

Exploring the Italian research landscape on Digital Library in the Conference IRCDL

Eleonora Bernasconi^{1,†}, Andrea Mannocci^{2,†} and Anna Maria Tammaro^{3,†}

¹ *Università di Bari, Dipartimento di Informatica*

² *CNR-ISTI, Pisa*

³ *Università di Parma, Parco Area delle Scienze 181/A 43124 Parma*

Abstract

This study aims to explore the structure of knowledge around digital libraries embedded in IRCDL Conference presentations and examine research trends over time. It also analysed the published articles' subject, the authors, their affiliations and provenance and the collaboration network in IRCDL. We applied several bibliometric techniques, including productivity visualisation, authorship network analysis, and subject analysis.

Keywords

Digital Library Research, Italian Research Conference on Digital Library (IRCDL), Bibliometrics

1. Introduction

The Digital Library was born in the USA from 1994 to 1998, when three US government agencies - the National Science Foundation, the Defense Advanced Research Projects Agency and the National Aeronautics and Space Administration - funded six projects in the first phase of the Digital Library Initiative. In 1999, the three founding institutions were joined for the second phase of the Digital Libraries Initiative by the National Library of Medicine, the Library of Congress and the National Endowment for the Humanities, with the participation of the National Archives and the Smithsonian Institution.

In Europe, similar initiatives emerged in those years, such as the UK Electronic Library Program and the European Network of Excellence on Digital Libraries, known as DELOS. All these Digital Library initiatives have encouraged scientists, engineers, and librarians to explore research problems on Digital libraries together.

[†] These authors contributed equally.

andrea.mannocci@isti.cnr.it (A. Mannocci); eleonora.bernasconi@uniba.it (E. Bernasconi);

annamaria.tammaro@unipr.it (A. M. Tammaro)

[0000-0002-5193-7851](https://orcid.org/0000-0002-5193-7851) (A. Mannocci); [0000-0003-3142-3084](https://orcid.org/0000-0003-3142-3084) (E. Bernasconi); [0000-0002-9205-2435](https://orcid.org/0000-0002-9205-2435) (A. M. Tammaro)

The Italian Research Conference on Digital Library (IRCDL) was born in the context of the European Network of Excellence DELOS, starting in 2005, and for the last 20 years, the Conference has been held without interruption. In Italy, there has not been a cooperative project similar to the DL initiatives in the USA; however, since its inception, IRCDL has sought to stimulate the sharing and collaboration of digital library research. IRCDL has become a key venue for the exchange of experiences and knowledge between researchers and professionals engaged in the creation, organization, management and research of digital libraries. The conference was essential for the digital library experience in Italy, and the analysis of research presented at the Conference over the last 20 years (2005–2023) can describe the evolution of theory and practice of digital libraries.

“Digital libraries may be viewed as a new form of information institution or as an extension of the services libraries currently provide”. This is the first definition given by IRCDL, which historically approached “Digital libraries” research embracing the field at large and comprehending three key areas of interest that can be synthesized as

- scholarly communication (e.g. research data, research software, digital experiments, digital libraries),
- e-science/computationally-intense research (e.g. scientific workflows, Virtual Research Environments, reproducibility), and
- library, archive and information science (e.g. governance, policies, Open Access, Open Science).

IRCDL's focus is on emphasizing the multidisciplinary nature of research on digital libraries, which not only goes from computer science to humanities but also crosses areas in the same field, ranging, for example, from archival to librarian sciences or from information management systems to new knowledge environments. Representatives from academia, government, industry, research communities, and others were invited to participate in this annual conference. In her preface to IRCDL 2014 [2] Agosti affirms: “The conference draws from a broad and multidisciplinary array of research areas including computer science, information science, librarianship, archival science and practice, museum studies and practice, technology, social sciences, cultural heritage and humanities, and scientific communities”.

The second major focus of IRCDL [2] is on the profound change that is happening in the world of scholarly communication, where the object of scientific communication is no longer a linear text, although digital, but an object-centric network that consists of text, data, images, videos, blogs, and so on.

This study analyses the subject, the authors, their affiliations and their collaboration network of the papers published in the IRCDL conference. The findings of this study aims to understand the digital library knowledge in Italy and relevant interdisciplinary and transdisciplinary areas of research and collaborative networks in this field. Another

objective is to explore the digital libraries concepts embedded in IRCDL Conference presentations and examine research trends over time.

2. Data and methods

For our purposes, we needed bibliographic metadata about the papers presented at the IRCDL throughout the years, such as title, abstract, authors and affiliations, keywords and subjects, and citations. Unfortunately, a unique source for all such a wealth of information is not available. In fact, several sources, overlapping and mutually completing, can be found. Therefore, we opted for collecting data from multiple sources and selecting the most appropriate based on the task at hand.

