

The Role of Informatics in Software Engineering: Literature Reviews, Agenda and Software Informatics

Ira Monarch
Software Engineering Institute
4500 Fifth Avenue
Pittsburgh, PA 15213
001-412-268-7070
iam@sei.cmu.edu

Sheila Rosenthal
Software Engineering Institute
4500 Fifth Avenue
Pittsburgh, PA 15213
001-412-268-7846
slr@sei.cmu.edu

Rachel Callison
Software Engineering Institute
4500 Fifth Avenue
Pittsburgh, PA 15213
001-412-268-7725
callison@sei.cmu.edu

ABSTRACT

Literature reviews have been and are increasingly being merged with semi-automatic versions of bibliometrics and text analytics to extend library capabilities in the direction of performing informatics. Such extended capabilities allow special libraries to move in the direction of supporting and even performing informatic functions for bioinformatics or medical informatics. Recently a new subfield of informatics, software informatics, has been discussed that opens informatic opportunities for software engineering special libraries. The paper discusses a software engineering library that is using informatic techniques to support characterizations of customer demand landscapes that inform software engineering agenda. Sources analyzed go beyond published periodical literature to include organizational reports like budget justifications and, potentially, use of web harvesting. It is proposed that this use of informatic techniques is part of software informatics and because of the potential impact on how software is developed and used, may be part of software engineering as well.

Categories and Subject Descriptors

D. Software; H.3.1 Content Analysis and Indexing; K.6.3 Software Management.

General Terms

Management, Measurement, Documentation, Design, Human Factors.

Keywords

Software engineering, software informatics, bibliometrics, text analysis, informatics, special libraries.

1. INTRODUCTION

There are many characterizations of informatics. The paper characterizes informatics as the study, use and communication of information including its analysis and organization. Sub-domains of informatics focus on specific domains like chemical-, bio-, medical- geo-, social- business-, library- and recently software-informatics. The sub-domains all employ information technology to manage, process and analyze data and information pertinent to a given domain. Informatic sub-areas also approach information from various perspectives, individual, social, economic and cultural. Employing information technology from these perspectives can help to provide support for constructing and carrying out agendas of various disciplines (e.g., chemical engineering, bio-medicine, business re-engineering and software engineering) in a given domain. Nevertheless, work in informatic sub-domains is not always thought to be part of these engineering,

medical and business disciplines, even if the work is thought to be part of a discipline. Moreover, reviews of textual sources and their analyses, and library informatics more generally, are not always thought to be part of the informatics of sub-domains, and even less likely to be thought of as part of the disciplines corresponding to an informatic sub-domain. To open the possibility for extending software informatics, the paper proposes, as a subject for discussion, the idea of using informatic techniques to support learning software engineering agenda and informing those who could influence these agenda. The paper further proposes, again as a subject for discussion, that this extended version of software informatics be considered part of software engineering.

In order to discuss these proposals, it is important to get clear on the notion of agenda. According to Michael S. Mahoney, an historian of mathematics who also wrote extensively on the history of software engineering and computation, an agenda lies at the heart of a discipline [5]. Here Mahoney is referring to disciplines such as applied mathematics, computer science or software engineering. An agenda for Mahoney is

a shared sense among its practitioners of what is to be done: the questions to be answered, the problems to be solved, the priorities among them, the difficulties they pose, and where the answers will lead. When one practitioner asks another, "What are you working on?" it is their shared agenda that gives professional meaning to both the question and the response. Learning a discipline means learning its agenda and how to address it, and one acquires standing by solving problems of recognized importance and ultimately by adding fruitful questions and problems to the agenda. Indeed, the most significant solutions are precisely those that expand the agenda by posing new questions ... Most important for present purposes, the agendas of different disciplines may intersect on what come to be recognized as common problems, viewed at first from different perspectives and addressed by different methods ... [5]

This notion of agenda provides a backdrop for much of the discussion in the paper. However, it has typically been the case that only practitioners of a discipline, those on the inside, get to determine the discipline's agenda. Customers or users of what the discipline produces according to its agenda are almost never considered. This would seem to make most sense in pure mathematics whose products are papers containing mathematical proofs, read for the most part, by other pure mathematicians. However, this is not always the case. Even pure mathematics often gets applied and sometimes different disciplines mutually influence each others' agendas as has been the case for abstract algebra and computer science [5 and 7]. Nevertheless, on Mahoney's view, disciplinary agenda are influenced from the

outside only when the outside is another discipline part of whose agenda can be borrowed by the receiving discipline. What is not considered is an outside that is articulated in terms of what a user needs or demands from the discipline. A classic articulation of this view for software engineering is from Dijkstra who believed that software engineering is not predicated on user demands but rather on mathematics [3]. This paper will consider a very different view discussing how laypersons can influence specialist agendas [8] through their document representatives [9].

