

Data-Driven or Background Knowledge Ontology Development

EMMA TONKIN

*UKOLN, University of Bath
Bath, BA2 7AY, UK
E-mail: e.tonkin@ukoln.ac.uk*

HEATHER D. PFEIFFER

*Klipsch School of ECE, New Mexico State University, Box 30001/MSC 3-O
Las Cruces, New Mexico, 88003-8001, USA
E-mail: hpfeiffe@nmsu.edu*

The development of ontologies for various purposes is now a relatively commonplace process. A number of different approaches towards this aim are evident; empirical methodologies, giving rise to data-driven procedures or self-reflective (innate) methodologies, resulting in artifacts that are based on intellectual background understanding. In this paper, we compare and contrast these approaches through two practical examples, one from a descriptive metadata domain and one from the area of physical computing. Both examples are chosen from domains in which automated extraction of information is a significant use case for the resulting ontology. We identify a relationship within the ontology development process that allies empirical evidence and user judgement to develop user-centred ontologies, either on an individual or collaboratively-focused basis. A qualitative treatment of the characteristics of this type of 'language game' is identified as an ongoing research goal.

1. Introduction

Ontologies, defined in the computer science area as “agreement about a shared, formal, explicit and partial account of a conceptualisation” (Spyns, 2002) are increasingly visible in various disciplines, particularly in the area of knowledge management, where the encoding of static domain knowledge is a key process (Aldea et al, 2003). Ontologies are generally applied for a number of purposes, including the following:

“to share common understanding of the structure of information among people or software agents; to enable reuse of domain knowledge; to make domain assumptions explicit; to separate domain knowledge from the operational knowledge; and to analyze domain knowledge (Noy and McGuinness, 2001).”

A variety of methods for ontology generation have been identified: these could be described collaborative development through introspection and discussion with joint use of expert knowledge (see Valle et al, 2008); data-driven or corpus-driven means (i.e. unsupervised methods of ontology generation); and in general as introspection, or self-reflection (externalisation of an expert's background knowledge). Linking two or more of these methods together is also possible; for example, Carvalheira and Gomi (2007) describe a method that makes use of automated ontology generation to fuel a semi-automatic, or 'hybrid', overall process.

In many ways, the debate surrounding the process of ontology building has interesting parallels to be drawn with other forms of knowledge organisation, and other historical

discussions on the topic of eliciting information about knowledge, language and the structures that underlie our everyday activities. In certain areas, the choice between data-driven and background knowledge approaches is one that defines the shape of the discipline. When studying grammar in human languages, an old and bitter argument rages between corpus-driven and introspective approaches. A similar gulf is perceived by some between structured taxonomy/vocabulary for classification and the use of unstructured or very loosely-structured approaches, such as social tagging.

We will discuss an approach that links together a data-first approach with a collaborative discovery. Reports synthesised from expert knowledge have the advantage of very closely approaching the individual's own viewpoint; if they are then bolstered by discussion with other individuals and groups, the result is intended to approach a consensus viewpoint. However, such an approach does not take into account the visibility or availability of features within the data. Under some circumstances, this approach is not a defect for an ontology. If the resulting knowledge collection is to be used only for purposes that involve human judgement, it matters only that they closely approximate consensus of opinion. However, if it is to be used for a data-driven or highly data-dependent application, such as a system that classifies documents within an ontology, or a context-aware wearable device that makes use of a set of sensor signals to characterise and perhaps identify the current context in which the user stands, it is advantageous for elements within that ontology to have a visible presence within the data.

There are various reasons to expect some features of human judgment to have little analogue in the data domain. Many judgments depend on knowledge and experience unavailable to the machine. We do not suggest that characteristics that are not directly visible within observable datasets have no role within an ontology. Rather, we simply note that the background knowledge's ability to be effectively tracked is limited, and that therefore engineering for use cases that limit this necessity may be preferable. Equally, constructions that are closely tied to a dataset, and are useful in that context, may have no English-language equivalence. In this paper we examine two example ontology developments with the aim of seeking consensus between a group of users and a set of machine-generated features describing the artifact (environment) in question. The first example involves the description of a set of documents, which have already been marked-up over time by a user population using an evolving ad hoc keyword set. Our second example is taken from the wearable computing domain, and examines an approach to collaborative ontology development for a small set of physical contexts, taking into account various sources; a set of geotagged photographs and sensor trails taken from the area itself. With these examples, we examine the enrichment of an ontology development effort by merging information from several sources; the result is a set of design proposals that we believe may be applicable for supporting the development of future collaborative or multiple-source ontology development platforms.