The first method involved manual data collection from the IRCDL series main website². We recorded all the relevant information for each presentation at the IRCDL Conference since its inception in 2005. As fine-grained topics were not available here, we initially relied on the IRCDL Conference's own session classifications, except for 2012, where session numbers were used instead.

The second method utilized a (semi-)automated approach based on the data available in DBLP³. The metadata available in DBLP is highly reliable, alas, the information is limited to a handful of useful fields such as DOI, relevant URLs, title, authors, the number of pages, and year. Author affiliation, keywords and subjects, and citations are unfortunately not available in this database. Also, no metadata is available for the first edition in 2005.

Therefore, to complement DBLP information, we enriched the data with information from Google Scholar⁴, Semantic Scholar⁵ and the OpenAIRE Graph [4], adding details on citation counts. Finally, following the URLs present in DBLP, PDFs were downloaded and processed with the Llama2⁶ large language model to extract keywords, affiliations, and abstract summaries. As a failsafe, author affiliations were also extracted from PDFs with Grobid⁷, which seemingly offered a better basis for affiliation analysis, yet not perfect.

This combination of manual and automated methods provided a comprehensive dataset for analysis, offering insights into the evolution of topics and trends within IRCDL over two decades.

² <http://ims.dei.unipd.it/websites/ircdl/home.html>

³ <https://dblp.org/db/conf/ircdl/index.html>

⁴ <https://scholar.google.com>

⁵ <https://www.semanticscholar.org>

⁶ <https://ollama.com/library/llama2>

⁷ <https://github.com/kermitt2/grobid>

3. Results

The results described in the following are about productivity, impact (as proxied by the number of accrued citations), author and institution collaboration networks, and subject trends over time.

3.1. Productivity

The IRCDL Conference has always been a productive conference, with an average of about 20 papers per year, as illustrated in Fig. 1. While the long papers have been constantly represented throughout the whole conference, short papers appear to have declined over the editions, reaching a minimum in 2017 and steadily increasing ever after. The inlay reports the length distribution; most of them are in the 10–12 pages range, highlighting a preference for full contributions.

The top authors in terms of the number of papers presented in IRCDL are reported in Fig. 2. In Fig. 2c and 2d, we filtered the authors by removing the authors present in the steering committee and advisory board and selecting the remaining ones with 4 or more papers presented overall.

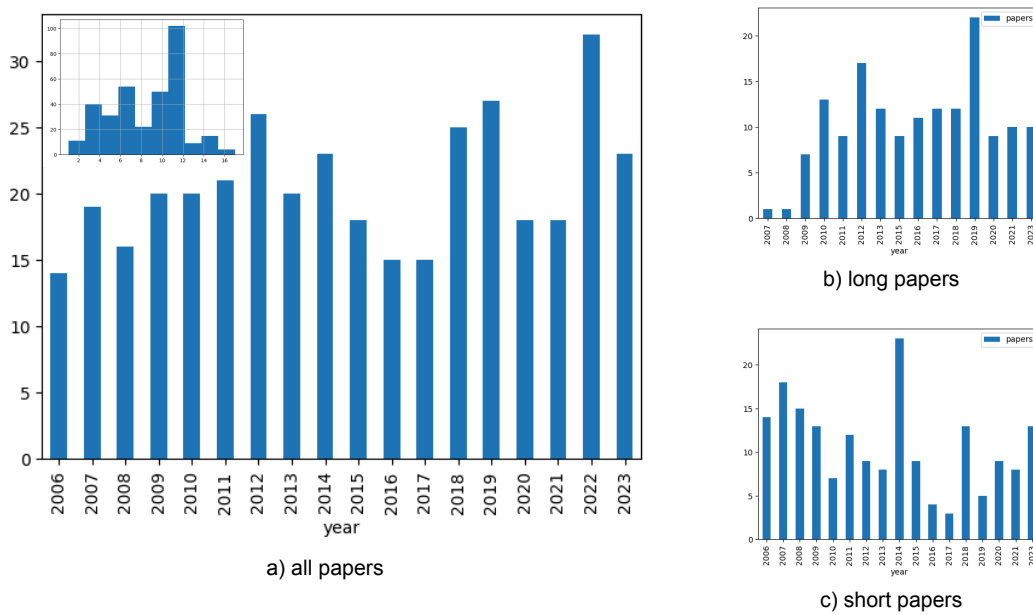
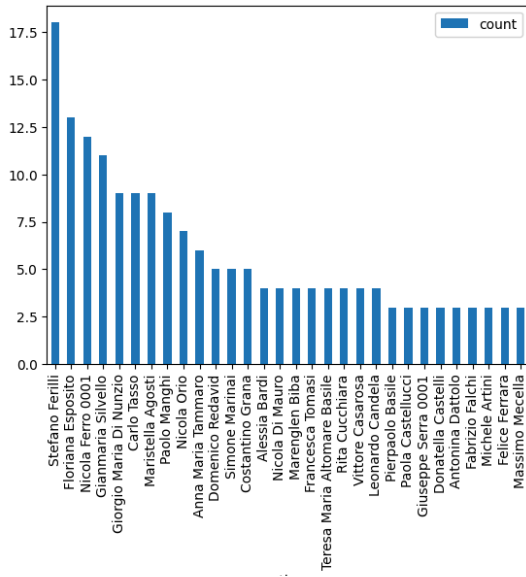
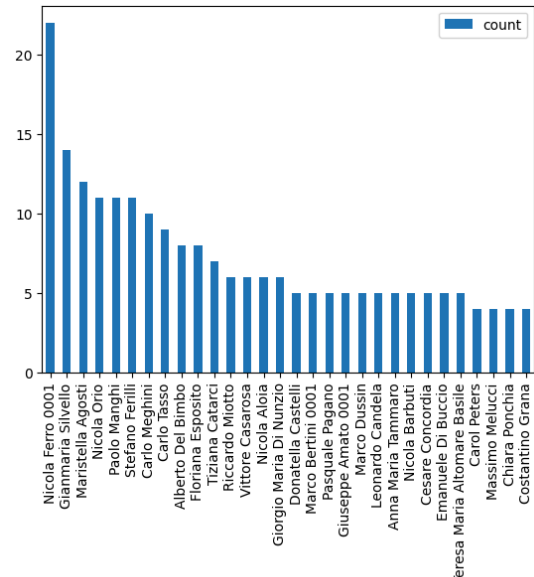


