# The Science Library: A Controlled Natural Language driven Science Gateway

Eunjin Lee*, Dave Braines*†

*Emerging Technology, IBM Research, UK,
†Crime and Security Research Institute, Cardiff University, UK

*Abstract*—The International Technology Alliance in Distributed Analytics and Information Science (DAIS-ITA) is a research program conducting fundamental research across a consortium of organizations from academia, industry and government in the US and UK. A key output of the program are the academic publications, and from these a rich social and topical network between authors and organizations emerges. To capture and convey this we have created the publicly available Science Library as a user-centric, interactive portal. The Science Library is a Controlled Natural Language (CNL) driven research gateway, which allows the community to explore and query the publications and networks through an open, web-based application. The data is represented through interactive visualizations along with the ability for users to query the model using a natural language conversational interface. The CNL based approach models the data through concepts, properties and relationships which are defined using CNL and are therefore both human readable and directly machine processable. This captures complex semantics in a simple format and enables non-technical users to participate in the continuous improvement of the data model behind the Science Library application. This paper presents the features, implementation and design considerations of the Science Library and the underlying CNL implementation.

*Keywords*—Controlled English; Human-Computer Interaction; Science Gateways; Science Library; Visualization.

## I. Introduction

When publishing a broad body of scientific research, we believe that it is important to retain some of the contextual information about the relationships between the papers, authors and topics. This helps the reader to understand the "bigger picture" behind each paper, beyond the traditional information captured in the references and citations. The International Technology Alliance in Distributed Analytics and Information Science (DAIS-ITA) [1] is a research program conducting fundamental research across a consortium of organizations from academia, industry and government in the US and UK. For this program, and the previous Network and Information Science (NIS) ITA [2] we created a large long-running multi-disciplinary research team from multiple organizations who are encouraged to pursue collaborative research across organizational, national and disciplinary boundaries. In this context, recording the wider organizational, topical and social context of the papers is a key consideration for us in measuring additional impacts and effects of the work.

In this paper we describe the Science Library[1] application,

---

[1]The live Science Library for the DAIS ITA research program can be accessed at http://sl.dais-ita.org/science-library

and give examples of the various interactive visualizations that enable the visitor to explore the rich knowledge graph of research publication information. We believe that the approach we have taken supports the idea of Cognitively Mediated Research Discovery [3] both in terms of the interactive visualization and search capabilities provided, but also in the way in which the underlying information is modelled, captured and shared. By this we mean that we are exposing multiple perspectives by which a reader might explore a particular research domain; able to follow different pathways between papers and authors to navigate the cognitive space.

To achieve this in a realistic timeframe we used the advantages available when developing such semantic knowledge graphs using a Controlled Natural Language (CNL), and we describe this CNL-based implementation of the Science Library in this paper. We also described the high-level architecture of the Science Library system, and the technique we use to solicit publication information from the DAIS ITA community for publishing.

The Science Library solution is built using technology from our earlier research (Controlled English) but the Science Library solution itself has not been reported in detail in any of our publications to date. The previously published work introduces the Science Library at a high level while this paper explores how this solution provides flexibility through the knowledge graph and the ability to explore complex relationships. Hence this paper describes the implementation of the Science Library and the potential power of the approach when using a Controlled English based solution. The Science Library currently handles 9,000 nodes, with approximately 40,000 links and 60,000 property values. It contains over 250 papers, approximately 300 authors and over 60 organizations. The site is publicly available and gets up to 50 active users per day.

