



HAL
open science

Knowledge Cartography and social network representation Application to collaborative platforms in scientific area

Michel Plantié, Pierre-Michel Riccio

► **To cite this version:**

Michel Plantié, Pierre-Michel Riccio. Knowledge Cartography and social network representation Application to collaborative platforms in scientific area. Sixth International Conference on Signal-Image Technology and Internet Based Systems, Dec 2010, KUALA LUMPUR, Malaysia. pp.100-110, 10.1109/SITIS.2010.45 . hal-00807943

HAL Id: hal-00807943

<https://hal.science/hal-00807943v1>

Submitted on 4 Apr 2013

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Knowledge Cartography and social network representation

Application to collaborative platforms in scientific area

Michel Plantié, Pierre-Michel Riccio,
Laboratory LGI2P, Ecole des Mines d'Ales EMA,
Parc Georges Besse, Nimes, France
Michel.plantie@mines-ales.fr, pierre-michel.riccio@mines-ales.fr,

Abstract—Our research revolves around collaborative platforms entirely dedicated to research activities for several scientific organizations. Here, researchers from different domains interact and exchange information using our platforms as the common ground involving new concepts, methods and services to encourage collaborative work for their research activities.

The work is based on the co-operation and collaboration between scientific specialists and engineers building the platform. In our proposal we attempt to capture the individual information and then build social network representations to build a group representation of knowledge, attracting and encouraging people to participate in collaborative tasks.

Keywords-Frameworks and Methodologies for Collaboration, Tools for Collaborative Environments, Coordination, Cooperation and Collaboration, Collaboration Enabling Technologies, Platforms for Collaboration, Collaborative Knowledge Management.

I. INTRODUCTION

Since several years we see organizations trying to consolidate their efforts and enhance collective efficiency in human and computer systems. We have launched several scientific projects gathering collaborative efforts in a multidisciplinary approach to cope with scientific and societal issues. We already have some results in using computer based collaborative platforms [10] which are very encouraging to go several steps further and enhance social network information to strengthen the projects efficiency.

Collaboration literally means an action or a work completed in common with two or several persons. It is a group activity where individuals unite to form alliance or unions with an intention to attain an objective. We can find evidence of group activity in many living beings ensuring early completion of tasks and better security against possible dangers. Every member in the group experiences better results when tasks are being accomplished in a co-ordinated group than attaining it individually.

This meaning of collaboration as a form of collective intelligence has been stated [1] as a hypothesis relative to the capacity of a group of cognitive actors and artificial agents to reach a higher level of performance than the addition of individual ones. Based on this definition, collaboration appears as a form of coordination [2] which is managing dependencies between activities in a group.

Collaboration and sharing is relatively a developing area in research introducing a methodology for the planned capture and re-use of organizational knowledge. Successful application of collaboration practices involves the understanding and constructive use of organizational learning and information flows within the organization. The concept of collaboration is becoming more important in the evolving context of global network, thus placing the user at the centre of a collective device. Collaborative work can either be of the nature, where each group member is involved in every activity with the work being highly interactive or where each group member is given an individual task.

In our proposal, we start by building a user friendly platform with built-in applications and features to guide the users in efficiently using the platform for information retrieval and management and also encourage them to share information about their research activities with the rest of the registered members belonging to different research communities, on the platform. In this paper we mainly focus on the utilization of knowledge cartography which acts as a mirror of the group activity carried out on the platform both in terms of research teams and members involved.

The following sections of this paper are organized as follows: In section 2 we give an introduction about the creation and use of our application platforms. In section 3, we propose novel approaches for information retrieval and management to boost collaborative thinking in group oriented working environment. This is basically a knowledge representation model called the Extended Semantic Network (ESN). In section 4, we introduce the concept of knowledge map, we have implemented on our platforms. In section 5 we derive social networks from

knowledge maps. In section 6 we discuss about additional tools helping the use of social network information. Finally, in the last section we conclude this paper with the summary of the work completed, advantages and results studied and direction for future work.

