# Towards Multi-Dimensional Adaptation of Digitised Historical Content

**Eleanor O'Neill[1], Mark Sweetnam[2], Owen Conlan[1], Séamus Lawless[1],**

**Alexander O'Connor[1], Micheál Ó'Siochrú[2], Jane Ohlmeyer[2], Vincent Wade[1]**

[1]Knowledge and Data Engineering Group,
School of Computer Science and Statistics,
Trinity College Dublin,
Ireland
{Eleanor.ONeill, Owen.Conlan, Seamus.Lawless,
Alex.OConnor, Vincent.Wade}@scss.tcd.ie

[2]School of Histories and Humanities,
Trinity College Dublin,
Ireland
{sweetnms, osiochrm, ohlmeyej}@tcd.ie

## Abstract

Traditional keyword-based search approaches are limited in their retrieval of digital content. In this paper we discuss progress towards multi-dimensional adaptation of digitised content in the context of the European Commission-funded CULTURA project. The supporting technology draws on adaptive hypermedia (AH) and adaptive web systems with the aim of enhancing user queries with contextual and user-specific information. Ultimately, the CULTURA project aims to provide a technological environment for the humanities, through which it becomes possible to address implicit questions that a researcher has not explicitly articulated. This is manifested in two particular aspects of CULTURA's objectives. The first is in assisting the user in finding content that meets their needs and in discovering other content in which they will be interested. The second is to present that content in the context of other material with which the user is familiar, as well as supplemental resources and community generated commentary.

## Introduction

There has been substantial effort in the area of digitisation and cultural heritage preservation. Much of this work has, until recently, been focused on the creation of digital representations of cultural artefacts, and the creation of metadata and documentation associated with this. The result of this effort is that there is a growing collection of content available to digital humanists, in text, images and other representations. Additionally, online availability means that many important cultural heritage collections are becoming accessible to the global research community and interested public for the first time.

However, the full value of these heritage treasures is not being realised. After digitisation, these collections are typically monolithic, difficult to navigate and can contain text, which because of the noisy nature of the primary source, is of variable quality in terms of language, spelling, punctuation and consistency of terminology and naming. As a result, they often fail to attract and sustain broad user engagement and so have only limited communities of interest. While there have been recent attempts to use Adaptive Hypermedia (AH) techniques to improve engagement by supporting personalised retrieval, interrogation and presentation of cultural heritage content collections, these efforts to-date have not been widespread.

The CULTURA digital humanities project aims to address this by delivering innovative personalisation which takes into account (i) individual user intent, motivation and diversity of use, (ii) awareness of the activities and interests of the community to which the user belongs and (iii) in-depth analysis of the structure and features of digital humanities artefacts and collections. In this way, the adaptive system can influence the pathway presented by CULTURA through the cultural heritage resources to address the (explicit and implicit) information need of the end user.

In order to create successful adaptive experiences, the adaptation mechanism must reconcile numerous pieces of information before generating an appropriate adaptive offering. This information can include user preferences, task context, pedagogical principles and evolving knowledge of the user. However, CULTURA must also be sensitive to the role or status of the individual for whom the experience is being generated. Five initial broad categories of user have been identified who would require very different forms and levels of adaptivity including Professional Researcher, Non-Domain Professional Researcher, Apprentice Investigator, Non-Professional Researcher, and General Public.

It is through this combination of function that CULTURA will be able address not only explicit queries such as keyword searches, but also implicit questions, i.e. those questions that the end-user has not articulated, but which may be of interest to the end-user. For example, this might facilitate an end-user to explore the contextual situation of people, places or events, and without explicit query to look at questions such as 'Who is the key influencer or protagonist, where is he, what is he doing and who is he with?'.

In this paper we present our progress from the CULTURA project to support personalised access to digitised historical documents. We discuss the features of the initial baseline technical implementation that support the CULTURA Environment with emphasis on the enabling technology and the features that it provides for research in the humanities.

