

Automatic Video Editing

AMIDA technology package description

Pavel Žák, Kubíček Radek

*Graph@FIT
Brno University of Technology
Faculty of Information Technology
Božetěchova 2
612 66 Brno, Czech Republic*



Technology developers:
Stanislav Sumec, Pavel Zemčik, Radek Kubíček,
Pavel Žák, Michal Hradiš

Contact persons:
Zemčik Pavel, zemcik@fit.vutbr.cz, responsible for AMIDA project
Herout Adam, herout@fit.vutbr.cz, responsible for Graph@FIT group



Contents

1 Purpose of the technology	3
2 Features	3
3 Technical description	4
3.1 Video Editing process	4
3.2 Editing algorithm	5
3.3 Input Data	7
3.3.1 Offline	7
3.3.2 Online	9
3.4 Output	9
4 Limitations	9
5 Technical specifications	10
6 Package content	10

1 Purpose of the technology

This technology provides automatic video editing, that allows automatic producing of single video stream from more cameras by choosing the one that contains most of relevant information about what is happening in the observed scene and also preserves several aesthetic aspects in the final video that makes it more attractive to the viewer.

It can be used for example in meeting environments when the meeting room is observed by several cameras. Here the technology enables meeting summarization, distant meeting participating or just for storing only important informations in distilled form (instead of all recorded streams).

2 Features

This technology has been developed primarily for meeting scenarios and is capable of processing either online or offline data. Its core forms versatile rule-based editing algorithm that can be configured for specific scenarios.

Online editing can be done in any meeting room without specific camera setup and its design enables customizing for specific scenarios. No special hardware is needed, the basic setup for realtime processing contains personal computer with several firewire cameras(see Technical specifications).



Figure 1: Online video editing system setup with three firewire cameras.

Offline editing can be done from meeting transcriptions or annotations,

which enable more accurate output. 20

The final package we propose contains implementation of this technology 21
together with its additional more detailed description. The implementation 22
is ready in configurable application for either online or offline processing. 23
Demonstration video of online editing capabilities can be seen in [1]. 24

3 Technical description 25

This section summarizes how the technology works. More detailed technical 26
information can be found in references[2][3]. 27

3.1 Video Editing process 28

The function of the proposed algorithm can be formulated as a problem of 29
one camera or of several combined cameras selection in each time point of 30
the recorded meeting. The image from the selected camera has to preferably 31
represent what is happening in the meeting room according to different (user 32
specified) aspects. Satisfaction of these aspects warrants that produced video 33
contains as much of the relevant information as possible. 34

The input of the editing process is a set of facts - results of audio/video 35
stream processing(and/or video annotations), set of rules - the predefined 36
general video editing rules and a dynamically changeable set of rules based 37
on the details of the audio/video processing, user requirements, etc. The 38
output of the video editing decision process is a text description of the editing 39
decisions that is afterwards interpreted in a video editing engine that actually 40
processes the audio/video data. The whole algorithm works in the following 41
way captured also in figure 2: 42

1. Source video streams of all cameras are simultaneously processed from 43
the beginning to the end of the meeting. 44
2. Specific set of events (or features) is detected in each stream(or just 45
loaded from annotations). 46
3. Based on detected events (which specify each stream information im- 47
portance) and aesthetic editing rules the output stream is chosen for 48
specific time. 49

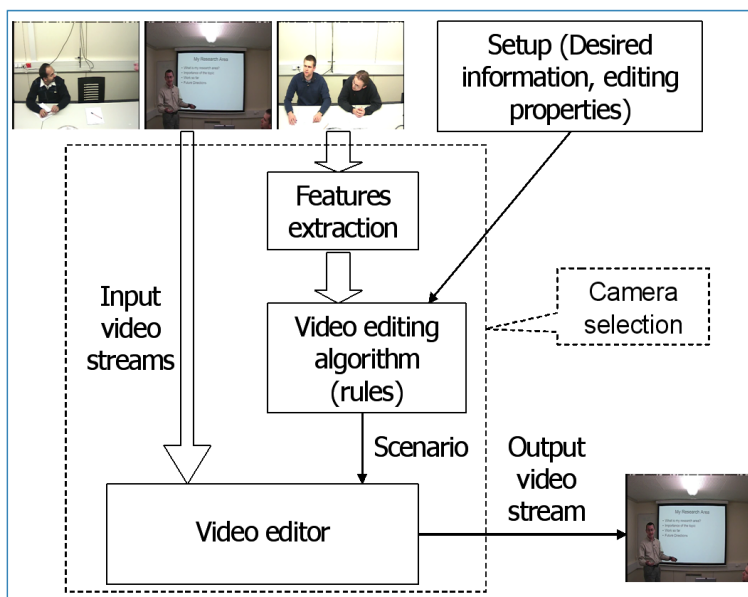


Figure 2: High level scheme of whole editing process

3.2 Editing algorithm

50

An image from the selected camera has to preferably represent what is happening in the meeting room according to different (user specified) aspects. Mainly, technical and aesthetical aspects warrant that produced video contains as much of the relevant information as possible.

