**AMIDA**

Augmented Multi-party Interaction with Distance Access

`http://www.amidaproject.org/`

Integrated Project IST–033812

Funded under the 6th FWP (Sixth Framework Programme)

Action Line: IST-2005-2.5.7 Multimodal interfaces

# Deliverable D6.5: Final Report on Integrated Applications and their Evaluation

**Due date:** 30/09/2009            **Submission date:** 15/10/2009
**Project start date:** 1/10/2006       **Duration:** 39 months
**Lead Contractor**: IDIAP Research Institute   **Revision:** 1

| Project co-funded by the European Commission in the 6th Framework Programme (2002-2006) | | |
|---|---|---|
| **Dissemination Level** | | |
| PU | Public | ✓ |
| PP | Restricted to other programme participants (including the Commission Services) | |
| RE | Restricted to a group specified by the consortium (including the Commission Services) | |
| CO | Confidential, only for members of the consortium (including the Commission Services) | |

# D6.5: Final Report on Integrated Applications and their Evaluation

**Abstract:**

This document reports on progress made in AMIDA WP6 – Applications: Interfaces, Integration and Evaluation – from **October 2008** to **September 2009**, the third year of the AMIDA project. The two integrative applications that are the main target of the WP include core technologies of the AMIDA partners in a user-oriented setting.

The Automatic Content Linking Device (ACLD) is a just-in-time speech-based document retrieval system for meetings, and the User Engagement and Floor Control demonstrator (UEFC) is an assistant for improving remote participation in meetings, in particular using addressee detection. A significant amount of work was devoted within WP6 to the ACLD and UEFC, which are collaborative achievements, and have reached the level of mature research prototypes, accompanied by user-friendly interfaces, as described in this document. Other meeting browsing/assistance applications were also completed, for meeting catchup, spatialisation of remote participants, and information access. The applications were built around the Hub client-server architecture for annotation exchange, which serves as an integration infrastructure for the AMIDA technologies.

This final WP6 public report follows up on previous descriptions of WP6 status: D6.7 from March 2008 (*AMIDA Proof-of-concept System Architecture*, public), D6.1a from June 2008 (*WP6 Planning - ACLD and UEFC*, confidential), and the 2-year progress report D6.4 from September (*WP6 Progress and Planning*, confidential).

## Contents

# 1 Shared Support Infrastructure: the Hub

**Improvements to the Hub architecture.** The Hub client/server architecture for annotation exchange has been extensively used to integrate modules for AMIDA demonstrators, mainly the ACLD and UEFC described below. The Hub[1], already documented in the AMIDA Deliverable D7.2 (AMI Consortium, 2007), has been improved based on the received feedback. As a result, version 2.2 of the Hub was released in autumn 2008. This is now robust to dead, inactive, or very slow clients. A C++ client class has been written, in order to enable C++ based applications to produce data for the Hub. Bugs were fixed and overall performance was improved, in particular by creating a compressed communication channel between the client modules and the Hub server.

**A new API to the Hub: the ezHub.** The latest development for the Hub architecture is a new API known as the *ezHub*[2]. The ezHub is a new, simple, elegant and efficient Java API to produce and consume Hub data, with no need to know about the intricacies of connection, data formats, error handling, complex threading issues, etc.

The ezHub provides a sophisticated object-oriented model of the data associated with a meeting, as a collection of attributed objects. The objects are distributed, so changes to an object's attribute on one machine are automatically distributed, via the Hub, to other clients observing the same object. Similarly, objects are related to other objects, and updating a relationship is again automatically distributed to interested parties, transparently. The net effect is that each Hub client appears to have a shared, distributed object model on which to operate. Producers of information simply update "their" model, and the consumers notice "their" model being updated, automatically.

Such models have existed before, but the final twist with the ezHub is that every attribute or relationship can be queried either now, or for any time in the past, or for future values (via call-backs), in the same API. For example, an attribute might indicate a person's name, which is (fairly) constant, or another might indicate what the person is saying – a fleeting ephemeral value.

The ezHub API was released for use by the AMIDA community in June 2009, currently at `http://www.idiap.ch/~flynn/Hub`. This web site includes a thorough tutorial in using the ezHub, complete code documentation (in JavaDoc), download facilities and example data. A simple web interface for exploring AMIDA example data, available by arrangement, was also designed.

# 2 Automatic Content Linking Device (ACLD)

The ACLD is a meeting assistant that provides just-in-time access to potentially relevant documents or past recorded meetings, based on speech from ongoing discussions. Participants in meetings often mention such documents, but do not usually have the time to search for them, therefore, the ACLD retrieves them automatically.