1. The challenge: linking concepts and data

Conceptualism takes the view that a word and its referents (the entities to which it refers) are linked by an intermediate mental entity - a concept (Sowa, 1984). Conceptualism may apply as easily to a realist or nominalist position - after all, the world is as real to us whether we construct our categories in a manner consistent to sensory experience.

However, for a device intended to operate within the range of a human conceptual model, the question suddenly becomes a great deal more than a sideline. Designing devices intended to operate in such a manner has as a prerequisite that language is not arbitrary with respect to sensory data. To what extent is this case? To quote John Taylor (1995):

"To the extent that a language is a conventionalized symbolic system, it is indeed the case that a language imposes a set of categories on its users. Conventionalized, however, does not necessarily imply arbitrary. The categories encoded in a language are motivated, to varying degrees, by a number of factors – by actually existing discontinuities in the world, by the manner in which human beings interact, in a given culture, with the world, and by general cognitive processes of concept formation."

That is, we may assume that the categories underlying language are indeed grounded in the actual shape and pattern of the world around us (our background knowledge), as well as the system in which it operates and the processes that give rise to concept development in the mind. The development of concepts, or categories, is not random - but neither is it a process bound entirely by a simple interaction between the environment and the individual.

Our task is not to come to an understanding of the system as a whole; it is merely to look for a means that simplifies the problem of eliciting categories, that may be encoded as concepts, which are sufficiently close to the world around us that they may be of some use to automated processes that must work through the filter of a mesh of sensors, or a series of features elicited from an image or a text file - the latter, of course, being quite a different problem to the former - but are also sufficiently close to the user's own understanding of the world to be of some relevance to the user.

This is scarcely a new way of looking at the problem; consider for example Jorgensen's (2007) work on reducing what is referred to in the literature of image access as "the semantic gap", a conceptual term which (Smeulders et al., 2000) identify as originating in computer science. Jorgenson notes that the term is still used in computer science literature "to refer to the difference between two descriptions of an object using different languages, specifically the difference between a human-readable description and a computational representation". The problem of relating symbols to their meanings, *symbol grounding* (Harnad, 1990), has itself recently been linked to the topic of conceptual graphs (Delugach & Rochowiak, 2008), and hence potentially to ontology development.

2.1. Developing an ontology for papers in the CS domain

There exist, of course, many excellent ontologies covering some portion of the problem area. For example, the computer science department ontology (Heflin and Hendler, 2000) provides an overview of roles, activities and types of dissemination within a computer science department. To fulfil our aim of harmonising the observed (extracted) characteristics of the dataset as far as possible with observed characteristics during the design phase, we suggest that one source of data input into the ontology model is output from unsupervised methods of examination of the data and of its characteristics. A second

source of data is from individual experts' recommendations. A third source results from collaborative development from several experts. We may then compare the result with existing ontologies such as the cs-dept-ontology, to see whether any features can be attributed to that model.

We take a document set from the Computer Science department at the University of Bristol, and begin by examining the full-text information available in the departmental repository. There is also an existing set of keywords. By extracting a series of features from the full-text data and performing an unsupervised clustering of the papers themselves, we gain something of an overview about how these terms are applied - whether some terms represent duplicates or subsets of other terms.

2.1.1. Individual ontology development

We began using a piece of cross-platform software, CharGer, a conceptual graph editor for research and education purposes. Conceptual Graphs (CG), as defined by John Sowa (1984), represent a logical formalism allowing the description of classes, relations, individuals and quantifiers. The software provides a simple visualisation method for conceptual graphs. The conceptual graph is the logical formalism used throughout the knowledge eliciting process; the ontology itself is contained within it, yet is broken out as a hierarchy of relationship that is used when exploring the knowledge relationships. For this example, we only used the ontology of the data and did not use the knowledge base of conceptual graphs.