II. COLLABORATIVE PLATFORM IN LARGE SCIENTIFIC PROJECTS

Several years ago several large scientific projects wanted to gather scientific knowledge in a collaborative platform. This was particularly the case with the French health research organization INSERM which wanted to have an overall view of all the research domain and resources done in the country and asked for a framework overall agreement to be able to share very different scientific knowledge.

This was also the case for the CARNOT-MINES association in charge of gathering and enhancing scientific exchanges between researchers of partner scientific organizations to impulse new scientific areas. For example the scientific domain of psycho-sensitive material is a very new research domain requesting scientific knowledge in various disciplines which were until now very separate: environment questions, micro and nano-technologies, micro-mechanics, color technologies, sounds and electronics, etc.

Finally, in 2001, a multi-field inciting program was initiated by the CEA (Commissariat à l'Energie Atomique in France) to stimulate the emergence of a community of experts and young researchers around a stake, mainly touching the public health and the environment. It mainly handled the question of understanding the mechanisms of actions of heavy metals and radio nuclides on the various levels of organization of the living beings. This research fundamental program extended in 2004 to four research organization partners (CEA, CNRS, Inra and Inserm) involving some of well established French research laboratories is a multi-field project which involved a great number of researchers from different research disciplines like biology, chemistry, medicine and physics.

Program management team main tasks were to manage and provide all necessary tools and applications for easy interaction among the vast community of researchers involved in the program. Thus at length favoring and supporting communication leading to information exchange between actors (researchers) of the program: grounded on a collaborative platform.

The idea of this platform rose from several questions on nuclear toxicology. Answering these questions eventually made it possible to have a thorough knowledge of the impact of the anthropical activities on human health and its environment. The recent studies and observations made on impact of toxicology on mankind are very few in France as well as abroad. Some of the field and

methodologies used scarcely integrate the projections of the revolutionary techniques proposed by genomic and biotechnology domains.

Contradictory to these observations, research in biology and genetics is developing at a vertiginous speed and all the resources of post-genomic available to renew the field of toxicology are highly neglected in biology. In order to contribute to this society and human health related questions, fresh impulse was given to this research within the framework of an inciting multi-field Program heading now "Environmental Nuclear Toxicology: ToxNuc" [3].

A *Scientific Objectives of the Program:*

The program mainly focuses on the question of including the mechanisms of actions of heavy metals and radio nuclides on the various levels of organization of living organisms (molecular, cellular, bodies and fabrics, whole organizations) in order to propose preventive technical solutions, provisions of effective monitoring and solutions to decontaminate these elements distributed in certain compartments of the tropic chain.

These chemical elements were primarily identified in dialogue with various actors involved in nuclear die, in industry and in research, and a list of interesting elements were identified and brought out. These elements are listed as follows: tritium, beryllium, boron, carbon, cobalt, selenium, strontium, technetium, cadmium, iodine, cesium, lead, uranium, plutonium, americium, zinc, copper and nickel.

The state of the art on this domain was very weak to be used in identifying these elements considering the fact that very little work been carried out in this field. This resulted as a primary factor in focusing the studies on two fields called environmental toxicology and human toxicology. In these two fields, it is a question of being interested: with the biological effects as of these substances and the molecular and cellular mechanisms of transport, of toxicity and de-toxicity. This subsequently leads to the issue of proper co-ordination between researchers from the two different fields.

- In environmental toxicology, studying transfer mechanisms from geo-sphere towards biosphere by bacteria and plants means helps to imagine applications to decontaminate terrestrial or watery environments.
- In human toxicology, imagining applications for contamination treatment by targeting studies on uranium and plutonium helps a lot. Organizations on which these studies are focused are preferentially those whose genome is sequenced i.e.: bacteria, yeast, arabidopsis, human cells,

mouse, and rats. This approach allows massive use of genomic methods.

