# A Bibliometric Analysis of The Term DataOps

Antonius Sony Eko Nugroho[1], Wella[2]

[1,2] Information System Department, Faculty of Engineering and Informatics, Universitas Multimedia Nusantara, Tangerang, Indonesia

antonius.sony@lecturer.umn.ac.id, wella@umn.ac.id

*Abstract* – **This paper uses bibliometric studies to provide a rigorous analysis of research trends related to the phrase "DataOps". Data were initially taken from the Scopus database from 2018 to 2023 which resulted 34 documents, while in other hands we found more than 1000 documents in the Google Scholar database collected using Publish or Perish. Five years of data is collected because the scientific publications on term DataOps started to increase since 2018. We have implemented the "bibliometrix 3.0" software package, along with the r-package and VOSviewer, to examine important elements of the literature. Observing the emergence of the phrase DataOps, one can note a substantial impact on its definition and actual implementations. DataOps enthusiasts work to provide a standardised approach for consistently implementing the technique across different data operation contexts. DataOps will continue to evolve into an effective and reliable data management system. The results of this study will assist researchers in the field to identify the prevailing research patterns in worldwide DataOps research and propose potential avenues for future research.**

*Index Terms* – *Bibliometric, Bibliometrix, DataOps, VOSViewer.*

## I. INTRODUCTION

DataOps is a new discipline that focuses on efficiently implementing data science on a large scale, inspired by the operational strategies of successful firms like LinkedIn, and eBay. Organizations require more than just cutting-edge AI algorithms, emerging technologies, and skilled individuals to transform data into actionable insights and valuable analytical data products [1]. DataOps involves continuous incremental use of data, utilizing complex systems orchestration techniques which in turns to data insight [2]. It aims to harness the potential of data while addressing the challenges of data deluge.

DataOps can be applied in various domains, such as government financial data analytics [3], drought mitigation in high-risk areas [4], data quality discovery [5], and healthcare analytics [6]. It offers solutions for data preparation, feature selection, and machine learning algorithms. By implementing DataOps, organizations can accellerate the efficiency and effectiveness of their data science and analytics projects, leading to improved business value.

DataOps enables the development of adaptable capabilities that can assist companies in responding swiftly to the constantly evolving digital landscape. It also facilitates the identification of essential prerequisites for organisations to effectively use DataOps processes, using Leavitt's Diamond Model as a basis [7]. DataOps offers a comprehensive and structured method for digital business transformation, and implementing DataOps leads to enhanced organisational performance in digital business transformation [8]. However, we are just at the beginning of data-driven transformation and understanding of the optimal techniques to derive our desired outcomes from raw data [1], in the next few years we'll see a revolution in data science, machine learning, and deep learning.

### A. Research Question

This paper examines the research patterns of the term DataOps between 2018 and 2023. It addresses five specific research questions:

a. What is the output profile of DataOps articles from 2018 to 2023 and how geographically distribution of articles worldwide?

b. What is the publication frequency of the term DataOps between 2018 and 2023?

c. What are the findings of the research theme cluster visual analysis conducted on the term DataOps?

### B. Data Collection for Literature Review

Bibliometrics is a technique for exploring scholarly material in research articles, books, conference proceedings, and reports from several well-established databases [9]. It is the process of extracting, organizing, and evaluating data in order to make strategic decisions [10]. Bibliometric information for this study is collected in two steps. Firstly, we selected the Scopus database for information collection purposes. The database contains high-impact factors and prestigious research publications. Then, in the second step, we narrow down the subject category and conduct the search query for more holistic data collection. We have applied specific filters to get the desired results. The databases were searched by using the specific related keywords from the period of 2018-

2023 at the date of November 19, 2023. The Final search query is TITLE-ABS-KEY ( dataops ) AND PUBYEAR > 2017 AND PUBYEAR < 2024. It is expected that the results may change in the future when more papers related to the issue are included.