---

[1] Developed by Mike Flynn at IDIAP, with Alex Nanchen (IDIAP) and Jonathan Kilgour (UEDIN).
[2] Developed by Mike Flynn at IDIAP.

The ACLD serves as a demonstration application for a number of AMIDA core technologies. The very first version of the ACLD was demonstrated at the AMIDA review in February 2008, and a second version was demonstrated at the following review, in December 2008. The ACLD is developed by a multi-site team (IDIAP, DFKI, UEDIN, with TNO, BUT and UT), which coordinates through weekly conference calls, the minutes of which are gathered on the AMIDA Wiki. Following an initial description paper at MLMI 2008 (Popescu-Belis et al., 2008), subsequent improvements to the ACLD were described in a book chapter (Popescu-Belis et al., 2009a) and a demo paper at ICMI-MLMI 2009 (Popescu-Belis et al., 2009b). Other publications are under submission or in preparation.

**Overview.** The main achievements of the reporting period (Oct. 2008 – Sept. 2009) were the following, and details for some of them appear below.

1. The User Interface (UI) has been completely redesigned, following in particular the reviewers' recommendations for better usability, as a modular architecture allowing considerable flexibility in the display. The design was first sketched and discussed extensively by the ACLD team, then implemented in a Java-based framework, upon which several functionalities were added (see below).

2. A new System Controller was designed, offering a user-friendly interface to execute the preparatory steps for running the ACLD, allowing the use of a local or a remote Hub and the selection of the execution mode.

3. The repositories to which the ACLD gives access were extended (in addition to AMI Corpus documents, and to websites via Google): the user has now the possibility to add any local file or folder to the file repository indexed by Lucene, or can use their entire local disk thanks to a module which uses Google Desktop search.

4. Integration with real-time input processing modules from AMIDA partners was improved and made more robust. The AMIDA real-time automatic speech recognition (ASR) system (Garner et al., 2009) was coupled with the ACLD in several settings, such as meetings but also monologues (i.e. talks), and was also experimented with in a remote meeting setting, with Adobe Connect. The keyword spotting module (KWS) (Szoke et al., 2005) was also used, using in particular the possibility to add new keywords to recognize in real-time. The broadcast of audio to these modules (and to the user, in demo mode) was made possible by the use of the UEFC Media Server from UT (see Section 3). Several versions of the ACLD, with various hardware and OS settings, were thus tested.

5. Feedback from potential users attending demo sessions was collected, and two pilot evaluation sessions were organized at UEDIN with English-speaking subjects.

**System Controller.** As the ACLD application grew more complex, it was decided to move the ACLD control functions, and in particular the repository management commands, into a separate System Controller, shown in Figure 1. In addition to document conversion and indexing, the controller allows the selection of the running mode or "scenario" (live meeting or demonstration on past meetings from the AMI Corpus, with ASR
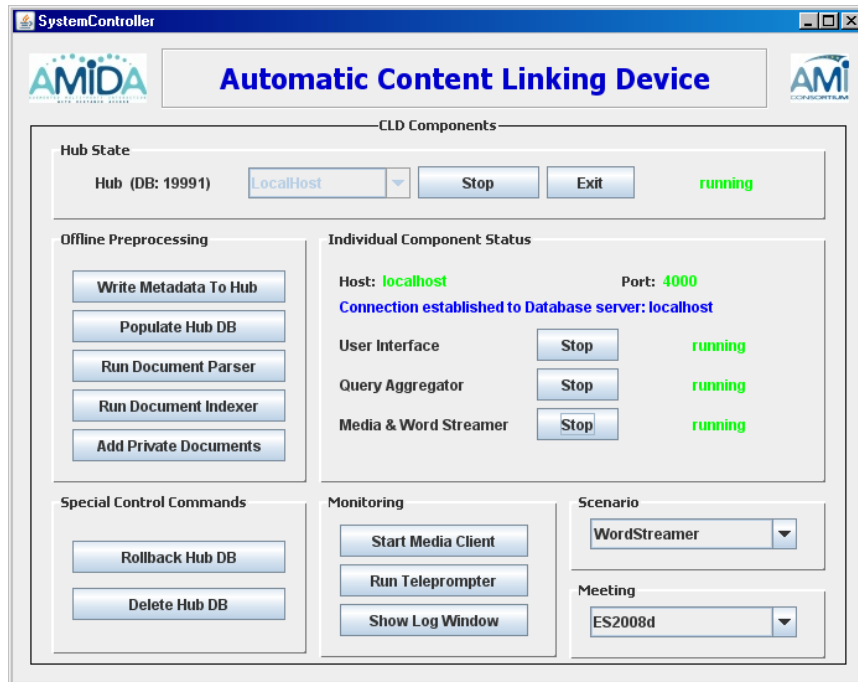
Figure 1: Snapshot of the ACLD System Controller, separated from the main UI.

or KWS), the selection of the meeting to replay for a demonstration, and includes commands for starting and resetting the various components and displaying logs.