In the first instance, users were simply asked to build an ontology showing all the concepts that they felt to be relevant to the domain, starting from the list of keywords previously applied to the document. The validity of the resulting ontology structure was not considered at this time. Each user gave quite different responses at this early stage. The first versions are shown in Figure 1.

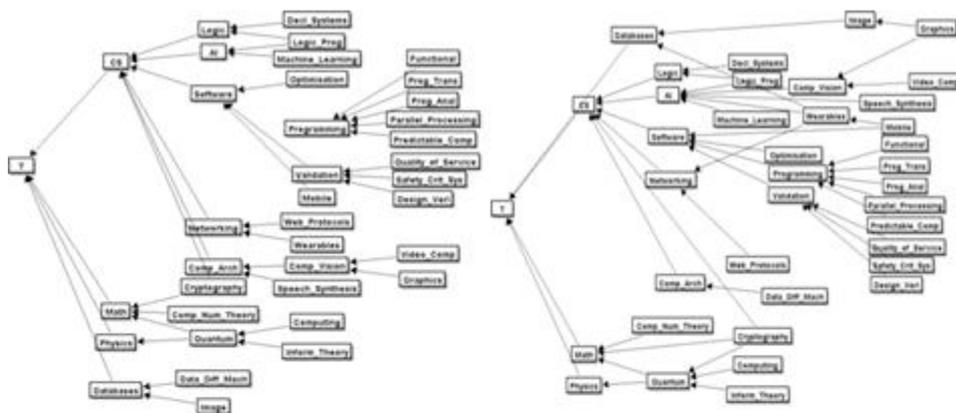


Figure 1: The manual ontology created by the initial two participants

2.1.2. Collaborative ontology development

Once an initial pair of participants hand-developed ontologies, these were compared and via an iterative cycle of editing, a jointly built 'strawman' was developed for later

revision. We then widened the scope of the process, inviting participants who were familiar with the departmental repository to make changes to the strawman model. Following a round of changes, the results were collated together manually and participant opinions were sought on the result. At this point, we found that the results were beginning to show signs of consensus. The general structure of the ontology was changing less radically than it had previously.

This was in part the result of increased user satisfaction with its present structure. One possible concern in this process, as with many other collaborative efforts, was the possibility that individuals were simply choosing to defer to the judgment of others; the majority of participants in the process worked together on a daily basis and were therefore enmeshed within an existing social structure, as well as a hierarchy of authority. In an attempt to mitigate this, we did not identify the authors of previous changes; however, it was notable that participants could at times correctly identify the participant who had made a given suggestion. As the ontology under development effectively related individuals' research topics into a hierarchy of terms, it is perhaps unsurprising that suggestions and recommendations could often be traced back to individual background knowledge.

2.1.3. Characterising a document set from observed features

We then took the results of our automated metadata extraction process, and sought to compare these with the characteristics of the emerging ontology. This includes features of the data itself; the departmental database that we are examining contains a series of eprints, and therefore our primary object of interest is extracted formal metadata (for example: title, author, date, noun phrases and characteristic phrases contained within the text). This included Dublin Core; however, it is worth bearing in mind that the majority of metadata featured in each of these classes is not present in the document itself, but is *extrinsic*, 'added-value' knowledge stored along with the document surrogate. In other cases, we might instead make use of metadata that is intrinsic to the object and may therefore be extracted. In certain cases this would include information such as the existence of features used for content-based image retrieval (Jorgensen, 2007), or of the types of features that are mentioned as candidates for the Open Text Mining Interface initiative by Nature (<http://opentextmining.org/>) - that is, information that is recognised as useful or characteristic of the document in some way. The primary advantage of this is the fact that such interfaces are typically built on top of services that are able to access a great deal of topical information held in knowledge bases or extracted data stores, and are therefore able to give an 'informed' judgment, in the sense that the suggestions given by the system are the result of inference or statistical analysis across a very large number of data points.