51
52
53
54

The main idea is that a methodology of the video editing can be described through a set of various rules. An application of these rules frame by frame to the whole meeting produces a scenario that can be used for generation of the final video. All rules can be divided into two basic classes. The rules of the first class can be used in realtime processing while the rules of the second class process data from the whole meeting and produce more accurate output in offline editing. The goal of the rules application is assignment of weight to every camera in the meeting room. After the weight of all the cameras is known, the camera with the highest weight is selected.

55
56
57
58
59
60
61
62
63

The main editing algorithm works as follows (see figure 3):

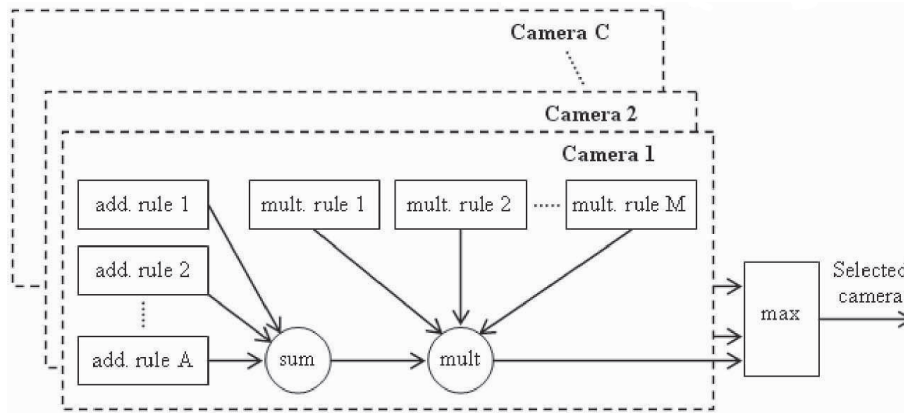
64

1. Input events detected in each camera streams are simultaneously processed frame by frame through the meeting.
2. Rules are used for evaluation of every camera in given time point.

65
66
67

3. The image from the camera with the highest weight is selected and presented as output in the given time. 68
69

Figure 3: Video editing algorithm, evaluation of different rules



Technical aspects of video editing are mainly represented by rules that evaluate e.g. activity of meeting participants or other interesting things such as slides projecting. The resulting weight is computed according to the activity and visibility of given person on certain camera. Other set of rules is important for satisfaction of the aesthetical aspects of video editing.

The first significant aspect of the participants activity is the information about whether the participant is speaking or not. The source data for these rules is obtained from meeting transcription or from automatic speaker identification. More important person is that person who starts speaking the first, also more important is that person who is speaking longer. It is also more important to show particular participant if he/she is gesturing by his head or hands.

Further, the rules can be also used to simulate some aesthetical aspects. Rules for handling periodic alternating of cameras add a little weight to cameras, which were not selected during long time period. This causes changes of cameras if no other activity is detected. Other rules are designed to assure that certain camera is selected at least for minimum given time period and the selection does not exceed maximal time period. These rules avoid quick camera changes that are not acceptable for the viewer and guarantee an interest of the produced video.

See [3] for details of proposed algorithm. There is also explained an approach for designing all kind of rules mentioned above and their detailed description.

3.3 Input Data

The input of the audio-video editing system is constituted by the multimedia audio-visual data themselves with metadata describing their properties and structure. As an input serve also events appearing in the data. These events – that are potentially somehow important for the video editing – are either manually annotated by a tool or are extracted automatically. The input data can be divided into two main parts – offline data and online data.

3.3.1 Offline

Offline data are used during the video editing process, when we already have collected audio and video records. In these records are some important events for the video editing process and it is necessary to have them annotated. In case of manual annotation, all the data is collected before the editing process starts, the automatic event extraction can happen during the video processing, whose output is the edited video sequence. All the manually annotated and automatically extracted events need to be stored for later use and for the editing process.

After all the necessary events are annotated, the offline input data are ready for use by the video editing process. It must be made an XML file describing these events and the used audio and video streams. This file serves as an input of the automatized video editing process.

The main disadvantage of this concept is in the fact, that is needed to have all the input data and detected events stored along, in most cases locally, and with a large number of data it is really space consuming and ineffective. So that concept of using annotated XML files and input data streams is usable especially in the case of processing a small amount of events and multimedia data, one time video editing or if the reusing of outputs is needless.