**The Modular User Interface.** A new interface was designed and implemented in the past year. The requirements highlighted the need for a flexible interface, to be used either as an unobtrusive "widget" or as a more informative "wide screen UI". The proposed solution, the new Modular UI, contains the following four widgets. A fifth widget, under development at the time of writing, will display the results of Google Desktop search in a similar way to the first two ones from the list below.

1. Labels of the relevant documents and past meeting snippets found in the meeting index (with an appropriate icon corresponding to the document type). The position of a label in the result list, as well as the font size and emphasis, indicate its relevance within the query result.

2. Labels of relevant web links found within a Web domain that can be specified from the menu, with relevance indicated in a similar way.

3. Keywords recognized in the respective time interval, represented as a tag-cloud coding for recency and overall frequency of mention.

4. All words recognized by ASR with highlighted keywords (this tab is activated in Figure 3).

The widgets can be enabled, disabled, and arranged at will. The 'View' menu allows for widget activation and deactivation, and arrangement in the window can be done simply by
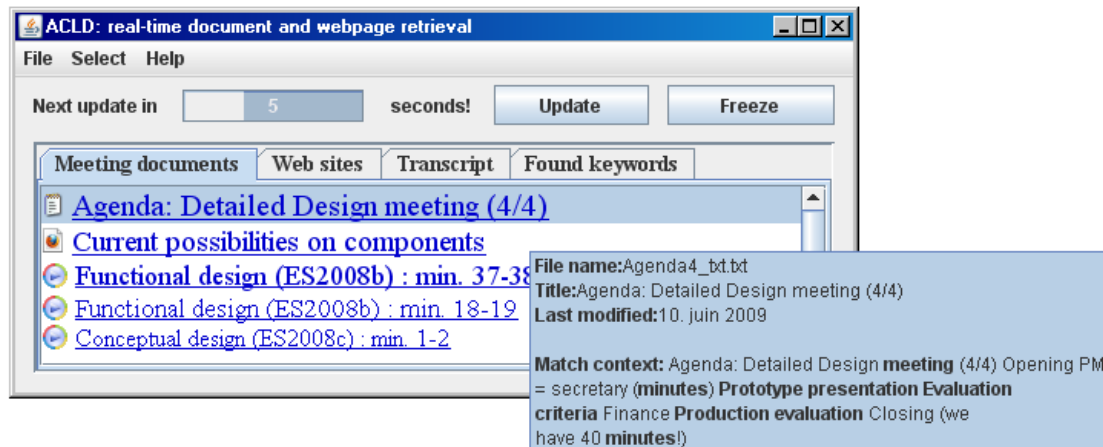
Figure 2: The ACLD UI with four superposed tabs, showing the list of relevant documents at a given moment, with explicit labels. Hovering over a label displays the metadata associated with the document, as well as excerpts where the keywords were found, with keywords in boldface.

dragging and dropping them with the mouse. Figure 2 shows a snapshot of the Modular UI with all widgets superposed as tabs, while Figure 3 shows the Modular UI with all four widgets represented as separate tabs, occupying the four quarters of the window.

The superposed arrangement of the ACLD UI (Figure 2) shows exclusively content that is actually delivered by the ACLD at runtime, with the other tabs in the background or de-activated. Hovering over a result link (document, meeting snippet, web link) provides metadata about it in a pop-up window, including most importantly the match context (as shown in Figure 2), which shows excerpts of the document that match keywords and words detected from speech, with surrounding words. Clicking on a document name opens the respective document using an appropriate viewing program – respectively, a native editor, the JFerret meeting browser (Wellner et al., 2004), or a Web browser.

**Personalization of document repositories.** A new module named Document Bank Parser was designed in order to insert into the Hub's database the information regarding all the documents related to a meeting from the AMI Corpus, or from a user's own folder, regardless of their native format (MS Office: Word, PowerPoint, Excel, and Visio; plus HTML and TXT files). The module also inserts into the Hub metadata for meeting documents, indicating in particular which documents from the AMI Corpus are available for each meeting of the corpus (documents presented prior to or at the meeting).