*C. Bibliometric Tools*

Bibliometrix is a package used for bibliometric analysis in research to gain insights and coordinate research efforts in the related fields [11]. It provides capabilities for analyzing performance indicators, identifying trends, and mapping technological developments. It can be used to analyze interdisciplinary research in fields like conversation and aphasia [12], patent analysis in areas like virtual worlds [13] or topic like agile IT Governance [14]. The analysis conducted using Bibliometrix helps in identifying predominant themes, emerging concepts, impactful authors, and sources, as well as making theoretical contributions and providing future research directions [15], [16].

VOSviewer is a freely available software utilised for conducting bibliometric analysis and visualising networks. It helps in understanding the structure and evolution of knowledge in scientific disciplines [17]. VOSviewer can be used to build bibliometric maps or networks based on different types of relationships, such as co-authorship, co-occurrence, and co-citation [18]. Researchers can utilise this tool to locate clusters of interconnected articles, determine precise keywords for effective searching, uncover potential collaboration partners, recognise influential papers, and pinpoint areas of knowledge that require further exploration [19], [20].

## II. RESEARCH METHOD

The purpose of this study is to establish a basis for future researchers and to outline previous research for scholars to expand upon in order to advance knowledge in the future. In order to accomplish this, a comprehensive literature review is carried out systematically, and the resulting data is examined through bibliometric analysis, keyword analysis, and citation analysis. The data can be summarised using many criteria, including top authors, journals, institutions, keywords, citations, publishing nations, and publication years. This organized data can assist individuals better understand the term of DataOps.

The research commenced by performing an online search in November 2023 within the Scopus database, which is renowned among the academic world for its collection of articles and conference proceedings considered highly pertinent. The entire procedure is depicted in Figure 1, adhering to the sequential process of doing bibliometric analysis. The 'Biblioshiny' is a specialised R package called 'Bibliometrix 3.0' that is designed for web use. It utilises analytical tools based on Bradford's law to describe selected texts. These tools include global citations, h-index, g-index, and m-index.
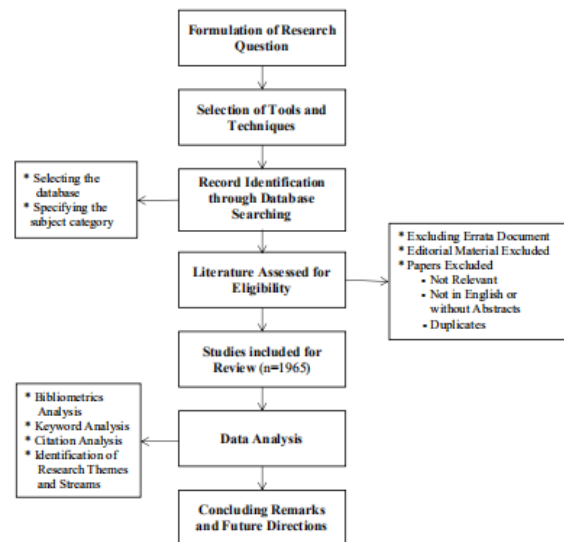


Fig. 1. The process of bibliometrics analysis

The next steps involve identifying important research themes and streams using scientific mapping approaches that analyse the conceptual framework and holistic keywords as input data. Upon completion of these analysis techniques, we will be capable of offering a comprehensive examination of the data and proposing a future research agenda.

The present study on bibliometrics analysis will provide information pertaining to the leading journals, authors, publications, universities, and countries. The study obtained thorough information regarding DataOps terminology through a comprehensive literature review. Table 1 provides a brief explanation of the features of the concepts under research, which is essential for understanding the attributes of the selected literature before proceeding to the analysis section.

