

# A Bibliometric Approach to Understanding Insider Threat Scholarship

Luke J. Osterritter<sup>1</sup>[0000-0002-4447-9928] and Kathleen M. Carley<sup>1</sup>[0000-0002-6356-0238]

<sup>1</sup> Carnegie Mellon University, Pittsburgh PA 15213, USA  
losterritter@cmu.edu, kathleen.carley@cs.cmu.edu

**Abstract.** The study of insider threat and insider risk is a relatively nascent academic field. It is at once a problem not only of computer science, information science, and cybersecurity, but also one of business, management, organizational theory, psychology, sociology, human behavior, and others. Given this highly interdisciplinary nature, there is not yet a cohesive understanding around which scholars and institutions are involved in this field, nor of what their central foci are. This work represents a preliminary step in collecting the body of scholarly work around insider threat across many academic fields and performing analyses through bibliometric and network analytical methods. We present a look at who the most prolific and influential scholars are, how they are working with one another, and what topics their research centers around. We further lay the groundwork for continuing these analyses to discover where and how these researchers are publishing, which methods they employ, the topics they are focusing on, and where researchers ought to set their sights for future research on insider threat.

**Keywords:** insider threat, bibliometrics, network analysis

## 1 Introduction

The problem of insider threat – the threat that a current or former employee, contractor, or trusted business partner could misuse their authorized access either maliciously or unwittingly to bring harm to their organization [1], [2] – is one faced by organizations of all stripes, whether in industry, governmental entities, or non-profits. Indeed, if there exists information of any importance within an organization and a workforce that is granted access to such data, then the possibility of insider threat exists.

Alongside these concerns, a growing body of research is beginning to coalesce around both what is called “insider threat” and, increasingly, “insider risk” – which according to the CERT division of the Software Engineering Institute refers specifically to the impact and likelihood of a realized instance of an insider threat[1]. Researchers and practitioners alike are beginning to shift the conversation from one of threat hunting to instead one of risk mitigation. Work in this field has also been referred to by slightly different terms in its history and has some slightly different nomenclature depending on the type of industry from which the work sprouts.

Couple this with the fact that insider threat, perhaps mostly assumed to be a problem of technology, is truly a very socio-technical issue that while mediated by technology

and its associated advances is in fact perpetrated by and centered squarely on the human element. Given the need to understand and reason about this inherently interdisciplinary, relatively nascent, and ostensibly decentralized academic field, bibliometric approaches afford us powerful tools with which to wrap our arms around the study of insider threat and finally begin to see how it stands as a body of scholarly practice.

In this initial work, we seek to answer three main questions: who are the most prolific authors in the insider threat space? In what way are they working with one another? What are the key topics and themes present in this academic field – what are these scholars studying, and are we able to identify any gaps that could benefit from further research?

### **1.1 Related Work**

Bibliometric methods have long been used to gain a sense of understanding around areas of research, the nature of academic disciplines, and derive scholarly trends. Donthu et al. [3] presents an overview and guidelines for performing bibliometric analyses within the field of business and management, one of the major fields which concerns itself with the problem of insider threat. This complements prior work in the management and organization sciences by Zupic and Čater which introduces “science mapping”, or the use of bibliometric methods “to examine how disciplines, fields, specialties, and individual papers are related to one another”[4]. A relatively recent example of these techniques being used to quantify a nascent or emerging field can be performed around Social Cybersecurity [5]. Studies have also been performed to characterize subdisciplines, such as those within computer science, another field closely related to insider threat [6]. Other studies have been done to gauge whether and to how great a degree a field of study is interdisciplinary [7]. These techniques have also been used to quantitatively examine the trajectory of specific conferences or journals [8], [9] or of entire databases of scholarship in a specific discipline [10].

Specific bibliometric techniques have been used in scholarship to great effect. Co-authorship analysis has been used in on example to understand the network of authors working in the emerging field of forest management an entrepreneurship [11]. Co-topic network analysis, as well as co-authorship networks and citation analysis, have been used to study the fields of language, linguistics, and computational linguistics [12], [13].

