# Aggregating Digital Traces into a Semantic-enriched Data Cloud for Informal Learning

Vania Dimitrova
University of Leeds, UK
v.g.dimitrova@leeds.ac.uk

Dhaval Thakker
University of Leeds, UK
d.thakker@leeds.ac.uk

Lydia Lau
University of Leeds, UK
l.m.s.lau@leeds.ac.uk

## ABSTRACT

Modern informal learning models require linking experiences in the training environments with experiences in the real-world. However, data about real-world experiences is notoriously hard to collect. Social spaces bring new opportunities to tackle this challenge, supplying digital traces where people talk about their real-world experiences, which can become valuable resource, especially in ill-defined domains which embed multiple interpretations. The paper presents a unique approach to aggregate content from social spaces into a semantic-rich data cloud to facilitate informal learning in ill-defined domains. This pioneers a new way to exploit digital traces about real-world experiences as authentic examples in informal learning contexts. The research effort to date allows us to make some observations about potential of technology and the overall approach, as well as to draw a roadmap with issues for further consideration.

## Categories and Subject Descriptors

K 3 [**Computers and Education**]: Computer Uses in Education.

## General Terms

Human Factors, Standardization.

## Keywords

Social Semantic Web, Linked Data, Semantic Augmentation, Adult Informal Learning.

## 1. INTRODUCTION

There is a radical change of the contemporary technology-enhanced learning environments responding to three main drivers. The first driver is the demand for new skills. There are widespread calls to emphasize the development of new skills, e.g. soft skills or self-regulation skills, that learners "will be required to have in order to be effective workers and citizens in the knowledge society of the 21st century" (Ananiadou & Claro, 2009). This brings forth the need to consider *ill-defined domains*, which are hard to specify including often multiple interpretations and viewpoints; and thus require highly contextualized and individualized pedagogical interventions (Lynch et al., 2009). Although learning environments in such domains are emerging,

the topic still awaits major breakthrough that should step on accumulated body of research with examples of potential technological solutions and lessons learnt from case studies in illustrative domains. Such a case study will be presented here.

The second driver is the new generation of learners. Today's learners, specifically young adults, are independent thinkers, self-regulated, ICT savvy, accustomed to social media and pervasiveness of technologies, motivated by self-realization, and shaped by experience. Hence, modern learning environments should be tailored to adult learning models (Knowles et al, 2005), which brings forth the need to create *learning situations linked to real world experiences* and opens a new set of challenges – how to capture, retrieve, aggregate, and utilize real world experiences in virtual learning environments.

To address these challenges, the third driver for the radical change in today's learning landscape – the emerging technological advancements and paradigms – should be taken into account. Specifically, we consider the *blurring boundaries between digital, physical and social worlds*, and the new opportunities for informal learning this can offer. For instance, people write reviews about their real-world experiences with staff (e.g. in hotels) or services (e.g. offered by companies), share their personal stories (e.g. in blogs or videos), leave comments pointing at situations they have experienced (e.g. when watching other people's stories or videos), etc. Although the massive popularity of social spaces and their availability via ubiquitous devices allows real-world experiences to be shared with a wide audience, using content for learning purposes is yet to be addressed, as recommended by the EU Joint Research Center (Redecker & Punie, 2010).

A major step in this new direction for exploiting social content as a source for experiential learning is being investigated within the EU funded project ImREAL[1] (Immersive Reflective Experience-based Adaptive Learning). The project focuses on a specific group of virtual learning environments – simulated dialogic situations for interpersonal communication skills (e.g. interviewing, advising, mentoring, and patient diagnosing). The key challenge addressed in ImREAL is: *how to effectively align the learning experience in virtual learning environment with the real world context and the day-to-day job practice where the skills are deployed* (Hetzner et al, 2011). The prime sources of real-world experiences are social spaces where the learners, or other users, leave traces about their 'real-world' experiences. A suite of intelligent services is being developed to augment intelligent

---

[1] http://www.imreal-project.eu

learning environments using digital traces from social media. The services for intelligent content assembly will be presented here.

