

An Authoring System for Interactive Mobile TV Services Using MPEG-7 and MPEG-4 LAsER

Albert Hofmann, Andreas Kriechbaum, Werner Bailer
Institute of Information Systems
JOANNEUM RESEARCH Forschungsgesellschaft mbH
Steyrergasse 17, 8010 Graz
Austria
{firstName.lastName}@joanneum.at



Journal of Digital
Information Management

ABSTRACT: Consuming mobile TV services is becoming increasingly popular. As hand held devices are personal and used in a “lean-forward” mode, they are well suited for interactive TV. The porTiVity system provides authoring tools for mobile interactive content that allow attaching interactivity to moving objects in a TV program. We describe automatic and semi-automatic annotation tools that create content descriptions represented as MPEG-7 documents. The descriptions of moving objects are then transformed into MPEG-4 LAsER, which is a suitable lightweight format for visualization and interactivity on the mobile device. Since mobile devices are well suited for consuming personalized services we describe different personalization aspects. We show the use of the authoring components of porTiVity in two example scenarios.

Categories and Subject Descriptors

H.5.1[MultimediaInformationSystems]:H.4.3[Communications Applications]:

General Terms: Authoring, Interactive TV, Mobile services, MPEG-7, MPEG-4 LAsER.

Keywords: Mobile Applications, Multimedia applications, User Evaluation

Received: 19 July 2009; Revised 18 August 2009; Accepted 28 August 2009

1. Introduction

1.1 Motivation

With the introduction of the mobile broadcast standard DVB-H [1] mobile TV is becoming increasingly popular. The number of mobile TV users is increasing at an enormous rate and consumers are showing their interest to enjoy innovative mobile applications. Market analysis predicts growth rates of 100% and more for the number of users in the coming years¹ and within the last year the average daily time a user spent watching mobile TV has nearly doubled². Mobile interactive TV can be one of these emerging applications. Results from previous interactive TV projects such as GMF4iTV [2] show that hand held devices are better suited for interactivity than classic TV sets as they are personal and used in “lean-forward” mode.

This paper presents results from the porTiVity project [3] which aims to provide interaction with objects in video shown on mobile devices. The interactive experience for the end-user is achieved by simply clicking on highlighted objects overlaid on the video on the touch screen of the mobile device. This interaction triggers e.g. downloading of additional content from the web, participating in a quiz, in mobile network based games or purchasing a chosen product. In a live scenario such as a football game interaction on the soccer players or via menu buttons allows to retrieve content from a player’s website, game statistics, watching replays of interesting events etc.

The efficient semantic annotation of audiovisual content is an important part of this workflow. MPEG-4 LAsER annotation tools supporting moving objects and interaction with those objects were not available when the project started. Relevant annotation tools for audiovisual content are among others M-OntoMat-Annotizer [4], developed by the aceMedia project and the IBM VideoAnnEx Tool [5]. However, these tools lack functionalities for efficient annotation support and are not applicable in live use cases.

We describe in this paper an authoring system for creating interactive TV content with interactivity based on moving objects. Section 2 describes the automatic and semi-automatic annotation of content described using MPEG-7 [6]. In Section 3 we discuss how this description is transformed into a suitable representation for the mobile device, namely MPEG-4 LAsER [7]. The integration of these authoring steps into the porTiVity system and the end-user experience in two scenarios is described in Section 4. Section 5 concludes this paper.

1.2 Authoring for mobile interactive TV

The architecture of the end-to-end platform for providing rich media interactive TV services for portable and mobile devices developed in the porTiVity project is shown in Figure 1 and described in Section 4 in more detail. In the authoring phase the main video content is analyzed and annotated to link it to additional content (images, web pages, etc.) and the possible user interactions are defined. The tools support workflows for both offline and online scenarios. The result of the authoring phase is a scene description in MPEG-4 LAsER which is multiplexed into an MXF [8] stream together with the main and additional content. MPEG-4 LAsER is a standard for rich media scene description on mobile and embedded devices. In our work it is used to link additional content to the main video stream which can be simple text, HTML pages, pictures, audio and video clips. In order to reduce the bandwidth and processing

