions. Combining these approaches allowed us to find groups that are close in content but not in context, which indicate potential opportunities for collaboration.

We plan to apply this framework to particularly entangled fields such as artificial intelligence and machine learning, and add a temporal dimension to our analysis. For example, can we recognize changes in citation network and prominent topics over time? Can we detect shifts in rhetoric and composition? Another direction is to examine the connection between content information and the citation network structure directly. If a link prediction model trained on the text accurately predicts citation links between papers, this would be evidence of interdependence between two forms of data.

Acknowledgments and Disclosure of Funding

The initial research was conducted at the UCLA Computational Applied Mathematics (CAM) Research Experiences for Undergraduates (REU), June 14-August 4, 2021. This work was supported by grant TWCF0333 from the Templeton World Charity Foundation.

References

- [1] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008, 2008.
- [2] Hung-Hsuan Chen, Liang Gou, Xiaolong Zhang, and Clyde Lee Giles. Collabseer: A search engine for collaboration discovery. In *Proceedings of the 11th Annual International ACM/IEEE Joint Conference on Digital Libraries*, JCDL '11, page 231–240, New York, NY, USA, 2011. Association for Computing Machinery.
- [3] James Evans and Jacob Foster. Metaknowledge. *Science (New York, N.Y.)*, 331:721–5, 02 2011.
- [4] Santo Fortunato, Carl T. Bergstrom, Katy Börner, James A. Evans, Dirk Helbing, Staša Milojević, Alexander M. Petersen, Filippo Radicchi, Roberta Sinatra, Brian Uzzi, Alessandro Vespignani, Ludo Waltman, Dashun Wang, and Albert-László Barabási. Science of science. *Science*, 359(6379):eaa0185, March 2018.
- [5] Rachel Grotheer, Yihuan Huang, Pengyu Li, Elizaveta Rebrova, Deanna Needell, Longxiu Huang, Alona Kryshchenko, Xia Li, Kyung Ha, and Oleksandr Kryshchenko. Covid-19 literature topic-based search via hierarchical nmf. *arXiv preprint arXiv:2009.09074*, 2020.
- [6] Jian Huang, Ziming Zhuang, Jia Li, and C. Lee Giles. Collaboration over time: Characterizing and modeling network evolution. In *Proceedings of the 2008 International Conference on Web Search and Data Mining*, WSDM '08, page 107–116, New York, NY, USA, 2008. Association for Computing Machinery.
- [7] Geraldo J Pessoa Junior, Thiago MR Dias, Thiago HP Silva, and Alberto HF Laender. On interdisciplinary collaborations in scientific coauthorship networks: the case of the brazilian community. *Scientometrics*, 124(3):2341–2360, 2020.
- [8] Mario Krenn and Anton Zeilinger. Predicting research trends with semantic and neural networks with an application in quantum physics. *Proceedings of the National Academy of Sciences*, 117(4):1910–1916, 2020.
- [9] Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
- [10] Rafael Núñez, Michael Allen, Richard Gao, Carson Miller Rigoli, Josephine Relaford-Doyle, and Arturs Semenuks. What happened to cognitive science? *Nature Human Behaviour*, 3(8):782–791, August 2019.
- [11] Venu Satuluri and Srinivasan Parthasarathy. Symmetrizations for clustering directed graphs. In *Proceedings of the 14th International Conference on Extending Database Technology*, pages 343–354, 2011.
- [12] Zhihong Shen, Hao Ma, and Kuansan Wang. A web-scale system for scientific knowledge exploration. In *Proceedings of ACL 2018, System Demonstrations*, pages 87–92, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [13] Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June Hsu, and Kuansan Wang. An overview of microsoft academic service (mas) and applications. In *Proceedings of the 24th international conference on world wide web*, pages 243–246, 2015.

A Data Acquisition and Preprocessing Details

Data Acquisition A total of 258,039 papers in the field “cognitive science” were obtained from the Microsoft Academic Graph [13], where the field tags of a paper are identified from its text and sometimes citations [12]. The papers are also given probabilities of being “important” as determined by [12, 13]. In addition, each paper is assigned a unique ID and include metadata such as title, author(s), journal of publication, year of publication, abstract text, and references.