TABLE I. DESCRIPTIVE CHARACTERISTICS OF TERM DATAOPS.

| Description | Results |
| --- | --- |
| **Main Information** | |
| Timespan | 2018:2023 |
| Sources (Journals, Books, etc) | 31 |
| Documents | 34 |
| Annual Growth Rate % | 37.97 |
| Document Average Age | 1.65 |
| Average citations per doc | 3.853 |
| References | 1 |
| **Document Contents** | |
| Keywords Plus (ID) | 299 |
| Author's Keywords (DE) | 106 |
| **Authors** | |
| Authors | 116 |

| Description | Results |
|---|---|
| Authors of single-authored docs | 7 |
| **Authors Collaboration** | |
| Single-authored docs | 7 |
| Co-Authors per Doc | 3.68 |
| International co-authorships % | 14.71 |
| **Document Types** | |
| article | 9 |
| book | 3 |
| book chapter | 1 |
| conference paper | 21 |

### III. RESULTS AND DISCUSSION

#### A. Publication Output and Document Sources

Based on data obtained from the Scopus database, it was found that the trend in the number of publications focus on DataOps research continued to increase from 2018 to November 2023. The number of publications decreased in 2019, but this downfall was compensated by a significant increase in 2020. The total publications during the 6 years were 34, with an increase of 37,97 % per year.

Fig. 2. Documents per year

More than half of the publications produced by the authors came from conference papers which reached 61.8 %, followed by publications in the form of articles at 26.5%. The remaining publications came from books by 8.8 % and book chapters by 2.9 %.
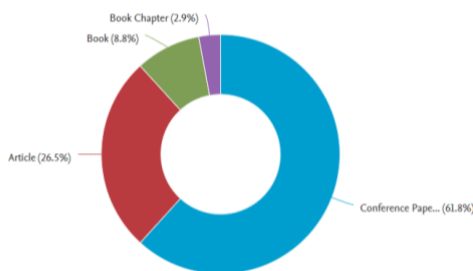
Fig. 3. Documents by type

It is evident that the research topic of DataOps is mostly focused on the field of Computer Science. The majority of DataOps research, approximately 55%, is conducted within the field of Computer Science. This is followed by Engineering and Mathematics, each contributing 9.6% of the research.
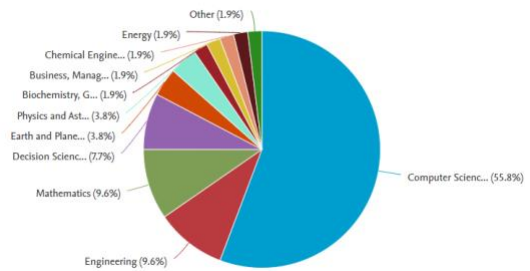
Fig. 4. Documents by subject area

During 2018-2023, DataOps-related publications come from 31 sources. Some journals even publish this theme consecutively for several years, these include: The Lecture Notes in Computer Science Including Subseries Bioinformatics (2018-2020), The CEUR Workshop Proceedings (2018,2022), and The ACM International Conference Proceedings Series (2020-2023). While other journals are not publishing research related to this regularly, for example The Communication in Computer and Information Science (2021) and The IEEE Transaction on Knowledge and Data Engineering (2023).
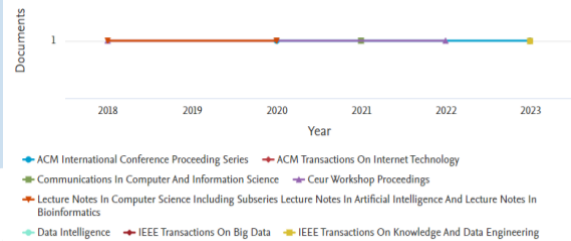
Fig. 5. Documents per year by source

In line with the increase in publications related to DataOps, the number of CiteScore also tends to increase significantly since 2018. The CiteScore of an academic publication is a metric that represents the average number of citations received by recent articles published in that journal on a yearly basis. This significant increase in citescore was experienced by The IEEE Transaction on Knowledge and The Data Engineering. Increased citations are also experienced by The ACM Transaction on Internet Technology, which experienced a decline in the 2012-2014 period, then increased again and achieved significant improvements in 2019-2022 period. Meanwhile, other citation sources such as The Lecture Notes in Computer Science and The Communication in Computer and Information Science tend to stagnate and do not experience a significant increase in citations.

The ACM International Conference Proceedings Series and The CEUR Workshop Proceedings has also been a source of citations since 2016. However, the number of citations is not much and tends to decline.