¹http://www.3gnewsroom.com/3g_news/apr_05/news_5811.shtml

²<http://futurezone.orf.at/stories/1500388/>

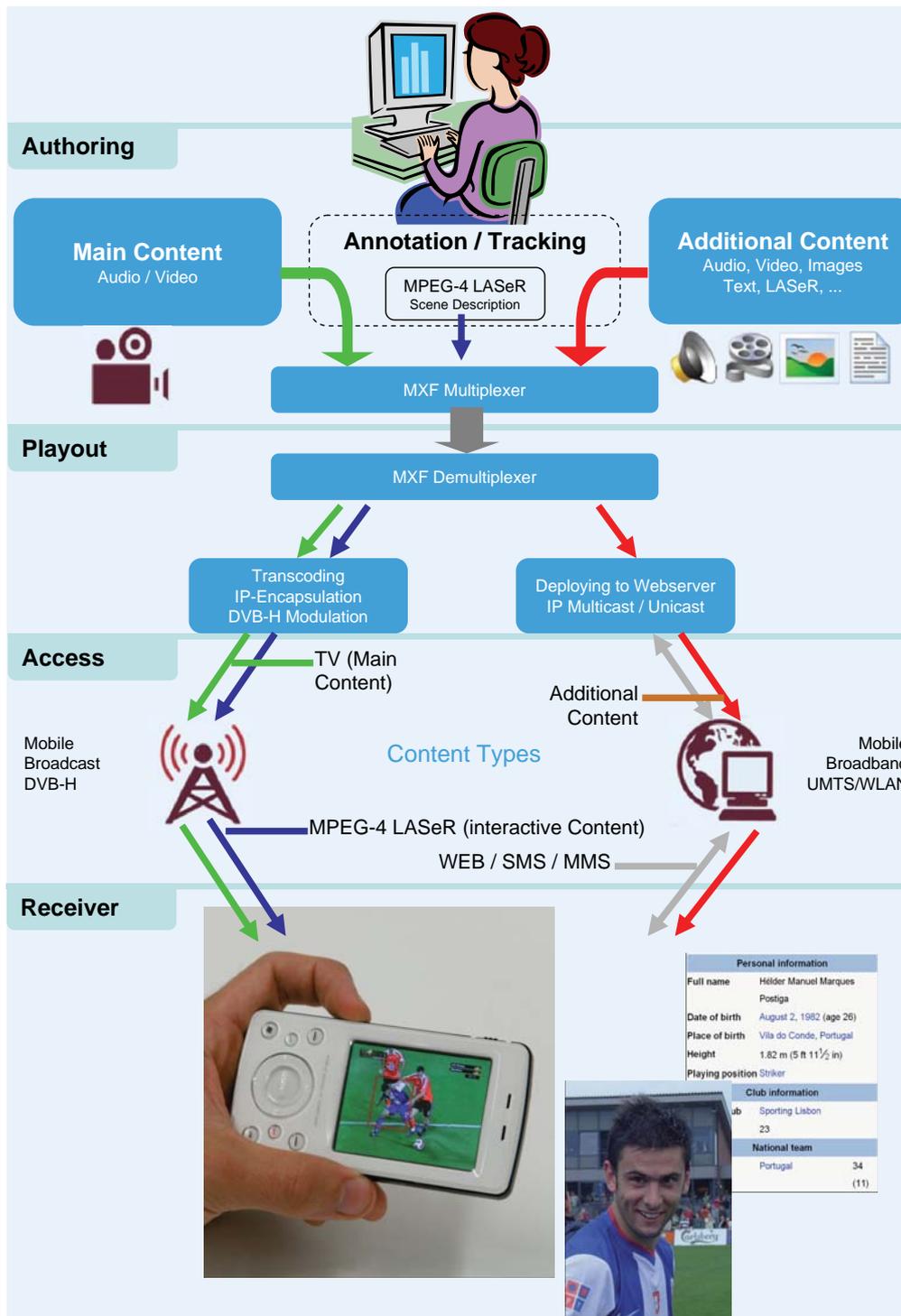


Figure 1. porTiVity system overview

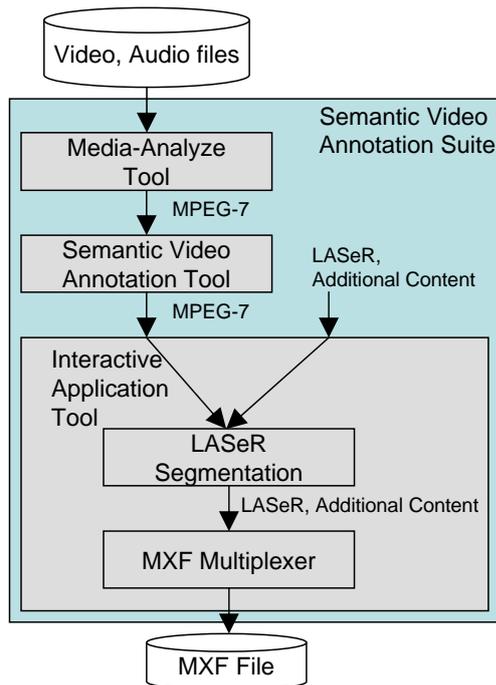
requirements it has been used as a more lightweight alternative to MHP [9]. MXF is a container format defined by SMPTE that can hold various types of audiovisual essence as well as metadata, with capabilities for multiplexing.

In Figure 2 the workflow between the tools of the proposed Semantic Video Annotation Suite in a typical authoring chain and the metadata formats used at the interfaces are visualized. Media-Analyze is an automatic video preprocessing tool (cf. Section 2.1). In the Semantic Video Annotation Tool (SVAT) the results from the automatic preprocessing are used to support the user in annotating the video by providing navigation aids, redetection and tracking of objects. The Interactive Application Tool (IAT) uses the MPEG-7 description of the annotated objects generated by the SVAT,

transforms them into MPEG-4 LAsER objects and adds interaction support on them. As a post-processing step it also performs the LAsER segmentation and packs all the required data into an MXF file suitable for the playout system.

2. Content Annotation

The content annotation is divided into two steps: an automatic preprocessing step done with the Media-Analyze tool and a semi-automatic annotation step performed with the Semantic Video Annotation Tool. Both tools are bundled together in the Semantic Video Annotation Suite [10] and are described here in the context of authoring interactive mobile TV services. The main task of the automatic preprocessing is to extract metadata



(Video, Audio, LASeR, Additional Content: Images, Videos, Audio, LASeR, ...)

Figure 2. Annotation chain in the Semantic Video Annotation Suite including the exchanged metadata formats

that aids the user later in the semiautomatic annotation step to perform interactive tasks more efficiently. The output of both steps is described using MPEG-7, a comprehensive standard for multimedia content description covering both low-level and high-level description. In particular, the Detailed Audiovisual Profile (DAMP) [11], an MPEG-7 subset aimed at fine-grained description of audiovisual content in media production and archiving, is used.

2.1. Automatic preprocessing

In a first step a video file is automatically analyzed by the Media-Analyze tool. We extract metadata for efficient navigation in videos and structuring the content. These are shot boundaries, key-frames and stripe images.

