

to adapt the acoustic models to the domain. The main issue is to process the imperfect transcription of the stenographs. We generate a word grammar which allows the insertion and deletion of words given in the stenography. In a second processing step, we apply confidence scoring to identify segments of speech which have been transcribed with a high degree of reliability and which can be used for retraining the acoustic models. In a first evaluation we achieved a recognition rate of 65% for the German parliament domain tested on 400 sentences[8].

MPEG-7 provides the SpokenContentDS to describe the content of speech material using an existing speech recognition engine, like Viavoice, HTK or ISIP. The SpokenContentDS contains a description of a word and a phoneme lattice to allow retrieval on a lattice. This descriptor scheme was proposed and developed by [9] and is now part of

[10]. A lattice is a graph containing recognition alternatives weighted by probabilities. Figure 2 shows a typical lattice with similar words:

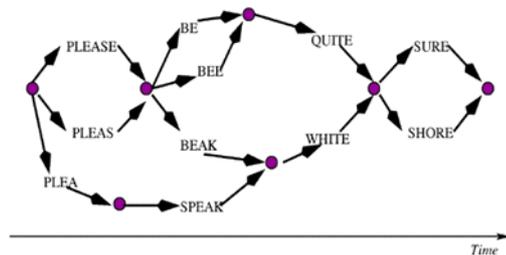


Figure 2: An example of a lattice for the SpokenContentDS (from [9])

Previous investigations on information retrieval have shown that the use of a lattice improves the performance of the system. Further, the phoneme lattice allows identification of words that are unknown to the recognition system, but that occur in retrieval queries. Proper names and new technical terms are particularly likely to be missing from the recognizer's vocabulary.

The SpokenContentDS is fully integrated in the MPEG-7 XM-system. Currently we are working on converting the ISIP recognition output to the format of the SpokenContentDS and combining this with the Transcriber toolkit [11] to annotate speech material. This system uses also XML and DTD to describe the segmentation and transcription of a spoken document.

4 Video Processing and FaceRecognitionDS

Although this work is in preliminary stage, we present our approach for video segmentation and face recognition. The video analysis is composed from three successive processing steps: temporal segmentation, face detection and tracking (figure 3), and face recognition. The temporal segmentation uses the difference image and the difference of the histograms of two consecutive frames to find the scene boundaries (cuts and dissolves) in the video sequence. A frontal face detector is used to find the faces in the scenes. The gaps between the detected frontal views of the face, where the person looks to the left or right, are filled using a simple histogram based tracking algorithm. The recognition of the detected faces is based on pseudo 2D Hidden Markov Models. One model trained on face images of each member of the German Bundestag is used to recognize the speaker at the parliament.



Figure 3: Face tracking result for a single speaker

5 AGMA System Design

The system design in AGMA is closely related to the design of the MPEG-7 XM (eXperimental Model). Because there are more than 100 different DS in MPEG-7 we have focused on those that are relevant to our goal of transcribing and annotating speeches of the German parliament. For example, the MelodyDS is not required for the German parliament domain. The most important DSs are the SpokenContentDS and the FaceRecognitionDSs. Both DSs are high-level descriptor schemes that contain semantic information already high level. Information on this level is important to do meaningful retrieval and have an effective application.

The components of the systems are shown in figure 4.

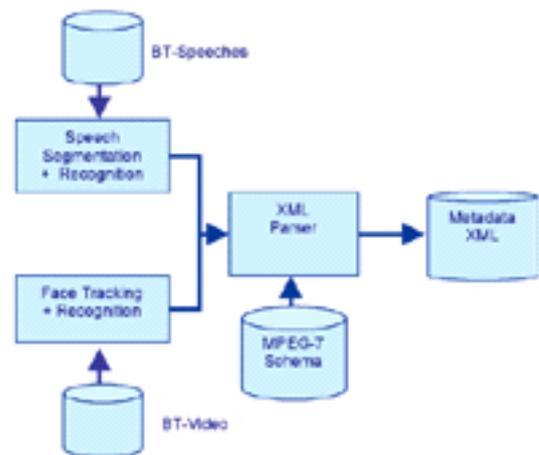


Figure 4: content extraction architecture in AGMA

The feature extraction process is divided in two stages. The speech recognition system generates a segmentation and transcription of the German parliament speeches. The output of the recognizer is a word and a phoneme lattice. In a simple version only the first best hypothesis is generated. The recognition engine is the ISIP speech recognition toolkit described in section 3. The output of the recognition process is passed to the XML parser. Here we use the Xerces parser from Apache, which is also used in the MPEG-7 XM system. The parser creates a DOM tree and is able to check the MPEG-7 schema. This step is an important stage in the process because here the recognition output format is validated. The validation guarantees conformance with the MPEG-7 standard and the interoperability with future applications.

The similar approach is carried out for the face recognition process. Before the speakers are recognized, a face finder locates the position of the face. Then the recognition process is initiated and the results are also passed to the XML parser. The results of the speech recognizer and the face recognition system are written to a XML file that contains the metadata of the audio-visual input file. This

