



ON AUTOMATIC ANNOTATION OF MEETING DATABASES

Daniel Gatica-Perez * Iain McCowan *
Mark Barnard * Samy Bengio *
Hervé Bourlard *

IDIAP-RR 03-06

JANUARY 2003

Dalle Molle Institute
for Perceptual Artificial
Intelligence • P.O.Box 592 •
Martigny • Valais • Switzerland

phone +41 - 27 - 721 77 11
fax +41 - 27 - 721 77 12
e-mail secretariat@idiap.ch
internet <http://www.idiap.ch>

* IDIAP, Martigny, Switzerland

1 Introduction

Multimedia content analysis addresses, among many others, the task of automatically annotating audio-visual material with labels relevant to browsing and retrieval [12]. Annotation is a rich domain; here we focus on single labels that name semantic entities. The labels could consist of relatively low-level events or concepts, such as words, identities and specific objects, as well as higher-level semantic concepts, such as a weather report within news, a goal within a football match, or the genre of a documentary.

Automatic video annotation has mostly been focused on a small number of application domains, including broadcast news, sports videos, and documentaries. The data in these applications is highly produced, and thus has a strongly imposed structure due to shot cuts that segment a video into a coherent ‘story’. For this reason, most automatic annotation approaches to date have used shots as the basis for event segmentation and classification. In the more general case, however, produced content cannot be assured, and so this reliance on shots as the fundamental unit for processing and recognition is a limiting assumption. Meetings [7, 11, 2] provide a counter-example in which the input media naturally consists of raw audio and visual streams.

In order to progress from low-level to more semantic annotations, statistical models are commonly used to infer high-level events from lower-level visual and audio features. Hidden Markov models (HMMs) are sequence models that have been used for a variety of visual and audio processing tasks. For the task of video annotation, HMMs have been used to model the content of broadcast material such as news [3], documentaries [5], and sports [13, 14].

Two distinct approaches to such event-based semantic annotation exist. Perhaps the most common approach is to consider a set of events that occur sporadically within the data stream. Commonly, systems following this approach classify each segment (usually a shot) according to the presence/absence of each event using decision thresholding. A second approach is to consider that the data stream consists of a continuous sequence of events, and in this case continuous decoding strategies can be employed, alleviating the need for a pre-segmentation. In this paper we investigate such an approach, in which meetings are decomposed into a sequence of *meeting actions*, such as discussions, presentations and note-taking. The annotation task is then clearly defined as recognising the correct sequence of these meeting actions.

As well as having a well-defined annotation task, there is a need for standardised measures by which the quality of the annotations can be assessed. One response to this need is the NIST TREC video track project [10]. The metrics used in the NIST TREC 2001 evaluation were *recall* and *precision*, which relate mainly to retrieval and two-class classification problems. However, to assess the quality of video annotations involving more than two classes, performance measures are still non-standard. In this paper, we advocate the use of the word error rate commonly used in the speech recognition domain. If the video annotation task is defined as the recognition of a continuous sequence of events (as discussed above), and when shot-cut boundaries are not present (or are irrelevant to semantic events), the word (or event) error rate presents a natural and effective metric for system performance.

The paper is organised as follows. Section 2 discusses semantic annotation of meetings in the context of multimedia content analysis. Section 3 describes our approach in detail, including event definition, the performance evaluation protocol, and discusses both its applicability to other domains and its limitations. Section 4 provides some final remarks.

2 Meetings as multimedia data

2.1 The “nature” of meeting data

Meetings depict people interaction, and occur in reasonably constrained yet challenging conditions. As a source of multimedia information, meetings consist of unedited streams of audio and video, captured with multiple cameras (covering participants and workspace areas, including whiteboards or projector screens) and microphones. A possible production model for real-time communication could merge AV

streams focused on the current speaker(s) [2], possibly using motorised cameras. However, in many real settings (including the one described in this paper) cameras are fixed and AV data are archived in raw form.

In this setting, the typical concepts of shots and scenes are absent. Cameras and microphones continuously capture people engaged in discussions, gesticulating, and making/listening to presentations. The continuous nature of the data renders methods for discovering syntactic rules (commonly studied in multimedia content analysis) of relatively little relevance.