In Section II we explore the various ways to interact with the complex knowledge graph through visualizations accompanied by the motivations behind them and in Section III we explain the CNL basis for the application, and the conceptual model schema that underpins the work. In Section IV we discuss existing work related to using CNL to drive Science Gateways and in Section V we talk about our future plans for the environment as the DAIS ITA program continues. Section VI concludes the paper.
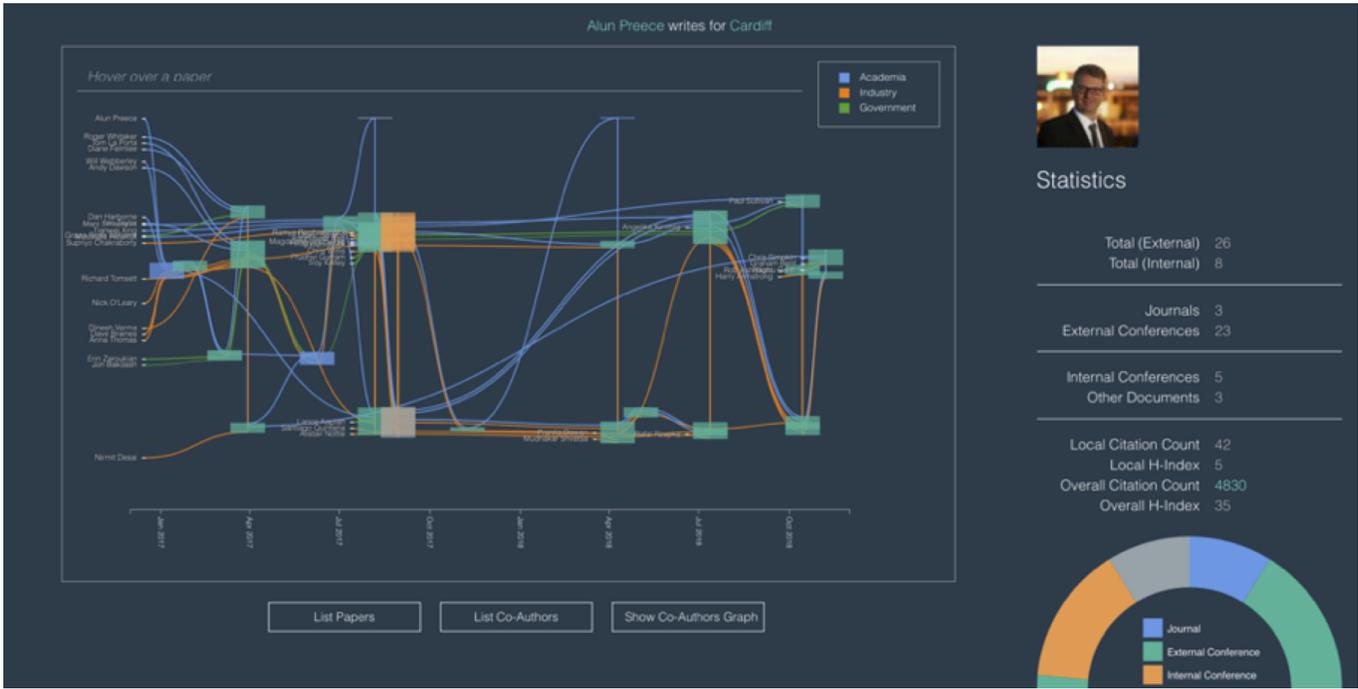
Fig. 1. Author narrative timeline

## II. INTERACTION WITH THE KNOWLEDGE GRAPH

Like many research groups or programs, one key output are the academic publications created as a result of the research. From these a rich social and topical network between authors and organizations emerges, and since the DAIS ITA is focused in particular on fostering a deep collaborative mind-set it is important for us to understand and explore this aspect of our work. To capture and convey this body of published research, for example as shown in the Author narrative timeline in Figure 1, we have created the publicly available Science Library[2] as a user-centric, interactive portal. The idea for the Science Library originally came from the earlier NIS ITA research program, and was originally developed and deployed[3] for that purpose, but has been reused and extended for this later DAIS ITA research program. The Science Library is a Controlled Natural Language (CNL) driven research gateway, which allows the community to explore and query the publications and networks through an open, web-based application built on a complex semantic knowledge graph. The data is represented through interactive visualizations along with the ability for users to query the model using a natural language query interface. The CNL based approach models the data through concepts, properties and relationships which are defined using CNL and are therefore both human readable and directly machine processable. This captures complex semantics in a simple format and enables non-technical users to participate in the continuous improvement of the data model behind the Science Library application. This paper presents the features, implementation and design considerations of the Science Library and the underlying CNL implementation.
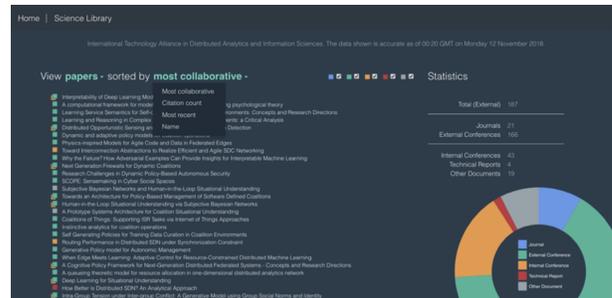