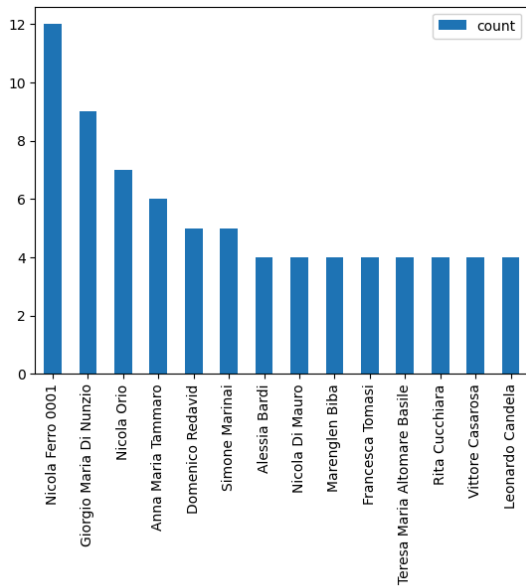
Fig. 1 Papers presented per year



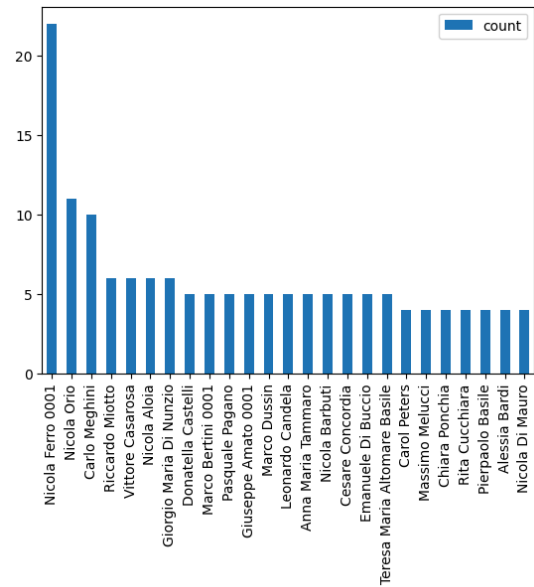
a) Top contributors of long papers



b) Top contributors of short papers



c) Filtered top contributors of long papers



d) Filtered top contributors of short papers

Fig. 2 Top contributors

3.2. Citation analysis

In Fig. 3 we reported the results obtained by analysing citational data. The top-10 most cited papers are reported in Fig. 3a and 3c, using data coming from different sources, namely Google Scholar and Semantic scholar in the first case, OpenAIRE in the second. Discrepancies are expected as different databases accrue citations with different criteria. Fig. 3b reports a list of the most cited authors, while Fig. 3d shows the years accounting for the highest number of citations. As can be noticed from the dates

of the top-10 most cited papers, it is interesting to see how the most recent editions (2017-2019) scored the highest impact in terms of citations.