B. Mobilization and Organization of the Program

Human means does exist but it's a question of mobilizing them on some clearly given scientific objectives. Program management team organised meetings in order to include some of the major researchers in the biological, chemical, physical and informative fields. Committees were organised and co-ordinators or heads for each research project were chosen. Several researchers geographically dispersed were brought into contact through this platform.

Registered members were over 700 researchers from diverse fields working on topics related to nuclear toxicology. In a very short period, vast information was collected on the platform. Now the problems like efficient data management, easy information retrieval and safety about sharing one's research results with other members of the platform known only through professional contact because of similar research interests, needed to be solved.

A positive response to this question would automatically encourage researchers on the platform to exchange information and discuss the research requirements and observation with other members of the community. Thus leading to a collaborative proceeding to resolve issues regarding nuclear toxicology. Precisely, for this requirement, we use knowledge cartography to provide a visual image of the collaborative work done on the platform. This helps members to strengthen confidence on each other and enhance co-ordination through utilization of our proposed tools on the platform. Using these tools will actually give a global knowledge view about research work carried out on the platform and as well provide information about researchers involved and their actual domain of interest. This will in turn boost confidence and encourage will for collaboration.

III. EXTENDED SEMANTIC NETWORK FOR EFFICIENT INFORMATION RETRIEVAL AND DOCUMENT CLASSIFICATION

Extended Semantic Network [4] is an innovative tool for knowledge representation and ontology construction, which looks for sets of associations between nodes semantically and proximally as opposed to present method of keyword association. Our goal is to achieve a semi-supervised knowledge representation technique with good accuracy and minimum human intervention, using heuristically developed information processing and integration methods. This model is built with information and research documents shared on the platform ToxNuc.

A. Hybrid Approach – Extended Semantic Network (ESN)

The basic idea of Extended Semantic Network is to identify an efficient knowledge representation and ontology construction method to overcome existing constraints in information retrieval and classification problems on ToxNuc platform. To realize this, we put our ideas into practice via a two phase approach. The first phase consists in processing large amount of textual information from the platform using mathematical models to make our proposal of automatic ontology scalable. The second phase consists in examining carefully and efficiently the various possibilities of integrating information obtained from our mathematical model with that of the manually developed mind model.

The first phase of our proposal is carried out by realising a lattice of words mathematically computed using different statistical and clustering algorithms, thus creating a proximity network computationally developed, essentially depending on word proximity in documents. The second phase is developed with a heuristically developed network extension method using outputs from the mathematical approach. This is achieved by considering the manually developed semantic mind model as the entry point of our concept network.

Here, the basic idea is to develop an innovative approach obtained by combining features of man and machine theory of concepts, whose results can be of enormous use in the latest knowledge representation, classification, retrieval, pattern matching and ontology development research fields. In this paper we discuss and highlight our methods used for information processing and integration to get a new visualising method for knowledge representation [5] and ontology construction [6]. This will help the ToxNuc researchers to easily retrieve information on the platform and will encourage information sharing.

1) Proximal Network for Efficient Data Processing

Proximity is the ability of a person or a thing to tell when it is near an object, or when something is near it. This sense keeps us from running into things and also can be used to measure the distance from one object to another object. The simplest proximity calculations can be used to calculate distance between entities thus avoiding a person from things he can hit. Proximity between entities is often believed to favour interactive learning, knowledge creation and innovation.

The basic theory of proximity is concerned with the arrangement or categorisation of entities related to one another. When a number of entities are close in proximity a relationship is implied and if entities are logically positioned; they connect to form a structural hierarchy.

This concept is largely used in medical fields to describe human anatomy with respect to positioning of organs.

Our Proximal Network Prototype model is built based on this structural hierarchy, of word proximity in documents. This approach is mainly employed to enable processing of large amount of data in a considerably small time. Another important aspect of this approach is its ability to automatically process input data into a network of concepts interconnected with mathematically [7] established relations.