In what follows we describe a specific case in which informatic-based text mining and analysis is performed with the goal of agenda learning/influencing based on characterizing parts of the customer demand landscape of software engineering. Interspersed in the description of this case, we further discuss the relationships between informatics, software informatics and software engineering.

The next section introduces the Software Engineering Library at the Software Engineering Institute (SEI) as a special library supporting Software Engineering and goes on to discuss its role in document searches and analyses for better understanding of the software engineering demand landscape. This is followed by a discussion of informatic techniques that can be, and to some extent are being, merged with special library functions. These techniques in the form of text analytics applied to the documents collected are used as a basis for understanding the demand landscape and informing the subsequent building of a software engineering roadmap. The final section summarizes the case for the new field of software informatics and its role in software engineering and makes a case for including special library functions, bibliometrics and text analytics in software informatics because of their role in building roadmaps and understanding/influencing agenda.

2. A Special Library Supporting Software Engineering

The definition of a Special library, according to the Online Dictionary for Library and Information Science, by Joan M. Reitz is “a library that is established and funded by a commercial firm, private association, government agency, nonprofit organization, or special interest group to meet the information needs of its employees, members, or staff in accordance with the organization’s mission and goals.” [6] The SEI Library fits the definition. Established and funded by the SEI in 1986, its mission is to efficiently provide timely and relevant information to the SEI by maintaining an expert library staff and a strong collection in software engineering, computer science and related disciplines.

The SEI Library is staffed by a library manager, a reference librarian, an archivist, and a paraprofessional. In addition, the library currently has the good fortune to have the assistance of a recent Carnegie Mellon art design graduate who has created the information visualization poster shown in Figure 1. This poster is prominently displayed on the wall behind the reference desk in the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

iConference 2010, February 3-6, 2010, University of Illinois at Urbana-Champaign, IL, USA.

Copyright 2004 ACM 1-58113-000-0/00/0004...\$5.00.

library, showcasing both the library and the archive. The SEI Archive was created in 2005 under the supervision of the library’s departmental Director. The archivist works closely with the librarians and many of the services overlap, as illustrated above. One of the major goals of the SEI Library is to increase its involvement in the research required by the Institute’s software engineers. Therefore, the library focuses on the following areas of commitment: encourage collaboration on SEI projects; develop beneficial programs for SEI staff; and provide special services that will enhance library support for their research.

In order to understand the Software Library’s mission, the SEI’s mission must also be stated. According to SEI’s website (<http://www.sei.cmu.edu/about/>) its mission is to advance software engineering and related disciplines, though with an eye to ensuring that the development and operation of systems is predictable with improved cost, schedule, and quality. This means, according to the terminology of this paper, that the SEI’s agenda for software engineering focuses on certain kinds of systems and certain of their attributes. The SEI Library by supporting SEI’s mission supports the discipline of software engineering in this sense. In the case to be described the SEI Library was able to participate in performing extended informatic functions that informed SEI agenda learning/influencing. What will be described in the rest of this section and in following sections are the role the library played in performing the informatic functions, whet these functions were and what findings were produced. Any use the ED project or the SEI made of the findings will not be described.

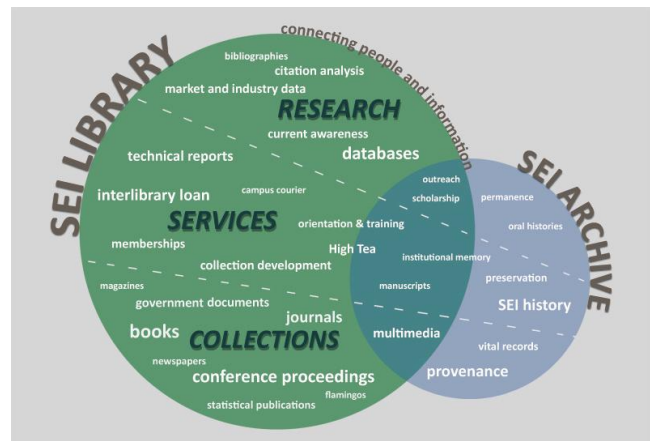


Figure 1: Work of the Software Engineering Institute Library

The SEI Library has participated in a project created by the new Executive Director of Interagency and Cyber Initiatives and head of the Acquisition Support Program (ED-ICI & H-ASP or ED for short) when she joined the SEI Staff in May 2009. Not long after the ED arrived, she informed the library staff that she was interested in everything *cyber* and was starting a new project¹ to locate all current information on this topic. The SEI Library participated in the project through one of its main functions, providing *current awareness*, including *table of contents* and *keyword alert* services.