The dataset can be created on the basis of a series of feature extractions, in particular, of a series of common terms from each paper following removal of stop words. A form of principle component analysis was then applied to reduce the dimensionality of the resulting dataset for visualisation (the x and y dimensions are not meaningful; they are automatically set to optimise the visualisation). Some clusters in the dataset can then be manually outlined - in general, increased distance between keywords in this visualisation suggests that there is little evidence from the examined feature set that the two keywords

are applied to similar sets of papers. Proximity suggests that, from the content of the papers themselves, the two keywords may be related. Of course, there are many kinds of relation, including a superficial similarity in terms of subject matter, and evidence of superficial relations can be found.

2.1.4. *Linking these two datasets*

Because of the number of changes of concept name, hierarchy and title (some of which are outlined) it is difficult to apply the information from the keyword chart directly by string matching. However, it is possible to see from Figure 2 that the evidence from the existing dataset directly supports a number of the relations outlined. The ontology is now clearer than it was previously and contains a stronger and more generalised hierarchy of concepts. Further revision would be required to create a totally valid ontology from this source.

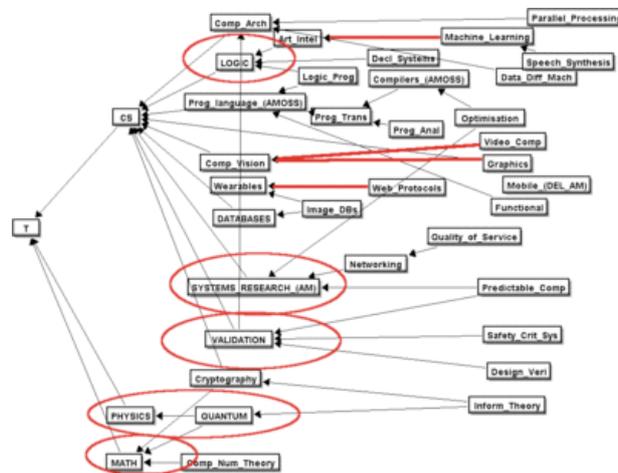


Figure 2: A composite dataset

2.1.5. *Structural stability*

It would seem reasonable to expect a composite ontology to trend towards structural stability. Review of changes shows us that this occurs by examining Figure 3.

2.2. *Developing an ontology for physical contexts in the wearable computing domain*

Research in the wearable computing domain includes a focus on knowledge management, with a significant research strand focusing on ontology development. The abstract concepts that are defined and encoded within ontologies are at some distance from the sort of information that is immediately available through sensor information, so most work focuses more on lower-level data analysis and use. However, there exist a number of ontology developments designed for the description and encoding of various contexts in the wearable computing domain.

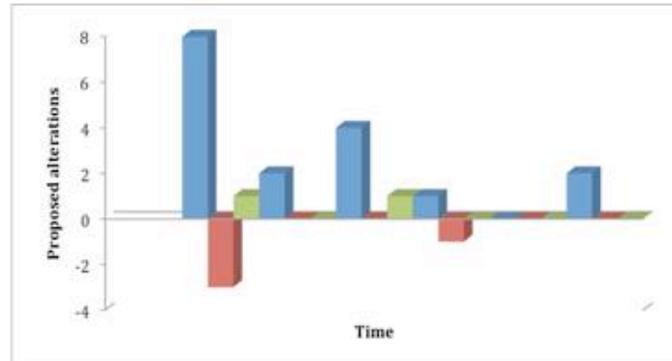


Figure 3: Data suggest stability with additions, pruning & structural changes

The wearable computing domain shares use cases with other areas in which ontologies are commonly used. Consider for example AT&T's 'Spirit' project Circa (Addlesee, et al, 2001), which coined the term sentient computing to describe a possible endpoint: using sensors and resource status data to maintain a model of the world which is shared between users and applications., so named because of the aspiration that the users and system should share a compatible model. Users were able to observe and act directly within the environment; applications could observe and act via an intermediate model (world model). This depended on a realistic intermediate model, which was achieved via an ultrasonic location system and a set of programmatic objects which were intended to 'correspond to real-world objects'. Interestingly, the project's web page notes that: "if the terms used by the model are natural enough, then people can interpret their perceptions of the world in terms of the model, and it appears to them as though they and the computer programs are sharing a perception of the real world."