To build this prototype we systematically process three phases for identifying our data and build the final network. We first start with a set of documents related to 3 major fields out of the 15 fields in the nuclear environmental toxicology domain, coming from the program ToxNuc. The documents obtained are first converted into simple text format using an external converter and are use as input into our first stage called the pre-treatment process. This process is carried out in 2 different approaches, one for identifying the most significant words and the other for eliminating hollow terms. Here, in this process the input document is processed in several stages and a word frequency matrix is created with rows representing words and columns representing document.

This program is primarily concerned with physical distance that separates words. Currently, we have successfully processed around 3423 words computing their actual physical occurrence. We have been able to successfully build a proximity network of 50,000 word pairs. Each of these word pairs is related by using the value obtained from the prototype and is visualised using the simple UML link of association.

This data processing method in itself can be independently used for data processing and representing knowledge in various domains. The small time taken for processing huge amounts of data is an important aspect to get scalable processes for constructing ontologies representing multiple domains.

2) Semantic Network Prototype

Semantic Network [8] is basically a labelled, directed graph allowing use of generic rules, inheritance, and object-oriented programming. It is often used as a form of knowledge representation. It's a directed graph where vertices represent concepts and edges represent semantic relations between concepts. The most recent language to express semantic networks is KL-ONE [9].

Nodes are labeled and edges are single labeled relationships between semantic nodes. Further, more than one relationship may co-exist between a single pair of connected words: for instance the relationship is not necessarily symmetrical and there can be relationship between nodes through other indirect paths. Technically a semantic network is a node- and edge-labeled directed graph, and it is frequently depicted this way.

The scope of the semantic network is broad, allowing semantic categorization of a wide range of terminology in multiple domains. Major groupings of semantic types include organisms, anatomic structures, biologic functions, chemicals, events, physical objects, and concepts or ideas. Links between semantic types provide the structure for the network and show important relationships.

In our semantic network prototype we reuse documents pertaining to each field of research in the program ToxNuc and then choose a set of concepts most significant to the field in consideration. This has been achieved with help of researches of ToxNuc who helped us in identifying them. This list of concepts pertaining to each field is given to each specialists who in turn rate each concept with respect to its importance in representing the field.

We then choose the first 50 concepts most representing the field from the above list and resented this list to people who were either specialists or people possessing good level of knowledge in each of these study area accompanied with our relational links. All the links used in connecting a node are based on the UML links, consisting of four different types of links for connecting these concepts. Thus the concept network is built based on the meaning each concept pair shares.

They have been currently chosen on an experimental basis, after proper consideration and analyzing the requirements of our approach. We start with our domain name representing the super class in our approach. The super class is then connected to its subclasses based on the category of the relation they share, which can be chosen from four links representing the simple UML links of association, composition, instantiation and inheritance.

the user requirements thus helping the user in efficient data search, retrieval, management, and sharing.

Some major points we hope to achieve through this method of knowledge representation network are:

- Make construction of semantic based concept networks cost effective by campaigning minimum human intervention. In turn reducing the construction time using mathematical models.
- Identify a good balance between mind and mathematical models to develop better knowledge representing networks with good precision and high recall.

IV. KNOWLEDGE DOCUMENT GRAPH AND ITS APPLICATION ON COLLABORATIVE PLATFORMS

As the size of digital libraries (more specifically in our case; research documents shared on the platform) grows in size/ number, the lack of adapted tools for browsing in these libraries appears more and more crucial. For instance, in the context of ToxNuc project, such a tool should be able to retrieve from the library not only the documents the user expects but also documents as suggested by the tool, unexpected for the user like information coming from other disciplines involved in the project, but dealing with the same scientific concepts.

Such an information retrieval tool is clearly semantic-oriented and the one we have proposed is based on knowledge maps. Building a knowledge map consists in representing contextual information, i.e. knowledge, using a visual metaphor chosen according to the future usage of this map. Here we want to represent scientific documents in the context of the nuclear toxicology domain. Hence, ontology of this domain will be represented on a map, linked to documents.