The latest developments in 2020-2021, The IEEE Transaction on Big Data and The Data Intelligence began to become a source of citations, even the number of citations from The IEEE Transaction on Big Data in 2021-2022 has exceeded the number of citations in The ACM Transaction on Internet Technology.
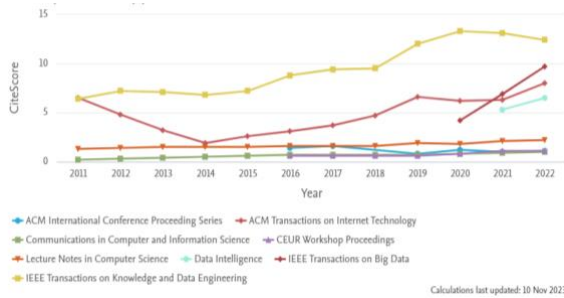


Fig. 6. CiteScore publication by year of the source

When discussing the legitimacy of publication sources, it is worth noting that certain journals have a high SCImago Journal Rank (SJR). The SJR indicator is a metric that quantifies the prestige of academic journals by considering both the number of citations received by a journal and the prestige of the journals from which the citations originate. For 22 years, The Lecture Notes in Computer Science, and The CEUR Workshop Proceedings have always been included in the list of Q1 journals. The ACM International Conference Proceedings Series was continuously on the Q1 list from 2003 to 2022, while The Communication in Computer and Information Science and The ACM Transaction on Internet Technology has been there since 2008-2009.

The IEEE Transaction on Big Data journal is still on the Q1-Q2 list. Meanwhile The IEEE Transaction on Knowledge and Data Engineering as the main source of publication with a focus on DataOps is still on the list of Q2-Q3 journals.
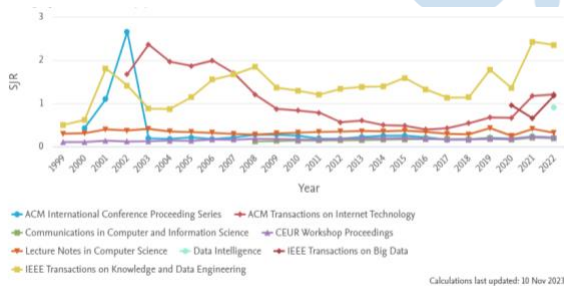


Fig. 7. SCImago journal rank by year of the source

Research productivity of an organization or university research reflects the relative position of the institution among others regarding a specific research interest and affects its ability to raise funds. Therefore, generally many organizations and institutions compete to produce research that will attract public attention. Although DataOps can be applied in various sectors, it is still relatively new.

Research related to DataOps is mostly carried out by universities and a small number of companies. Most publications are affiliated with universities, but the top affiliation comes from companies, namely with Microsoft Research as many as 7 documents, and Ubitech Limited as many as 6 documents.
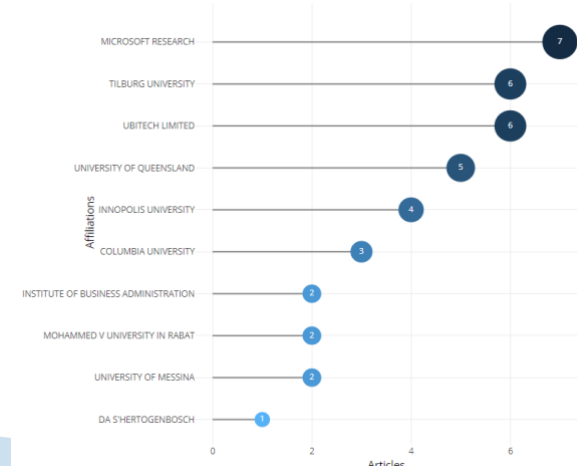


Fig. 8. Most relevant affiliations

Tilburg University has the highest number of author affiliations with 6, followed by the University of Queensland with 5 and Innopolis University with 4. Additional affiliates originate from universities in the Americas, Europe, and Africa.