First, we discard 58,039 papers with the lowest probabilities of being “important” because 1) $\sim 0\%$ of them have abstracts, 2) $\sim 0\%$ have references, 3) none are published in recent years, and 4) the probability is significantly lower than the rest. We then remove papers published prior to 1950 in order to limit the scope to the modern notion of cognitive science from the 1950s [10]. Next, we keep only the papers that contain references, and whose abstracts are between 30 and 500 words long. We found that many exceedingly short abstracts are actually titles and publication information, while exceedingly long abstracts tend to contain extraneous text such as table of contents or the text of the entire first page of the paper. Finally, after removing all papers with duplicate abstracts, we have a dataset of 59,384 papers for analysis.

Bag-of-Words Matrix Construction We first lemmatize the abstracts; remove numbers and punctuation; remove English stop words, and stop words specific to abstracts (e.g. “et al”, “this paper”). We then construct the data matrix using the bag-of-words model and term frequency-inverse document frequency (tf-idf) weighting, including tri-grams and excluding words that appear in more than 80% or less than 0.05% of abstracts. This yields a word-by-abstract matrix \mathbf{X} of size $9,106 \times 59,384$.

B Additional Figures

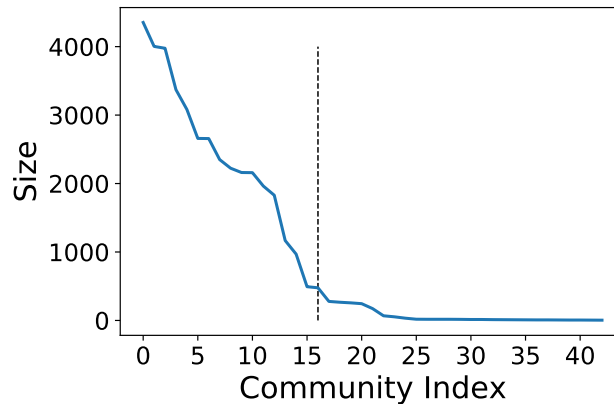
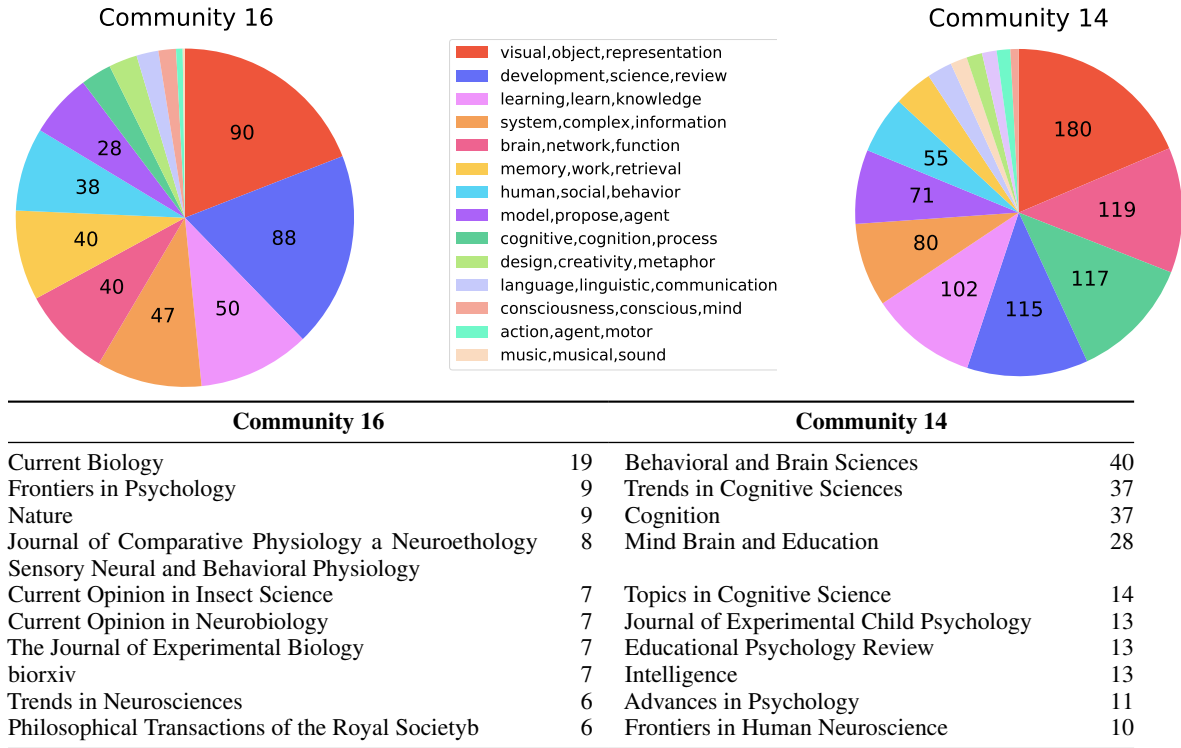
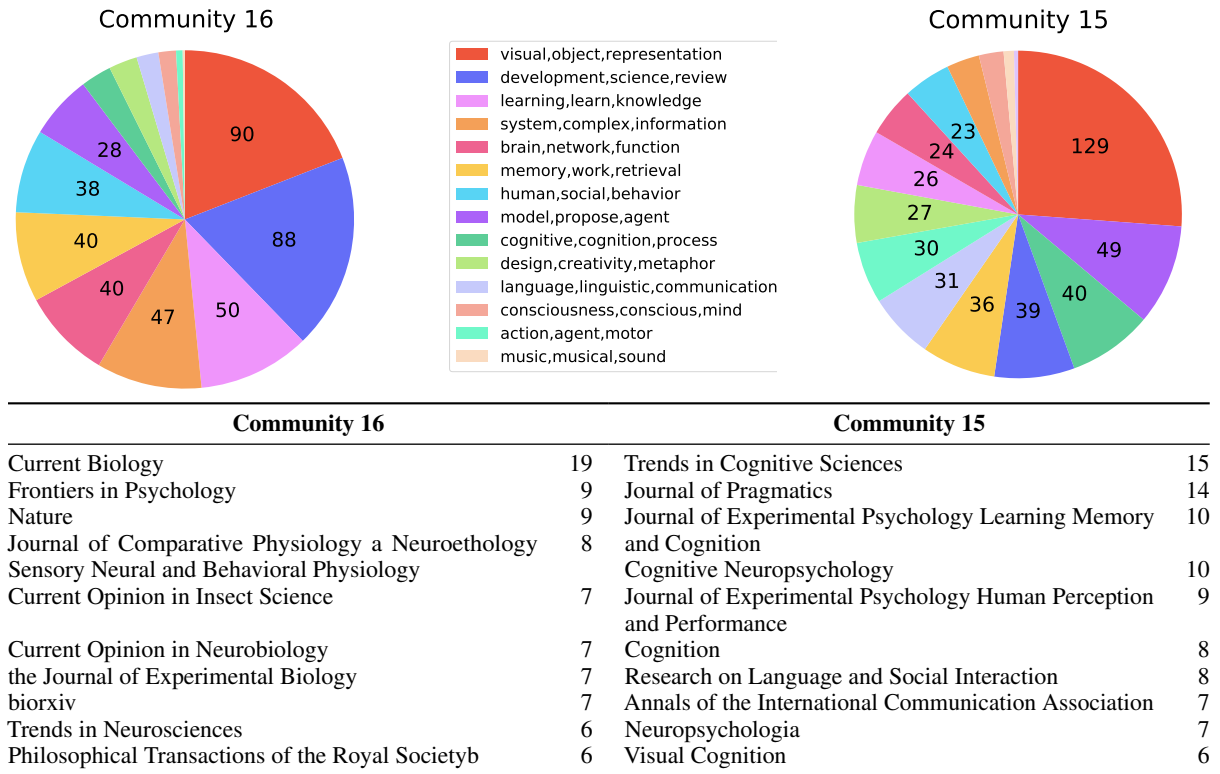