## Related Research

The MultimediaN N9C Eculture project (MultimediaN) aims to provide multimedia access to distributed collections of cultural heritage objects. It is an aim of the project to support the generation of various types of personalised and context-dependent presentations of cultural material. However, the current system provides static semantic search across entities in manually

annotated content collections. The CHIP project (CHIP) aims to provide personalised presentation and navigation of the Rijksmuseum cultural resources. The Artwork Recommender supports the rating of artworks/topics to generate a user profile, which is then used to drive future artwork recommendations. The Tour Wizard is a web-based tool, which uses the user profile to semi-automatically generate personalised museum tours.

Additionally, in parallel, the rise of 'i', 'me' and 'my' prefixes for various web portals (e.g. iGoogle) and web services (e.g. MobileMe) are intended to give the impression of some form of personal adaptation of content and service to an individual user's needs, preferences or history to enhance the individuals experience. Typically however, such services tend to focus on the following:

- identification and ranking of relevant content (web pages) or services (Agichtein, 2006; Teevan, 2005; Dou, 2007);
- simplistic 'personalisation' of the content presentation by inclusion of the user's name, historical information/recently used resources; or
- simple augmentation of screen layout (Ankolekar, 2007).

To effectively empower communities of researchers with personalised mechanisms, which support the collaborative exploration, interrogation and interpretation of complex digital cultural artefacts, the adaptivity provided in CULTURA will need to be more integrated and intelligent than in the portals described above. Such next generation adaptivity will need to support the dynamic adaptive composition and presentation of digital cultural heritage resources. It will also need to support personalised visualisation of, and interaction with, social networks, which are identified in both the content collections and the related research communities.

In terms of content collections, there are many potential sources to which CULTURA could be applied. One example is the Europeana project, which represents metadata from collections across many EU member states. While Europeana does not directly host content, it is a large repository of metadata, which could be processed, alongside a specific collection's content to fuel the CULTURA experience. This offers a significant opportunity to make use of linked data approaches, for example to support collection cross-walking.

## Adaptive Content Retrieval for Historical Documents

The 1641 Depositions, which have been digitised and transcribed, are being used as the baseline historical content to validate the techniques developed within CULTURA. The 1641 Depositions are seventeenth-century manuscripts that comprise of over 8000 (or 20,000 pages) depositions or witness statements, examinations and associated materials in which mainly Protestant men and women of all classes and from all over Ireland told of their experiences following the outbreak of the rebellion by the Catholic Irish in October 1641.

From a historical perspective, this body of material is unparalleled anywhere else in early modern Europe and provides a unique source of information for the causes and events surrounding the 1641 rebellion and for the social, economic, cultural, religious, and political history of seventeenth-century Ireland, England and Scotland. From a technological perspective, the 1641 Depositions represent a textually-rich digital humanities collection, which is characterised by noisy text, inconsistent sentence structure, grammar and spelling.

These cultural artefacts have important similarities to the user-generated content found on the world-wide-web (WWW) today. They are inconsistent in almost every aspect, including spelling, punctuation, case and language. This means it is possible to draw on state of the art approaches in Adaptive Hypermedia and Adaptive Web systems research. Adaptive Hypermedia and Adaptive Web systems research is concerned with improving the retrieval and composition of information based on an individual's needs and interests (Brusilovsky, 2007).

An extract from the 1641 Depositions is included below to illustrate the difficulties and challenges that exist for applying text analysis to these documents. Due to the form of English that appears in the Depositions, along with human error evident in the texts, state of the art text correction and normalisation techniques must be applied to the text before applying modern information retrieval techniques. For example, during the transcription process, 47 variant spellings were uncovered for the barony of Kinalmeaky, which, for the purpose of computer-based text analysis, must be reconcilable to a single normalised form or a single unique identifier. The normalisation process is currently underway by our project partners at Sofia University.

In the interim however, we have been able to make use of metadata recorded for each of the depositions. As part of the 1641 Depositions digitization process, the transcription team captured a set of metadata about each deposition document, in addition to undertaking the transcription process. The metadata documented for the depositions is extensive, including the following categories: *occupation, person type, gender, religion, nationality* and *nature of deposition (crime/event)*. This provides us with some insight into the content of a deposition, prior to performing full text analysis. However, the limitation of this is that we cannot determine the significance, or frequency of occurrence, of any given term for a particular deposition. This means, using the metadata alone, we cannot determine if a crime or event is mentioned as a major or minor topic of the deposition.