## B. Publication Distribution Across Countries

Academic publications are thought to be indicative of a country's scientific prosperity, both in terms of quantity and influence. The Netherlands and the United States are the primary countries that publish documents on DataOps, based on the document's country of origin. Each country released 7 documents. Following that is Germany and Spain, both released 3 documents. Australia, along with Italy, Norway, the Russian Federation, and Sweden, each published a total of 2 documents. Austria is represented by 1 document.
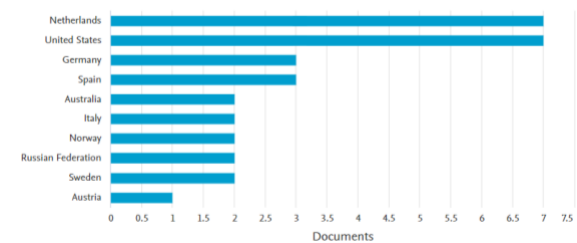


Fig. 9. Documents by country or territory

## C. Top Authors in The Research Term DataOps

Tamburi, D.A. is the most prolific author in DataOps-related research, with subsequent active authors producing only half of what Thamburi produces.
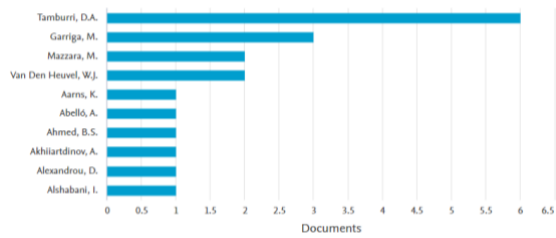
Fig. 10. Documents by author

Fig. 11. Thamburi released a total of 6 documents. The following authors have the highest number of active documents: Garriga (3), Mazzara (2), and Van Den Heuvel (2). Additional authors who have each published one document include Aarns, Abello, Ahmed, Akhilarddinov, Alexandrou, and Alshabani. Figure 11 shows that Thamburi D.A., as a sole author, receives the highest number of citations in comparison to other writers. The article titled Sustainable MLOPs: Trends and Challenges garnered 51 citations. In addition, Thamburi's collaborative work with Heuvel and Garriga has garnered 12 citations.

Fig. 12. The subsequent item that garnered significant interest was a publication authored by Munappy, Mattos, Bosch, Olsson, and Dakkak titled "From Ad-Hoc Data Analysis to DataOps," which amassed 23 citations. Ereth's work titled "DataOps: Towards Definition" got third place with a total of 16 citations. As of November 2023, documents authored by either a single individual or numerous authors receive fewer than 10 citations.



Fig. 13. Documents cited per author

## D. Publication Pattern

The researchers used several keywords to describe their publications as we can see in figure 12, in DataOps related publications, the most popular keyword used is DataOps itself (9), followed by Data Analytics (6), information management (5), Life cycle (5), decision making (4), develops (4), machine learning (4). The rest use big data keywords (3), codes (3) and data handling (3).
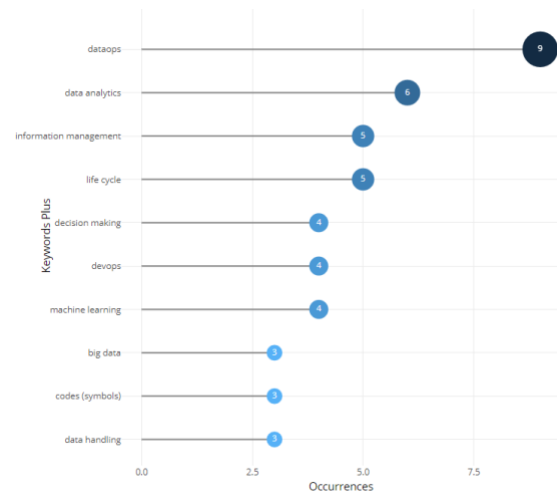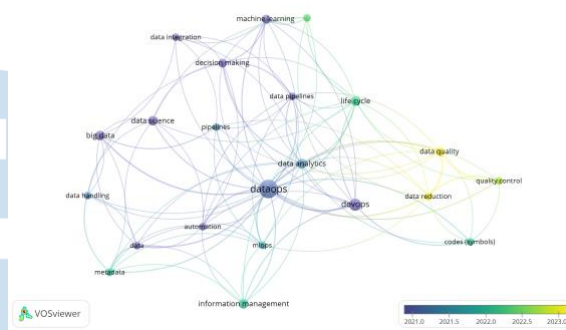