A “reference frames of knowledge” type collaborative work platform is intended to help a scientific community to develop its collective processes: presentation of researchers and teams, presentation of program, capitalization of information and results, sharing knowledge, internal communication, filing institutional documents, joint workspaces, exchange forums, and diffusion of information to general public. Advanced functions of content dynamic cartography (evolutionary trees and matrices) make it possible to follow evolution of data bases.

Each researcher registered in the program is a contributor authorized to deposit documents, consult filed documents and communicate with the other researchers. A follow-up and a report at all stages are done every six months. We have accumulated feedbacks from the platform users on our provided tools and have analyzed its influence on researchers to co-ordinate and work collaboratively.

The scientific assessment of the period 2001-2003 is as follows: 79 publications with a factor of average impact of 4,17 including 6 publications with a factor of impact higher than 10 (average: 14,70); 4 articles of synthesis in referred works; 8 cards of synthesis summarizing the principal results by chemical element studied and 4 patents deposited. The human assessment is more difficult to quantify however in bringing together every Net between biologists and chemists, in the broad sense, has been achieved.

This one led for example to the appropriation of the analytical step by the biologists and of taking into account the complexity of living organisms world by chemists. Without the detailed attention of all the scientific components, technical and administrative support of the CEA, this program could not have become a success. On this same set of themes, and with the same rigor in selection of projects and their follow-up, we will continue our adventure with this project. The program “Environmental Nuclear Toxicology” was initiated in 2004 for three years duration by selecting best teams of the four organizations of research partners.

Recently, we have conceived another new type of graph able to represent distance between documents and able to take into account any document from the platform whatever the application domain. To compute the distance between two documents we first use natural language processing techniques to “lemmatize” each word and documents, to be able to filter by POS tags words used for distance computations. All documents types are taken into account: pdf, text, docs ... All documents are represented in our system as bags of words. Only common nouns are considered to compute distances between documents since mainly nouns express the general meaning in a sentence. This filtering on POS tags is a very important feature to have a very good computation time. The distance then is computed by using the Jaccard distance [11]: the number of words present in both documents divided by the number of total words for the two documents. This distance is quite simple but gives very good results. We have checked the results by asking users the relevance of these new knowledge graphs. It has very good impact on users. For example in our corpus we have a domain related to a book edition on nuclear toxicology. With this type of graph we were able to show that each one of the book chapters were very different semantically speaking as the documents related to each chapter appeared very much gathered by chapter except for one chapter which was dedicated to a transversal topic and appeared linked to all others.

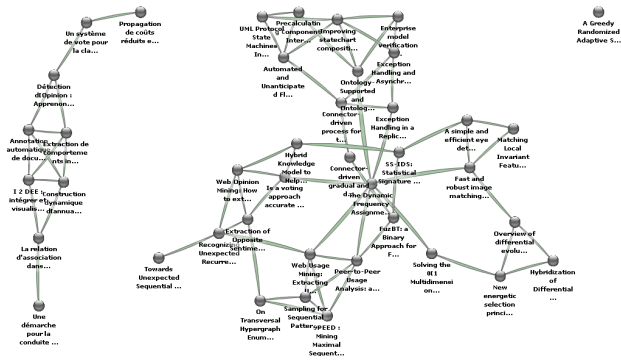


Figure 4. A document knowledge graph computed with the Jaccard distance

In figure 4 you can see a graphic copy of such a graph showing all documents of a domain. This graph is used to show documents graph for several domains at the same time. This graph reveals some features that do not appear clearly in other views of the platform. For example for a particular document, we see all its neighbors and among them we can perceive that older versions of the same document are present on the platform. Even if this document was visible in the platform, it was not so obvious to see the redundancy. It is then possible to decide to eliminate or put in some different folder such old documents.