The journals the ED selected for her table of contents current awareness included: *Defense News Early Bird* and *Daily News*

¹ This will be called the ED Project in the rest of the paper.

Roundup, The Economist, and the Wall Street Journal. In addition to these publications, the library staff provided her Tech America Headlines from Infoition, which are scanned every morning for information on *cyber*. This news service includes such publications as the *Washington Post, Information Week, NextGov, Wall Street Journal, Federal Computer Week, Business Week, Associated Press, New York Times, Technology Review, Computer World, CNET News, Government Computer News, and Time*, providing very thorough daily coverage on this topic.

Keyword alerts for the ED are run in *Academic One File*, which is a service provided by Gale Publications; Elsevier's *E.I. Engineering Village*, which contains two files, *Compendex and Inspec*; and UMI's (University Microfilm International's) *ProQuest Database*. The retrieval from these databases comes to the ED's email address directly, with many providing full text options. The ED shares important key and high impact articles with other members of her new project who include representatives from all SEI Programs, Services, and Functions. As certain specific phrases using the term *cyber* gathered more impact or focus through the media, e.g., the phrase, *cyber command*, the ED would request targeted searches to run on these specific phrases.

After being in her new position for only a few months, The ED created the concept for another type of library alert service. She asked if the library could, on a regular basis, search certain websites for the following specific terms: cyber warfare, cyber assurance, cyber component commander, 24th Air Force, and AF cyber security. In addition and also on a regular basis, the ED requested a similar set of alerts from the OSD Homepages for OSD NII, OSD USD-Intelligence, Office of the Director of National Intelligence, IARPA and DARPA to start.

One of the teams working on the ED Project devised specific literature search strategies for retrieving *cyber* related information focusing in on DOD, civil agencies and the business sector. The team wanted to identify and quantify future business opportunities for software reliant systems in *cyber environments*; establish the dollars to be spent; forecast business risk over the next three to five years for software reliant systems in *cyber environments*; and identify budgetary focus on such items as training, software, infrastructure, systems and platforms. They identified specific budget topics or keywords to query in the context of software reliant systems, software intelligent systems, and systems of systems in *cyber environments*. This team's dual role of identifying and quantifying business opportunities builds a bridge, in effect, between the special library function to support literature searches and the extended library function of using bibliometrics and text analytics to help the ED project interpret the results of the literature searches and make forecasts. The text analysis described in the section after next is applied to some of the documents found in the searches described above.

3. Bibliometrics and Text Analytics

A characterization of informatics was provided in the introductory section. Bibliometrics and text analytics will be introduced in this section.

Bibliometrics is a set of methods, usually statistical or mathematical, used to study or measure the attributes and relations of collections of bibliographic references or documents. Bibliometric methods are most often used in the field of library and information science. In this sense bibliometric methods can be seen as part of informatics. Citation analysis and content analysis

are commonly used bibliometric methods, though historically bibliometrics has devoted more attention to citation analysis. Bibliometric methods can be used to explore the impact of research fields on one another, the impact of a set of researchers or the impact of a particular paper. In each case, impact on practice is determined by bibliometric analysis of documents made available in other ways than through publishing vehicles such as journals.

Citation analysis has been used to evaluate the importance of a research article, the researchers themselves through their published work or a whole department in a university. Citation analysis is often used to identify the watershed publications in particular disciplines through quantitative means and the interrelationships between authors from different institutions and schools of thought. Such citation analysis depends on citation indices, such as Institute for Scientific Information's (now Thomson Reuters') Web of Science.

Some limitations on the value of citation data have been pointed out, especially in regard to quality ratings. For example, a low correlation has been found between peer evaluation of computer science groups and citation impact indicators of their papers [2]. In addition, the h-index (A scholar has an index of h if he or she has published h papers each of which has been cited by others) can be misleading [2]. These points are reinforced by a report that strongly cautions against over-reliance on citation statistics such as the impact factor and h-index [4].