Fig. 2. List of publications

One primary purpose of the Science Library is the exploration of publications from a given research program, such as the DAIS ITA. Figure 2 shows a simple list of publications that can be sorted by various criteria including: age, name and citation count. Additional information is available on this page, such as an overall summary of the publications for the program, broken down into publication type (journal, external conference, internal conference, technical report or patent). Each publication has a colored icon to indicate the paper type, and any paper *variants* are indicated by stacked icons. For example, it is often the case that an original journal or conference publication will be summarized or reworked into an internal conference publication for presentation to the DAIS ITA research community at the annual meeting of the alliance.

Fig. 3. An example publication

Throughout the Science Library application all entities can be clicked on to take the user to a page that shows the details about that particular entity. In Figure 3 we show the publication details page which provides a large thumbnail image of the paper and the ability to download all relevant information regarding that publication (e.g. the paper PDF, any presentation slides, poster, or other supporting material). We also display the author information, the title and abstract and any venue information, including the location at which the paper was presented (for conference papers). All of this information is available in the semantic knowledge graph that underpins the system and is available for searching or exploration of the graph.

Figure 1 shows a *narrative timeline* for the publications of any author. The visualization shows the development of research ideas and topics over time, along with co-authorship sub-groups. Users can read this visualization from left to right, with each node on the network being a particular publication at a point in time and the links between them showing co-authorship flows between the papers. This can be particularly useful to determine when new threads or topics form part-way through the research, and how they manifest in terms of publications over time.



Fig. 4. Authorship social network

Building on this is the more traditional co-authorship chart shown in Figure 4. The nodes represent authors, indicated by their initial (e.g. "DB" for Dave Braines), and links to other authors who they have co-authored with. The weight of the links indicate the number of times that co-authorship has occurred, and this also drives the proximity of the nodes. In the interactive visualization, users can hover over particular links to have them highlighted, and have information about the author in question appear over the node. This is particularly useful for highlighting common co-authorship subgroups, which often indicate particular topics or threads of work that subsets of authors regularly publish. The colors of the nodes indicate whether an author is from academia, industry, or government.
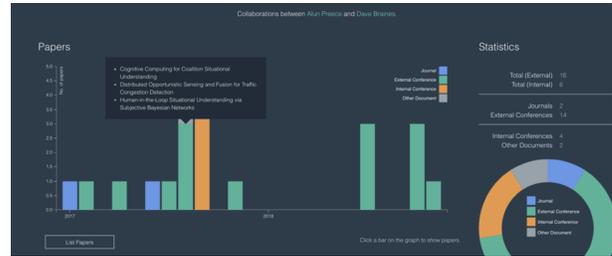


Fig. 5. Co-authored publications

Sometimes it is useful to see the specific publications for a pair of co-authors, i.e. only those papers for which both authors are involved. This is shown in Figure 5 and is represented as a stacked bar chart over time, with the bars representing the papers of a particular type. Hovering over any bar shows the individual papers for that period, enabling the user to click through to see their details. This is a very similar visualization to the one used for displaying the publications for any given organization and seems to give good insight for this particular perspective. It is worth noting however that the visualizations chosen for each view are purely subjective based on our informal opinion as to what "works best" for the given situation and is not according to any scientific basis.