The paper proposes an original, semantic-rich, and linked data inspired approach to support methodologically and technologically the development of a framework for assembling content which includes digital traces about real world experiences on a particular interpersonal communication activity. The framework is generic and can be applied to a range of activities associated with soft skills, e.g. in ImREAL it is being instantiated for job interviews, intercultural mentoring, and medical interviews. The illustrations in the paper are from the job interview activity, showing the aggregation of digital traces with personal stories (in a blog-like format), videos with examples and training materials from YouTube, and comments on these videos, in a semantic-enriched data cloud to support informal learning.

The paper will first outline the two main components of the content assembly framework – ontological underpinning (Section 2) and semantic services (Section 3). Then, an intelligent content assembly workbench will be presented in Section 4, illustrating the use of the data cloud with digital traces on job interview in a learning scenario. Finally, Section 5 will conclude by outlining lessons learnt from the ImREAL experience to date.

## 2. ACTIVITY MODEL ONTOLOGY

In order to aggregate digital traces on an activity, it is important to have a semantic model that describes the key aspects of that activity in the form of an ontology. Hereafter, we call this activity model ontology, in short AMOn. In contrast to domains for which a number of ontologies are available, there are no ontologies which describe activities in ill-defined domains. It is therefore required to develop corresponding activity model ontologies to identify the backbone of the semantic underpinning (the key activity aspects). This ontology can then be combined with other ontologies to capture additional aspects. We will illustrate below how this is being done in ImREAL.
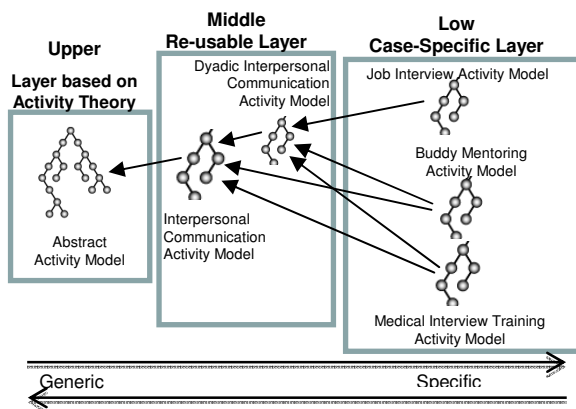


**Figure 1. Multi-layered activity model ontology (AMOn) capturing the three use cases in the ImREAL project.**

AMOn is being developed in ImREAL by a multi-disciplinary team of computer scientists and social scientists. We follow an established social science approach of applying the Activity Theory to model human purposeful activities (Thakker et al, 2011). This provides a heuristic framework to identify the key activity aspects as experienced in the 'real-world'. The Activity Theory is used as the theoretical framework for conceptualizing

the activity model ontology. This allows modularization by identifying domain specific versus reusable versus abstract modules (layers). Figure 1 outlines the three main layers in AMOn.

The first layer is upper ontology layer that covers base concepts describing activity system as outlined by Activity Theory. Activity theory provides base layer of concepts such as: `Activity`, `Tools`, `Action`, `Operation`, `Motivation`, `Outcome`, `Subject` and `Rules and Norms`, see Figure 2.



**Figure2. Upper ontology layer in AMOn.**

For the middle layer, we define concepts and respective modules that are used across use cases. This layer specializes the concepts for concepts related to interpersonal communication activities. For example, the abstract activity model module has concept of Tool as "something that is used by Subject in an Activity to achieve an Object". To specialize this concept for interpersonal communication, we expand concept of Tool to include Mental Tools, Interpersonal Skills, Non Verbal Communication Skills and Body Language. These concepts are presented below.

```
InterPersonalCommunicationActivity ⊑ Activity

InterPersonalCommunicationActivity ⊑
hasCultureNorm.CultureNorm

⊤ ⊑ ∀ hasCultureNorm. CultureNorm

InterPersonalCommunicationActivity ⊑
hasDisplayRule.DisplayRule

⊤ ⊑ ∀ hasDisplayRule.RuleNorm

MentalTool ⊑ Tool
InterpersonalSkill ⊑ MentalTool
NonVerbalCommunication ⊑ Mental Tool
BodyLanguage ⊑ NVC
BodyLanguage ≡ Kinesic
```