In order to aid the user in efficient semi-automatic object annotation we provide an object redetection functionality which is based on SIFT descriptors [12] extracted around interest points. The steps for building the SIFT descriptors in the automatic preprocessing as well as during the matching in the semi-automatic object annotation are described in the following.

a) *Build scale-space*: The input image is represented at different resolutions in a Gaussian pyramid by convolution with Gaussian smoothing filters and by sub-sampling. A second pyramid of edge-detection images, called Difference of Gaussian (DoG) pyramid, is produced by subtraction of adjacent images from the Gaussian pyramid. To build this scale-space we have used a full-octave approach, as proposed in [13].

b) *Extract key-points*: After the scale-space of an input image is built, DoG points are extracted from the DoG pyramid. Therefore, local intensity minima and maxima are detected by comparing each of the image points with their neighbor in scale-space. Each detected key-point is initially composed of a position (x and y coordinates) and a scale.

c) *Assign orientation*: Next, the dominant orientation of each key-point is computed with an orientation histogram that is generated from a small image region around the key-point. If several directions have similar likelihood, multiple orientations can be assigned to a key-point. This is only useful for keypoints from model images.

d) *Generate features*: The last step of this process is the generation of SIFT descriptors for each key-point and each orientation. For each orientation of a key point a separate set of SIFT descriptors is generated, which is composed of low-dimensional and high-dimensional features. Different features are used to achieve efficient feature matching and storage of the model database on disk.

To enable a short search time for detecting similar objects and shots we use an approach similar to Video Google proposed in [14]. The SIFT descriptors are extracted for each key frame in the video. In order to reduce noise only robust descriptors over time are used. Any feature point that does not exist for at least six frames is rejected as unstable point. The estimate of the descriptor is computed by averaging the descriptors of the tracked points.

The next step is to setup a visual vocabulary. First, we cluster the extracted feature vectors. Each of these clusters represents a visual word. All extracted descriptors are assigned (matched) to these visual words. The visual words and the matched occurrences (image position, scale, and an orientation value) are stored in a k-means cluster tree as proposed in [15]. This avoids the problems of large list-like databases for storing the visual vocabulary.

The tree query is bound to a certain tree level, i.e. all queried descriptors are propagated in the tree down to a maximum level. If the descriptor cannot be propagated to this level, the query is discarded. The tree nodes are enumerated relatively to the tree level, thus yielding to the resulting index of the query. If a tree level is not fully populated, there will be gaps in the node enumeration. Figure 3 illustrates a query in the tree.

When the user identifies single object occurrences, these occurrences need to be tracked within the same shot. Since this step is interactive it requires a very efficient tracking algorithm. For this purpose feature points and descriptors are extracted

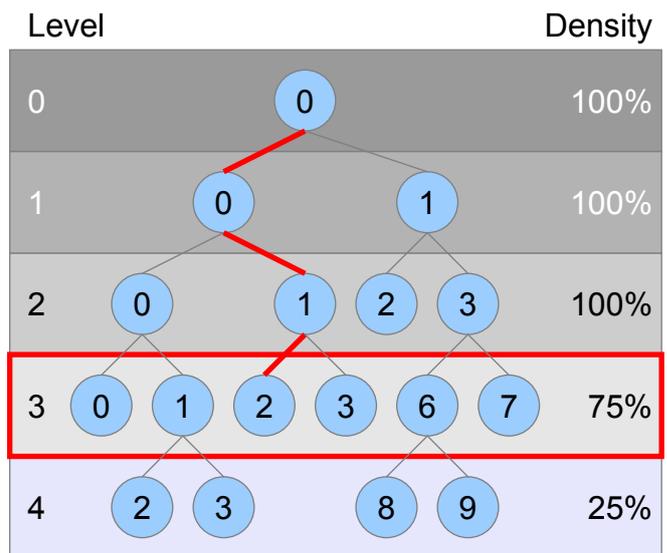


Figure 3. Query on a hierarchical k-means cluster tree, with k = 2. The thick red path is a successful query up to the maximum level 3

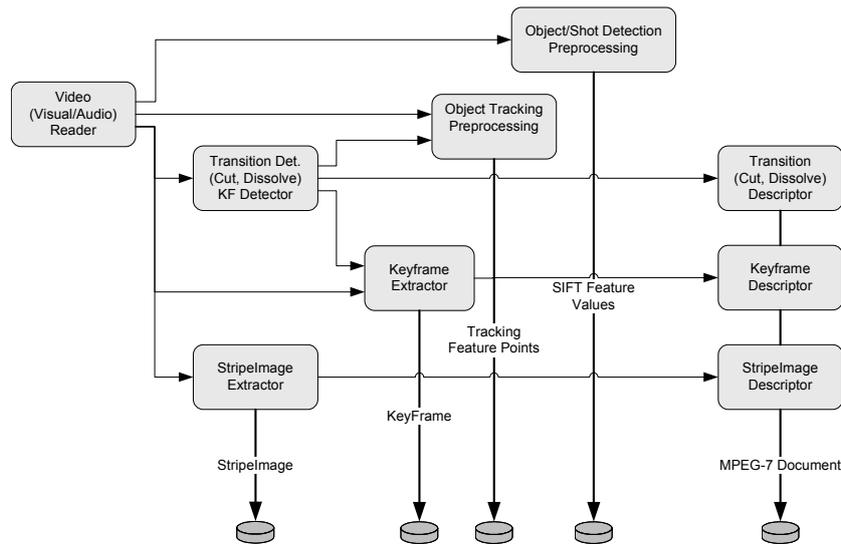


Figure 4. Analysis Module Graph used in Automatic Preprocessing

within this automatic preprocessing step [16, 17]. The descriptors for object tracking are stored in a binary format for efficiency reasons.