↓ cites_scholar	Authors	Title	Year
0	98	F Ferrara, N Pudota, C Tasso	2,011
1	80	M Basaldella, E Antolli, G Serra, C Tasso	2,018
2	34	L Candela, D Castelli, P Pagano	2,009
3	33	J Allprantis, M Konstantakis, R Nikopoulou...	2,019
4	30	S Peroni, F Tomasi, F Vitali	2,013
5	24	G Castellano, G Vessio	2,020
6	23	R Jankovic	2,019
7	21	E Vocaturo, E Zumpano, L Caroprese...	2,019
8	20	AL Gentile, Z Zhang, L Xia, J Iria	2,010
10	19	F Bolelli	2,017

a) Top-10 most cited papers (source Google scholar + Semantic scholar)

Author	↓ Total Citations
C Tasso	213
N Pudota	110
F Ferrara	98
G Serra	93
S Ferilli	82
M Basaldella	81
E Antolli	80
M Agosti	79
N Ferro	75
F Esposito	57
L Candela	53
D Castelli	53
P Pagano	50
C Meghini	49
F Tomasi	48

b) Most cited authors

doi	title	authors	year	count(1)
10.1007/978-3-319-73165-0_18	Bidirectional LSTM Recurrent Neural Network fo...	Basaldella, Marco,Antolli, Elisa,Serra, Giusep...	2017	40
10.1007/978-3-642-27302-5_2	A Keyphrase-Based Paper Recommender System,A K...	Ferrara, Felice,Pudota, Nirmala Mary,Tasso, Carlo	2011	37
10.1007/978-3-319-68130-6_4	Indexing of Historical Document Images: Ad Hoc...	Bolelli, Federico	2017	19
10.1007/978-3-030-39905-4_11	Towards a Tool for Visual Link Retrieval and K...	Giovanna Castellano,Gennaro Vessio	2020	17
10.1007/978-3-642-15850-6_14	Semantic Relatedness Approach for Named Entity...	Ziqi Zhang,José Iria,Anna Lisa Gentile,Lei Xia	2010	17
10.1007/978-3-319-73165-0_15	XDOCS: An Application to Index Historical Docu...	BOLELLI, FEDERICO,BORGHI, GUIDO,GRANA, Costantino	2017	16
10.1007/978-3-319-68130-6_15	The Use of Hashtags in the Promotion of Art Ex...	Furini, Marco,Mandreoli, Federica,Martoglia, R...	2017	12
10.1007/978-3-030-11226-4_11	OpenAIRE's DOIBoost - Boosting Crossref for Re...	La Bruzzo, Sandro,Manghi, Paolo,Mannocei, Andrea	2019	12
10.1007/978-3-642-35834-0_23	Reflecting on the Europeana Data Model,Reflect...	PERONI, SILVIO,TOMASI, FRANCESCA,VITALI, FABIO	2013	11
10.1007/978-3-642-15850-6_8	A New Domain Independent Keyphrase Extraction ...	Pudota, N,Dattolo, Antonina,Baruzzo, A,Tasso, ...	2010	11

c) Top-10 most cites papers (source OpenAIRE)

Year	↓ cites_scholar
2,019	212
2,018	165
2,017	151
2,011	151
2,013	108

d) Most cited years

Fig. 3 Citation analysis

3.3. Network analysis of authors and organisations

As author name are unique by constriction within DBLP, we used them as identifiers in order to build a network using the library igraph⁸ for Python. The resulting network is reported in Fig. 4a where nodes represent authors, and the presence of an edge represents the existence of a co-authored paper. The node size represents the degree of an author (i.e., the total number of co-authors), while the weight of an edge represents the number of papers co-authored by the two authors connected. The resulting network is obviously disconnected and counts several connected components. The bigger connected components are formed thanks to authors with a higher degree who were capable of establishing fruitful collaborations across different organisations (see Fig. 4c), and can be seen as IRCDL powerhouse and core

⁸ <https://igraph.org>

community. Smaller clusters represent teams that contributed to IRCDL with one or more papers without establishing other connections with the core IRCDL research community. Tab. 1 reports the most productive author couples encountered.