Our modeling approach allows us to utilize relevant ontologies. Data to enhance the coverage of the concepts in the AMOn modules. For example, to improve the coverage of the `Body Language` concept we utilize an ontology developed from external resources (Despotakis et al., 2011). Using this ontology, we are able to further specify body language into various body language signals such as arms, eyes, and head (see Figure 3).



**Figure3. Utilizing external ontologies to specialize Body Language concepts in the middle layer (left: before expansion, right: after).**

Application-specific layer, one for each ImREAL use case, is being derived to capture the specificity of interpersonal communication activities in the use cases (e.g. job interview).

The activity system shows a converged model of the activity from the perspective of the interviewer including the important concepts and the relationships between them.
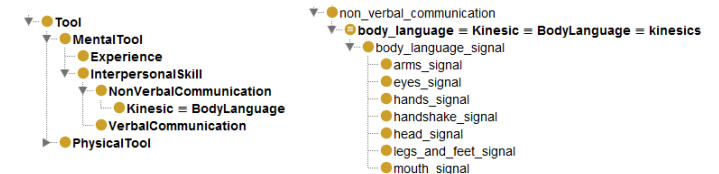
As there are a number of formats for a job interview i.e. group interview, panel interview, telephone interview, one-on-one interview and so on, we narrow our activity to the dyadic interview, conducted in a face-to-face setting. We focus only on the specific requirements of the Job interview case as envisaged in the ImREAL project. The `job interview` specific layer contains specialized concepts for dyadic job interviews by extending the interpersonal communication concepts from the `middle` layer. For example, Figure 4 outlines the specialization of `Mental Tool` concepts from the abstract layer.
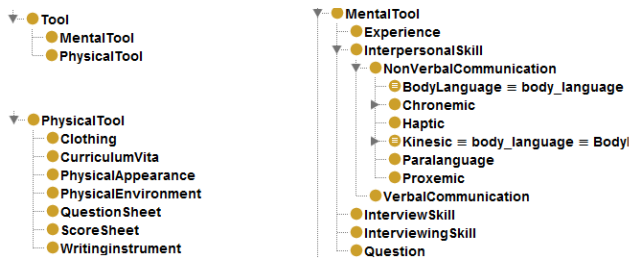


**Figure 4. Specialising Mental Tools for Job Interview**

The semantic underpinning for aggregation of digital traces on an activity in ImREAL utilizes AMOn, Linked Data Cloud datasets (mainly DBPedia ) and public ontologies (WordNet-Affect and SentiWordNet ). These ontologies are used by intelligent services for semantic content augmentation and semantic query.

## 3. SEMANTIC SERVICES
The main goals of the content assembly framework are to aggregate digital traces with learning and training experiences from the contributors with different levels of competence; and to retrieve relevant content for specific learning/training needs. It utilizes semantics (a) to annotate content contribution denoting its relevance to interpersonal communications across a range of potential learning activities (e.g. job interviews, mentoring, or medical interviews) and (b) to guide users through the search process and increase their awareness of relevant content. This is provided with corresponding semantic services: semantic augmentation service and semantic query service.

### 3.1 Semantic Augmentation Service
Semantic augmentation of content from diverse data sources using AMOn makes it possible to enrich unstructured or semi-structured data with an activity modeling context. This context is further linked to structured knowledge about these activities. Figure 5 illustrates the architecture of the semantic augmentation service. It takes any text and outputs semantically augmented text.