Fig. 6. Geospatial contributions

In Figure 6 we show the visual aspects of a particular event mainly to showcase that we have the geo-spatial data, rather than it being a particularly useful view. We show the event venue as a large yellow circle, and the authors of the publications at this event (color coded according to their organization type) are shown in the geographic location associated with their organization.
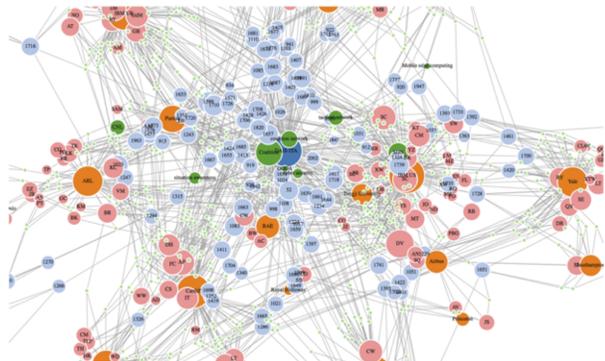
Fig. 7. The underlying complex semantic knowledge graph

Each page also has a simple search bar at the bottom of the screen. This will accept any search terms from the user and performs a reasonably complex semantic search into the semantic knowledge graph. If the terms that are searched for result in matches to specific named instances then the result shown to the user is the relevant page on which this information can be located. For example, the search term "*Dave Braines co-authors*" will return a result which is the author page for Dave Braines in the co-author graph mode. In the event that multiple matches are made, or matches are made to which there is no pre-defined visualization type, then a simple list of clickable links (to papers, authors, organizations, topics etc) are shown. The search is based on a combination of fuzzy keyword searching (to match title and abstract terms) as well as semantic graph searching. This capability is underpinned by the "Hudson" natural language search APIs[4] which provide a rich JSON response indicating which concepts, properties or instances were matched in the semantic knowledge graph based on the query terms used. There is a lot of additional power that can be gained from using this approach which is not yet implemented into the system. For example the ability to ask questions such as "*list authors who have written papers with someone from industry on the topic of IoT before 2018*". Today the interpreter API will match all of the terms in this question to the corresponding features in the semantic graph, but we have not yet developed a suitable "answerer" capability to take this interpretation and show the result to the user in a meaningful way.

In all of these visualizations we have designed a very simple user experience and have tried to showcase a variety of visual techniques to communicate certain aspects of the information in the underlying semantic knowledge graph. There is far more information available in the graph than we are making use of (See Figure 7), and it would be possible in the future to create a much richer set of visualizations if time permits. See Section V for some discussion about our plans for possible future extensions to this work.

## III. BUILDING THE SCIENCE LIBRARY

All of the visualizations and interaction features described in the preceding section are implemented using traditional web development technologies using HTML, CSS and Javascript (client and server side),and they could be integrated with any suitable back-end database-like system. The real novelty in our approach comes in the technology used to implement the semantic knowledge graph, which drive the visualizations and allow users to easily contribute to, and extend, the model. To achieve this, we use a Controlled Natural Language (CNL) technology named Controlled English (CE), which we defined and developed in previous research [4] in order to provide agile knowledge representation and reasoning capabilities like those available in the Semantic Web stack. CE is a constrained form of English that is easy for non-technical human users to read, and it is possible for them to write it too (it's harder to write, but still less technical than traditional alternative machine formats such as RDF/OWL). CE is a language that is both the human friendly language and the directly machine-processable language **at the same time**. There is no need to convert CE into some technical form (such as RDF or JSON) for the machine to process it. It is directly processed in its original readable English form. This is a key differentiating factor to our approach, and is why it is so agile as an ontology development and semantic graph building language.

The following CE examples taken from the Science Library illustrate the structure and readability of the sentences used to rapidly build up the concepts and relationships in the knowledge graph:

```
there is a person named 'dave_braines' that
  has 'Dave' as forename and
  has 'Braines' as surname and
  has 'Dave Braines' as full name and
  has the organisation 'IBM UK' as default
  organisation and
  has 'https:/.../dave_braines.jpg' as profile
  picture.

the person 'dave_braines' has the person
    'alun_preece' as co-author.

there is a document named 'doc-1009' that
  is a external conference paper and
  has 'A Generative Model For Predicting
  Terrorist Incidents' as title and
  applies to the programme 'DAIS-ITA' and
  has the event 'SPIE DSS 2017' as venue and
  has the project 'IPP16P4' as project.

the published person 'dave_braines'
  wrote the document 'doc-1276'.
```