## **2 Methods**

### **2.1 Scoping and Search Criteria**

Collecting scholarship on insider threat proposed several interesting challenges. Our initial question was one of scope: what constitutes research on “insider threat”? After all, insider threat is itself scoped into several categories: theft of intellectual property, fraud, sabotage, espionage, and negligence or otherwise non-malicious activities. More recently, workplace violence has been noted as having similar precursors and outcomes

to what had traditionally been considered insider threat [14], and there are other adjacent phenomena such as research integrity which deals specifically of the largely adversarial nation-state threat to academic research institutional data [15]. This is further confounded by the existence of research that, by topic, relates either wholly or in part under the umbrella of insider threat and/or risk, but does not specifically use those terms within the text. Further, while we could hand-curate a corpus of scholarship, it is desirable and more thorough to use programmatic means.

Ultimately, we decided to use a rather simple search criteria: “insider AND (threat OR risk)” within a text’s title, abstract, and keywords (whether those keywords are author-supplied or added by the publisher or database). Searching on these two phrases allows us to reduce the number of out-of-scope results and affords authors the ability to “self-select” into the discipline of insider threat and risk, focusing solely on work that is squarely aimed at this specific problem rather than work which, while perhaps strongly related, belongs to its own distinct discipline. We further decided that we would include any scholarship from conferences, conference workshops, journals, or other collected volumes that are themed around insider threat or risk, even if the individual papers did not themselves meet the criteria of having those terms in their titles, abstracts or keywords.

## 2.2 Data Sources

Our choices for data sources fall into three broad categories: interdisciplinary citation databases, discipline-specific databases, and organization-specific databases. Pulling from multiple sources with differing scopes allowed us to both cast a wide enough net to capture research across disciplines, while also allowing us to go deep enough to ensure we captured the bulk of work done in the field.

We identified Web of Science, Scopus, and Dimensions as our primary citation databases, as each has an extensive breadth of disciplines and literature, while also affording powerful and consistent searching capability. We further identified IEEE Xplore and the ACM Digital library for their thorough collection of literature within computer and information sciences as well as cybersecurity, and the INFORMS for its collection of scholarship around business and management sciences. The citation data in our corpus dates from 1977 to 2023, though the corpus does not contain the entirety of scholarship published in 2023 due to the collection cut off time.

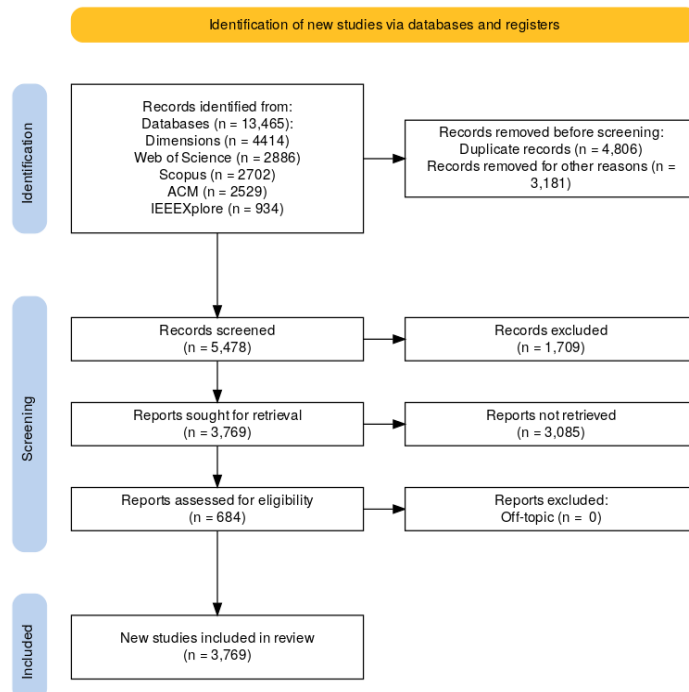
Through expert elicitation via research staff with the Insider Threat team at the CERT division of Carnegie Mellon University’s Software Engineering Institute, we collected a list of scholars in the insider threat domain who are prolific and important in the field. We were able to use this list to validate that the scholarship collected through these databases were returning the results that contained the work of recognized experts in the field. In doing this, we were able to identify a third category of sources that were necessary to capture the “grey literature” – work that is made available outside the traditional academic publishing venues. This includes documents such as industry reports, white papers, and government reports. The sources we used for these types of documents included the SEI Digital Library, MITRE News & Insight, public reports from the RAND corporation, the websites for several over federally funded research