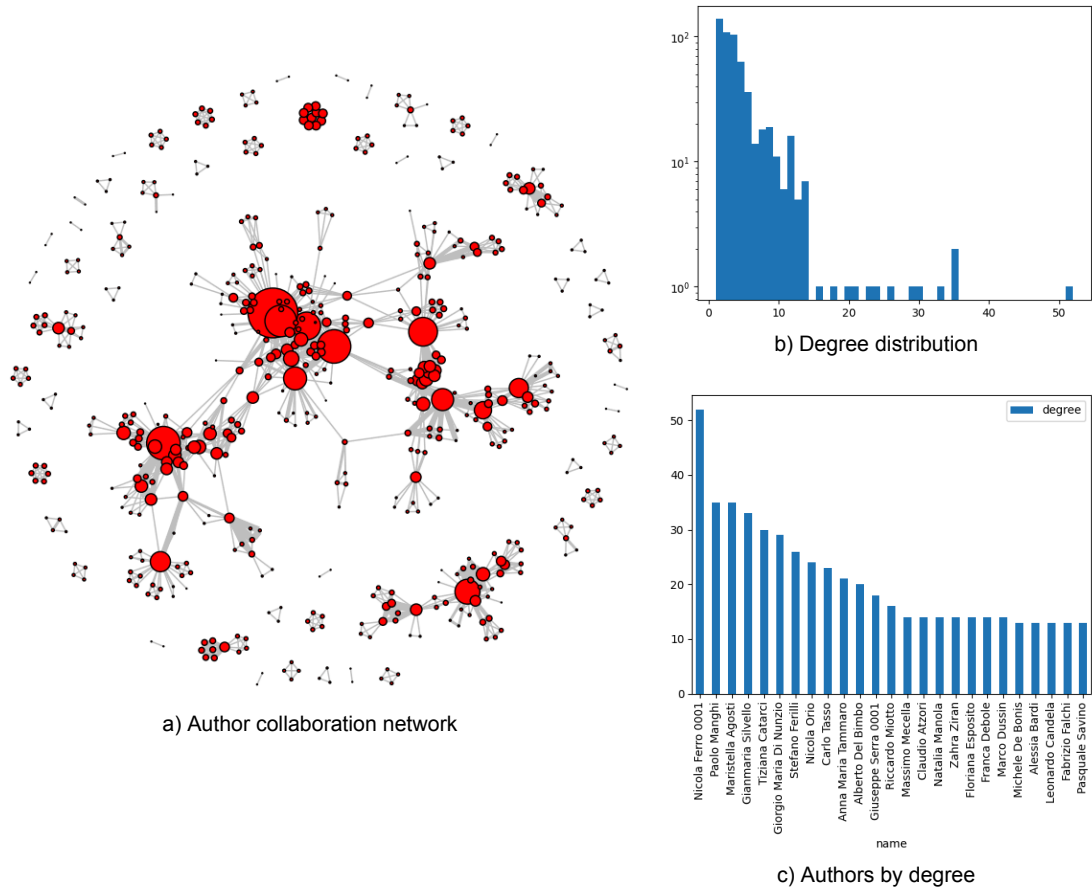


Fig. 4 Author collaboration network

Tab. 1 Most productive author couples

Couple	Co-authored papers
Stefano Ferilli + Floriana Esposito	17
Nicola Ferro + Gianmaria Silvello	17
Stefano Ferilli + Domenico Redavid	9
Leonardo Candela + Donatella Castelli	9
Floriana Esposito + Teresa Maria Altomare Basile	9
Paolo Manghi + Alessia Bardi	8

Stefano Ferilli + Teresa Maria Altomare Basile	8
Floriana Esposito + Nicola Di Mauro	8
Nicola Ferro + Maristella Agosti	7
Leonardo Candela + Pasquale Pagano	7
Costantino Grana + Rita Cucchiara	7
Andrea Mannocci + Paolo Manghi	6
Stefano Ferilli + Marenglen Biba	6
Giorgio Maria Di Nunzio + Nicola Ferro	6
Giorgio Maria Di Nunzio + Maristella Agosti	6
Gianmaria Silvello + Maristella Agosti	6
Nicola Orio + Riccardo Miotto	6
Donatella Castelli + Pasquale Pagano	6
Carlo Meghini + Nicola Aloia	6
Floriana Esposito + Marenglen Biba	6

Fig. 5 represents instead the collaboration network from the organisation standpoint. Similarly, nodes represent organisations; the bigger the node, the higher the degree, while edge weights represent the number of papers the two organisations co-participated to. The most profitable collaborations between institutions can be found between the University of Bari "Aldo Moro" and Artificial Brain S.r.l (8 papers), the Sapienza University of Rome and the University of Bari "Aldo Moro" (5), and ISTI-CNR + University of Parma (5). The University of Bari "Aldo Moro" is the organisation with more affiliates (46), followed by the University of Padua (44) and CNR-ISTI⁹ (34).

Despite the national character of IRCDL, the conference has been able to attract contributions from several foreign organisations, mainly from Europe, as reported in Tab. 2. Information professionals have collaborated almost every year, involving the most important cultural institutions that have been building digital libraries in Italy, as reported in Tab. 3. This is one of the unique features of the IRCDL Conference.

⁹ This is actually a lower-bound estimation as, in some circumstances, CNR-ISTI authors deliberately used the generic CNR affiliation rather than the most specific one.