We operationalize the CE language in the open source ce-store[5] library which is a simple Java-based web-services component to allow developers to easily integrate with CE knowledge graphs using REST APIs that are familiar to them

---

[4]Hudson natural language apis for the ce-store, https://github.com/ce-store/ce-store/wiki#hudson

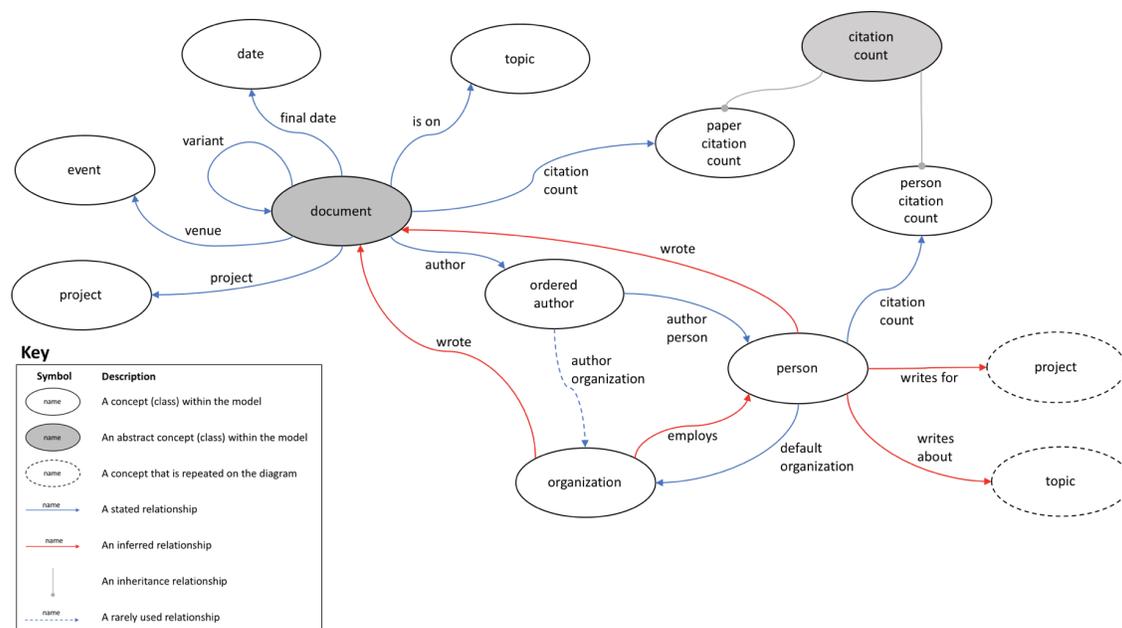[5]ITA Controlled English store (ce-store), https://github.com/ce-store/ce-store

Fig. 8. Conceptual model for the Science Library data

from using other technologies. For the developer we want the environment to be as familiar and easy to use as possible, even if they don't need to know that the underlying implementation is CE. In our previous research we have written extensively about the use of CE and ce-store in a variety of settings directly relevant to the usage in this paper. In [5] we define a conversational mechanism using this approach which enables users to converse with the underlying knowledge graph using natural language queries, and this is what led to the development of the "Hudson" APIs (mentioned in the Section II) which underpin the Science Library semantic search capability.

The ability to rapidly define conceptual models in real-time to support the capture of complex semantic data came from our field trial work on various exercises, for example as reported in [6]. This evolved in various directions, one of which was the highly extensible capture of biographical and social information to aid human decision-making [7] and this led us directly to the idea for building the Science Library to capture our own contextual and social network for our research publications.

In addition to the core ce-store engine being available as open source, we have also open sourced the Science Library user interface[6] and the Controlled English conceptual models that underpin it[7]. All the components and functionality presented in this paper can be used freely by other solutions; the entire approach is open sourced and available for anyone else to adopt as a separate, stand-alone instance. Full details and tutorials on how to set up a similar Science Gateway using CE can be found on the ce-store wiki page[8].