**Pilot evaluation sessions.** The evaluation setting was the same as in the former Task-Based Evaluation (TBE) experiments (AMI Consortium, 2006), in order to be able to compare results to control groups from these experiments. Four subjects were asked to complete the design of a remote control that was started by another group during three previous meetings – ES2008a, b, c from the AMI Corpus (Carletta et al., 2006). The subjects had a limited time to consult the meetings using JFerret and the documents, but were given the possibility to use the ACLD during their meeting. The ACLD was
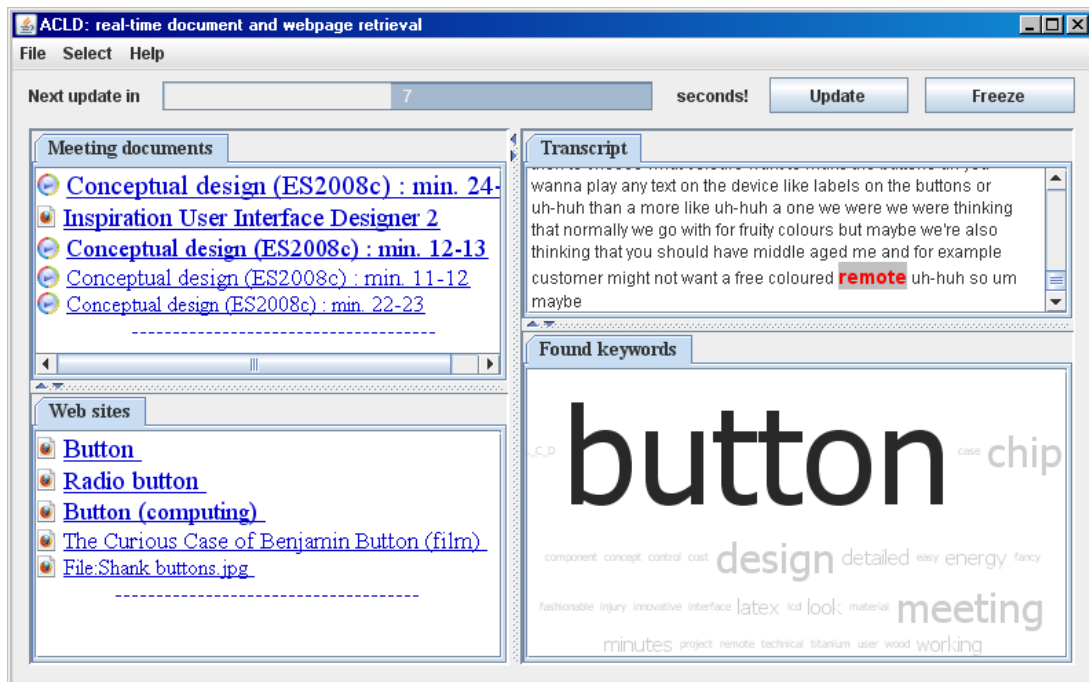
Figure 3: The ACLD Modular UI with the four widgets displayed side-by-side ("wide screen" mode).

configured to run on a central computer, using ASR over the signal obtained from mixing input from all four headset microphones, but each user had an instance of the UI on their laptop.

The results from two sessions showed that in this setting, users seemed too busy with carrying out the task to devote their attention to the ACLD. Although some users did consult some of the documents found by ACLD, some others minimized the UI, and then ignored its results. The average consultation, 3–5 documents per user per session, did not provide a sufficient quantity of data to perform user-oriented evaluation, and therefore this protocol was not considered a realistic way to acquire data for evaluation, given our limited resources in terms of number of sessions that can be set up. However, one of the lessons learned from these experiments was that preliminary instruction on how to use the ACLD is essential to obtain a sufficient amount of user-data. To this purpose, a presentation video was produced[3], which improved usage of the ACLD in the second group with respect to the first one.

## 3   User Engagement and Floor Control (UEFC) Demonstrator

The goal of the User Engagement and Floor Control (UEFC) system is to improve integration into meetings of the remote participants, by detecting and displaying significant communicative events in real-time, and in particular alerting the remote participants on when they are being addressed, and whether words of interest are pronounced (op den Akker et al., 2009). In the final year of the AMIDA project, the following progress has