## Sample of Deposition Content: *Deposition of Elizabeth Gough*

**fol**. 2r

192

Elizabeth Gough late of Bellamenagh in the County of Cavan spinster deposeth as followeth beeing duly sworne & examined

That at Lismore in the County aforesaid at the house of Philip mc Mulmore o Reily this deponent beeing in the company of Cahil o Reily [ ] <a> Cahir ô Reily & Thomas mc Encor of Lismore aforesaid (the servants & warders of the said Philips house) about the 25 of November last she demaunded of Cahil O Reily the reason of these outrages against the English above others: the said Cahil answering that it was pitty that all the English in England & Ireland were not hangd drawne and quartered before now. this
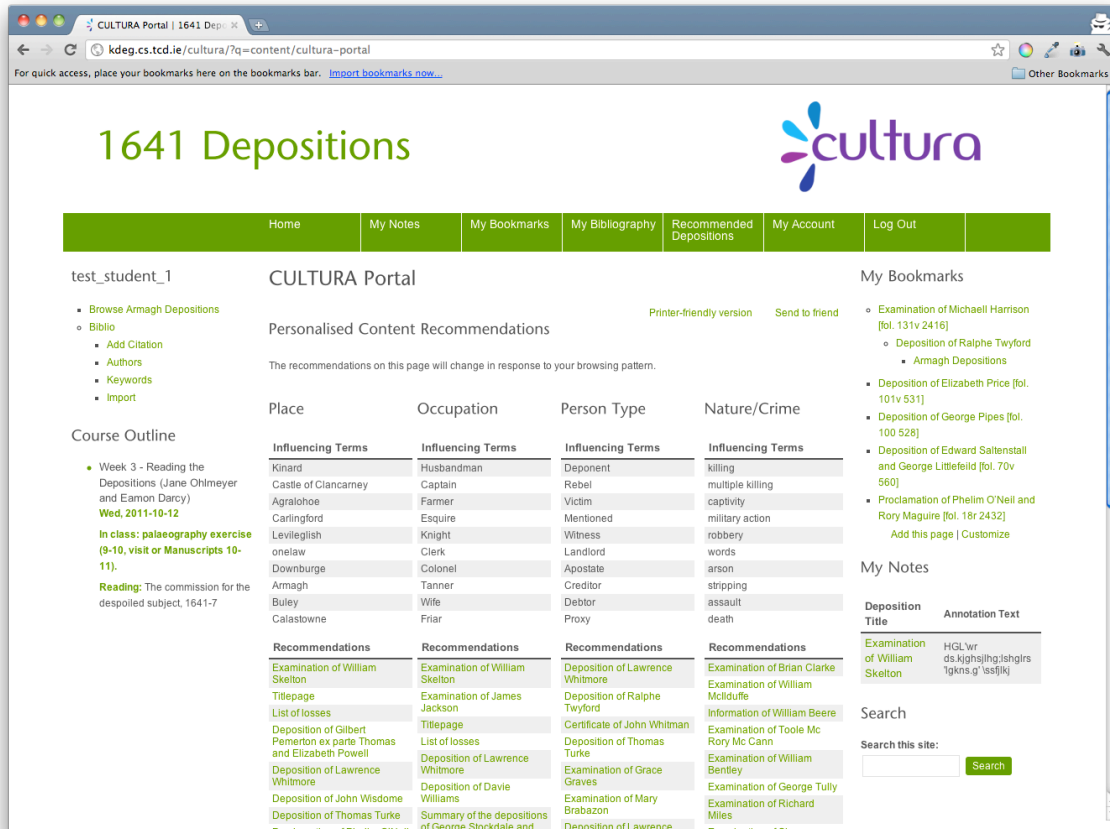
**Figure 1. CULTURA Web Portal**

deponent demaunding the reason he replied, ther they had hanged & quartered the Queenes priest in her presence: & had put gunpowder in her sadle to blowe her up: *the said English* calling her whore & her children bastards: whereupon she ~~was glad~~ *laboured* to flee to her brother into France, having first sent letters into Ireland to kill all the English men women & children. the said Cahill adding that the said English had favour that they lost all & escaped vnslaine he further said that the Irish purposed to have a king of their owne in Ireland, and that Sir Phelim o Neale should be he.