It is this ability to rapidly define the underlying semantic conceptual models and then populate them with corresponding data that enabled us to build the Science Library application in a very short amount of time. We were able to develop simple models quickly, populate a small amount of data and then build some example visualizations in a very rapid manner. The flexibility to modify and extend the concepts, properties or instances in the underlying graph gave us a very fluid environment in which we could rapidly explore new ideas and new conceptualizations before finalizing on those that worked. In Figure 8 we show the high-level conceptual model for the Science Library. Key concepts are *document* (an academic publication), *ordered author* (an authorship of a paper), person (an individual who may author multiple papers), organization, event and so on. Each of these concepts has one or more named *relationships* to other concepts, or textual *attributes* (such as title or abstract) which are collectively known as properties. These concepts and properties which are the equivalent of a conceptual graph [8] and the semantics are added through the definition of inheritance and the specification of logical inference rules. We use numerous rules in our conceptual model to infer additional information. Many times, this is simply to more fully connect the graph to enable the development of the visualizations to be more straightforward, but in many cases they also add significant additional "logical value" to the model. For example we infer co-authorship from paper authorship: If two authors write the same paper then they are inferred to be co-authors.

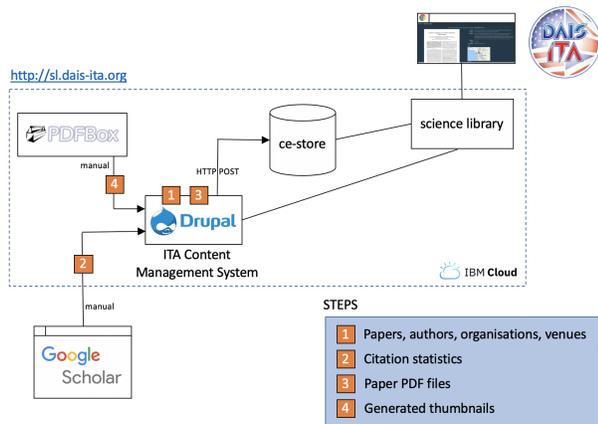In Figure 9 we show the implementation of the system,

---

Fig. 9. System outline for the Science Library solution

and how the information flows from our private content management system into the public Science Library application. We have designed a system to place minimum administrative burden on our consortium of paper authors, enabling them to simply provide a PDF of their publication, and we as administrators manually extract all of the required information (authors, venue, organizations etc) in order to publish this on their behalf.

Once the paper authors upload a PDF to the Drupal based content management system, the administrators manually create an instance of a Science Library publication and add the related properties. A Jpeg thumbnail of the paper is generated using the PDFBox tool and the citation statistics are manually pulled from Google Scholar and also added to the Science Library publication instance.

The ce-store engine is run to ingest and generate CE for the new publications and to update the existing instances and relationships. Finally, the Science Library makes calls to the ce-store semantic knowledge graphs via simple REST APIs to generate the visualizations of relationships and present the data to the user on the interactive website.

## IV. RELATED WORK

There is limited literature around the use of CNL to drive Science Gateways and interactive visualizations, however methods to allow users to quickly and easily contribute to ontologies using CNL are known [9] [10]. There are also several techniques which help users to explore and query ontologies using other forms of Controlled English. These methods allow users to query a knowledge graph using CNL, such as Attempto Controlled English [11]. Semantic Wikis capture knowledge about the data within pages and the relationship between pages, which allow semantic querying of the data. Some semantic wikis export the ontologies of the wiki as RDF or OWL, and parse it to a CNL to enable queries [12] [13].

Domain representation in general aims to effectively present the relevant information to users from a complex and vast set of data. There are various techniques and workflows to visualize digital networks. A key aspect of domain representation lies in the ability to allow the user to explore the data and understand the relationships between entities in an intuitive way. An example technique is visually presenting the semantic relationships between entities alongside an ordered list to give the user a wider perspective of the data [14] [15].

## V. FUTURE PLANS

We are continually improving the Science Library and the underlying technologies which underpin the project. One such improvement is to enhance the user interface of the website to make it accessible from any device, including mobile phones and tablets.