---

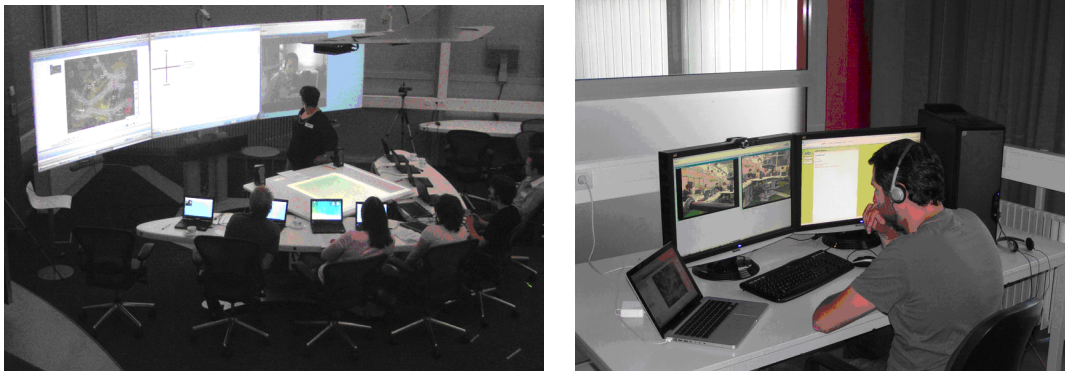[3]http://www.idiap.ch/technology-transfer/demonstrations/content-linking-device

Figure 4: Meeting room and remote participant in the TXchange project.

been made on the software for media streaming (HMI Media Server), and the UEFC demonstrator. Below is a list of the main achievements over the past year.

**Website and Deliverables.**   The website for the HMI Meeting Room software with Web Start versions of the UEFC application has been improved and contains user manuals. The software described below is available for project partners on the AMIDA internal CVS repository at UEDIN.

**More versatile configurations for media streaming.**   The first version of the UEFC demo used a fixed configuration for media streaming (op den Akker et al., 2009). In that version, the configuration was fixed for a scenario with one remote participant, one overview camera, and one or more local participants. In the current version, it is possible to make configurations for any topology of audio and video streams, with variable numbers of broadcasters and receivers. Several configurations have been tested. The media streaming package is used in the technology-transfer AMIDA mini-project between the U. of Twente and TXchange, and shown in Figure 4.

**Recovery of lost connections.**   The current streaming software is much more robust in unreliable network environments, and allows restarting separate components rather than the whole distributed system.

**Access to the raw (uncompressed RGB) video stream on both the broadcasting and the receiving ends.**   In the previous version of the UEFC demo an extra camera was needed for capturing the video data for the visual focus of attention (VFOA) recognizer. This new feature enables direct streaming to VFOA (no separate camera needed), and software rendering, for example on 3D surfaces (see 3D User Interface below).

**Synchronisation.**   Experiments showed that DirectShow (on which our media framework is built) does not synchronise streams even if they are in one graph. The current version automatically synchronises a video stream to the audio stream, but it's also possible to synchronise streams to any clock. In the current UEFC demo we broadcast a

clock over the Hub and use it to synchronise all streams, even when they are captured on different computers.

**Integration of the speech recogniser.**    The real-time ASR is integrated in the UEFC demo.
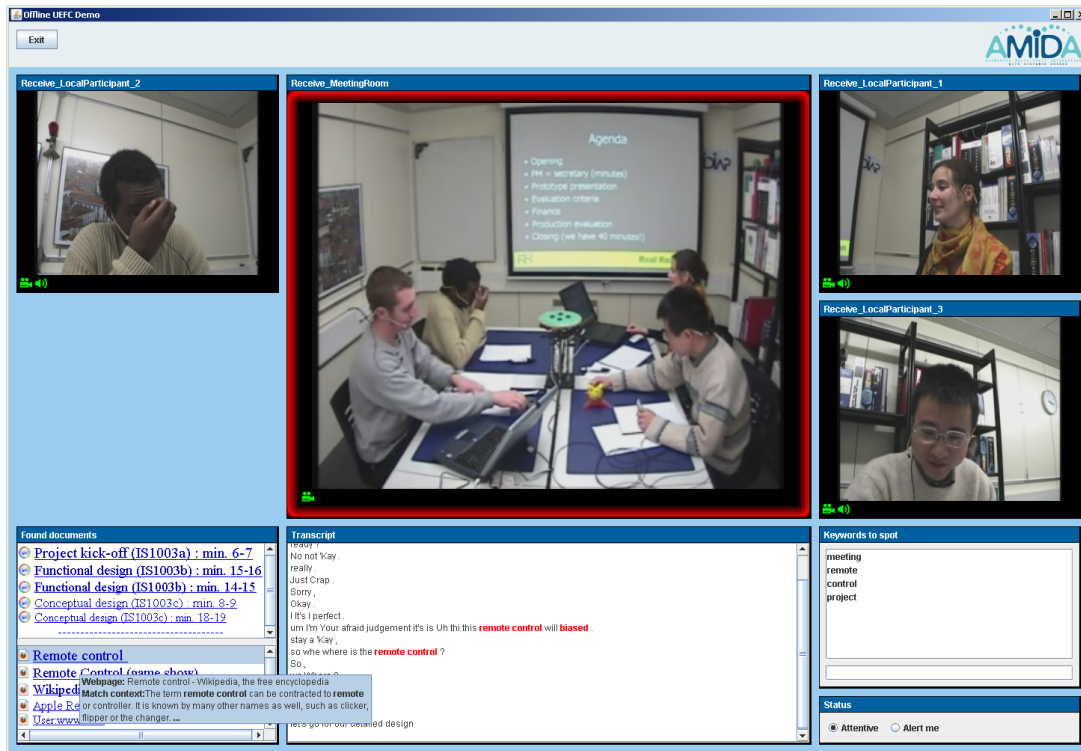