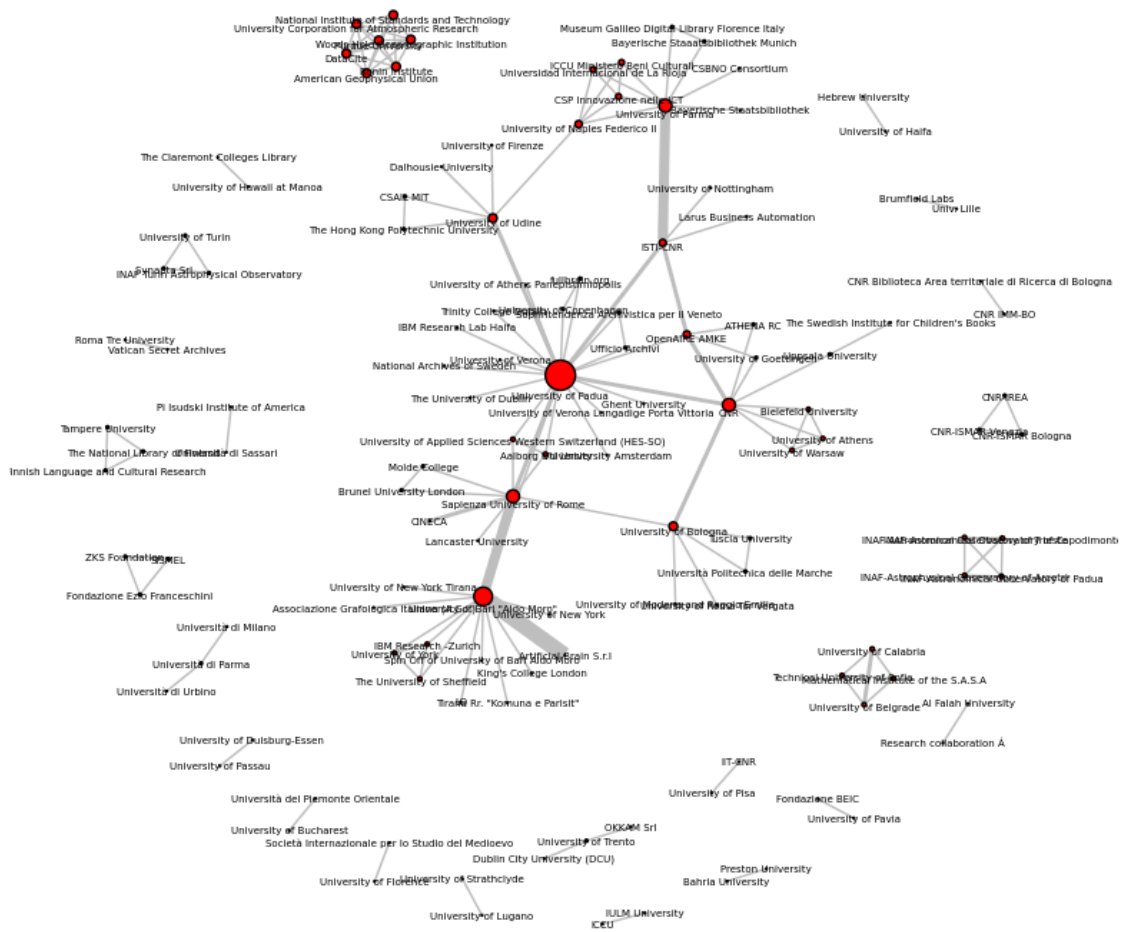


Fig. 5 Organisation collaboration network

Tab. 2 Authorship and international networking

Institutions	Authors
Uni-Basel	G. Brettlecker, P. Ranaldi, and H. Schuldt
Uni Lugano	R. Bache, F. Crestani
Uni Athens	Y. Ioannidis, E. Stamatogiannakis, M.L. Triantafyllidi, M. Vayanou, N. Manola, Y. Foufoulas, H. Dimitropoulos, T. Giannakopoulos
Uni Sheffield	A.L. Gentile, Z. Zhang, L. Xia, J. Iria
Uni Brunel London	Kai Olsen
Uni Ghent	T. Vets, M. Leman
Yonsei Univ Korea	Y. Chu, R. Allen

Tab. 3 IRCDL Conference collaboration with information professionals

Institutions	Authors
ICCU	R. Caffo, M.T. Natale, S. Di Giorgi
SISMEL	E. Degli Innocenti, A. Cosco
Biblioteca Europea	C. Consonni
Archivio	Vitali
Astrofisica	M. Gargano, A. Gasperini, E. Olostro Cirella, Riccardo Smareglia, V. Zanini
BNCF	Bergamin, Messina, A. Lucarelli, E. Viti
BNM	O. Braides, E. Sciarra
FRD	M. Lunghi
FAO	C. Caracciolo
ESA	Y. Coene, P.G. Marchetti, S. Smolders

3.4. Subjects analysis

Following the IRCDL session classifications, the authors identified two distinct trends in the subjects of presentations: one spanning from the beginning of IRCDL until 2014 and another from 2015 to the present day, as listed in Tab. 4 and Tab. 5.