and development centers (FFRDCs) as well as the national labs, university affiliated research centers (UARCS), and both the Defense Technical Information Center (DTIC) to capture work from and related to the US Department of Defense and the U.S. Department of Energy’s Office of Scientific and Technical Information (OSTI) database. Note that often these works are attributed to an organization rather than one or more individual authors – whenever this occurs in our collection, we treat the organization as the author node for the purposes of network analysis.

It should be noted that although there is a great deal of research that occurs on insider threat within the classified space, such work, along with other non-fundamental research, cannot be and is not included in our analysis and is as well out of scope for our purposes.

### 2.3 Deduplication and Screening

The citation data was collected from each database as an exported BibTeX file, which was then imported into our reference management software, Zotero [16]. Within Zotero we performed initial data cleaning, collected missing metadata and collection of references that were not otherwise available in BibTeX format (via the Zotero web browser plug in on DOI URL targets, as well as utilizing Google Scholar as necessary), and performed data deduplication as there was some significant overlap between results from each source database. At this stage we made sure to populate abstracts and full text whenever possible to ensure their availability for the screening process.



**Fig. 1.** PRISMA Diagram Detailing the Screening Process

We used a web-based solution called Rayyan [17] to perform inclusion screening. This software allowed us to quickly scan titles, abstracts, and keywords to assess whether they were appropriately on-topic. One illustrative example of the importance of the screening process is that a significant number of papers in the initial corpus referred to the “risk” of “insider trading”, a related but ultimately out-of-scope phenomenon. The Rayyan software afforded us the ability to work through the corpus in many separate sessions, assign reasoning to our inclusion criteria, mark certain papers for further review, and also computed a score as to how likely it was for a paper to be marked as included which sped up the process somewhat. Figure 1 shows a Preferred Reporting Items for Systematic reviews and Meta-Analyses (PRISMA) diagram that provides a flow chart of the screening process [18]. From an initial total of 13,465 bibliographic entries, after deduplication and screening for appropriateness, 3,769 were included in our corpus for this stage.

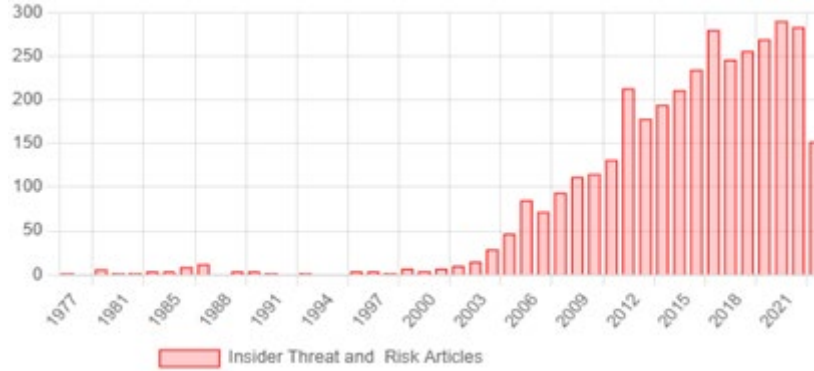
At the conclusion of screening process, we exported the resulting citation data as a BibTeX from Rayyan, then re-imported this into Zotero. Here we performed some further meta data cleanup, added any final citation data, and then exported the corpus, again as a BibTeX file, to import into the ORA-PRO [19] software. The ORA-PRO software is where the bibliometric and network analyses were performed.

### 3 Analysis

An analysis of papers published per year, as shown in Figure 2, reveals a linear increase in work in the field from the early 2000s which continues to increase. Note that the corpus does not include the entirety of all work performed in the year 2023 due to the collection cut-off. There is some historical context to keep in mind that may help to explain the periods of growth depicted in this graph.