The various modules for extracting metadata are plugged into an analysis graph (Figure 4) that is executed in the Media-Analyze tool, considering the interdependencies between the modules.

2.2 Semi-automatic object annotation

After this automatic annotation step objects of interest can be annotated with interaction. This is done with the Semantic Video Annotation Tool and its support for spatio-temporal object annotation.

As operators often wish to assign the same additional content to similar objects we provide an object redetection feature. The

input for the object redetection is a region of interest identified in a video player component. This can be done by simply marking a bounding box around the region, or by using a color segmentation functionality in order to get the exact border of the object. By using the visual vocabulary from the preprocessing step we can search for appearances of similar SIFT descriptors as extracted from the query region. For this application scenario higher recall is preferred since false positives in the result list can be easily deleted.

In the next step these single occurrences can be tracked over time within the boundaries of the shot. The input for the tracking algorithm can be a bounding box as well as a complex polygon around an object. The resulting moving region is always described as bounding box. The tool also supports manual

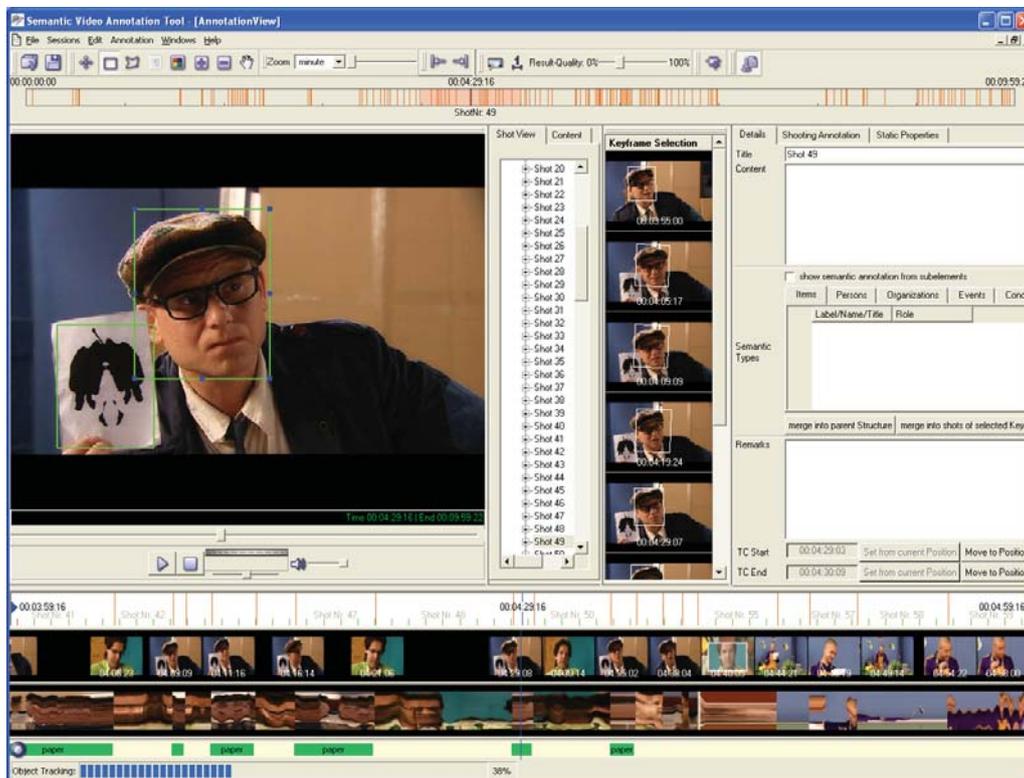


Figure 5. Semantic Video Annotation Tool for object annotation

tracking in case the automatic tracking fails. In this case the operator specifies key-regions over time that will be connected with a spline interpolation algorithm.

The output of the Semantic Video Annotation Tool is an MPEG-7 description which includes the metadata from the automatic analysis step enriched with the annotation of moving objects. An example of such a moving region is shown in Listing 1. The moving region shows an object with name *bottle* (*SemanticRef* element) and a bounding box that changes over time (*ParameterTrajectory* elements as defined by [6, part 3]). The values inside the *Coords* element define the three points that build a closed bounding box.

```
<MovingRegion>
  <SemanticRef idref="bottle" />
  <SpatioTemporalLocator>
    <ParameterTrajectory motionModel="still">
      <MediaTime>
        <MediaTimePoint>T00:01:23:1F25</MediaTimePoint>
      </MediaTime>
      <MediaDuration>P0DT0H0M0S2N25F</MediaDuration>
    </ParameterTrajectory>
    <InitialRegion>
      <Polygon>
        <Coords dim="2 4">421 96 0 -96 300 0 151
      </Coords>
      </Polygon>
    </InitialRegion>
  </ParameterTrajectory>
  <ParameterTrajectory motionModel="still">
    <MediaTime>
      <MediaTimePoint>T00:01:23:3F25</MediaTimePoint>
    </MediaTime>
    <MediaDuration>P0DT0H0M0S1N25F</MediaDuration>
  </ParameterTrajectory>
  <InitialRegion>
    <Polygon>
      <Coords dim="2 4">422 95 0 -95 301 0 151
    </Coords>
    </Polygon>
  </InitialRegion>
  </ParameterTrajectory>
  ...
</SpatioTemporalLocator>
</MovingRegion>
```

Listing 1. Moving object described in MPEG-7.

3. MPEG-7 to MPEG-4 LAsER Transformation

3.1. Transforming moving object descriptions

In this step the representation of moving regions described in MPEG-7 format are transformed to MPEG-4 LAsER, a suitable delivery format for the mobile device. MPEG-7 is intended for content description, but does not provide visualization and interaction possibilities. In addition, a large part of the MPEG-7 description is not needed for the service at the mobile device. We present an approach for transforming shape descriptions of moving objects (like polygons and rectangles) to MPEG-4 LAsER.