The **information extraction** component is implemented using General Architecture for Text Engineering (GATE)[2] and outputs Annotation Sets of extracted entities with offset, ontology URI and type information. There are three sub-components: gazetteers, Java Annotation Patterns Engine (JAPE) grammar and disambiguation. The gazetteers are the list of known entities, in

the form of lists, that the system utilizes during the initialization process. The ontologies stored in the semantic repository provides the context for the gazetteers in terms of what information needs to be stored in these gazetteers and also provides known entities as a set of lists. The JAPE grammar component contains linguistic filtering rules that allow detecting additional entities and at the same time confirming the entities detected by the gazetteers. Disambiguation and summarization component are under implementation. The disambiguation component will deal with any ambiguity generated by previous components. Here, the ambiguity refers to the cases where the same piece of text is either given two class labels or where two or more entity identifiers are assigned to the same piece of text.



**Figure 5. Semantic augmentation service architecture.**

The **semantic indexer** converts the annotation sets in semantic format (i.e. set of triples), checks the existing index for the content and stores new or updated indexes into the semantic repository. The semantic indexer component is built using the Sesame SPARQL API[3] and extensively utilizes SPARQL queries.

A **semantic repository** is a tool, which combines the functionality of an RDF-based DBMS and an inference engine and can store data and evaluate queries, regarding the semantics of ontologies and metadata schemata. As a result, semantic repositories offer easier integration of diverse data and more analytical power. We utilise OWLIM[4] as the semantic repository in our implementation In general, any semantic repository that implements the providers for Sesame SPARQL API can be used instead of OWLIM.

### 3.2 Semantic Querying Service
This service provides a mechanism for querying and browsing semantically augmented content. The service takes term(s) or concept(s) as keywords and outputs information relating to the matching concept(s) and content(s). Figure 6 illustrates the architecture of the query service.

The **query processor** implements various concept/content lookup functionalities to find related and relevant concept(s) or content(s) from the semantic repository. The query processor component is built using the Sesame SPARQL API and extensively utilizes the SPARQL queries in a variety of forms such as SELECT, ASK.

---

The query processor takes the concept URI(s) as input and outputs the result sets in SPARQL query results XML format.



**Figure 6. Semantic query service architecture.**

The **semantic repository** was described in Section 3.1.

The REST-based implementation of the semantic augmentation and semantic query services allows their utilization in various applications, one of which is presented below[5].

# 4. INTELLIGENT CONTENT ASSEMBLY

An Intelligent Content Assembly Workbench (I-CAW) is developed in ImREAL to enable learners and tutors to contribute relevant content - videos and comments related to an activity (e.g. job interview videos and comments from YouTube, as well as personal stories in a free-text form. In both cases, textual content (stories, comments) and meta data (when exists, e.g. the metadata from YouTube) is instantly tagged using onologies (AMOn, DBPedia, WordNet affect and SentiWordNet) and the semantic augmentation service. Hence, I-CAW enables the composing of a semantic-rich data cloud on the chosen interpersonal communication activity. I-CAW also utilizes the ontologies, semantic augmented content, and the semantic query service to provide intelligent browsing and a dialogue for goal setting to enable the use of the local data cloud for informal learning.

The architecture of I-CAW is presented in Figure 7. The user categories for I-CAW are[6]:

- *Trainers* - who manage the training and may use I-CAW to contribute or retrieve content linking the experience in the virtual learning environments with real-world experience;

- *Simulator scriptwriters* or *subject matter experts* – who plan the simulation script and may use I-CAW to retrieve authentic examples useful for designing the simulation situations;

- *Learners* – who are being trained to develop interpersonal communication skills, e.g. a job interviewer or an applicant,

---

[5] There is also a storyboarding environment which enables simulator developers to plan the simulation script and use example authentic digital traces as content for the simulation, see ImREAL web site for more detail about this application.

[6] See demonstration videos for the use of I-CAW by each of the user categories:
http://imash.leeds.ac.uk/imreal/wp3/icaw.html#demonstration

and may use I-CAW to become aware of relevant activity aspects and to set personal goals.



**Figure 7. I-CAW architecture.**

In the remaining part of this section, we will illustrate, with the help of scenarios, the use of I-CAW by *learners*.