Here we see that one focus within wearable and ubiquitous computing research is ensuring that terminology used is 'natural'- is used as the user would expect - and that the user has the impression that the computer shares their perspective of the world. Inaccuracy in a document search function would only result in limited relevance of search results has more significant effects in the real world. If the user's romantic dinner engagement with their partner is interrupted by a software alarm, due to the software mistakenly concluding that the current context is one in which reviewing voicemail is an appropriate task, then the user is undoubtedly and very reasonably irritated by this. At its best, the effect of getting it slightly wrong might lead to an effect not unlike that predicted by the 'uncanny valley' hypothesis, which states that when robots look and act almost human, it causes a negative response among human observers; when an interface is able to converse about concepts in a manner that is almost, but not quite equivalent to that expected by humans, small errors are seen as far more significant than they would be if presented differently.

The problems posed by development of an ontology in this area are many. Whilst the most common use of ubiquitous computing technology is to identify location -- which is to say, beginning from a concept of 'position', we generate further information via sensor

data in order to gain information about how that position might be described by the user. The authors' respective offices, for example, may be pinpointed by two sets of GPS coordinates - but the GPS coordinates are not a sufficient description of the nature of the office as perceived, as a social construction, by each user. Without this information it is merely a building. With this information, an office becomes a place whose parameters are known - is it a place for quiet work and reflection? Should the telephone be allowed to ring loudly? Is it reasonable to play music, and under what circumstances? For the student who works in the same office, are the rules the same? The relevance of a position to an individual user represents a complex interplay of many variables.

Consider the example of the Context Broker Architecture, or CoBrA (Chen et al, 2004), which was designed to support context-aware systems in smart spaces, such as intelligent meeting rooms, smart homes and smart vehicles. Within this project, a collection of ontologies called COBRA-ONT (expressed in OWL, the Web Ontology Language) were defined for modelling the context in an intelligent meeting room environment. Within the Chen paper, the authors describe previous systems (such as the Intelligent Room, Cooltown and Context Toolkit) as suffering from weak support for knowledge sharing and context reasoning, in significant part due to the fact that they 'are not built on a foundation of common ontologies with explicit semantic representation'. They provide three reasons why ontologies are key requirements for building context-aware systems:

- a common ontology enables knowledge sharing in an open and dynamic distributed system
- ontologies with well defined declarative semantics provide a means for intelligent agents to reason about contextual information
- explicitly represented ontologies allow devices and agents not expressly designed to work together to interoperate, achieving "serendipitous interoperability"

One is struck by the simplicity of the representation. Considering the problem from the gender studies perspective, user gender is encoded into a property entitled 'Gender' (the distinction between gender and sex receives no comment). Given that this is the only characteristic of the user encoded into the ontology, it is an interesting first choice, perhaps considered necessary in order to encode the distinction between the 'LadiesRoom' and the 'MensRoom', and enable practical use of the property 'AccessRestrictedToGender'. There are unhandled exceptions to many such rules; indeed, access restrictions according to gender are related to other facets of the user's identity in a manner more complex than has been represented here.

The problem of defining and encoding an ontology able to contain users' own views of their context is susceptible to alteration through users' own perceptions of their current task or occupation, status, and other physical, social and cultural factors.

A practical example: Views of Clifton Bridge

Clifton Bridge in Bristol, UK, is a well-known landmark and one that is often visited and photographed. As a result, there exist several sources of data about the bridge and its environs. Here we examine the question of what an ontology describing the area might appear, how a number of users carrying wearable devices -- as with the well-trodden

problem of development of a mobile tour guide, see Abowd (1997) -- might train a system to recognise each context - and develop labels for the contexts that are shared between devices. Initially, we take an approach very similar to the one developed in our first example above. We seek out all the available sources of data; Google Maps, tag sets from flickr and panoramio, and actual sensor information from sessions at each location.