However, meetings are strongly structured data in semantic terms. Events at different semantic levels, ranging from low-level actions and gestures (entering or leaving the room, standing to make a presentation, raising a hand) to high level actions (discussing, doing monologues, making presentations) to very high-level notions (planning, negotiating, making decisions) are all common in meetings, and induce semantic hierarchies in the data. In many cases, these events represent meaningful annotations and could be directly used as queries in a retrieval system.

As multimedia data, meetings have common features with other data types generated by “looking at people”, like surveillance and instructional videos (raw data, multiple cameras, with an obvious difference in the number and quality of the audio sources), and also share characteristics with some highly produced content, like news and interviews programs (where speakers play a leading role, and the audio track represents a very strong cue). As a result, many of the problems in analysing meetings are shared by other domains.

2.2 The relevance of analysing meetings

The amount and relevance of semantic labels that can be potentially extracted from meetings for annotation are considerable, both from what is said and from what is done. At the individual meeting level, annotations could improve collaborative work by helping people quickly retrieve information from a meeting archive without having to listen/view entire recordings. At the database level, important high-level trends could be discovered and attached to the database as labels, useful for organisational management tasks.

Paraphrasing [1], meeting databases have three clear highlights in terms of research relevance and applicability. In the first place, meetings -as raw data- are suitable for the automatic generation of metadata not available from production. Although identities of participants -and to some extent the basic semantic structure if an agenda was available- could be potentially extracted by automatic means at the time of acquisition, peoples’ statements and actions are natural and cannot be generated at production time without human intervention. In the second place, off-line annotation of meetings is a task for which humans are not good/fast at generating. In the third place, meetings occur regularly and are often generated in large numbers. As the number of meetings increases, however, their individual value tends to decrease, mainly in terms of novelty. Analysing a meeting database would increase the value of the raw data if annotations related to management tasks (like the progress of a project over a period of time) could be produced.

Needless to say, the research problems at hand are ambitious. It is also clear that, even though many of the very high-levels events cannot be recognised by current state-of-the-art means (computer vision, speech processing, data fusion, language modeling, text retrieval), the current technology (as witnessed by other domains in multimedia content analysis) should allow for the labeling of low- and some high-level semantic events.

2.3 Semantic annotations in meetings

Although far from perfect, identifying meeting participants, transcribing what they say, and partially inferring what they do, are becoming feasible. However, given the large list of potential events, what specific events should we (or can we) annotate? We can briefly state four broad categories, each generating annotations of distinct (but possibly complementary) nature and complexity.

Speech transcription-based annotation. The analysis of ASR transcriptions by text retrieval and language modeling techniques is expected to generate the highest-level concepts for annotation, varying from specific key-word detection and recognition of participant identities, to topic and subtopic detection, etc. [11, 7].

Audio-based annotation. When rooms are equipped with microphone arrays, the approximate location of participants can be robustly inferred with from the audio streams [4], so location-based events can be identified, including monologues, turn-taking, or presentations. However, the amount of (non-speech) events that can be extracted from audio-only is limited.

Video-based annotation. Recognition of people and some of their actions can currently be addressed by computer vision algorithms to perform person identification in individual meetings (face detection/recognition) and across meetings (face clustering), gesture recognition, facial expression recognition, etc. Video-based annotation faces two main challenges : robustness and usefulness. Many of the low-level semantic labels that can be generated (usually related to recognition of sparse, low-level actions performed by individuals) do not constitute annotations directly useful for indexing or retrieval (nobody needs to query a system looking for people standing up from their seats, or pointing to the whiteboard). In other words, the problem of mapping low-level features and gestures to semantically meaningful concepts remains open, as in all other multimedia content analysis domains. These events, however, can be the building blocks towards recognising high-level semantic events. Note that the definition of high-level events admits multiple dictionaries and different levels of semantic granularity.

Multi-modal annotation implies the development of principled frameworks for the integration of multiple data streams of different nature and frame rate (audio, text from speech transcripts, and video in this paper) to detect events. In meetings, events are inherently multimodal, but the involved modalities have complex relations (they might be asynchronous, and contain significantly distinct amounts of relevant information related to the event). The general goal is to combine low-level features and events provided by the individual modalities into high-level event recognisers.