V. SOCIAL NETWORKS

Each document is related to one or more authors. So we derive from document graph a new very interesting graph: the author graph or the social network of the domain! To draw a graph between authors we consider that each author is associated with his documents (the documents the author has written). To compute distance between two authors we simply consider the minimum distances between documents of the two authors respectively. This distance reveals clearly
 This social network graph is very useful to have hints on who is working on the same topics in the platform. It reveals obvious situations like persons working on the same topics. But it shows also unexpected results such as persons from different domains that are using the same concepts and topics in their job, even if they are considered as working on different areas. It is used to make recommendations on what documents are useful for one person and propose new avenues of interest, or suggest some persons to share their common experiences to have a better efficiency. In figure 5 we see a copy of such a social scientific network representing semantic links between researchers based on shared platform documents.

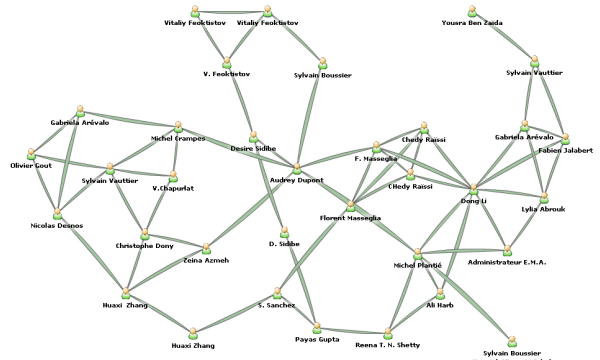


Figure 5: A social network graph derived from the document graph with the Jaccard distance

VI. OTHER COLLABORATIVE TOOLS

Social network information may be derived from several tools helping users to share collaborative knowledge.

A. Full text research engine

Our social network and document knowledge map is associated with an elaborated research engine able to filter information for semantic graph visualization. Users ask for a particular request on the collaborative platform.

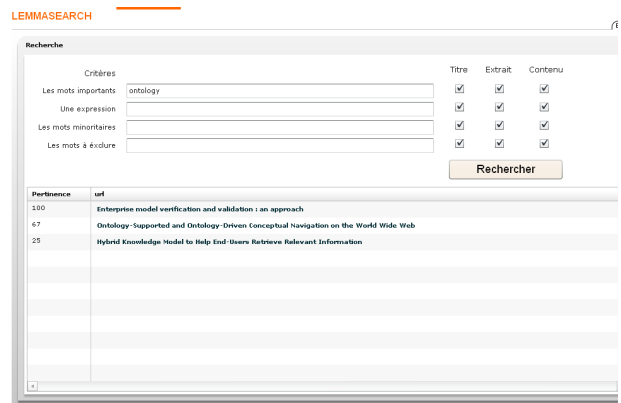


Figure 6. Research engine view

Answer is shown in a list (see figure 6). Answers may be used to draw a new document knowledge map showing graphically answers to the request. This is a personalized knowledge map corresponding to users' choices.

B. Personal workspace

We consider in our platforms that each member may customize his workspace and especially choose among all documents, those who are to be memorized in their personal space. This personal information is very useful to modify user profile and alter the document knowledge

map according to the personal preferences. In the same logic, the social network graph is modified to take into account the users preferences and show relations between people more relevant for the user.

C. Geographic and statistic maps

In order to know who is working in which geographical area, we have developed a new type of map indicating the number of each member or scientist of the platform working in a particular geographic location. This map also allows filtering geographic map by scientific domain and gives statistic data on scientific domains and research.

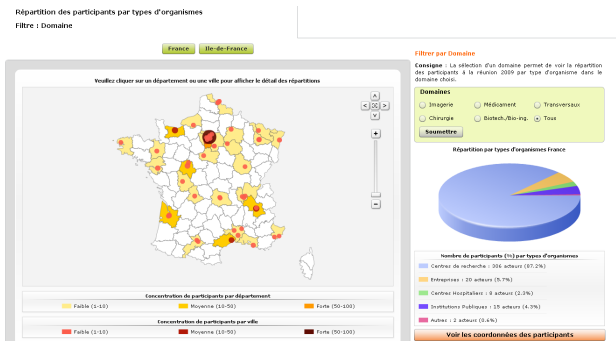


Figure 7. Geographic graph

VII. RESULTS AND FUTURE WORK

In CARNOT-MINES project, more than 100 scientists use the platform and 5 scientific collaborative projects have been launched to share very different scientific skills and knowledge. These new projects were made possible only by sharing in a efficient way knowledge of several different scientific domains.