A better solution to link directly to an external citation source that can return the citation data in a structured form, to allow computation and processing within the Science Library will be developed to replace the current manual system of collecting data from Google Scholar.

We are also planning to improve the way the new instances in the ce-store model get created as well as updated, to improve the efficiency of the publishing process. A Node.js based tool will enable generation of new instances and link to existing instances (e.g. people, organisations, projects). This will enable us to publish to the Science Library more effectively by automatically generating the required CE.

## VI. CONCLUSION

The Science Library holds the publication information of the DAIS-ITA research program and provides the community with an interactive, visual web portal to view this data. It presents contextual information about the relationships between papers, authors and topics alongside the traditional information captured in the references and citations. Users are able to interact with the rich body of information in different forms including lists, relationship graphs and timelines to better understand the wider context of the publication, authors and topics.

ITA Controlled English was used to capture and generate the semantic knowledge graph which drives the Science Library. Since CE is language that is both machine processable and human readable at the same time, it has been the key in how we were able to implement and rapidly develop this Science Gateway. We have open sourced the core ce-store engine as well as the entire Science Library user interface to share this flexible framework with the wider scientific community.

REFERENCES

[1] T. Pham, G. Cirincione, A. Swami, G. Pearson, and C. Williams, "Distributed analytics and information science," in *2015 18th International Conference on Information Fusion (Fusion)*. IEEE, 2015, pp. 245–252.

[2] A. Preece and W. R. Sieck, "The international technology alliance in network and information sciences," *IEEE Intelligent Systems*, vol. 22, no. 5, pp. 18–19, 2007.

[3] G. de Mel, D. Braines, A. Thomas, T. Pham, and W. Dron, "Cognitively mediated research discovery: A context-aware rich visualized knowledge graph co-created by humans and machines using a common language," in *Workshop on Hybrid Human-Machine Computing (HHMC), Guildford, UK, 20-21 Sep 2017*, 2017.

[4] D. Mott, "Summary of ita controlled english," *ITA Technical Paper, http://nis-ita.org/science-library/paper/doc-1411a (Visited on 27-Nov-2017)*, 2010.

[5] A. Preece, D. Braines, D. Pizzocaro, and C. Parizas, "Human-machine conversations to support multi-agency missions," *ACM SIGMOBILE Mobile Computing and Communications Review*, vol. 18, no. 1, pp. 75–84, 2014.

[6] D. Braines, D. Mott, S. Laws, G. de Mel, and T. Pham, "Controlled english to facilitate human/machine analytical processing," *SPIE Defense, Security, and Sensing*, pp. 875 808–875 808, 2013.

[7] D. Braines, J. Ibbotson, D. Shaw, and A. Preece, "Building a living database for human-machine intelligence analysis," in *Information Fusion (Fusion), 2015 18th International Conference on*. IEEE, 2015, pp. 1977–1984.

[8] M.-L. Mugnier and M. Chein, "Conceptual graphs: Fundamental notions," *Revue dintelligence artificielle*, vol. 6, no. 4, pp. 365–406, 1992.

[9] V. Tablan, T. Polajnar, H. Cunningham, and K. Bontcheva, "User-friendly ontology authoring using a controlled language," in *LREC*, 2006.

[10] G. Hart, M. Johnson, and C. Dolbear, "Rabbit: Developing a control natural language for authoring ontologies," in *ESWC*, 2008.

[11] A. Bernstein, E. Kaufmann, A. Göhring, and C. Kiefer, "Querying ontologies: A controlled english interface for end-users," in *International Semantic Web Conference*, 2005.

[12] S. Schaffert, "Ikewiki: A semantic wiki for collaborative knowledge management," *15th IEEE International Workshops on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE'06)*, pp. 388–396, 2006.

[13] M. Buffa and F. Gandon, "Sweetwiki: A semantic wiki," *Journal of Web Semantics*, vol. 6, pp. 84–97, 2008.

[14] K. Börner, C. Chen, and K. W. Boyack, "Visualizing knowledge domains," 2003.

[15] K. Börner and C. Chen, "Visual interfaces to digital libraries: Motivation, utilization, and socio-technical challenges," in *Visual Interfaces to Digital Libraries*, 2002.