Tab. 4 IRCDL Classification of Sessions 2005-2014

Year	IRCDL Classification of Sessions
2005	General presentation-Focused research presentation
2006	3D e Audio-System architecture-Personal information management-Video
2007	Annotations-Intelligent services for user-DL System architectures-Access-Video e 3D Data
2008	Multimedia and Multilingual Digital Libraries-Personalization and Preservation Clustering and Classification in Digital Libraries-Digital Library Architecture-Scientific Digital Libraries
2009	Models for Digital Libraries-Content description-Information access-Relevant project
2010	System Interoperability and Data Integration-Infrastructures, Metadata Creation and Management-Representation, Indexing and Retrieval in DL-Handling Audio/Visual and Non-traditional Objects

2011	Information Extraction and Access -Digital library models and systems
2012	The sessions are only numbered. No name is given.
2013	Information access-DL Architecture-DL Projects-Semantics and DL-Models and evaluation of DL-DL applications-Discussing DL perspectives
2014	DL projects-DL Models and modeling-Future trends in DL and scholarly communication-DL Infrastructures

Tab. 5 IRCDL Classification of Sessions 2015-2023

Year	IRCDL Classification
2015	Semantic modelling-Projects-Models and applications-Content analysis-Infrastructures
2016	Practices-Multimedia-Semantics-Collection management-Evaluation-Layout
2017	Bibliometrics and education-Multimedia-Data Management and presentation-Cultural heritage-Applications
2018	Models and applications-Cultural Heritage-Digital library architecture-Content analysis and text mining-Multimedia
2019	Open Science and open access-Open Science publishing and scientific workflows-OpenAIRE Workshop on Open Science Publishing Practices and Prospects-Text mining
2020	Information retrieval-Big data and data science in DL-Cultural heritage-Open science
2021	Data and platforms-Data access and monitoring
2022	Text recognition and multilinguality-NLP and AI-Digital edition and preservation
2023	Ontologies and knowledge graphs-Linguistics-Education-NLP and knowledge extraction-Projects-Document and data processing-Open Linked data

To facilitate the exploration and visualization of these findings, we have developed a Streamlit¹⁰ web application. This application allows users to interactively explore the results of crawling data from DBLP, Google Scholar and Semantics and testing an LLM model (LLAMA2) to extract knowledge and analyse recurrent keywords, topics, and affiliations.

¹⁰ <https://streamlit.io>

By employing Latent Dirichlet Allocation (LDA) topic modelling [3], applied to titles and abstracts (where available) within the crawled data, six distinct thematic clusters emerge.

Topic 1: Digital Library Ecosystem. This topic revolves around the digital library ecosystem, covering aspects such as system architecture, research projects, access mechanisms, and user services. Notable terms include "digital," "library," "system," "research," "information," "approach," "access," and "management." The trend associated with this topic has been decreasing over the years.

Topic 2: User-Centric Approach. Topic 2 emphasizes a user-centric approach to digital libraries, focusing on projects, tools, and resources aimed at enhancing user experience and access. Key terms include "user," "system," "research," "content," "library," "project," "provide," and "multimedia." The trend for this topic is also decreasing.

Topic 3: Project Management and Semantic Annotation. This topic delves into project management within digital libraries, with a particular emphasis on semantic annotation, research initiatives, and collaborative efforts. Noteworthy terms include "system," "document," "project," "group," "automatic," "digital," "biographical," and "dictionary." The trend for this topic is decreasing over time.

Topic 4: Cultural Heritage Preservation. Topic 4 focuses on cultural heritage preservation efforts within digital libraries, covering aspects such as information retrieval, semantic annotation, and content management. Key terms include "digital," "library," "information," "document," "user," "retrieval," "open," "cultural," and "semantic." This topic shows a decreasing trend over the years.

Topic 5: Knowledge Dissemination and Retrieval. The final topic highlights knowledge dissemination and retrieval mechanisms within digital libraries, focusing on services, resources, and tools for accessing and managing information. Important terms include "library," "digital," "retrieval," "information," "system," "search," "present," "service," and "user." This topic also exhibits a decreasing trend.

The evolution of topics over the years (Fig. 6) demonstrates varying trends in relevance and prominence. While some topics exhibit a consistent decrease in importance, others show fluctuations or even an increase in relevance. Notably, Topic 1, focusing on the digital library ecosystem, shows an increasing trend, indicating sustained interest and developments in this area. Conversely, Topic 2, which emphasizes a user-centric approach, displays a decreasing trend, suggesting a potential shift in research focus or priorities. Overall, these trends underscore the dynamic nature of digital library research and its continual evolution to meet changing needs and priorities.