The first obvious source is the pre-existing information retrievable from services such as Google Maps, which is able to provide a series of constraints defining which geographical coordinates (latitude, longitude) should be considered as within the area of the bridge itself. It is also able to provide the landmark's name. In that sense, therefore, one might have thought that the problem is already solved, because we know where and what the bridge is. Crowd source data, however, presents a more complex picture of user opinion. We harvested social tags from the web sites Flickr and Panoramio, and then compared reported locations to the coordinates of the bridge - the result is shown in Figure 4 – and compared user-defined labels for the landmark.

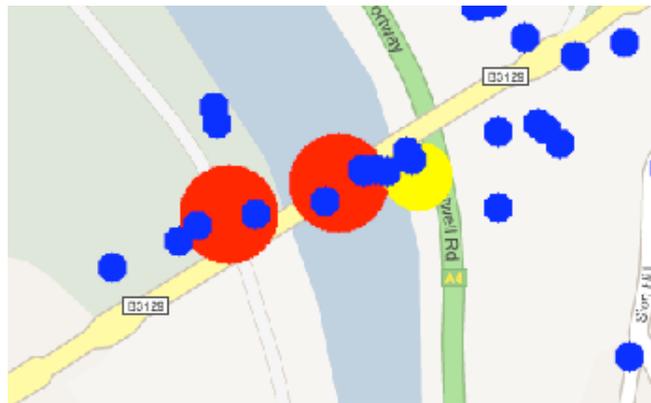


Figure 4: The distribution of positions tagged as 'Clifton Bridge' in a data set harvested from the image collection websites Flickr and Panoramio

One of the major foci of these positions (the red points on the map) are centred on the bridge itself, but most are not. Only around a third of the map positions that are identified as the bridge are located on the bridge. Another third is scattered around the periphery of the bridge, whilst the rest are scattered around the area, mostly representing vantage points from which the bridge has been photographed. But in looking at the other tags applied to the area, we see that the area can be described in any number of different ways by other users.

Assuming that users have made an attempt to place tags for the Clifton Bridge, Figure 5, shows several features that suggest task-and context-based location labels. These concepts are being based on the background knowledge of the users.

As a result, in this instance, we do not build a collaborative ontology with a common core of features, as was described in the previous section. Rather, we build a collaborative meta-ontology that maps between individual sets of concepts, which themselves are linked to a series of constraints on a shared variable.

During the first case study, participants' actions suggested various improvements that could be made, such as annotations. Even though the collaborative ontology was difficult to build, it resulted in a clearer and more refined ontology than the initial individual effort.

This suggests that communication between participants not only clarifies conceptual ideas but also generalizes the ontology making it more abstract. The second case study also showed how game theory can be brought into the process of the development of the ontology; therefore, giving a way to add meta-data to the ontology which broadens its generalities.

4. Conclusion

In this paper, we have explored ways of linking ostensibly dissimilar methods of designing, and informing the design of, an ontology. We have discussed two case studies examining problems of collaborative grounded ontology development in different domains; document set analysis and description of physical contexts in wearable computing. Our results show the problem of bringing together background knowledge and data-driven approaches can be modeled as a language game. This approach, which links empirical evidence and user judgement, is of interest in domains in which the resulting ontology is likely to be applied by an automated process, and in which successful automation of classification is the primary issue of importance. In this paper, we have qualitatively discussed some characteristics of this form of language game; a qualitative treatment of this process is an ongoing research goal.

We evaluate design methodologies; support of user annotation, where applicable; support of the collaborative design process; the ability to accept and handle input that leads to what is formally an invalid ontology. Future work in the area of grounded concept map and ontology development, particularly work that showcases practical approaches to the problem, is of interest.