**Reflecting on personal experience and linking with similar experiences**

*Jane has recently had her first experiences as a job interviewer. She uses I-CAW to reflect on her experience and learn from experience of other people. Jane logs to I-CAW and goes to the personal stories option. She is then presented a screen to type a personal story in a free text. She types her story and uploads to I-CAW – this provides a digital trace with her personal experience which is added to the pool of digital traces aggregated in I-CAW. Note that Jane could have also entered key words in the semantic search option (which maps key words to concepts and retrieves relevant content).*

*The semantic augmentation service tags Jane's digital trace (personal story) with relevant concepts from the ontologies used in I-CAW, which enables retrieving other digital traces linked to the tagged concepts or other concepts (see Figure 8).*



**Figure 8. The learner is shown key concepts in her personal story and presented with relevant content related to her story.**

*In this case, Jane has dealt with an anxious applicant and has noticed some eye contact body language signals. The concepts linked to her story are given at the top of the window, while the bottom part of the window shows relevant digital traces, which refer to similar concepts. Jane can click on the links with videos and watch some examples with using eye contact in interpersonal communication. She clicks on the video with body language at work. This enables her to make a link between the body language at job interview (which she has recently experienced) and general use of body language in interpersonal communication at work. This can help her see generic patterns and become aware of the broader application of the skills she has practiced.*

*Jane also notices the concept 'anxious' in the concept cloud at the top and clicks on it. She is presented with the window shown in Figure 9 offering definitions of anxious (extracted from the ontologies used in I-CAW), as well as digital traces with personal stories shared by other users (learners or tutors).*



**Figure 9. Exploring the semantic space for 'anxious'.**

*Jane clicks on the first story and reads somebody else's experience as a job interviewer. This story refers to anxiousness and includes eye contact (which Jane has experienced), but also refers to related concepts (such as voice) which Jane was not aware of. She can now link her personal experience and her future experience.*

*At the end of the interaction, Jane is encouraged to set personal goals and to indicate job interview aspects which she may wish to study further (for instance, by using a simulator where she can go thorough several different job interview situations).*

This scenario illustrates the potential of I-CAW to support informal learning via intelligent browsing through a semantic-rich data cloud on an interpersonal communication activity. We envisage that interactions in I-CAW will help learners become aware of a broader range of activity aspects, reflect on their real-world experience, and set personal goals for further training and development. These expectations drive the evaluation studies, which are ongoing at the time of writing of this paper.

## 5. DISCUSSION AND CONCLUSIONS

The work presented here contributes to a recent strand of innovative learning environments that are "open-ended and exploratory in nature, allowing learners to question and enhance their understanding about areas of knowledge in which they are motivated to learn" (Woolf, 2010). Our approach to this vision

capitalizes on social content and adopts techniques from building data clouds with heterogeneous content. We believe that with the recent proven success of semantic web and ontologies and the pervasiveness of social spaces, the learning technologies are ready to take on the challenges offered by ill-defined domains, and to better understand the benefits semantics brings (e.g. reasoning, aggregation, automation) for informative learning applications in these domains.

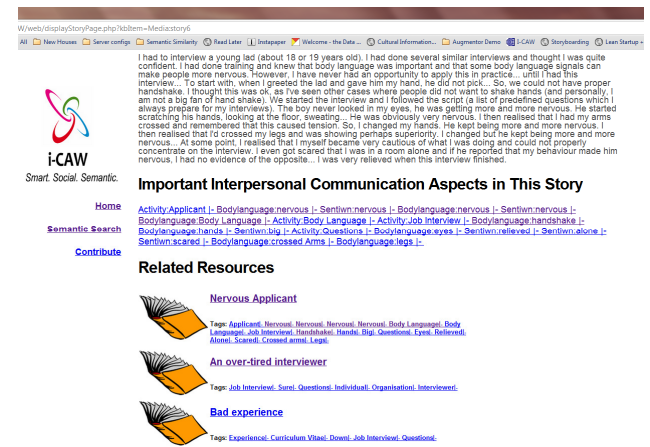The linked learning initiative is gaining a momentum with notable successes for sharing content between educational practitioners and augment learning, such as sharing videos across higher education institutions (Fernandez et al, 2011) or linked open resources for specific subject areas, e.g. medical education (Mitsopoulou et al, 2011) or organic agriculture (Sicilia et al, 2011). There is a notable move towards exploitation of linked data in educational context, see review in (Fernandez et al, 2011).