MPEG-4 LAsER is based on SVG³ Tiny 1.2 and is thus very powerful regarding shape representation. We have focused

³Scalable Vector Graphics (SVG) is a W3C recommendation for the XML based representation of 2D vector graphics (<http://www.w3.org/Graphics/>)

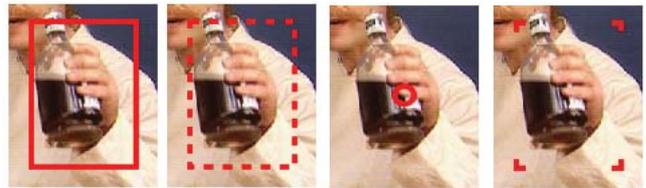


Figure 6. Styles for highlighting moving objects: A rectangle with solid or dashed line style, a contact point in the middle of the area and rectangle with only corner indicators

on transforming polygons and simple bounding boxes to SVG taking different personalization options for the end user in consideration. Currently three options are supported which are also reflected in the MPEG-4 LAsER snippet in Listing 2 and the images in Figure 6: a rectangle with solid or dashed line style, a contact point in the middle of the area and rectangle with only corner indicators.

```
<g pointer-events="all" fill="none" stroke="red" id="
bottle0" stroke-width="2">
  <title>bottle</title>
  <desc>Moving Object</desc>
  <rect width="96" x="421" y="300" height="151">
    <animate values="421;426;432;439;447;..."
      fill="freeze" begin="83.040s" dur="9.680s"
      attributeName="x" />
    <animate values="300;305;302;277;229;192;..."
      fill="freeze" begin="83.040s" dur="9.680s"
      attributeName="y" />
    <animate values="96;94;93;95;99;..."
      fill="freeze" begin="83.040s" dur="9.680s"
      attributeName="width" />
    <animate values="151;149;148;149;153;..."
      fill="freeze" begin="83.040s" dur="9.680s"
      attributeName="height" />
  </rect>
  <circle cx="469" cy="375" display="none" r="10">
    <animate values="375;379;376;351;305;..."
      fill="freeze" begin="83.040s" dur="9.680s"
      attributeName="cy" />
    <animate values="469;473;478;486;496;..."
      fill="freeze" begin="83.040s" dur="9.680s"
      attributeName="cx" />
  </circle>
  <g display="none">
    <desc>corner rectangle</desc>
    <polyline points="421,310 421,300 431,300">
      <animate values="421,310 421,300
        431,300;426,315 426,305 436,305;..."
        fill="freeze" begin="83.040s" dur="9.680s"
        attributeName="points" />
    </polyline>
    <polyline points="507,300 517,300 517,310">
      <animate values="507,300 517,300
        517,310;510,305 520,305 520,315;..."
        fill="freeze" begin="83.040s" dur="9.680s"
        attributeName="points" />
    </polyline>
    <polyline points="517,441 517,451 507,451">
      <animate values="517,441 517,451
        507,451;520,444 520,454 510,454;..."
        fill="freeze" begin="83.040s" dur="9.680s"
        attributeName="points" />
    </polyline>
    <polyline points="431,451 421,451 421,441">
```

```

    <animate values="431,451 421,451
421,441;436,454 426,454 426,444;..."
    fill="freeze" begin="83.040s" dur="9.680s"
attributeName="points" />
</polyline>
</g>
<set begin="0s" attributeType="XML" to="none"
attributeName="display" />
<set begin="83.040s" attributeType="XML"
to="block" attributeName="display" />
<set begin="92.720s" attributeType="XML"
to="none" attributeName="display" />
</g>

```

Listing 2. Moving object described in MPEG-4 LAsER

For the SVG description of the animation of the moving object over time we considered two options: specifying the coordinates in a fixed interval for the whole time range, or using a variable interval depending on the movement of the object and specifying those time points in the *keyTimes* attribute. Using a variable interval has advantages in terms of description size for moving objects that hardly move or keep their position for a longer period, but increases the complexity of the synchronized broadcast (see LAsER Segmentation in Section 3.2). So we have chosen the approach using a fixed interval (e.g. 0.2 sec) that can be configured by the user.

The result of the conversion from the moving region described in the MPEG-7 snippet in Listing 1 to MPEG-4 LAsER is shown in Listing 2. The id from the MPEG-7 description is reflected in the *title* element (unchanged) and also in the *id* attribute, taking into account that XML requires unique ids. Also when accessing specific elements during MPEG-4 LAsER command processing we require unique ids when adding, modifying or removing objects.

The graphic object is embedded inside the SVG document of the main MPEG-4 LAsER scene. In order to reduce the complexity of MPEG-4 LAsER handling within the authoring tools, no LAsER commands are used in this authoring step for time related updates of the scene. This is done in a separate post processing step – the LAsER segmentation – which is described in Section 3.2. The animation of objects over time is specified for the whole object time range. Also the visibility of the object is modeled by setting the SVG *display* attribute to 'none' and 'block' for the respective time points. So up to this step the produced MPEG-4 LAsER document just contains one scene with one SVG document containing the whole description for the whole video time range. Also no interaction possibilities are added up to this point.