3 Annotating Meetings as a Sequence of Events

In this section, we consider meetings as continuous sequences of AV events with natural transitions. If a list of possible events is defined, the task of annotation then consists of finding the sequence of events that constitutes a particular meeting. Given such a definition, the video annotation task becomes analogous to that of speech recognition, and so a similar training, decoding and assessment methodology can be employed.

3.1 Definition of Events

Many different sets of events could be defined for the task of meeting annotation. In [6], we proposed a list of events characterised by group behaviour of meeting participants. This list included monologues (by participant), discussions, consensus, disagreement, presentations, white-boards and note-taking. These are all natural actions in which participants play and exchange similar, opposite, or complementary roles. As these events are based on group interactions, we refer to them as *meeting actions*.

The definition of such a lexicon of meeting actions is interesting from a research perspective, as recognition of group interactions could be approached from at least two distinct angles. In a first case, the actions of individual participants could be recognised, and then these responses fused at a higher level to recognise the interaction. Such an approach, however, overlooks the fact that the behaviour of individuals in meetings is somehow constrained by the behaviour of the other participants. A second approach (taken here) is to model the interactions directly, by integrating all observations into a unique probabilistic model and learning the constraints from the data. If the group as a whole provides enough evidence for the performed action, recognition of personal actions could be bypassed altogether, potentially increasing robustness to imperfect feature extraction and measurement processes.

From the retrieval viewpoint, defining the events based on group actions has the benefit of attaching a single semantic annotation to all audio-visual streams. In contrast, annotations based on individual actions would result in different annotations for each camera or microphone. Also, individual actions would tend to be sparse in nature, while the above list of meeting actions can be treated as a continuous sequence.

3.2 Methodology

To annotate meetings as a sequence of events, we use statistical generative models based on HMMs [8]. HMMs have been successfully used to recognise speech, visual and audio-visual sequences. When the video annotation problem can be posed as recognising a continuous sequence of events, techniques and assessment metrics can be borrowed from these other tasks.

To use HMMs to annotate meetings, we require an event lexicon (as described above) and feature vectors appropriate for measuring the defined events. Given a training sequence of feature vectors with the corresponding labelling (but not necessarily the precise alignment), HMMs can be trained using the classical embedded training method based on Expectation-Maximisation (EM). Recognition then simply involves application of the Viterbi decoding algorithm to find the most likely sequence of meeting actions.

3.3 Meeting Database

A corpus of meetings was recently recorded in the IDIAP smart meeting room [6]. Meetings were recorded using 3 cameras and 12 microphones, with all channels fully synchronised. Currently the database contains 60 meetings (30 train, 30 test), where each meeting consists of 4 participants and lasts approximately 5 minutes. The meetings are loosely scripted in terms of the type and schedule of the high-level actions, but otherwise the content is natural. The corpus is fully described in [6] and is being expanded and made available for public distribution [15].

3.4 Performance Evaluation

Speech recognition is often quoted as an example of a processing domain where research has been greatly aided by the use of standard performance metrics, facilitating comparisons between different systems. While performance measures for retrieval have been largely standardised, methods of assessing the accuracy of video annotations are still largely system dependent [3, 13, 14]. A major benefit of posing the annotation problem as described above, is that standard performance metrics, such as the word error rate used in speech recognition, may be employed. This was acknowledged long ago in computer vision for gesture recognition [9].

The above methodology for annotating meetings was applied in [6] using our meeting corpus. A feature vector of 19 audio-visual features was extracted from the input channels at a rate of 5 Hz. From 2 cameras looking at people at the table, GMM models of skin/background colours in RGB space were used to extract head blobs. Skin/background pixel classification and morphological post-processing were performed inside image regions enclosing typical head locations. For each person, the detected head blob was represented by the vertical position of its (normalised) centroid. From a wide-view camera capturing the presentation screen and white-board area, moving blobs were detected by background subtraction and represented by their (quantised) horizontal position. Audio features were extracted to measure the speech activity of different locations, as well as the occurrence of a set of positive and negative keywords. These features were used to train HMMs for meeting actions using the train set.