Fig. 14. Most frequent words



Fig. 15. Overlay Visualization

The topic of publication related to DataOps is also likely to shift from 2021 to 2023 as we can see in the visualization of the words in figure 13. In 2021, many publications discuss keywords related to defining DataOps, which include data, automation, integration, data science, big data, data pipelines, decision-making, and machine learning. In 2022, topics discussed include metadata, information management, mlops, codes, and life cycle. Shifting in 2023, more is discussed related to quality, data reduction, and quality control. We also able to see the trends of the word showed in figure 14 density visualization, DataOps and DevOps are one of the word mentioned a lot in the articles.
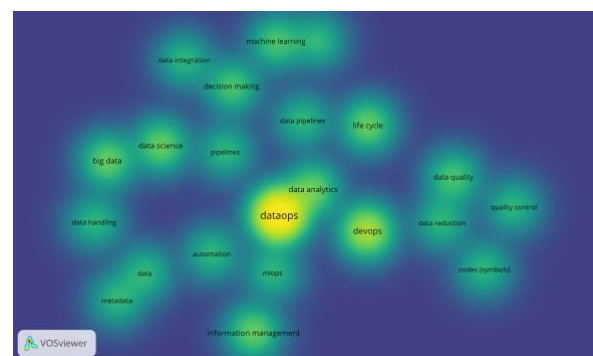


Fig. 16. Density visualization

## E. Theme Cluster Visual Analysis

Bibliometrix can display the mentioned word in the articles by displaying as wordcloud as in figure 15, so we can observe that words like data analytics, information management, devops, etc are following the trends of the term DataOps. While we can use Vosviewer to see the cluster network visualization shown in figure 16 Researchers from around the world generated four main clusters, shown by the colours yellow, red, green, and blue, according to the Scopus database. The initial cluster, shown by its yellow colour, is formed based on phrases such as bigdata, pipelines, and data integration, which ultimately leads to the subject of DataOps. The second cluster, represented by the colour green, is formed based on terms such as devops, lifecycle, decision making, and machine learning. This cluster primarily focuses on the topic of DevOps, which is closely related to DataOps articles. The third cluster, represented by the colour blue, is formed based on terms related to data analytics, data quality, data reduction, and quality control. This cluster revolves on the issue of data analytics, which is an integral part of dataops. The final cluster (red) is categorised based on keywords such as data science, data handling, metadata, automation, MLOps, and information management, which signify the management of data operations.
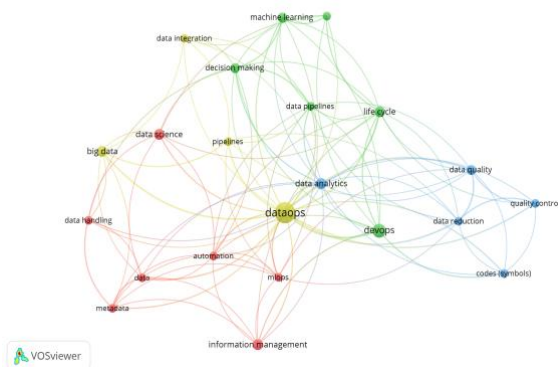


Fig. 17. DataOps terms wordcloud



Fig. 18. Cluster Network Visualization

Now, we review the bibliometric result from bibliometrix in figure 17 and 18. It is found that a wide variety of topics are covered in the journal during the study period. All the research themes are classified in order from highest amount of coverage to least to find out most preferred areas.
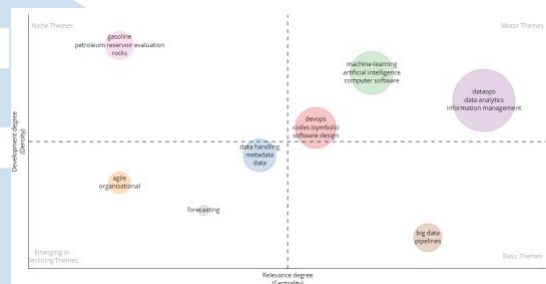