<b> As for her losses she deposeth that there were taken from her <46-2-00 b> ~~husband~~ five English Cowes worth 20 li. one horse worth 5 li. corne in stack 5 li. houshold goodes 10 li. one fowling peece 32 s. quailenets, larkenets *an* [hay?] & ~~a~~ plovernets with there appurtenances to the value of 5 li. the moste of these beeing taken away by commaund ~~from~~ *of* Mulmore ô Rely of Cavet in the County aforesaid esquire & by others whose names this deponent knoweth not.
[her marke]
Elizabeth [mark] Goughs
Deposed before vs feb. 8. 1641
Hen: Jones
John Sterne
[Parte of ~~there~~ Smith]
fol. 2v
[Copy at MS 832, fols 119r-v]
fol. 3r

| Term Category | Deposition Metadata |
|---|---|
| Nature/Crime | Killing, Robbery, Words |
| Occupation | Spinster, Knight |
| Type of Person | Deponent, Mentioned, Rebel |
| Place | Bellamenagh, Cavan, Ireland |

**Table 1. Metadata recorded for the Deposition of Elizabeth Gough**

Table 1 documents the terms recorded in the metadata for the Deposition of Elizabeth Gough. For example, *spinster* and *deponent* refer to Elizabeth Gough, while *knight* and *rebel* refer to Sir Phelim o Neale. *Killing*, *robbery* and *words* refer to the crimes documented in the deposition, for example the killing of the Queen's priest and the list of losses that Elizabeth Gough deposes were taken from her.

## Case Study

To provide an example of how the CULTURA Environment might be used from an end-users' perspective, a use case is described. Sir Phelim O'Neill was one of the leaders of the Irish rebellion, and as such, his name appears repeatedly throughout the depositions in a variety of contexts. He is often mentioned as the leader of the rebellion or is accused of being directly involved in events happening on the ground.

An expert researcher may wish to identify the full extent of the allegations made against him by the deponents. These allegations, however, are distributed over such a wide number of manuscript resources that it is almost impossible to ascertain with any degree of certainty, the

extent or the veracity of the charges. The completion of such a task through the manual analysis of manuscript resources would take months, if not years, and therefore the full extent of his actual involvement has never been accurately identified.

CULTURA is addressing this by using natural language processing to normalise the text so that network analysis techniques can be applied to the documents as a whole in order to examine the relationships that exist between people, places and events over the documented time period. By making such traditionally difficult tasks achievable by a variety of user communities and by facilitating enhanced adaptive investigations, CULTURA is working to add real benefit to the area of cultural heritage and enhance the interpretation of culturally influential resources.

## Personalised Web Portal for Digitised Historical Documents

The initial prototype of the CULTURA Environment has been developed on top of an existing and popular open source content management system (CMS) called Drupal. Drupal provides secure, flexible, extensible, standard Web 2.0 functionality such as user profile pages and discussion forums, and is richly supported with over 7 million installations worldwide. Drupal also provides an extension framework, which has allowed us to develop custom functionality specifically to support accessibility and personalised presentation of content within the CULTURA Environment. In the baseline recommendation engine, we make use of the following 1641 metadata, occupation, person type, nature of deposition, and place, to broadly make recommendations about people, events and places. These categories in the metadata were chosen specifically to align the system with coursework for an MPhil class at Trinity College, under the supervision of Jane Ohlmeyer, who are trialling the environment during the Michealmas (Winter) Term, 2011. The current interface for the CULTURA environment is shown in a screenshot in Figure 1.

The central columns visible on this webpage each provide recommendations with respect to one of chosen metadata categories, i.e. place, occupation, person type and nature/crime. At the top of each column the Environment provides a list of the key influencing terms that are driving the personalised recommendations for the end-user. This supports an element of scrutability for the end-user by allowing the user to have some understanding about why depositions are appearing as recommendations. A secondary aspect of displaying the ranking of terms is that it provides the end-user with insight into their own interaction and browsing activity with the depositions in the environment.