In 2001, the CERT Division of the Software Engineering Institute formed the National Insider Threat Center, which coincides with when the research in the field began to increase. In 2011, U.S. president Barack Obama issued Executive Order 13587, “Structural Reforms to Improve the Security of Classified Networks and the Responsible Sharing and Safeguarding of Classified Information” [20] which called for all federal agencies to establish insider threat programs and as well the National Insider Threat Task Force (NITTF).

Note also that there is some scholarship that predates the point in the early 2000s where the insider threat definition was formalized. There are two main reasons for this; the first is that this scholarship was collected from entities that performed work on insider threat by another name (such as “internal threat” or “the insider problem”), or that the work has been tagged specifically as “insider threat” retroactively when digitally archived.



**Fig. 2.** Articles Published per Year.

A list of the top authors by total degree centrality gives us an idea of the most prolific scholars in the field. We queried the Crossref API to derive author affiliations, as they are not usually included in bibliographic citation data and used data from Google Scholar to supplement and corroborate. This does not give us an idea of where the work for each paper was performed but does give us a sense of the most current affiliation of each author – we believe this trade off was the correct one for our purposes in order to get a sense of the state of the field as it is at the time of this publication.

From this list in Table 2, we can begin to get an idea of not only the top authors, but the top institutions as well – Carnegie Mellon University appears five times, and the University of Oxford appears three times. We also see that the top centrality authors largely hail from traditional academic institutions of which five are located in the United State, four are in Europe, and two in Asia.