Fig. 19. DataOps terms treemap



Fig. 20. Thematic map

The most frequently discussed topics in research related to DataOps include DataOps itself 98%, data analytics 65%, information management 54%, and life cycle 54%. These topics are related to the next sub-topic which is also often discussed including: decision making 43%, data handling 33%, data quality 33%, Data reduction 33%, machine learning 33% and so on. A very fewer numbers of research paper focused on codes 3%, computer software 3%, organization 2%, and others which account less than 2% of the total topic categories.

## F. Limitations and Future Research Agenda

Currently, the process of devising and assessing research is of utmost importance for researchers. Utilising the references provided in article metadata from the Scopus database can be considered a highly reliable method for identifying taxonomies of research fields. The correctness of topic-level taxonomies can be evaluated by comparing document grouping methods such as direct citations, bibliographic coupling, and co-citation. Efforts of investigation should focus on the disparities in value across consistent taxonomic subjects and integrate historical records with fields of study that are prone to rapid changes. When employing more precise techniques to identify subjects, it is imperative for research papers, researchers,

organisations, and governments to enhance the accuracy of their innovation. Utilising direct citations results in a highly precise taxonomy, which is advised as a suitable foundation for decision-making.We exclusively utilised the bibliometric analysis technique on published documents alone from the Scopus database, without incorporating any additional databases. In addition, the study exclusively concentrated on the evaluation of digital collections, digital resources, and e-books. This evaluation was conducted using the bibliometric analysis method, which involved analysing published documents indexed in Scopus from 2018 to 2023.

## IV. CONCLUSION

The term DataOps has gained prominence and has made substantial contributions to its definition and practical uses. DataOps enthusiasts work to provide a standardised approach for consistently implementing the technique across different data operation contexts. Despite the extensive efforts made, the diverse nature of the data analysis process gives rise to numerous uncertainties regarding the use of DataOps. Data analysis is a vast field with different techniques, approaches, and technologies that might produce the same results. DataOps, on the other hand, provides a solution by combining a comprehensive and process-focused approach to data with automation and methodologies from agile software engineering and DevOps. This approach aims to enhance quality, velocity, and cooperation, while also promoting a culture of continuous development. DataOps will continue to evolve into an effective and reliable data management system.

## REFERENCES

[1] H. Atwal, *Practical DataOps*. Berkeley, CA: Apress, 2020. doi: 10.1007/978-1-4842-5104-1.

[2] A. R. Munappy, D. I. Mattos, J. Bosch, H. H. Olsson, and A. Dakkak, "From Ad-Hoc Data Analytics to DataOps," in *Proceedings of the International Conference on Software and System Processes*, New York, NY, USA: ACM, Jun. 2020, pp. 165–174. doi: 10.1145/3379177.3388909.

[3] A. Darono, "Dataops dalam Analitika Data Keuangan Negara: Studi Eksploratif," *Indonesian Treasury Review Jurnal Perbendaharaan Keuangan Negara dan Kebijakan Publik*, vol. 8, no. 2, pp. 125–136, Jun. 2023, doi: 10.33105/itrev.v8i2.545.

[4] D. A. Tamburri, V. R. van Mierlo, and W.-J. van den Heuvel, "Big Data for the Social Good: The Drought Early-Warning Experience Report," *IEEE Trans Big Data*, vol. 9, no. 3, pp. 773–791, Jun. 2023, doi: 10.1109/TBDATA.2022.3191749.

[5] S. Yu, T. Chen, L. Han, G. Demartini, and S. Sadiq, "DataOps-4G: On Supporting Generalists in Data Quality Discovery," *IEEE Trans Knowl Data Eng*, pp. 1–1, 2022, doi: 10.1109/TKDE.2022.3151605.

[6] S. Bahaa, A. Z. Ghalwash, and H. Harb, "DataOps Lifecycle with a Case Study in Healthcare," *International Journal of Advanced Computer Science and Applications*, vol. 14, no. 1, 2023, doi: 10.14569/IJACSA.2023.0140115.