In addition to the recommendation system, the CULTURA Portal provides a set of features to facilitate use of the site in an educational and research capacity. These features include private note-taking, bookmarks and bibliography, each of which are entered and saved by an individual user. Additionally, the environment provides supporting features, such as the class schedule, reading list and deadlines, which are important to the class as a group. The baseline implementation also supports explicit search queries using both basic and advanced search options.

A listener module (Drupal extension) developed for the CULTURA Portal captures the activity of each individual user as evidence of user interest. In the baseline implementation we focus on the actions of viewing a deposition, book-marking a deposition and making notes on a deposition, with each action respectively considered of increasing importance as indicators of user interest. Based on these observations, the Environment calculates a term matrix using a vector of decaying interest for each individual user. This provides a representation of user interest that is aggregated over multiple sessions. When a user performs a relevant action on a deposition, each of the terms recorded in the deposition's metadata are increased in importance in the user model, all other terms are decreased in importance. This allows repeatedly recurring terms to float to the top, while other terms sink to the bottom. Four matrices of interest are maintained, one each for occupation, person type, place and nature of deposition. Depositions are recommended based on these weightings, which are normalised according to the number of terms influencing the recommendation.

Although the metadata driven approach that we have described is acknowledged to be relatively simplistic, it is beneficial in providing the baseline comparison point. In our continued development of the CULTURA Environment, we can leverage normalised text, entity and relationship analysis of the text, as well as annotations made by users, each of which are technologies provided by partners within the CULTURA Consortium. Through these additional features, we have the opportunity to generate a user experience based on the relationship between people, places and events in the underlying material, to compliment an individual's interest model.
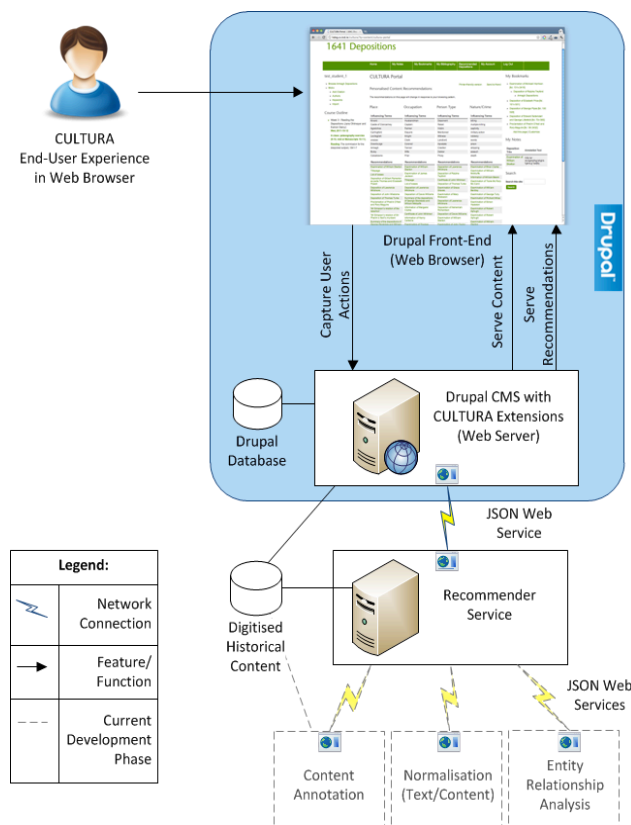


**Figure 2. CULTURA Technical Architecture**

# CULTURA Technical Architecture

The development of the CULTURA Portal using the existing content management system, Drupal, has provided many ready-made features and functionality. As a result, this has enabled us to commence trials early in the project life-cycle. However, the potential requirement for heavy-weight processing, analysis and calculations as the Environment matures, means that we have integrated Drupal with a web service architecture. This enables us to move some of the processing required in the Environment outside Drupal to more appropriate technologies.

Figure 2 provides an architectural overview of the development work that supports this. The integration over a web service architecture allows for loose integration of the CULTURA Environment, which is beneficial both in terms of speed of integration and minimising workload at integration and maintenance. The Drupal install makes use of the existing Services and JSON Server modules (extensions) to support this integration approach.