Figure 5: GUI of the UEFC demonstrator, shown here over a pre-recorded meeting. The left lower frame contains links to documents – retrieved by the document retrieval system of the integrated ACLD – in which key phrases occur that the user has specified in the right lower frame.

**Real-time Addressee Detection.**    In the first version of the online UEFC demo, the module for detecting if the remote participant is being addressed by the speaker was based on the VFOA module.  A new addressee detection module has been developed and connected to the Hub; the module uses dialogue act segmentation and labeling as inputs.  This module uses a classifier trained using machine learning, as described in (op den Akker and op den Akker, 2009). Improvement of the quality of the new addressee detection module using output of the real-time ASR and DA detection remains to be assessed.

**Integration with the ACLD.**    The first version of the UEFC demo was using an older version of the Hub.  The UEFC demo has been upgraded to the same Hub version as the ACLD, and it conforms to the same Hub message formats as the ACLD. This means that the offline UEFC and ACLD demos can run on the same instance of the Hub and share

the same annotations and keyword lists. A component of the ACLD interface showing relevant document links, has been integrated into the UEFC interface. Figure 5 shows the GUI of the integrated system. This means that integration of ACLD and UEFC demo as was proposed by the reviewers has been achieved by integrating the central part of the ACLD functionality into the GUI of the UEFC. The resulting application allows remote participants to benefit from the documents suggested by the ACLD, but also to better sense the current scope of the meeting, or even be alerted when certain documents of interest are found by the ACLD (not necessarily as mentioned by participants, but related to the discussion).

**3D User Interface.** A new 3D user interface for the UEFC demo is being developed. It allows the presentation of a number of remote participants as if they were in the same room. We are in particular interested to see how such a 3D interface affects the feeling of presence and the floor control of participants. User evaluations are planned in end of this year.

## 4   Other Applications and Demonstrations[4]

### 4.1   Audio Spatialisation

We continued our work on audio spatialisation – i.e., creating a 3D-ambience for a remote meeting participant with a different position for each speaker – and expanded our interface to allow users to freely place speakers in two dimensional acoustic space. The basic GUI is shown in Figure 6. We explored the use of the interface and how listeners placed speakers in this space when given a simple keyword spotting task (Wrigley et al., 2009). We also built a demonstrator application which allowed listeners to freely place speakers in 2D space and also automatically placed speakers according to the amount of time they had been speaking for. Thus the interface allows meeting participants to identify speakers who are not fully participating and ensures that those speakers will be heard over the rest of the meeting participants when they start speaking.

### 4.2   Audio Catchup

We completed our prototype implementation of audio catchup and carried out a user evaluation of the technology (Tucker et al., 2010). The catchup principle is shown in Figure 7. Normally late arrivals in meetings have to either permanently miss the information which was discussed when they were not present, or the meeting must be disrupted in order to summarise it for them. With Catchup the late participant sacrifices some of the meeting time in order to catch up with the missed content which they do under temporal compression (Tucker and Whittaker, 2008). Once caught up they join the meeting in real time. Our evaluations showed that this catchup period was effective and that the sacrifice is worthwhile: Catchup participants had an increased understanding of the meeting in