Content analysis using semi-automated text analytic techniques can avoid some of the issues with citation analysis. It can identify the role that authors play in originating and/or contributing to the thematic structure of a discipline. The emphasis can therefore be on who introduces new ideas and applies them, rather than on sheer citation counts. Moreover, by focusing on the thematic structure of a discipline or an ultra large collection of documents, content analysis, or text analytics, can address questions concerning the *content* of the discipline or document collection that citation analysis cannot. This paper concentrates on the role of text analytics and its application to software engineering literature reviews, road maps and to other parts of software informatics because its value is underappreciated. However, citation analysis used critically either alone or in conjunction with content analysis also has value.

4. Applying Text Analytics to uncover a Software Engineering Demand Landscape

This section describes the application of semi-automated text analysis to two sets of documents being collected in a literature search, partly described above, to gain insights into the current software engineering demand landscape. One set of documents focuses on budget and financial summaries from the DOD, DARPA, DHS, DOE and US Military Services etc. that provide justifications for their technology and service acquisitions; the other focuses on descriptions of technologies and services that are seen as attractive or promising to actual or potential SEI customers. So far, 10,000 pages of budget and financial summaries and over 1000 pages describing promising technologies have been analyzed. The aim of these and other analyses is to support (1) understanding of the demand landscape of potential and actual SEI customers, (2) determining the extent to which current SEI services align with the demand landscape

and (3) making suggestions for better SEI alignment with demand and vice versa. This paper will focus mostly on 1 and a little on 3.

Two main questions guided the text analysis. (1) How well are software related concepts integrated into budget justifications of promising technologies that support cyber environments? (2) Do cyber concepts play a mediating role between software related concepts on the one hand and concepts involved in budget justifications and promising technologies on the other? In both cases the short answer, based on the text analysis and interpretation of concept maps, turned out to be only in a limited way. Further, in answer to the second question, the text analytic findings also indicated that the concepts of **data** and **networks** already do play a more important mediating role than **cyber** concepts and in the future could play an even larger mediating role. Approaching engineering agenda in terms of mapping the conceptual space of what is engineered is important because design is a primary activity in engineering, and the design of software, like the design of any artifact, is based on the meaning of an artifact negotiated by all design stakeholders including users and customers as well as engineers [10]. Concept maps generated by automated text analysis techniques can be useful in this regard (given that the right text sources are selected for analysis) to provide a snapshot of the design or “drawing together” of an artifact in all its interrelationships and complexity [10].

Automated text analysis produces what are called concept maps.² A concept map generated for the budget and financial summaries is shown in Figure 2. The themes are the large circles with labels inside the circles capturing the theme. The circles representing themes also contain light grey dots or nodes representing concepts (not labeled in Figure 2) and lines representing links between the concepts. Links between concepts are determined on the basis of co-occurrence of terms.³ Because concepts with different meanings can frequently co-occur, concepts in different themes are often linked. The intensity or thickness of the lines representing the links indicates likelihood of co-occurrence. The themes with the most concepts and links are ranked more highly in terms of Connectivity and overall Relevance to the contents of the collection of documents being analyzed. The top theme in terms of Connectivity and Relevance is **technological**. It is set at a benchmark of 100% that other themes are measured against. A table showing the ranking of the themes is included under the concept map in Figure 2.

The nodes representing concepts are also labeled (see the concept **technological** in the upper left part of the circle containing the theme **technological**⁴ in Figure 3). Concepts stand for more than the literal words or phrases that label a node. They also stand for an affinity list of terms that strongly co-occur with the term standing for the concept. The term standing for the concept is automatically chosen on the basis of its more strongly co-occurring with all the other terms on the affinity list than any other term on the list. Themes are clusters of concepts that have more similar co-occurrence patterns with each other than concepts

in other thematic clusters have with each other. Themes are automatically named by the concept that has a co-occurrence pattern more similar to the co-occurrence patterns of all the other concepts in the cluster than does any other concept in that cluster.

The meaning of a theme derived from a collection of documents, for example **technological** in the budget and financial document collection, is determined by the concepts in the cluster that constitute the theme. In the case of **technological**, this list of concepts includes: **renewable** (technologies, e.g., **ethanol**), **engine**, **aircraft**, **hyperspectral/polarimetric** (remote sensing) **vehicles**, **amplifiers**, **weapon**, **biology**, **catalytic**, **chemicals**, **nuclear**, **medical**, **environmental**, **missile**, **algorithms**.