# User Focus Groups

A key characteristic of CULTURA's novel approach is that it is designed to situate a user's personalised experience within the community that surrounds them. The 1641 Depositions have active communities of interest because of the wider social and historical implications that transcend geographical and chronological boundaries and continue to shape opinions and values to this day. These communities are providing a range of end-user groups with whom we have been conducting focus groups and trials. Prior to commencing development of the baseline environment, focus groups and interviews were conducted with prospective users of the CULTURA Environment, some of whom have prior exposure to the original 1641 website at http://1641.tcd.ie.

Focus group discussions were carried out with representatives of a number of the user communities targeted by CULTURA. Detailed exchanges were carried out with Professional Researchers, Non-Domain Professional Researchers, and Apprentice Investigators, see Table 2 for definitions of these groups. Many of the professional researchers have had extensive experience of using the 1641 Depositions in their teaching and research. The non-domain professional researchers were expert on seventeenth century religious violence and seventeenth century Ireland more generally. The Apprentice researchers included undergraduates who are or have written dissertations on aspects of the 1641 Rebellion, and Ph.D. students who are working on the Rebellion. Some of these students have completed an earlier version of the M Phil class, and have an excellent insight on the challenges of working with the 1641 Depositions.

These discussions revealed a central core of functionality common to the requirements of each user community. Unsurprisingly, this core accurately reflected the core interests of historians. Thus, the requirement to trace interactions between people, places and events was at the top of the list for each community, and reinforces our choice of these three axes in our baseline recommender system. In particular, researchers who had prior experience of attempting to reconstruct the dynamics of the rebellion from the fragmentary evidence furnished by the Depositions stressed the importance of being presented with the relevant depositions necessary to construct a synoptic version of events. Users also identified visualisations and mappings as tools of significant potential. These tools are built on and represent essentially the same interactions between persons, places and events.

| Category | Example 1641 Depositions User |
|---|---|
| Professional Researcher | Scholars, Academics, Tutors and Curators in directly related domains. |
| Non-Domain Professional Researcher | Scholars, Academics, Tutors and Curators in directly related domains. |
| Apprentice Investigator | Undergraduate, Post-graduate and Post-Doctoral Students. |

**Table 2. Categories of participants in User Focus Groups.**

# User Trial: MPhil Students

The baseline is being trialled October-December 2011 with a class of MPhil students at Trinity College Dublin. The objective of this user trial is to investigate degree to which recommendations can be made, and adopted by end-users, using only metadata as the basis of the recommendation. During the course of the trial, MPhil students have been set tasks and assignments to complete using the CULTURA Environment. These tasks centre on the key metadata supported by the Recommender system, i.e. people, place and event. The coursework ask students to address a set of the following questions, in the form of an essay style answer:

**Person:** *Sir Pheilim O'Neill*
- Where is he?
- What is he doing?
- Who's he doing it with?
- What terms are used to describe O'Neill?

**Place:** *Armagh*
- What sort of events take place in Armagh? How many of each?
- How did the frequency of different sorts of events change over time?
- Who are the key people involved in perpetrating events [influencers]?
- To what extent are individuals associated with the city of Armagh – as inhabitants, refugees, soldiers, traders, or through family members?

**Event:** *Portadown Bridge*
- Who's there? Victims, perpetrators, witnesses?
- Who are seen as the organisers/influencers of this event?
- How many people were killed at Portadown Bridge?
- Is there any pattern in the events that are mentioned along with Portadown? [How is the event situated in witness testimony? Is it reported as standalone or as part of a series of atrocities?]

To overcome the issue of the cold start in the recommender system when a user first logs into the system, we are adopting a two-pronged approach. By seeding the terms weights with key terms from these questions presented to students, for example the terms

'killing' as an event and 'Portadown Bridge' as a place, the system does not present a flat, non-prioritised list of documents. The second part of warming the system is to instruct students to begin browsing the depositions by inputting some explicit search queries. The subsequent browsing of these search results allow each individual's experience to begin to deviate from the starting system state and commence their own personalised experience.