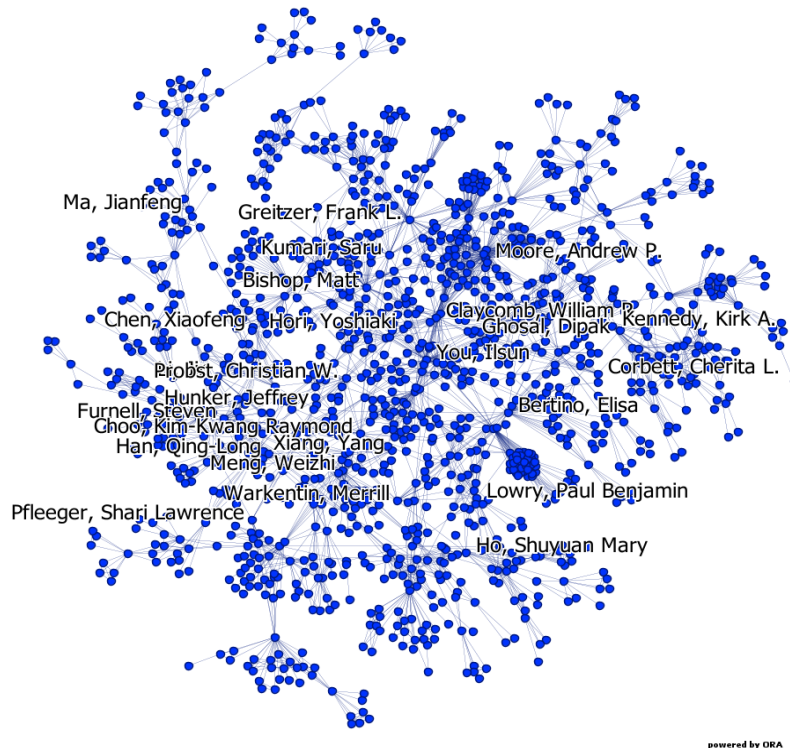
Top Authors by Total Degree Centrality

Rank	Author	Affiliation	Value	Unscaled
1	Moore, Andrew P.	Carnegie Mellon University	0.001	380
2	Cappelli, Dawn M.	Carnegie Mellon University	6.747e-04	214
3	Creese, Sadie	University Of Oxford	6.432e-04	204
4	Trzeciak, Randall F.	Carnegie Mellon University	5.990e-04	190
5	Goldsmith, Michael	University Of Oxford	5.675e-04	180
6	Greitzer, Frank L.	PsyberAnalytix	5.234e-04	166
7	Agrafiotis, Ioannis	University of Oxford	4.288e-04	136
8	Nurse, Jason R.C.	University of Kent	4.225e-04	134
8	You, Ilsun	Kookmin University	4.225e-04	134
9	Li, Wenjuan	Guangzhou University	4.162e-04	132
10	Legg, Philip A.	University of the West of Eng-land	4.099e-04	130
11	Bishop, Matt	University of California Davis	4.036e-04	128
11	Collins, Matthew L.	Carnegie Mellon University	4.036e-04	128
12	Jaros, Stephanie L.	University of Maryland	3.783e-04	120

13	Meng, Weizhi	Technical University of Denmark	3.657e-04	116
14	Stolfo, Salvatore J.	Columbia University	3.594e-04	114
15	Bertino, Elisa	Purdue University	3.531e-04	112
15	Costa, Daniel L.	Carnegie Mellon University	3.531e-04	112

From the co-authorship network we can derive some insights. There are a total of 3615 unique authors in the dataset. Of those, 392 are isolates, meaning that they have published solely on their own. There are 399 sets of dyads, where two authors only published with one another, and a set of 297 triads, wherein a set of three authors only published with each other. The largest network component contains 1374 authors, which represents our largest connected community.

We can see a visualization of this in a graph of this component shown in Figure 3, as well as the top 25 authors in terms of betweenness centrality, or those authors in the network who are in the most shortest paths between others in the network, and thus enjoy certain levels of information access. This graph shows us the largest network of authors in the field who have co-authored papers with one another; it's one way at which we can begin to reveal the “invisible college” of insider threat practitioners – those individuals who, while not sharing a formal organization, are socially connected with one another to produce scholarship and exchange ideas [21].







explore further the institutions that are furthering this work, the publications venues that bring the field together, and a closer look at the network of citations that can reveal the foundations upon which this field is being formed.

### Acknowledgements

This research was supported through the Minerva Research Initiative, in partnership with Office of Naval Research under grant #N000142114012, and by the Center for Computational Analysis of Social and Organizational Systems (CASOS) at Carnegie Mellon University. The views and conclusions are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Office of Naval Research or the U.S. Government.

### References

- [1] Software Engineering Institute, “Common Sense Guide to Managing Insider Threats, Seventh Edition,” Software Engineering Institute, Carnegie Mellon University, Pittsburgh, PA, White Paper, Sep. 2022. Accessed: Aug. 24, 2023. [Online]. Available: <https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=886874>
- [2] D. Costa, “CERT Definition of ‘Insider Threat’ - Updated.” Mar. 06, 2017. [Online]. Available: <https://insights.sei.cmu.edu/blog/cert-definition-of-insider-threat-updated/>
- [3] N. Donthu, S. Kumar, D. Mukherjee, N. Pandey, and W. M. Lim, “How to conduct a bibliometric analysis: An overview and guidelines,” *Journal of Business Research*, vol. 133, pp. 285–296, Sep. 2021, doi: 10.1016/j.jbusres.2021.04.070.
- [4] I. Zupic and T. Čater, “Bibliometric Methods in Management and Organization,” *Organizational Research Methods*, vol. 18, no. 3, pp. 429–472, Jul. 2015, doi: 10.1177/1094428114562629.
- [5] K. M. Carley, “Social cybersecurity: an emerging science,” *Comput Math Organ Theory*, vol. 26, no. 4, pp. 365–381, Dec. 2020, doi: 10.1007/s10588-020-09322-9.
- [6] S. Guha, S. Steinhardt, S. I. Ahmed, and C. Lagoze, “Following bibliometric footprints: the ACM digital library and the evolution of computer science,” in *Proceedings of the 13th ACM/IEEE-CS joint conference on Digital libraries*, in JCDL ’13. New York, NY, USA: Association for Computing Machinery, Jul. 2013, pp. 139–142. doi: 10.1145/2467696.2467732.
- [7] M. Youngblood and D. Lahti, “A bibliometric analysis of the interdisciplinary field of cultural evolution,” *Palgrave Commun*, vol. 4, no. 1, pp. 1–9, Oct. 2018, doi: 10.1057/s41599-018-0175-8.
- [8] M. Meyer, M. A. Zaggl, and K. M. Carley, “Measuring CMOT’s intellectual structure and its development,” *Comput Math Organ Theory*, vol. 17, no. 1, pp. 1–34, Mar. 2011, doi: 10.1007/s10588-010-9076-0.