The framework presented here contributes to the linked learning initiative by offering a unique approach to *aggregate content from social spaces into a semantic-rich data cloud to facilitate informal learning in ill-defined domains*. This pioneers a new way to exploit digital traces about real-world experiences as authentic examples in informal learning contexts.

The paper presents work in progress. The technical infrastructure has been built, the first application profiles have been implemented, and small user testing has been done. We are now in a process of conducting an experimental study with users.

Our research experiences to date allow us to make some observations about potential of technology and the overall approach, as well as to draw a roadmap with issues for further consideration.

### 5.1 Positive Experience from ImREAL

We will summarize here the key aspects which we found positive and worth preserving in further work exploring similar avenues.

**Ontological underpinning**. Addressing ill-defined domains requires tailor-made ontologies to capture the key aspects of the activity. Such ontologies create the backbone for the semantic augmentation and should be built in a modular way (Thakker et al., 2011) to enable extendibility and linking with linked data repositories or other available ontologies.

**Using social content.** Content from open social spaces can be used in closed social spaces to stimulate contributions, e.g. interesting videos/comments from social media can trigger useful comments and telling personal stories. Educational practitioners often use this method - starting with fictitious stories/videos which amplify certain activity aspects and trigger interactions and personal reflections. This can be exploited also for creating socially inspired data clouds.

**Robustness and scalability.** In general, with the rapid progress in scalable Semantic Web technologies (Van Harmeleen, 2011) it becomes more and more feasible to process and integrate vast quantity of heterogeneous data. We must note that the current size of the data we operate with did not allow us to test the scalability of the proposed approach. However, using robust and scalable techniques is crucial, and in this respect, it is important to follow results from benchmarking studies, e.g. (Thakker et al, 2010).

### 5.2 Issues Requiring Further Attention

During our a slightly over one year experience in ImREAL, we noted several issues which require attention if the proposed

approach is to be deployed in educational settings. Some of them have a wider implication for linked learning repositories.

**Noise/reliability from social media**. The issue with noise applies mainly to content taken from open social spaces. Since we exploit YouTube comments, we had to deal with noise filtration and developed a generic framework for this, e.g. (Ammari et al, 2012). However, it has to be noted that the filtration algorithms are dependent on the existence of reliable content or human expertise, which may not scale. Furthermore, reliability, especially when it comes to social content, but also for any shared content, is subjective, which makes it hard to come with universal tools. Our ongoing experimental studies are examining whether tutors and learners would find content from social media relevant and useful.

**Intelligent scaffolding.** Interaction with data clouds, especially if this interaction is to promote learning, requires appropriate assistance. Using semantics (as annotation and or reasoning) allows development of intelligent scaffolding techniques, but we still do not know what scaffolding is needed in this particular case to turn semantic browsing experience into informal learning experience. We are currently exploring a range of prompts and nudges which exploit semantic similarity and entity summarization, e.g. (Cheng et al, 2011).

**Using linked data.** Although significant benefits are observed of using an activity model ontology (as discussed in Section 5.1.), AMOn is a bespoke ontology, and the other ontologies used (with exception of DBPedia) are not from linked data. The approach relies on rather heavyweight ontologies, which enable semantic scaffolding. A more lightweight approach exploiting other Linked Data resources is also being considered. We are examining the sources available and their suitability for interpersonal communication training. A benchmarking study is being planned.