But in the case that is collected a huge amount of data and extracted events for them or when reusing of extracted events is needful, it is better to use an input data distributed system. To allow this, a component serving as the data storage of the annotations has been defined, that allows a unified way

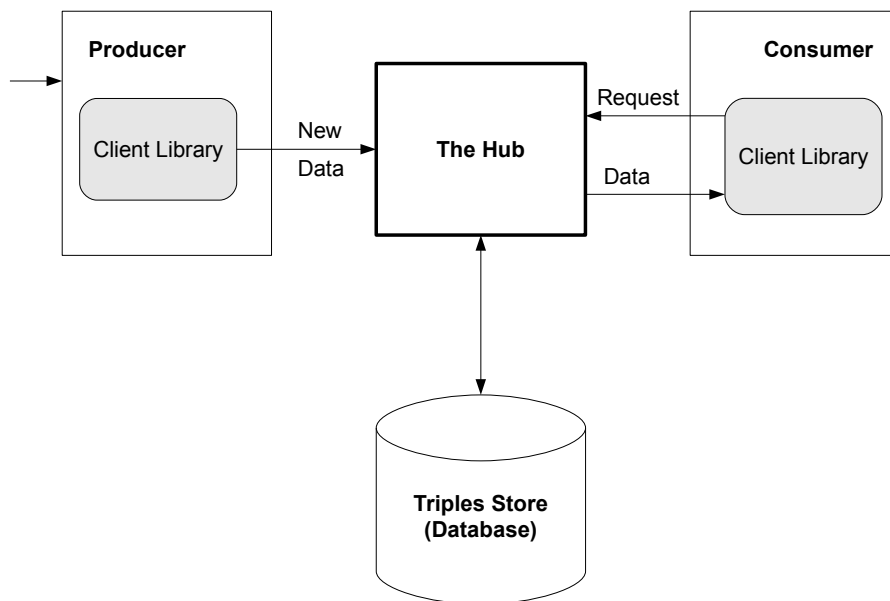


Figure 4: Producer-consumer scheme of Hub.

of storing and retrieving the data. This complete infrastructure is referred 123
to as the Hub. 124

The Hub is intended to provide all of the storage that a group or a 125
company needs for annotations about their archived meetings in one place. 126
Each such Hub contains a database, which keeps the stored annotations and 127
serves requests to the data. Along with the annotation/extracted data it 128
contains all related documents, presentation slides and notes shown during 129
the meeting. At the same moment when one event-extracting process inserts 130
data to the Hub, the consuming application can retrieve it. This allows the 131
Hub to be used both in real-time on a pending meeting to provide a newly 132
connected client with all available information about the meeting as well as 133
to post-processing the meeting data and production of summaries and similar 134
material. 135

Records in the Hub have to be stored in hard defined format. In the video 136
editing software it is necessary to have all needed data already stored in the 137

Hub before the editing process starts because of the editing process makes 138
only one shot data reading from the Hub and then works with result data. 139

3.3.2 Online 140

Sometimes the audio and video records are not reachable at the moment 141
and only data that can be processed by the video editing process is an online 142
stream of a few cameras. For these cases it was developed an possibility to use 143
an online (real-time) input data. The audio and video streams are acquired 144
by number of cameras and microphones, in the next step they are sent to the 145
editing process. During the editing the detectors are applied on the input 146
multimedia streams to check and obtain the important events, which are sent 147
into an editing process and according them the editing algorithm evaluates 148
the best possible shot in every moment. 149

This online option is still in the form of some technological demo and for 150
the real deployment is necessary to create the appropriate audio and video 151
detectors searching for important events in the audio and video streams. 152

3.4 Output 153

After the editing process is finished, the results can be divided into two types, 154
according to used mode of editing. At first, it could be the resulting set of 155
events – marking the shot boundaries – which could be saved into the file 156
or they could be saved to the Hub for future use. This type of output is 157
necessary for the re-editing process, which results into edited video movie. 158

Edited video is the other type of an editing process output. If an editing 159
process is launched in cutting mode, audio and video streams are available 160
and a set of camera cutting events is already created, the editing process 161
follows these camera events by choosing the proper camera and progressively 162
creates resulting movie. 163

4 Limitations 164

Crucial limitation of presented video editing tool is lack of wider variety of 165
real-time feature extractors, that could process the input video streams and 166
detect more distinct events. However the design of the tool is prepared for 167
implementing new extractors which makes it easily extendable. 168

5 Technical specifications 169

For the realtime processing computer with at least Quad-Core processor and 4GB of RAM is needed. The online setup is ready to use together with Unibrain firewire cameras. 170
171
172

- Processor: Intel Core 2 Quad or better 173
- Memory requirements: 4GB of RAM 174
- Disk space: 1.5GB (including demonstration data and configurations) 175
- Operating System: Windows 176
- Camera type: Unibrain Fire-i camera(s) 177
- Camera interface requirements: IEEE-1394a (FireWire) port 178

6 Package content 179

- Papers containing more detailed information 180
- Video editing software (both binary and source files) with additional tools 181
182
- Software using instructions 183
- Prepared demonstration configuration and data 184
- Demonstration videos 185

References 186

- [1] GraphAtFit Youtube channel. Automatic video editing demo, 2009. 187
- [2] Adam Herout, Radek Kubicek, Pavel Zak, and Pavel Zemcik. Automatic video editing for multimodal meetings. *Proceedings of International Conference on Computer Vision and Graphics*, 0, 2008. 188
189
190
- [3] Stanislav Sumec. Multi camera automatic video editing. In *Proceedings of ICCVG 2004*, pages 935–945. Kluwer Verlag, 2004. 191
192