References

- Abowd, G. D., Atkeson, C. G., Hong, J., Long, S., Kooper, R., and Pinkerton, M. (1997). *Cyberguide: A Mobile Context-Aware Tour Guide*, volume 3, Springer Berlin / Heidelberg, October. pp 421–433.
- Addlesee, M., Curwen, R., Hodges, S., Newman, J., Steggle, P., Ward, A. and Hopper, A. (2001). "Implementing a sentient computing system", *IEEE Computer*, August.
- Aldea, A., Bañares-alcántara, R., Bocio, J., Gramajo, J. and Isern, D. (2003). "An Ontology-Based Knowledge Management Platform", in *Proceedings of the Workshop on Information Integration on the Web (IIWeb-03)* at the 18th International Joint Conference on Artificial Intelligence. Retrieved Jan 20, 2009 from <http://www.isi.edu/info-agents/workshops/ijcai03/papers/DIsern-article-ijcai.pdf>.
- Carvalho, L.C., and Gomi, E.S. (2007). "A Method for Semi-automatic Creation of Ontologies Based on Texts", in *Advances in Conceptual Modeling – Foundations and Applications*. Springer.
- Chen, H., Finin, T., and Joshi, A. (2004). "An ontology for context-aware pervasive computing environments", in *The Knowledge Engineering Review*.
- Delugach, H. S. and Rochowiak, D. M. (2008). "Grounded Conceptual Graph Models",

- in *Conceptual Structures: Knowledge Visualization and Reasoning*, Lecture Notes in Computer Science Volume 5113/2008 Springer Berlin / Heidelberg, ISSN 0302-9743, DOI 10.1007/978-3-540-70596-3.
- Flanagan, J.A. (2006). "An Unsupervised Learning Paradigm for Peer to- Peer Labeling and Naming of Locations and Contexts", volume 3987 of Lecture Notes in Computer Science, Springer Berlin / Heidelberg. pp 204–221.
- Harnad, S. (1990). "The symbol grounding problem", in *Physica D* 42, pp. 335–346.
- Heflin, J., and Hendler, J.A. (2000). "Dynamic Ontologies on the Web", in *Procs. of the 7th National Conference on Artificial Intelligence (AAAI-2000)*. AAAI/MIT Press, Menlo Park, CA. pp. 443-449.
- Jorgensen, C. (2007). "Image Access, the Semantic Gap, and Social Tagging as a Paradigm Shift", in Lussky, J., Eds. *Proceedings 18th Workshop of the American Society for Information Science and Technology Special Interest Group in Classification Research*, Milwaukee, Wisconsin.
- Keeler, M. and Pfeiffer, H.D. (2005). "Games of Inquiry for Collaborative Concept Structuring", in F. Dau, M-L. Mugnier, and G. Stumme (Eds.): *Lecture Notes in Artificial Intelligence*, Vol 3596. Berlin: Springer, pp. 396-410.
- Noy, N. F. and McGuinness, D. L.(2001). *Ontology Development 101: A Guide to Creating Your First Ontology*, Stanford Knowledge Systems Laboratory Technical Report KSL-01-05 and Stanford Medical Informatics Technical Report SMI-2001-0880, March.
- Smeulders, A. W. M., Worring, M., Santini, S., Gupta, A., and Jain, R. (2000). "Content-based image retrieval at the end of the early years", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(12): pp. 1349–1380.
- Sowa, J.F. (1984) *Conceptual Structures: Information Processing in Mind and Machine*. Addison-Wesley, Reading, MA.
- Spyns, P., Meersman, R. and Jarrar, M. (2002). "Data modelling versus Ontology engineering". Retrieved Jan 20th, 2009 from <http://lsdis.cs.uga.edu/SemNSF/SIGMOD-Record-Dec02/Meersman.pdf>.
- Taylor, J. R. (1995). *Linguistic Categorization: Prototypes in Linguistic Theory*, Second Edition, Oxford University Press.
- Valle, E. D., Celino, I., and Cerizza, D. (2008). "Agreeing While Disagreeing, a Best Practice for Business Ontology Development", in *Business Information Systems, Lecture Notes in Business Information Processing*, Volume 7. ISSN 1865-1348 Springer Berlin Heidelberg.