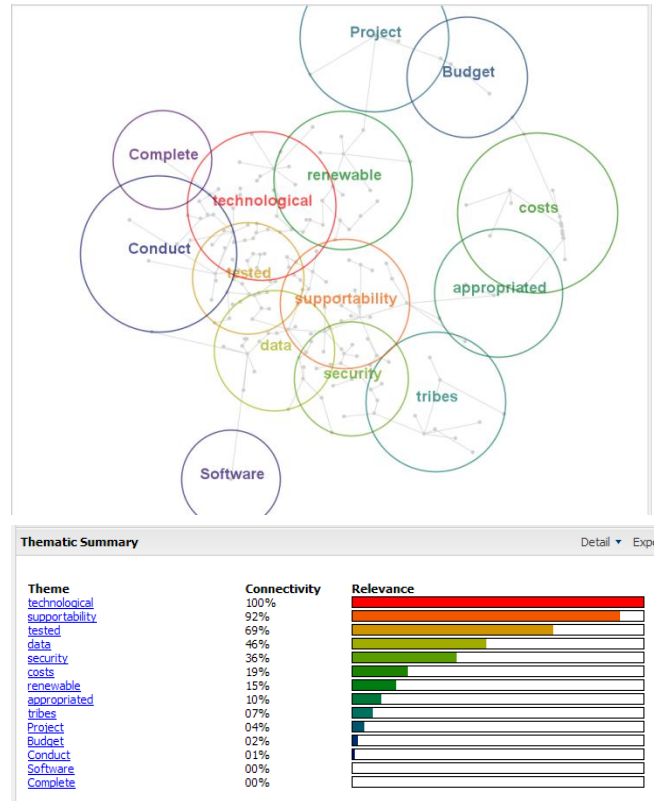


Figure 2: Concept Map and Rankings generated from Budget Documentation

None of the concepts seem to have much to do with software except **algorithms**. Though **Software**⁵ and **algorithms** occur nearly the same number of times, **Software** is not well connected with any other concepts, having a 4% or less likelihood of co-occurring with any of the top 300 concepts in the budget and financial collection, whereas **algorithms** has a likelihood of co-occurring of up to 21% with **polarimetric**, 14% with **hyperspectral** and 13% with **model-based**. Overall, **Software** has a 7% relevance to the other content in the document collection and the concept **algorithms** has an 8% relevance. In contexts outside of this document collection, **algorithms** can be understood to have a tight meaning relation with **Software**, e.g., an algorithm is basically the logic implemented in software by software developers. However, connections between **algorithms**

² The automated tool used in these analyses is called Leximancer. It has been available in an evolving form for almost 10 years.

³ Two terms co-occur if they both occur within a two sentence block. A terms consists of all word variants.

⁴ Note when a term stands for a concept is in bold font and when it stands for a theme it is bold italic font.

⁵ Software is capitalized whether as theme or concept because it more often appears that way than not in the budget texts.

third at 23% relevance of all the name-like concepts. On the other hand, it is part of a second tier theme, **computer**, that is just at 5% connectivity and relevance. This suggests that the concept of **Software** does not stand out from the concept of **computer** in discussions of promising technologies. What stands out more in these discussions are the themes **security**, ranked at 100% connectivity, **systems**, ranked at 45%, **technology**, ranked at 32%, **operations** at 29% and **Risk** at 24%.

The last two significantly overlap with **security**. **Cyber** is a concept in the **operations** theme with a rather high relevance ranking of 30%. **Cybersecurity** is a concept in the **security** theme but is ranked at only 9% relevance. **Cybersecurity** is a **policies** oriented concept concerned with **infrastructure** issues but with only very minimal connection to software related concepts, similarly with **cyber** (both are at 1% or less likelihood of co-occurrence with such concepts).