[7] I. Gur, F. Moller, M. Hupperz, D. Uzun, and B. Otto, "Requirements for DataOps to foster Dynamic Capabilities in Organizations - A mixed methods approach," in *2022 IEEE 24th Conference on Business Informatics (CBI)*, IEEE, Jun. 2022, pp. 166–175. doi: 10.1109/CBI54897.2022.00025.

[8] J. Xu, H. Naseer, S. Maynard, and J. Filippou, "Leveraging Data and Analytics for Digital Business Transformation through DataOps: An Information Processing Perspective," *Australasian Conference on Information Systems*, 2021.

[9] F. Narin, D. Olivastro, and K. A. Stevens, "Bibliometrics/Theory, Practice and Problems," *Eval Rev*, vol. 18, no. 1, pp. 65–76, Feb. 1994, doi: 10.1177/0193841X9401800107.

[10] A. Firdaus, M. F. A. Razak, A. Feizollah, I. A. T. Hashem, M. Hazim, and N. B. Anuar, "The rise of 'blockchain': bibliometric analysis of blockchain study," *Scientometrics*, vol. 120, no. 3, pp. 1289–1331, Sep. 2019, doi: 10.1007/s11192-019-03170-4.

[11] M. Aria and C. Cuccurullo, "bibliometrix : An R-tool for comprehensive science mapping analysis," *J Informetr*, vol. 11, no. 4, pp. 959–975, Nov. 2017, doi: 10.1016/j.joi.2017.08.007.

[12] W. Wei and Z. Jiang, "A bibliometrix-based visualization analysis of international studies on conversations of people with aphasia: Present and prospects," *Heliyon*, vol. 9, no. 6, p. e16839, Jun. 2023, doi: 10.1016/j.heliyon.2023.e16839.

[13] E. A. D. Moresi and isabel Pinho, "Patent Analysis: An Approach Using Bibliometrix," in *19th CONTECSI International Conference on Information Systems and Technology Management*, TECSI. doi: 10.5748/19CONTECSI/PSE/ITM/6962.

[14] A. S. E. Nugroho, V. U. Tjhin, W. Kosasih, and H. Prabowo, "Bibliometric Analysis of Research Trend on Agile IT Governance," *Business and Accounting Research (IJEBAR) Peer Reviewed-International Journal*, vol. 6, no. 1, 2022, doi: 10.29040/ijebar.v6i2.2976.

[15] S. BÜYÜKKIDIK, "A Bibliometric Analysis: A Tutorial for the Bibliometrix Package in R Using IRT Literature," *Egit Psikol Olcme Deger Derg*, vol. 13, no. 3, pp. 164–193, Sep. 2022, doi: 10.21031/epod.1069307.

[16] L. Senyu, "A review and insight of vosviewer-based China community renewal-related studies," in *Advances in Urban Engineering and Management Science Volume 1*, London: CRC Press, 2022, pp. 600–607. doi: 10.1201/9781003305026-80.

[17] R. De Jong and D. Bus, "VOSviewer: putting research into context," *Research Software Community Leiden*, Mar. 2023, doi: 10.21428/a1847950.acdc99d6.

[18] A. Kirby, "Exploratory Bibliometrics: Using VOSviewer as a Preliminary Research Tool," *Publications*, vol. 11, no. 1, p. 10, Feb. 2023, doi: 10.3390/publications11010010.

[19] A. S. E. Nugroho and M. Hamsal, "Research Trend of Digital Innovation in Banking: A Bibliometric Analysis," *Journal of Governance Risk Management Compliance and Sustainability*, vol. 1, no. 2, pp. 61–73, Oct. 2021, doi: 10.31098/jgrcs.v1i2.720.

[20] I. Shkola, M. Andriichuk, and A. Petruniok, "Using vosviewer to analyze articles, indexing in pubmed database, about emerging infections," *Ukrainian Scientific Medical Youth Journal*, vol. 134, no. 4, pp. 53–61, Dec. 2022, doi: 10.32345/USMYJ.4(134).2022.53-61.