Figure 4: Size distribution of citation network communities. Dotted line is at community 16.

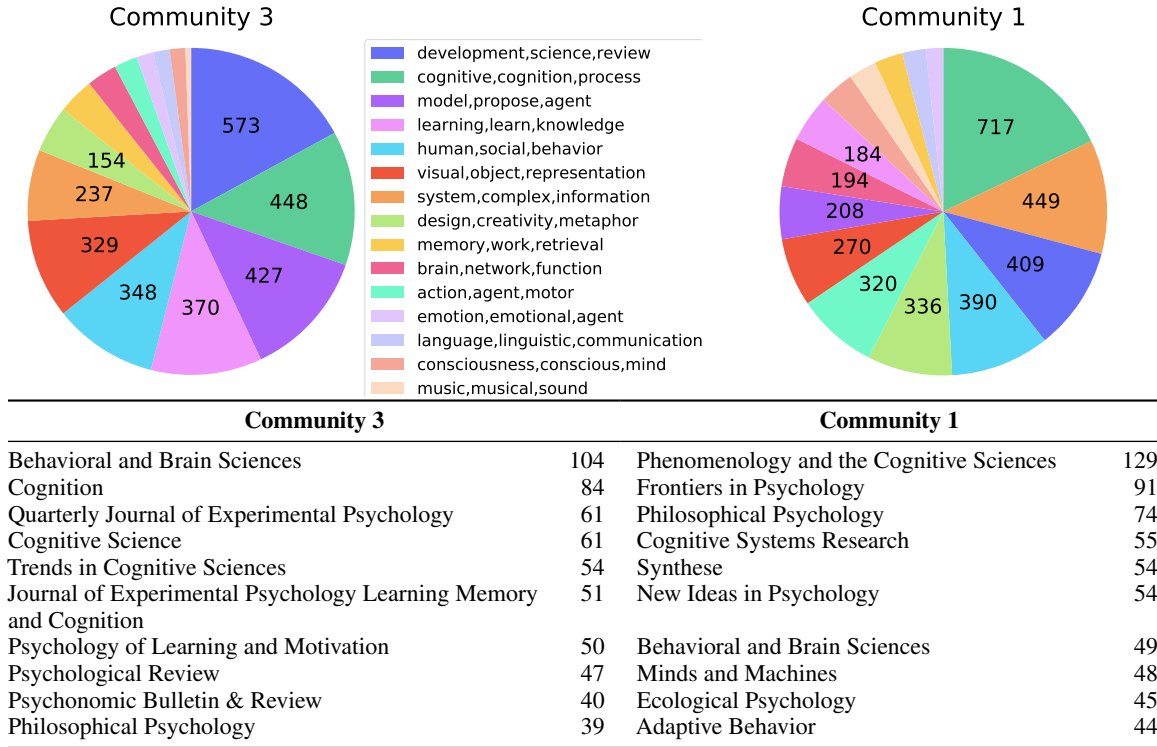
The six pairs of interest that were identified with community 4, 8 in Figure 2.
 $\rho = 0.8146, T = 0.9204, J : 0.2116.$



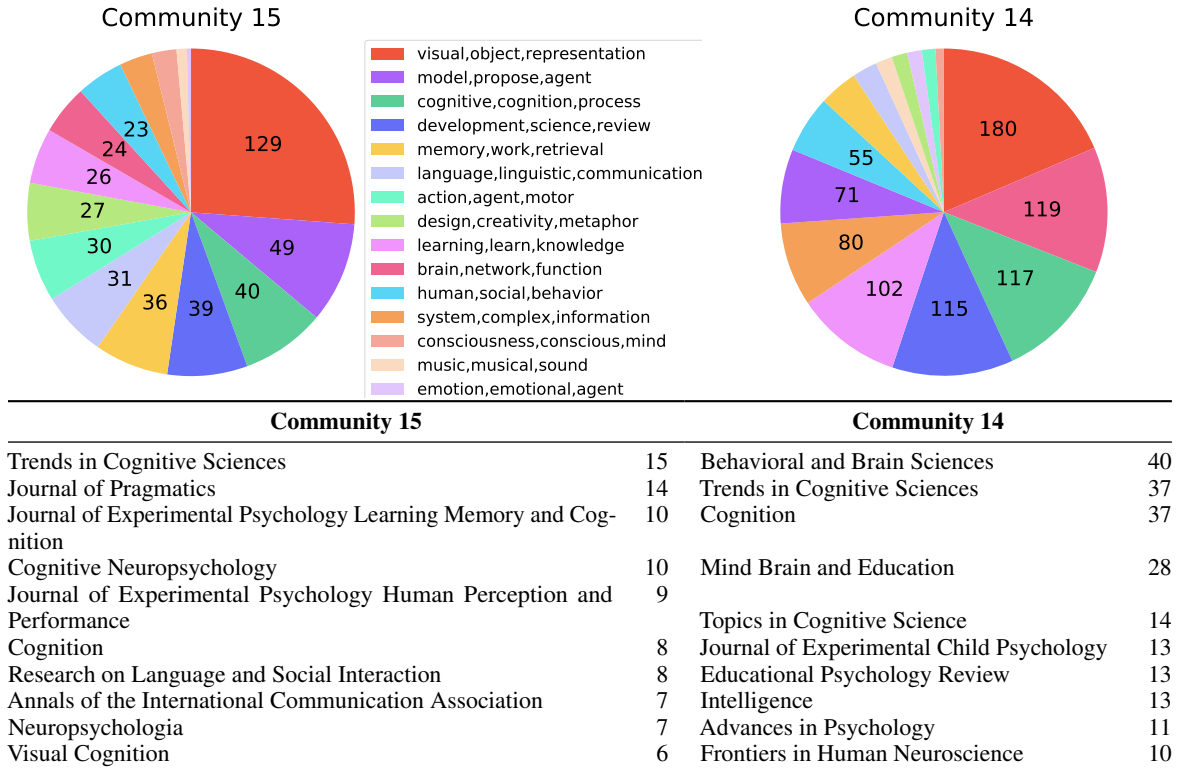
$\rho = 0.7486, T = 0.8445, J : 0.1919.$