Perhaps a better way to mediate operational and promising technological concepts to software related concepts in this collection is through the concept of **network**, at 41% relevance and **data** at 39% relevance. Both are in the **systems** theme. **Data** is related to **network** at 16%, **compiler** at 14%, **dead code** at 13%, **cloud** at 12%, **World Wide Web** at 8%, **software** at 7%, **cyber** at 6% & **cybersecurity** at 5%. **Network** is related to **GIG** at 33%, **computer** at 23%, **data** at 16%, **cyber** at 16%, **World Wide Web** at 8%, **cybersecurity** at 8%, **dead code** at 4%, **cloud** at 4%, **software** at 2%. **Data** and **network** are better integrated with software related concepts than **cyber** and even provide a bridge to the latter from the former. Although **Cyber** and **cybersecurity** are more highly ranked concepts in this collection than in budget descriptions, they do not do a good job of mediating or linking operational concepts and other promising technology concepts to software related concepts. **Data** and **network** do a somewhat better mediating job, but *again better strategies for negotiating the meaning of .software related concepts are needed.*

5. Conclusion: Software Informatics and Software Engineering

Recently, software informatics has been defined as the science of information, practice, and communication around software that studies the individual, collaborative, and social aspects of software production and use, spanning multiple representations of software from design, to source code and to application [1]. One of its key characteristics is to treat software in terms of the information flows that help produce it, including design specifications, source code, documentation, software libraries, applications and services, code repositories, revision histories, and more [1]. Treating software in terms of such information flows enables it to be studied using the methods discussed in this paper, most prominently text analytics. In addition this characterization of software informatics can be broadened to include using these techniques to study and influence how users and customers of current software products understand software and what they want to do with it. To begin this work, well over ten thousand pages of governmental budget justifications for purchasing large scale typically software intensive systems were analyzed. In addition, over a 1000 pages of published literature on promising software technologies were also analyzed. We found that development, operation, and maintenance of software for these large scale systems was not much considered in the budget justifications and even marginalized in the promising technologies collection. Our

basic recommendation to the ED project was to find better ways to find new ways of negotiating the meaning of software with the large varied group of software engineering stakeholders, including users and customers of large-scale software intensive systems, that were the sources of the documentation analyzed.

There are counter arguments to this approach to understanding and influencing software engineering agenda represented by the views of E. W. Dijkstra. He advanced a rather Platonic argument that progress in software engineering is not predicated on user demands but rather on mathematics (usually thought of as the pinnacle of user unfriendliness), i.e., on all the mathematical power and elegance that software can muster [3]. Dijkstra states that a *well documented* [author's emphasis] program is an object that is logically isomorphic with a constructive mathematical proof. This statement, even if not technically democratic, could very well be an hypothesis of software informatics. In other words, software informatics can study both the demand landscape of software use and software as a mathematical structure. Either way it seems to be critically relevant to software engineering.

6. REFERENCES

- [1] Jones, M. C. and Twidale, M. 2009. Software Informatics? Isociety Conference, University of North Carolina at Chapel Hill.
- [2] Mattern, F. Bibliometric Evaluation of Computer Science – Problems and Pitfalls. 2008. Presented at European Computer Science Summit – ECSS 2008, 9-10 Oct. 2008, Zurich.
- [3] Vlissingen, R. F. van. Interview Prof. Dr. Edsger W. Dijkstra, Austin, 04-03-1985. E. W. Dijkstra Archive the manuscripts of Edsger W. Dijkstra, 1930–2002.
- [4] Adler, R., Ewing, J. (Chair) and Taylor, P. Citation Statistics. Joint Committee on Quantitative Assessment of Research, A Report from the International Mathematical Union (IMU) in cooperation with the International Council of Industrial and Applied Mathematics (ICAM) and institute of Mathematical Statistics (IMS). June 2008.
- [5] Mahoney, M.S. The Structures of Computation and the Mathematical Structure of Nature. 16 June 2006. paper delivered at the 21st International Workshop on the History and Philosophy of Science, The Origins and Nature of Computation, Tel Aviv and Jerusalem.
- [6] Reitz, Joan M. ODLIS Online Dictionary for Library and Information Science. Special Library.
- [7] Eilenberg, S. Automata, Languages, and Machines (2 vols., NY: Columbia University Press, 1974), Vol. A, xiii.
- [8] Callon M, Lascoumes P, Barthes Y (2009) Acting in an uncertain world: an essay on technical democracy. The MIT Press, Cambridge (Inside Technology Series).
- [9] Callon, Michel, Law, J. and Rip, A. (eds) (1985), Texts and their Powers: Mapping the Dynamics of Science and Technology, Macmillan, London.
- [10] Latour, B. A Cautious Prometheus? A Few Steps Toward a Philosophy of Design (with Special Attention to Peter Sloterdijk), Keynote lecture for the Networks of Design* meeting of the Design History Society, Falmouth, Cornwall, 3rd September 2008.