In INSERM project, more than 200 scientists use the platform and 3 scientific collaborative projects have been launched to share very different scientific knowledge in health, medicine and biology.

In the organization of Tox-nuc program, twelve scientific projects have been selected. Each project controlled by one or more coordinators includes various specialists: biologists, chemists, doctors, physicists, pharmacists. The work flows involved in the program are 99 positions of statutory personnel primarily from CEA and also from the CNRS, Inra, Inserm, the combined laboratories of the CEA. The Program finances 30 post doctorates and 15 doctorates.

These 150 personnel correspond in fact to more than 250 researchers established in several areas and related to various operational directions of the CEA. These source data encouraged members of the program management team to install communication and management tools making it possible to create a community around the program. These tools are as follows: Newsletter - the Letter of the Nuclear Program Toxicology is a monthly

recto-back which is used as a bond between the researchers of the projects and allows a fast circulation of information useful to all. It is also an external tool of communication towards the directions of CEA and scientific and industrial partners.

Our platforms are efficient in sharing scientific knowledge. They provide collaborative help to share new ideas and results and we will continue in our future work to develop more social network statistics and analysis. Moreover we will develop tools to store, share knowledge and make new scientific knowledge to “bubble” from our platforms.

REFERENCES

- [1] JM Penalva, Monique Commandre, Typology of collaborative working variations around collectives in action, in Intelligence Collective, Les Presses Mines Paris, France.
- [2] Thomas W. Malone and Kevin Crowston, The Interdisciplinary Study of Coordination, ACM Computing Surveys, 1994 (March), 26 (1), 87-119, USA.
- [3] M T Ménager, Environmental Nuclear Toxicology program : how to federate a scientific community around a societal stake, Intelligence Collective Partage et Redistribution des Savoirs, Nîmes, France, septembre, 2004.
- [4] Reena T N Shetty, Pierre M Riccio and Joël Quinqueton. Hybrid method for knowledge processing, integration and representation. IEEE-IRI 2006 proceedings, September 2006, Hawaii, USA.
- [5] J.F Sowa, Knowledge Representation: Logical, Philosophical, and Computational Foundations, Brooks Cole Publishing Co., Pacific Grove, CA, 2000.
- [6] Natalya F. Noy and Deborah L. McGuinness, Ontology Development 101: A Guide to Creating Your First Ontology, Ontology Tutorial, Stanford University, Stanford, CA.
- [7] J.F Sowa, Conceptual structures: information processing in mind and machine, Addison-Wesley Longman Publishing Co., Inc, Boston, MA, 1984.
- [8] M.R Quillian, Semantic memory. M Minsky, Ed, Semantic Information Processing. pp.216-270. Cambridge, Massachusetts: MIT Press, 1968.
- [9] J Brachman, L Deborah, McGuinness, F Patel-Schneider, A Resnick Living with CLASSIC: When and How to Use a KL-ONE-Like Language, Special issue on implemented knowledge representation and reasoning systems Pages: 108 – 113, ACM Press, NY, USA, 1991.
- [10] Reena Shetty, Joël Quinqueton, Pierre-Michel Riccio, Jean-Michel Penalva, Jean Villerd, Collaborative Platform Using Knowledge Cartography - ToxNuc-E: CTS'07: Collaborative Technologies and Systems, Orlando, Floride, USA, 2007.
- [11] KAUFMAN, L., P. ROUSSEUW, et J. B. P. (1990). Finding groups in data: An introduction to cluster analysis. WILEY Interscience.