---

[4]These applications have been designed at USFD, as shown also by the references quoted for each of them. The last one is a collaboration between IDIAP and USFD.
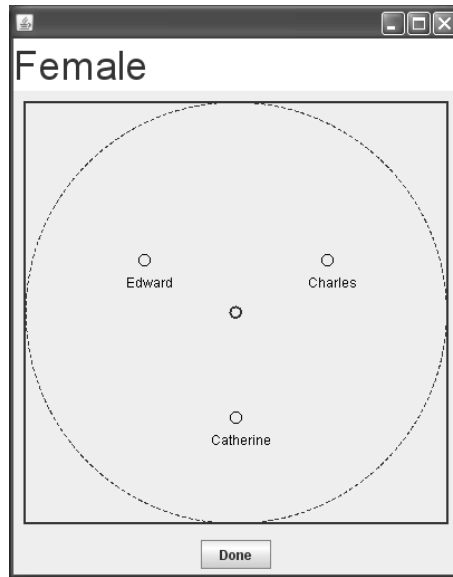
Figure 6: Interface for supporting the free placement of speakers in a two dimensional space.

comparison to those that missed the initial start of the meeting but started listening in real time earlier.
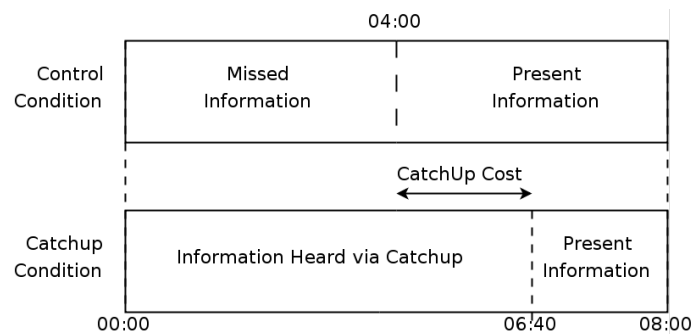


Figure 7: Representation of the benefits of audio catchup: the top portion shows typical result of arriving late to a meeting, the bottom portion shows how missing information can be recovered using Catchup.

## 4.3   NetNotes

NetNotes is a simple note taking device for both recording audio and co-indexing it with notes taken. The system records notes as individual key strokes and also broadcasts the time that the key was pressed (but not the actual key itself) to the Hub. This allows a network of meeting participants to record both their notes as text (and the time at which the notes were taken) and the times at which the other participants were taking notes. The interface to NetNotes is shown in Figure 8. The audio can then be combined with the notes in order to produce summaries of the relevant portions of the meeting as well as the personal notes of the user being used to locate information within the meeting.
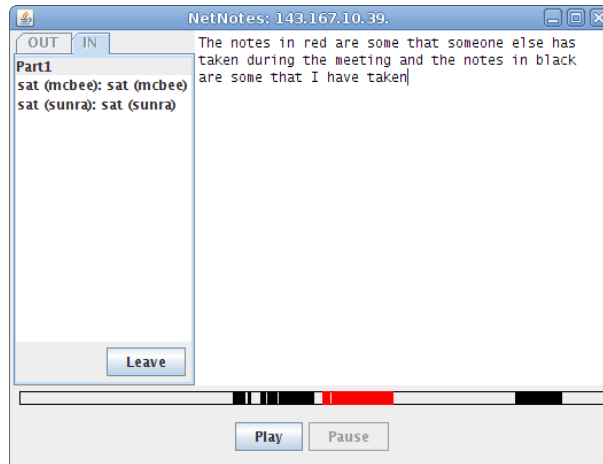
Figure 8: The NetNotes interface. On the left hand side is the list of participants within the meeting, the centre screen shows the personal notes taken by the user. The bottom panel shows when the user (black) and others(red) have been taking notes in the meeting.

## 4.4 Augmented Search

We have also been experimenting with a mechanism to augment search results[5], such as from Google or those presented in the ACLD demonstration. The augmented search results show discriminative terms, which help to indicate how each search result is *different* from the other search results. The mechanism is currently based on the well-known TF-IDF technique, but could easily be extended to more sophisticated techniques. A simple demonstration is available at: `http://mmm-web.idiap.ch:8180/tfidf/` and Figure 9 shows an example of result when searching the Web for 'meeting AMIDA'.

## 5 Conclusion

This report has shown that the focus of AMIDA.WP6 has been, during the past year, on the two main integrative demonstrations foreseen for the end of the AMIDA project, the ACLD and the UEFC, which have now reached acceptable maturity as research prototypes. Among the functionalities added during the past year, priority was given to those increasing acceptability, user-friendliness, and overall acceptability. Pilot trials have tended to demonstrate the validity of the concepts, and the suitability of the core technologies involved, as well as the support infrastructure.

## References

AMI Consortium (2006). Meeting browser evaluation. Deliverable 6.4, AMI (Augmented Multi-party Interaction) Integrated Project FP6-506811.

---

[5]A joint exploration between Mike Flynn at IDIAP and Steve Whittaker at USFD.