- [9] R. L. Abduljabbar and H. Dia, "A Bibliometric Overview of IEEE Transactions on Intelligent Transportation Systems (2000–2021)," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 14066–14087, Sep. 2022, doi: 10.1109/TITS.2021.3136215.
- [10] R. Raman, P. Singh, V. K. Singh, R. Vinuesa, and P. Nedungadi, "Understanding the Bibliometric Patterns of Publications in IEEE Access," *IEEE Access*, vol. 10, pp. 35561–35577, 2022, doi: 10.1109/ACCESS.2022.3161639.
- [11] P. R. Mourao and V. D. Martinho, "Forest entrepreneurship: A bibliometric analysis and a discussion about the co-authorship networks of an emerging scientific field," *Journal of Cleaner Production*, vol. 256, p. 120413, May 2020, doi: 10.1016/j.jclepro.2020.120413.
- [12] M. CheshmehSohrabi and A. Mashhadi, "Using Data Mining, Text Mining, and Bibliometric Techniques to the Research Trends and Gaps in the Field of Language and Linguistics," *J Psycholinguist Res*, vol. 52, no. 2, pp. 607–630, Apr. 2023, doi: 10.1007/s10936-022-09911-6.
- [13] D. R. Radev, M. T. Joseph, B. Gibson, and P. Muthukrishnan, "A bibliometric and network analysis of the field of computational linguistics," *Journal of the Association for Information Science and Technology*, vol. 67, no. 3, pp. 683–706, 2016, doi: 10.1002/asi.23394.
- [14] A. P. Moore, T. M. Cassidy, M. C. Theis, D. Bauer, D. M. Rousseau, and S. B. Moore, "Balancing Organizational Incentives to Counter Insider Threat," in *2018 IEEE Security and Privacy Workshops (SPW)*, May 2018, pp. 237–246. doi: 10.1109/SPW.2018.00039.
- [15] D. Halbert, "Intellectual property theft and national security: Agendas and assumptions," *The Information Society*, vol. 32, no. 4, pp. 256–268, Aug. 2016, doi: 10.1080/01972243.2016.1177762.
- [16] Zotero. (Aug. 26, 2024). Corporation for Digital Scholarship. [Online]. Available: <https://www.zotero.org/>
- [17] M. Ouzzani, H. Hammady, Z. Fedorowicz, and A. Elmagarmid, "Rayyan—a web and mobile app for systematic reviews," *Syst Rev*, vol. 5, no. 1, p. 210, Dec. 2016, doi: 10.1186/s13643-016-0384-4.
- [18] M. J. Page *et al.*, "The PRISMA 2020 statement: an updated guideline for reporting systematic reviews," *BMJ*, vol. 372, p. n71, Mar. 2021, doi: 10.1136/bmj.n71.
- [19] K. M. Carley, "ORA: A Toolkit for Dynamic Network Analysis and Visualization," in *Encyclopedia of Social Network Analysis and Mining*, R. Alhajj and J. Rokne, Eds., New York, NY: Springer, 2018, pp. 1693–1702. doi: 10.1007/978-1-4939-7131-2\_309.
- [20] "Executive Order 13587 -- Structural Reforms to Improve the Security of Classified Networks and the Responsible Sharing and Safeguarding of Classified Information," [whitehouse.gov](https://www.whitehouse.gov). Accessed: Jun. 02, 2024. [Online]. Available: <https://obamawhitehouse.archives.gov/the-press-office/2011/10/07/executive-order-13587-structural-reforms-improve-security-classified-net>
- [21] D. Crane, *Invisible colleges; diffusion of knowledge in scientific communities*. Chicago: University of Chicago Press, 1972.