**Different viewpoints.** Aggregating content from heterogeneous sources brings forth the issue about multiple viewpoints – different people see different aspects of an activity, in some cases these aspects are complimentary, in other they contradict. Making learners and tutors aware of different viewpoints is crucial in learning applications. Thus, appropriate techniques for capturing, aggregating, comparing viewpoints are highly needed. The semantics makes such techniques possible, but further work is needed to develop viewpoint frameworks applicable in learning environments. Some of our ongoing work exploits this issues, e.g. (Despotakis et al, 2011).

**New opportunities.** Using social content brings in new sources for learning. There are things which were not possible before, e.g. knowing more about the real-world experiences. Further work is needed to identify prospective application scenarios and conduct conclusive experimental studies to capitalize on the new opportunities brought by the broad availability of social content. Our approach is just a step in this direction, and we expect that there will be further research effort and studies which will exploit the use of social content for learning.

## ACKNOWLEDGEMENTS

## REFERENCES

Ammari, A., Lau, L., Dimitrova, V. Deriving Group Profiles from Social Media to Facilitate the Design of Simulated Environments for Learning, Proceedings of LAK2012.

Ananiadou, K., & Claro, M. (2009). 21st Century Skills and Competences for New Millennium Learners in OECD Countries. OECD Education Working Papers, No. 41, OECD Publishing.

Cheng, G., Tran, T., Qu, Z: RELIN: Relatedness and Informativeness-Based Centrality for Entity Summarization. International Semantic Web Conference (1) 2011: 114-129.

Despotakis, D., Lau, L., Dimitrova, V. A Semantic Approach to Extract Individual Viewpoints from User Comments on an Activity. Workshop on Augmented User Models at UMAP 2011, July 11-15 2011.

Fernandez, F., d'Aquin, M., Motta, E. Linking Data Across Universities: An Integrated Video Lectures Dataset, Proceedings of ISWC2011.

Hetzner, S., Steiner, Dimitrova, V., Brna, P., Conlan, O. Adult Self-regulated Learning through Linking Experience in Simulated and Real World: A Holistic Approach. In. Proceedings of 6th European Conference on Technology-Enhanced Learning, EC-TEL2011.

Knowles, M. S., Holton, E. F., Swanson, R. A. (2005). The Adult Learner: The Definitive Classic in Adult Education and Human Resource Development. Elsevier Science & Technology.

Lynch, C., Ashley, K. D., Pinkwart, N., & Aleven, V. (2009). Concepts, structures, and goals: Redefining ill-definedness. International Journal of Artificial Intelligence in Education, 19(3), 253–266.

Mitsopoulou, E., Taibi, D., Giordano, D., Dietze, S., Yu, Q., Bamidis, P., Bratsas C., Woodham, L. Connecting medical educational resources to the Linked Data cloud: the mEducator RDF Schema, store and API, In Proceedings of Linked Learning 2011.

Redecker, C., Punie, Y. Learning 2.0: A Study on the Impact of Web 2.0 Innovations on Education and Training in Europe, European Commission, Joint Research Centre, 2010.

Sicilia, M-A., Ebner, H., Sanchez-Alonso, S., Alvarez, F., Abián, A., Garcia-Barriocanal, E. Navigating learning resources through linked data: a preliminary report on the re-design of Organic.Edunet, In Proceedings of Linked Learning 2011.

Thakker, D., Dimitrova, V., Lau, L., Denaux, R., Karanasios, S., Yang-Turner, F. A Priori Ontology Modularisation in Ill-defined Domains. Proceedings of I-Semantics 2011.

Thakker, D., Osman, T., Gohil, S., Lakin, P.: A Pragmatic Approach to Semantic Repositories Benchmarking. In ESWC (1)(2010) 379-393.

Van Harmeleen, F. 10 Years of Semantic Web research: Searching for universal patterns, ISWC 2011 keynote talk.

Woolf, B. (2010). A Roadmap for Education Technology. GROE Available online at http://www.cra.org/ccc/docs/groe/GROE%20Roadmap%20for%20Education%20Technology%20Final%20Report.pdf