The system performance was assessed on the test set in terms of the *action error rate*, which is equivalent to the word error rate in speech recognition. The word (event or action) error rate is an appropriate metric where finding the correct sequence of annotations is more important than precisely

	mono1	mono2	mono3	mono4	white	note	cons	disc	pres	disa	DEL
mono1	10										1
mono2		9							1		
mono3			17								
mono4				10							1
white					18						
note						6					
cons								6			9
disc								45			
pres					1				12		1
disa								1			7
INS		1		1				1			

Table 1: Confusion Matrix of Recognised Meeting Actions, including monologues (mono1-4), whiteboards (white), note-taking (note), consensus (cons), discussions (disc), presentations (pres) and disagreements (disa). Zero values are represented as empty cells. Columns and rows show desired and obtained labels, respectively.

determining their temporal boundaries. This is often the case when the annotation labels are high-level semantic concepts. The word error rate is calculated as the number of substitution, insertion and deletion errors, divided by the correct number of words. Due to the inclusion of insertion and deletions in the error rate calculation, it is a more severe measure than classification accuracy. Video annotation systems are commonly designed as ‘shot classifiers’, in which case insertions and deletions do not occur, however the use of shots as a fundamental unit is often not appropriate, and the word error rate is a measure with more general application. The overall action error rate achieved in these experiments was 20.0% [6].

In addition to a standard performance measure, it is also necessary to analyse the results to determine common sources of errors. A useful analysis tool for a small vocabulary recognition task is the *confusion matrix*, which shows the distribution of recognised events according to events in the ground-truth (note that this differs from the standard multi-class confusion matrix, due to the lack of hard boundaries, consequently including deletions and insertions). The confusion matrix for the above task is shown in Table 1. Analysis of the confusion matrix is particularly useful in this case, as it shows that neither consensus and disagreement are recognised correctly, instead being commonly confused with discussion or deleted. These are examples of events for which features are clearly not discriminative enough and requires further research. This observation is discussed in detail in [6], and it is shown that if consensus and disagreement are removed from the lexicon by relabelling them as discussions for both training and testing purposes, then the action error rate decreases to 5.7%.

3.5 Applicability to other Video Annotation Domains

The above methodology for meeting video annotation is applicable to other domains where inherent structure exists such that the video can be considered as a continuous sequence of events. For example, in [3], such an approach is taken to annotate televised news broadcasts in terms of content classes. Sports videos and documentaries are other domains where such structure may exist.

While such an approach has been investigated across different annotation domains, the method of reporting results is still non-standard. Different methods of assessing event segment boundaries are used, and classification accuracies differ based on an assumed level of segmentation (frames, shots, scenes, etc). As discussed above, where the annotation labels are high-level semantic concepts (eg, presentations, discussions, interviews, shots at goal), often the concept of precise boundaries between segments has little relevance. Also, as multiple events can occur within a video shot, or conversely a single event could span multiple shots, shots are not always an appropriate unit for classification. In such a context, the ‘word’ error rate is a meaningful performance measure that could be adopted across different video annotation systems recognising such a continuous sequence of semantic events.

3.6 Limitations

This methodology for video annotation has a number of limitations. First, it excludes the co-occurrence of multiple events at a given time. Second, it cannot explicitly handle the case when events occur sporadically, and not as a continuous sequence. In some cases, the first limitation could be addressed by employing a hierarchical recognition scheme, in which recognised events are decomposed into a further sequence of sub-events. As a simple example of this in the context of meetings, we could handle the occurrence of note-taking during presentations by first recognising presentations, and then recognising this as a sequence of segments with or without note-taking. Clearly such a hierarchical system has limited application, and a need exists for a more general methodology allowing the joint occurrence of multiple events.

The second limitation could be addressed by introducing a ‘silence’ or ‘garbage’ event model to match periods where no explicit events occur. This is analogous to the approach taken in speech keyword spotting systems. In such an approach, however, selection of an appropriate garbage model is often a non-trivial task.