3.2. Adding interactivity

Next to the powerful shape representation capabilities of MPEG-4 LAsER, it provides support for interactivity to the end user. In our scenario we link additional content to moving objects like statistics and photos of soccer players, background information on actors etc. If the user clicks on a moving region this additional content should be displayed. Moreover, the additional content can be personalized, which means that different additional content items can be linked for different target groups (e.g. fans of the different teams in sports broadcasts, male and female viewers, ...). We can edit these actions in the Interactive Application Tool (IAT, see Figure 7). The resulting MPEG-4 LAsER snippet after adding an additional content is shown in Listing 3. For non-personalized additional content we place an *a* element around the shapes for highlighting the moving object. Accessing personalized additional content is described in Section 4.3.

```

<g pointer-events="all" fill="none" stroke="red"
id="bottle0" stroke-width="2">
  <title>bottle</title>
  <desc>Moving Object</desc>
  <a xlink:href="http://server.portivity.org/
PoEl_a.xsr">
    <rect width="96" x="421" y="300" height="151">
      <!-- animation of rect -->
    </rect>
    <circle cx="469" cy="375" display="none"
r="10">
      <!-- animation of circle -->
    </circle>
    <g display="none">
      <desc>corner rectangle</desc>
      <!-- polylines of corner rectangle -->
    </g>
  </a>
</g>

```

Listing 3. Interactivity described in MPEG-4 LAsER

Before sending the MPEG-4 LAsER description to the playout system, a post processing step is needed to make the service ready for transmission over a broadcast channel. Users can join the TV channel at any time, so it is necessary to transmit the LAsER content periodically. This is required for the general scene description containing the static objects (like menus), while moving objects just need to be transmitted during the time they are visible. This is done by the LAsER segmentation module. It takes the whole LAsER scene from the Interactive Application Tool as input and splits the objects into different access units. An example of how the MPEG-4 LAsER snippet presented in Listing 2 is segmented is shown in Listing 4.

```

<!-- Initial object description
in the main SVG document: -->
<g pointer-events="all" fill="none" stroke="red"
id="bottle0" stroke-width="2">
  <!-- shape representation cut here -->
  <set begin="0s" attributeType="XML"
to="none" attributeName="display" />
  <set begin="83.040s" attributeType="XML"
to="block" attributeName="display" />
  <set begin="92.720s" attributeType="XML"
to="none" attributeName="display" />
</g>
<!-- After the segmentation into access units
for insert and delete: -->
<saf:sceneUnit time="83040">
  <lsr:Insert ref="<parent_id>">
    <g id="bottle0">
      <!-- shape representation cut here -->
    </g>
  </lsr:Insert>
</saf:sceneUnit>
<saf:sceneUnit time="92720">
  <lsr>Delete ref="bottle0"/>
</saf:sceneUnit>

```

Listing 4. LAsER Segmentation into Access Units

Another task of this module is the periodic generation of *RefreshScenes*, which are access units containing a refresh of the scene as it should be at that precise moment. These access units are generated at a given interval – usually every two



Figure 7. Interactive Application Tool for creating MPEG-4 LAsER scenes

seconds – and they are ignored by the terminal if it already has this information (either from a *NewScene* or from a previous *RefreshScene*). These access units will contain the information on static and dynamic objects that become visible or invisible, and will also update the position of visible dynamic objects. The additional binary MPEG-4 LAsER stream only costs a small amount of bandwidth compared to the TV stream. An alternative system design to periodically poll a web server would lead to a traffic bottleneck due to the huge amount of requests to be handled by the web server.

The final step of the authoring is the MXF multiplexing. The video and audio content as well as all additional content files are packed into an MXF stream to ensure time synchronization.

4. porTiVity System

A detailed description of the porTiVity system is given in [18]. This section summarizes the steps after the authoring as shown in Figure 1.

4.1. Playout system

The playout system consists of several modules making the video stream and MPEG-4 LAsER content ready for transmission over DVB-H and the additional content accessible on a web server.

The MXF demultiplexer takes the produced MXF file from the authoring system as input, extracts all tracks and forwards them to the respective modules in the chain. The transcoder module converts the high-quality video and audio streams used in the production environment to the target formats on the mobile device: MPEG-4 AVC (H.264) and MPEG-4 AAC. Additional content files are also transcoded to the target format if necessary. The Metadata Processor Module receives LAsER access units from the MXF demultiplexer, encapsulates them over RTP

(Real-Time Transport Protocol), and sends them to the DVB-H module responsible for broadcasting.

4.2. Receiver

The mobile device receives the DVB-H stream consisting of an audio, video and LAsER MPEG-4 elementary stream encapsulated in RTP. These streams are synchronized on the terminal by RTCP.

The software responsible for rendering the porTiVity service is based on the Osmo player of the open source multimedia framework GPAC [19]. An implementation with the full support of the porTiVity features is available for Windows Mobile version 5 and 6, Windows XP and an experimental version for Symbian OS. The developed scenarios work on normal off-the-shelf mobile phones with DVB-H receivers and touch screens. Furthermore any kind of web access is required for downloading the additional content files requested by the user clicking on the annotated objects. For various demonstrations we used terminals with limited processing power like a Gigabyte GSmart T600, a QTec 9000 and an LG KS20. As high-end terminals an HTC Shift and a Samsung Q1 Ultra were used.

4.3 End-user experience

1) *iSports – interactive sports scenario*: The iSports scenario (see Figure 8) demonstrates how interactivity can be applied to live sport events. For this case a separate annotation tool – the Live Annotation Tool [20] – has been developed in the porTiVity project. With this tool it is possible to define objects of interest and link additional content, track the objects in realtime and directly send the results to the playout system – with an overall delay of 5 seconds to the input stream.

Next to the interaction possibility with the highlighted objects, a menu with buttons is shown in the top left corner. The menu,



Figure 8. Interactive sports scenario

which can be collapsed in order not to consume valuable display space, allows to access replays of interesting events like goals, to access a help page that explains the available features of the service and to configure personalization settings.