Figure 9: The Augmented Search application, showing search results with added discriminative terms.

AMI Consortium (2007). Commercial component definition. Deliverable 7.2, AMIDA (Augmented Multi-party Interaction with Distance Access) Integrated Project IST-033812.

Carletta, J., Ashby, S., Bourban, S., Flynn, M., Guillemot, M., Hain, T., Kadlec, J., Karaiskos, V., Kraaij, W., Kronenthal, M., Lathoud, G., Lincoln, M., Lisowska, A., McCowan, I., Post, W., Reidsma, D., and Wellner, P. (2006). The AMI Meeting Corpus: A pre-announcement. In Renals, S. and Bengio, S., editors, *Machine Learning for Multimodal Interaction II*, LNCS 3869, pages 28–39. Springer-Verlag, Berlin/Heidelberg.

Garner, P. N., Dines, J., Hain, T., El Hannani, A., Karafiat, M., Korchagin, D., Lincoln, M., Wan, V., and Zhang, L. (2009). Real-time ASR from meetings. In *Proceedings of Interspeech 2009*, pages 2119–2122, Brighton, UK.

op den Akker, H. and op den Akker, R. (2009). Are you being addressed? – real-time addressee detection to support remote participants in hybrid meetings. In *Proceedings*

*of SIGDIAL 2009 (10th Annual Meeting of the Special Interest Group in Discourse and Dialogue)*, pages 21–28.

op den Akker, H. J. A., Hofs, D. H. W., Hondorp, G. H. W., op den Akker, H., Zwiers, J., and Nijholt, A. (2009). Supporting engagement and floor control in hybrid meetings. In Esposito, A. and Vich, R., editors, *Cross-Modal Analysis of Speech, Gestures, Gaze and Facial Expressions, Prague*, volume 5641 of *Lecture Notes in Computer Science*, pages 276–290, Berlin/Heidelberg. Springer-Verlag.

Popescu-Belis, A., Boertjes, E., Kilgour, J., Poller, P., Castronovo, S., Wilson, T., Jaimes, A., and Carletta, J. (2008). The AMIDA Automatic Content Linking Device: Just-in-time document retrieval in meetings. In Popescu-Belis, A. and Stiefelhagen, R., editors, *Machine Learning for Multimodal Interaction V (Proceedings of MLMI 2008, Utrecht, 8-10 September 2008)*, LNCS 5237, pages 272–283. Springer-Verlag, Berlin/Heidelberg.

Popescu-Belis, A., Carletta, J., Kilgour, J., and Poller, P. (2009a). Accessing a large multimodal corpus using an Automatic Content Linking Device. In Kipp, M., Martin, J.-C., Paggio, P., and Heylen, D., editors, *Multimodal Corpora: From Models of Natural Interaction to Systems and Applications*, LNAI 5509, pages 189–206. Springer-Verlag, Berlin/Heidelberg.

Popescu-Belis, A., Poller, P., Kilgour, J., Boertjes, E., Carletta, J., Castronovo, S., Fapso, M., Flynn, M., Nanchen, A., Wilson, T., de Wit, J., and Yazdani, M. (2009b). A multimedia retrieval system using speech input. In *ICMI-MLMI 2009 (11th International Conference on Multimodal Interfaces and 6th Workshop on Machine Learning for Multimodal Interaction)*. ACM Press, Cambridge, MA.

Szoke, I., Schwarz, P., Matejka, P., Burget, L., Karafiat, M., Fapso, M., and Cernocky, J. (2005). Comparison of keyword spotting approaches for informal continuous speech. In *Eurospeech 2005 (9th European Conference on Speech Communication and Technology)*, pages 633–636, Lisbon.

Tucker, S., Ramamoorthy, A., Bergman, O., and Whittaker, S. (2010). Catchup: A useful application of time-travel in meetings. In *Proceedings of CSCW 2010*, Savannah, GA.

Tucker, S. and Whittaker, S. (2008). Temporal compression of speech: An evaluation. *IEEE Transactions on Audio, Speech and Language Processing*, 16(4):790–796.

Wellner, P., Flynn, M., and Guillemot, M. (2004). Browsing recorded meetings with Ferret. In Bengio, S. and Bourlard, H., editors, *Machine Learning for Multimodal Interaction I*, LNCS 3361, pages 12–21. Springer-Verlag, Berlin/Heidelberg.

Wrigley, S., Tucker, S., Brown, G., and Whittaker, S. (2009). Audio spatialisation strategies for multitasking during teleconferences. In *Proceedings of Interspeech 2009*, pages 2935–2938, Brighton, UK.