4 Conclusions

This paper has discussed meetings as a source of data for multimedia content analysis, specifically focusing on the task of automatic audio-visual event annotation. A methodology for treating meetings as a continuous sequence of events was proposed, leading to a well-defined annotation task and clear performance evaluation. As a case study, a system annotating a database of meetings as a sequence of meeting actions (monologues, presentations, discussions, white-boards, note-taking, consensus and disagreement) was presented and assessed in terms of the word (action) error rate. The advantages of our methodology, and its applicability to other types of multimedia data were discussed, along with potential limitations of the approach. In conclusion, we propose to the research community the use of this corpus [15], in general, and the particular task and evaluation measure used in this article (and [6]).

5 Acknowledgements

The authors thank the Swiss National Science Foundation for supporting this work through the National Centre of Competence in Research (NCCR) on “Interactive Multimodal Information Management (IM2)”. The work was also funded by the European project “M4: MultiModal Meeting Manager”, through the Swiss Federal Office for Education and Science (OFES). We also thank several of our colleagues at IDIAP for their collaboration in the smart meeting room project (Thierry Collado, Guillaume Lathoud, Sebastien Marcel, Olivier Masson, Florent Monay, Darren Moore, Jean-Marc Odobez, Pierre Wellner), and many others for their assistance during the database creation.

References

- [1] S.-F. Chang. The holy grail of content-based media analysis. *IEEE Multimedia Magazine*, 9(2):6–10, April-June 2002.
- [2] R. Cutler, Y. Rui, A. Gupta, JJ Cadiz, I. Tashev, L. He, A. Colburn, Z. Zhang, Z. Liu, and S. Silverberg. Distributed meetings: A meeting capture and broadcasting system. In *Proc. ACM MM Conf.*, 2002.
- [3] S. Eickeler and S. Müller. Content-based video indexing of TV broadcast news using HMMs. In *Proc. IEEE ICASSP*, Phoenix, 1999.
- [4] G. Lathoud and I. McCowan. Location based speaker segmentation. In *Proc. IEEE ICASSP (to appear)*, 2003.
- [5] T. Liu and J. R. Kender. A Hidden Markov Model approach to the structure of documentaries. In *IEEE Work. on CBAIVL*, 2000.

- [6] I. McCowan, S. Bengio, D. Gatica-Perez, G. Lathoud, F. Monay, D. Moore, P. Wellner, and H. Bourlard. Modeling human interaction in meetings. In *Proc. IEEE ICASSP (to appear)*, 2003.
- [7] N. Morgan, D. Baron, J. Edwards, D. Ellis, D. Gelbart, A. Janin, T. Pfau, E. Shriberg, and A. Stolcke. The meeting project at icsi. In *Proc. Human Lang. Techn. Conf.*, San Diego, March 2001.
- [8] L. Rabiner and B.H. Juang. *Fundamentals of Speech Recognition*. Prentice Hall, 1993.
- [9] T. Starner and A. Pentland. Visual Recognition of American Sign Language using HMMs. In *Proc. Int. Work. on Auto. Face and Gesture Recognition*, Zurich, 1995.
- [10] E. M. Voorhees and D. K. Harman, editors. *NIST Special Publication 500-250: The Tenth Text REtrieval Conference (TREC 2001)*.
- [11] A. Waibel, M. Bett, F. Metze, K. Ries, T. Schaaf, T. Schultz, H. Soltau, H. Yu, and K. Zechner. Advances in automatic meeting record creation and access. In *Proc. IEEE ICASSP*, Salt Lake City, 2001.
- [12] Y. Wang, Z. Liu, and J. Huang. Multimedia content analysis using both audio and visual clues. *IEEE Signal Processing Magazine*, 17(6):12–36, November 2000.
- [13] L. Xie, S.-F. Chang, A. Divakaran, and H. Sun. Structure analysis of soccer video with HMMs. In *Proc. IEEE ICASSP*, 2002.
- [14] G. Xu, Y.-F. Ma, H.-J. Zhang, and Shiqiang Yang. Motion based event recognition using HMM. In *Proc. ICPR*, Quebec, 2002.
- [15] IDIAP data distribution. <http://rhonedata.idiap.ch/>.