There is increasing demand for accessing information relevant and adapted to the user's needs. Mobile devices are typically personal items and thus well suited for consuming personalized services. The iSports service provides personalization support for normal TV programs based on user profile settings. The profile can be configured while watching the program (see Figure 9) and is stored in a local cookie file.

These user preferences have impact on three levels: the interaction options presented to the user (e.g. only provide interactivity for one team), the additional content assigned to objects (buying options for articles based on gender) and the behaviour of the service itself (e.g. enable time shift viewing). The first two issues need to be considered by the program creator during the authoring process. The profile categories that have impact on the interaction options are organized in corresponding group elements in the main SVG element of the main MPEG-4 LAsER scene. Scripting functionality on the terminal side sets the visibility for the respective groups based on the user profile. For customized additional content the operator assigns different content to the different target groups. Listing 5 shows how to represent and access personalized additional content using MPEG-4 LAsER. Instead of using an `a` element for non-personalized additional content, we install a listener on the object. A script on the terminal (`loadAdditionalContent`) evaluates the list of additional content links and activates best matching for the user profile.

2) *Spur & Partner* – interactive crime series for children: The *Spur & Partner* scenario (see Figure 10) is based on a successful interactive TV crime series for children produced by the public German broadcaster Rundfunk Berlin Brandenburg (RBB). Users can interactively collect hints by clicking on the objects annotated by the creator. They can also answer questions referring to those pieces of evidence and receive points for correct answers, which gives the TV program a little game characteristic.

The end-users can access background information on actors and learn from an interactive detective magazine more about the approach that the detective is using to identify the culprit. And



Figure 9. Personalization options for the interactive sports scenario

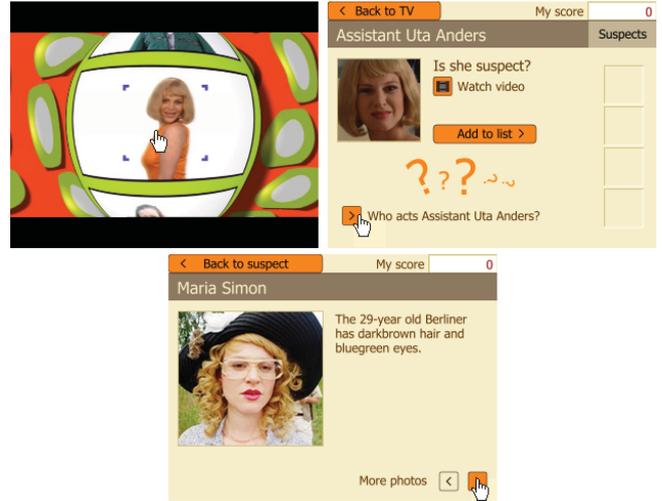


Figure 10. Interactive detective story Spur & Partner

finally at the end of the series the users have the chance to name the culprit by themselves and directly transmit their results.

```
function loadAdditionalContent(urls) {
  var argSize = urls.length;
  var profile = new Profile("personalization");
  /* get url according to profile */
  var url = profile.extractURL(urls);
  document.getElementById("saveURL")
    .setAttribute("xlink:href", url);
  document.save("lastPage", "saveURL",
    "xlink:href");
  xhr_object = new XMLHttpRequest();
  xhr_object.onreadystatechange = ready_
  AdditionalContent;
  xhr_object.open("GET", url, true);
  xhr_object.send("");
}
<g id="player_FCP_16_0">
  <!-- shape representation cut here -->
  <ev:listener event="click" handler="#hdl_pl_
  fcp_16_1"/>
</g>
<handler id="hdl_pl_fcp_16_1" type="application/
  ecmascript">
  <![CDATA[loadAdditionalContent(
    ['http://server/player_FCP_16_1_male.
  svg',
    'http://server/player_FCP_16_1_female.
  svg']]);]]>
</handler>
```

Listing 5. Personalized additional content described in MPEG-4 LAsER

Spur & Partner was evaluated by RBB with 6 children at the age of 9 to 11 years. They had normal computer skills and were familiar with the original TV series and the interactivity by its MHP application. Evaluation results have shown that recognizing the bounding box of the moving objects and further interaction with the additional content improve with the experience and are no problem after a little trying.

4.4 Evaluation of the authoring tool

The authoring system was evaluated by RBB with professional users from various departments. An evaluation of the tools has

shown that the time required to annotate objects is acceptable in a TV production environment. The system was tested on the interactive detective story with a video length of 8:40 minutes and consisting of 16 interactive elements. Not considering the automatic preprocessing step, the whole annotation process took 1:15 hours with the additional content already prepared. Adding shortcuts that were missing in the evaluation version for common task could further increase the annotation performance.

5. Conclusion

This work has shown that MPEG-7 is suitable as description format for video including the description of moving objects. We presented tools that aid an operator in efficient object annotation by using automatically extracted metadata in MPEG-7 format and enrich it with annotations for moving objects. On the delivery side, MPEG-4 LAsER has proved to be a powerful standard for enabling rich multimedia interactive services on mobile devices with limited processing power. Due to the fact that MPEG-4 LAsER is based on SVG it allows scalability for a wide range of target devices. Also personalized services are possible and have been demonstrated on different personalization levels. The chosen approach to use existing annotation tools based on MPEG-7 and adding MPEG-4 LAsER functionality combined the advantages of both standards and resulted in a full featured MPEG-4 LAsER authoring suite.

6. Acknowledgements

The authors would like to thank their colleagues at the Institute of Information Systems at JOANNEUM RESEARCH, especially Helmut Neuschmied, Severin Stampler, Georg Thallinger and Werner Haas, as well as the partners of the porTiVity project for the good collaboration.