A central motivation in conducting early focus groups and baseline trials is to capture requirements through user engagement in order to build user and community profiles as well as strategies for personalisation for the various types of user. The outputs from these focus groups and trials are key to guiding the technical development work. We make extensive use of logging to capture information surrounding visited links, referring pages, user identity and the timeline of user activity, to support our investigation of the paths that users take through the content. This has value in terms of assessing the Recommender system, for example to determine what proportion of a user's browsing activity draws on recommendations vs. explicit search queries. Additionally, it also enables us to draw comparisons between the browsing activity of experts and novices, which inform our strategies to support the communities that we identified in this paper. We also hope to gain insight into the key variables of the personalisation algorithm. For example, do boundaries exist for the set of influencing terms, which improve the recommender's precision in selecting depositions? These boundaries might include the size of the set or the relative significance or insignificance of a term within the set.

## Conclusions

The central objective of CULTURA is to improve the quality of access to collections, which are not exhibited physically, and to support a spectrum of users. CULTURA also offers substantial opportunity to introduce new methods of engagement and interaction for historical artefacts. It is not the aim of the CULTURA project to confine end-users within a defined subset of content, but rather to provide users with as much freedom and control as possible to browse a content collection.

In this paper we have described initial work to provide a form of implicit search by monitoring a user's activity within the CULTURA Environment in order to generate a model of the user's interest. We also outlined CULTURA's first user trial to investigate the effectiveness of the personalised recommender system, in terms of its adoption by end-users. The outcome of this trial will be used to inform the next iteration of development work for the CULTURA Environment.

The long term objectives over the course of the project are to produce a set of services that are generally applicable to a wide variety of cultural artefacts. This generality will be demonstrated by the fact that CULTURA will support both the IPSA (Imaginum Patavinae Scientiæ Archivum) and 1641 Deposition collections, which differ in morphology, language, modality and metadata.

### Acknowledgements

## References

1641 Depositions http://1641.tcd.ie (accessed 14th October 2011)

Agichtein, E., Brill, E., Dumais, S. (2006). 'Improving Web Search Ranking by Incorporating User Behaviour Information', Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval, Seattle, Washington, USA 2006, pp19-26.

Ankolekar, A. Krötzsch, M. Tran, T. Vrandečić, D. (2008). 'The two cultures: Mashing up Web 2.0 and the Semantic Web', in Web Semantics: Science, Services and Agents on the World Wide Web, Volume 6, Issue 1, Elsevier Science Publishing, ISSN:1570-8268.

Brusilovsky, P., Kobsa, A. & Nejdl, W. (eds.). (2007) 'The Adaptive Web: Methods and Strategies of Web Personalisation'. In The Adaptive Web, Lecture Notes in Computer Science, LNCS, vol. 4321, Berlin Heidelberg New York: Springer-Verlag.

Brusilovsky, P. (2007) 'Adaptive Navigation Support', In The Adaptive Web, Lecture Notes in Computer Science, LNCS, vol. 4321, P. Brusilovsky, A. Kobsa, W. Nejdl (eds.), Berlin Heidelberg New York: Springer-Verlag.

CHIP The Cultural Heritage Information Presentation project http://www.chip-project.org/ (accessed 15th October 2011)

CULTURA EU STREP, grant agreement no. 269973. http://www.cultura-strep.eu/ (accessed 14th October 2011)

Dariah http://www.dariah.eu (accessed 15th October 2011)

1641 Depositions http://1641.tcd.ie (accessed 15th October 2011)

Delos http://www.delos.info (accessed 15th October 2011)

Drupal Content Management System www.drupal.org (accessed 15th October 2011)

Dou, X. Song, R. Wen, J. (2007) 'A Large-scale Evaluation and Analysis of Personalized Search Strategies' Proceedings of 16th International Conference on World Wide Web (WWW16), Banff, Alberta, Canada.

Europeana http://www.europeana.eu/portal/ (accessed 15th October 2011)

iGoogle http://www.google.com/ig (access 14th October 2011)

IPSA (Imaginum Patavinae Scientiæ Archivum) http://www.ipsa-project.org/ (accessed 14th October 2011)

Mobile Me http://www.me.com (accessed 15th October 2011)

MultimediaN N9C Eculture project – http://e-culture.multimedian.nl/ (accessed 15th October 2011)

Teevan, J. S.T. Dumais, Horvitz, E. (2005). 'Personalizing search via automated analysis of interests and activities', 28th International ACMIR Conference on Research and Development in Information Retrieval, Salvador, Brazil.

QViz, http://www.qviz.eu/ (accessed 15th October 2011)