$\rho = 0.6629, T = 0.8754, J : 0.425.$

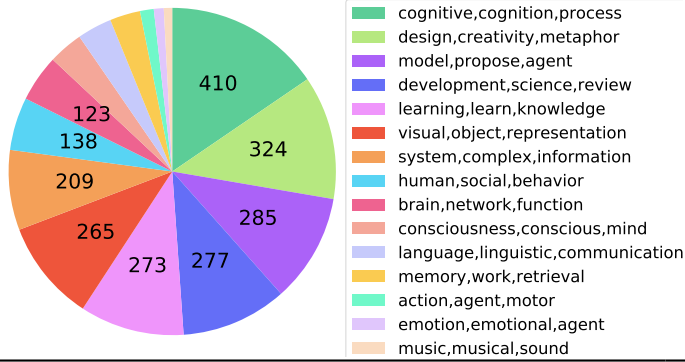


$\rho = 0.657, T = 0.8785, J : 0.443.$

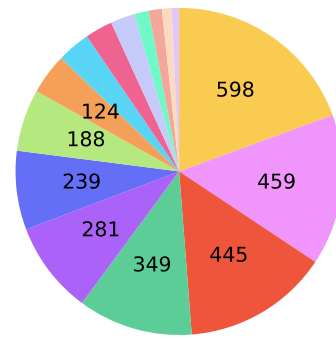


$\rho = 0.5968, T = 0.815, J : 0.4364.$

Community 6



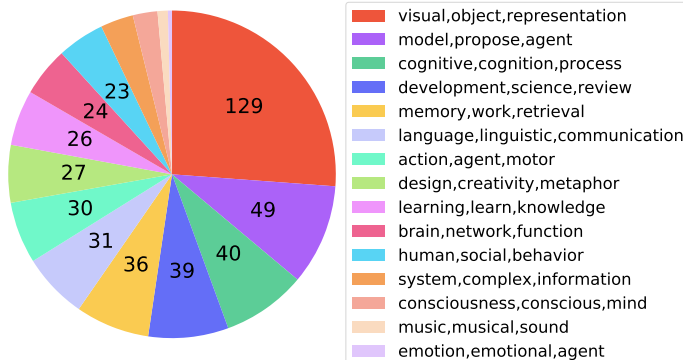
Community 4



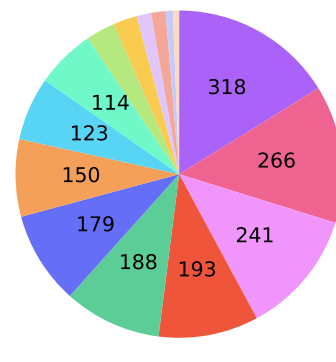
Community 6		Community 4	
Minds and Machines	98	Advances in Psychology	78
Cognitive Science	59	Memory & Cognition	66
Philosophical Psychology	56	Journal of Experimental Psychology	63
Behavioral and Brain Sciences	51	Applied Cognitive Psychology	61
Synthese	34	Educational Psychologist	52
Design Studies	31	Educational Psychology Review	44
Journal of Experimental and theoretical Artificial Intelligence	27	Psychology of Learning and Motivation	43
Advances in Psychology	26	Journal of Educational Psychology	35
Topics in Cognitive Science	23	Psychonomic Bulletin & Review	34
AI& Society	22	Memory	32

$\rho = 0.543, T = 0.7648, J : 0.4437.$

Community 15



Community 11



Community 15		Community 11	
Trends in Cognitive Sciences	15	Trends in Cognitive Sciences	80
Journal of Pragmatics	14	Behavioral and Brain Sciences	71
Journal of Experimental Psychology Learning Memory and Cognition	10	biorxiv	40
Cognitive Neuropsychology	10	Current Opinion in Behavioral Sciences	35
Journal of Experimental Psychology Human Perception and Performance	9	Neuron	34
Cognition	8	Frontiers in Psychology	32
Research on Language and Social Interaction	8	arxiv Artificial Intelligence	28
Annals of the International Communication Association	7	arxiv Neurons and Cognition	25
Neuropsychologia	7	Cognition	22
Visual Cognition	6	Philosophical Transactions of the Royal Societyb	22