The research leading to this paper has been partially supported by the European Commission under the contracts FP6-027098, "porTiVity – Portable Interactivity" (<http://www.portivity.org>), FP6-045032, "SEMEDIA - Search Environments for Media" (<http://www.semedia.org>) and FP7-231161, "PrestoPRIME" (<http://www.prestoprime.eu>).

References

- [1] Transmission System for Handheld Terminals (DVB-H), ETSI EN 302 304, Jun. 2004.
- [2] Cardoso, B. (2005). Hyperlinked Video with Moving Objects in Digital Television," in *International Conference on Multimedia and Expo*, 2005, pp. 486–489.
- [3] porTiVity web page, <http://www.portivity.org>
- [4] M-OntoMat-Annotizer web page, <http://www.acemedia.org/aceMedia/results/software/m-ontomat-annotizer.html>
- [5] IBM VideoAnnEx Tool web page, <http://www.research.ibm.com/VideoAnnEx>
- [6] Information Technology—Multimedia Content Description Interface. (2001). ISO/IEC standard 15938.
- [7] MPEG-4 LAsER, ISO/IEC 14496-20, 2006.
- [8] Material Exchange Format (MXF) – File Format Specification (Standard) (2004). SMPTE 377M.
- [9] MHP 1.2, ETSI 102 590, 2007.
- [10] Semantic Video Annotation Suite, <http://www.joanneum.at/en/fb2/iis/products-solutions-services/semantic-video-annotation.html>
- [11] Bailer, W., Schallauer, P. (2006). The Detailed Audiovisual Profile: Enabling Interoperability between MPEG-7 Based Systems, In: *Proceedings of 12th International Multi-Media Modeling Conference*, H. Feng, S. Yang, and Y. Zhuang, Eds., Beijing, CN, Jan. p. 217–224.
- [12] Lowe, D. (1999). Object Recognition from Local Scale-Invariant Features, In: *Proceedings of the International Conference on Computer Vision*, 1999, pp. 1150–1157.
- [13] Lowe, D. (2004). Distinctive image features from scale-invariant keypoints, *International Journal of Computer Vision*, 91–110.
- [14] Sivic, J., Zisserman, A. (2003). Video Google: A text retrieval approach to object matching in videos, In: *Proceedings of the International Conference on Computer Vision*, v. 2, Oct. 2003, p. 1470–1477.
- [15] Nistér, D., Stewnius, H. (2006). Scalable Recognition with a Vocabulary Tree, In: *Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. Washington, DC, USA: IEEE Computer Society, p. 2161–2168.
- [16] Trichet, R., Mérialdo, B (2007). Generic object tracking for fast video annotation, In: *VISAPP 2007, 2nd International Conference on Computer Vision Theory and Applications*, 8 - 11 March, 2007 Barcelona, Spain.
- [17] Neuschmied, H., Trichet, R., Mérialdo, B. (2007). Fast annotation of video objects for interactive TV In: *MM 2007, 15th international ACM conference on multimedia*, September 24-29, 2007, Augsburg, Germany.
- [18] Deigmöller, J., Fernandez, G., Lopez, A., Mérialdo, B., Neuschmied, H., Pinyol, F., Trichet, R., Wolf, P., Salgado, R Milagaia, F., Glaser, S., Duffy, A (2008). porTiVity : Object Based Interactive Mobile TV System Networked and Electronic Media, in *NEM SUMMIT 2008, Networked and Electronic Media*, October 13-15, 2008, Saint-Malo, France.
- [19] Feuvre, J. L., Concolato, C., Moissinac, J.-C. (2007). GPAC: open source multimedia framework, In: *MULTIMEDIA '07: Proceedings of the 15th international conference on Multimedia*. New York, NY, USA: ACM. p. 1009–1012.
- [20] Deigmöller, J., Fernandez, G., Kriechbaum, A., Lopez, A., Mérialdo, B. Neuschmied, H., Margalef, F. P., Trichet, R., Wolf, P., Salgado, R., Milagaia, F. (2009). Active Objects in Interactive Mobile Television, In: *15th International Conference on MultiMedia Modeling*, Sophia-Antipolis, FR.

Authors bibliography



Albert Hofmann studied Software Engineering at the Hagenberg University of Applied Sciences (Upper Austria) and finished in 2004 with the diploma thesis "Media Management as a Component for IMDAS-Pro". This work was performed at the Institute of Information Systems at JOANNEUM RESEARCH, where he is

working since 2003. He is involved in a number of national and European research projects in the area of audiovisual archiving, digital film restoration, quality analysis and visual surveillance. He has expertise in content based analysis and retrieval of audiovisual material, metadata modeling and software architectures.



Andreas Kriechbaum finished his study of Telematics at the University of Technology Graz in July 2007 with the master thesis "Moving Object Segmentation in Video and Film". This work was performed at the Institute of Information Systems at JOANNEUM RESEARCH, where he works since

2001. He is involved in a number of national and European research projects in the area of interactive TV and surveillance. His areas of interest and experience are content based analysis and retrieval of audiovisual information, and their application in the domains of audiovisual archives, video annotation and surveillance.



Werner Bailer received a degree in Media Technology and Design from the Hagenberg University of Applied Sciences (Upper Austria) in 2002 for his diploma thesis on "Motion Estimation and Segmentation for Film/Video Standards Conversion and Restoration". He is working at the Institute of Information Systems since 2001 and

has been involved in a number of national and European research projects in the area of audiovisual archiving, digital cinema production, digital film restoration and quality analysis and interactive TV. He has experience in development of image and video processing algorithms, metadata modelling and software architectures. Since 2007 he is working on a PhD thesis on the topic of multimedia content abstraction.