Further analysis reveals the top topics for each year along with their associated weights, indicating the relative importance of each topic (see World Cloud in Fig. 8). The topics vary across different years, reflecting the evolving landscape of digital library research and the emergence of new trends and priorities. Notably, certain topics show consistent relevance over multiple years, while others exhibit fluctuations or appear more prominently in specific time periods.

In LDAVis Topic Visualization (Fig. 7), there's an observed overlap between Topic 1 and Topic 2, indicating a potential overlap in the themes and terms captured by these topics. This overlap suggests interconnectedness or shared characteristics between these two thematic clusters within the dataset.

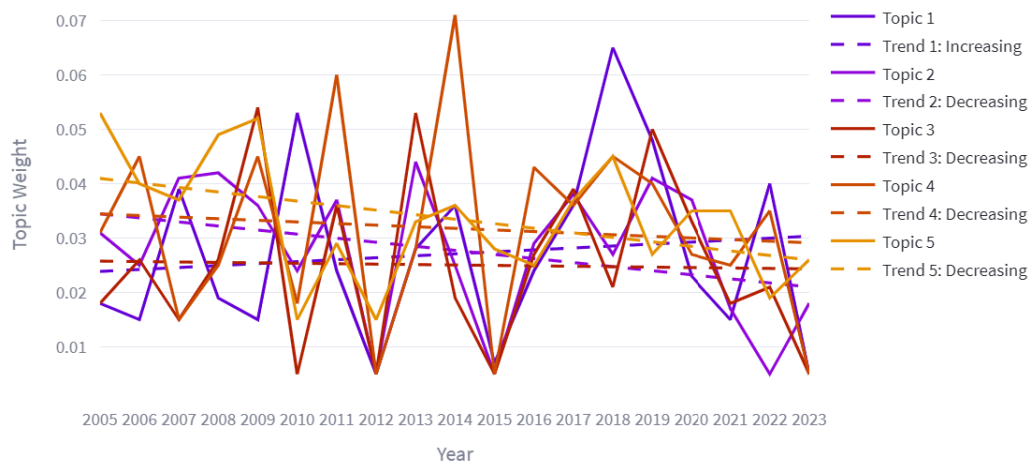


Fig. 6 Evolution of topics over the years

The trends (Fig. 6) associated with each topic indicate their evolving relevance over time. While Topic 0 shows an increasing trend, suggesting growing interest or research activity, Topics 1 through 5 exhibit decreasing trends, potentially indicating shifting research priorities or maturation of the respective fields. Notably, Topic 2 demonstrates very little overlap with Topic 3 (Fig. 7) in the LDA visualization, underscoring the distinct thematic boundaries between these two topics despite their related subject.



Fig. 7 Intertopic distance map



Fig. 8 Word cloud by year

4. Conclusions

The Italian Research Conference on Digital Library (IRCDL) Conference has become a key venue for the exchange of experiences and knowledge between researchers and information professionals engaged in the creation, organization, management and research of digital libraries. IRCDL was essential for the digital library

experience in Italy, and the analysis of research presented at the Conference over the last 20 years (2005–2023) can describe the evolution of the theory and practice of digital libraries.

The results of this study can facilitate further research to understand the knowledge structure of digital libraries in Italy in the context of Information science and relevant interdisciplinary and transdisciplinary research areas. This area of research has never been the subject of collaboration networks in this field, as demonstrated by IRCDL's twenty-year experience.

References

- [1] Lynch, C. (2005). Where do we go from here? The next decade for digital libraries. *D-Lib Magazine*, 11(7/8). www.dlib.org/dlib/july05/lynch/07lynch.html.
- [2] Agosti M., Catarci T., Esposito F. (2014) Pushing the Boundaries of the Digital Libraries Field Preface IRCDL 2014, Elsevier <http://www.dei.unipd.it/~agosti/papers/2014/2014-cover-and-preface-IRCDL2014.pdf>
- [3] Jelodar, H., Wang, Y., Yuan, C., Feng, X., Jiang, X., Li, Y., & Zhao, L. (2019). Latent Dirichlet allocation (LDA) and topic modeling: Models, applications, a survey. *Multimedia Tools and Applications*, 78(11), 15169–15211. <https://doi.org/10.1007/s11042-018-6894-4>
- [4] Manghi, P., Atzori, C., Bardi, A., Baglioni, M., Schirwagen, J., Dimitropoulos, H., La Bruzzo, S., Foufoulas, I., Mannocci, A., Horst, M., Czerniak, A., Iatropoulou, K., Kokogiannaki, A., De Bonis, M., Artini, M., Lempesis, A., Ioannidis, A., Manola, N., Principe, P., ... Pierrakos, D. (2023). OpenAIRE Graph Dataset (6.0.0) [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.8217359>