

iLex: Handling Multi-Camera Recordings

Thomas Hanke, Jakob Storz, Sven Wagner

Institute of German Sign Language and Communication of the Deaf, University of Hamburg
 {thomas.hanke,jakob.storz,sven.wagner}@sign-lang.uni-hamburg.de

Abstract

Until recently, sign language researchers were quite happy with just one or two views for each recording session. New corpus projects, however, offer the transcriber five or more camera views. This requires much more flexibility in the transcription environment for switching between different views in order to save network bandwidth, local CPU usage, and screen real estate. Here we present a user interface study within the iLex transcription environment that allows flexible switching between video layouts whenever the transcription focus changes. Switching (including zooming) may be initiated by the user at any point of time, or can be automated to depend on tagging such as tasks or turns. The user interface is backed up by a server infrastructure providing the videos in different spatial resolutions as needed for optimal display.

1. Introduction

More than 15 years ago, we introduced the first sign language transcription environment working with digital video (syncWRITER, cf. Hanke&Prillwitz 1995). However, back then digital video in very small spatial resolution was good enough to show the video in combination with the transcript, but not really to transcribe every detail from it. Rather, one had to use VCRs – either remote-controlled by the transcription environment or directly operated by the transcriber. In the following years, technological advances finally allowed to digitize video full-size SD and then to create digital video directly with the camera and to easily transfer the material to the computer. Now, processing speed and storage capacities would also allow HD videos to be used full-size in a transcription environment. However, even on very large screens, video competes with the space needed for a useful transcription layout. This is even more true so with material that has been shot with multiple cameras. Two of our projects, Dicta-Sign and DGS Corpus, use seven cameras to record a pair of informants, too much to be displayed full-size at the same time.

Sign language transcription environments such as ELAN (Crasborn&Sloetjes 2008) or iLex (Hanke&Storz 2008) have been designed at times when researchers were using digital video in the size of up to half SD (such as 320x240) and certainly need to be improved for the requirements of today's projects delivering multi-camera HD material.

ELAN allows the user to relate several media files to a transcript and to sync them. iLex just allows one single media container and relies on the container format, such as QuickTime, to group and sync several video streams into one container.

To save screen real estate, both systems allow the user to vary the display size from a fraction of the videos' spatial resolution to full size (and beyond) for all visible videos. iLex in addition allows the user to switch on or off individual tracks within the media container. This works quite fine with two or three different views grouped, but fails to provide an adequate solution when more camera views are available: A spatial layout of the tracks (defined in the container) that might be optimal when focussing on one

informant can be far from optimal in situations where both informants need to be watched in parallel.¹

In both systems, different display sizes for individual video views are not possible except by relying heavily on container formats to include one video in multiple sizes and the user switching one on and the others off as needed or to produce copies of one movie in several spatial resolutions.

Zooming onto specific parts of a video is also not possible except by providing the zoomed version as a separate movie (cf. Crasborn & Zwitterlood 2008).

Here we present a user interface study that promises to deliver the flexibility needed and at the same time to save transfer bandwidth and local processing power which even nowadays are an issue when dealing with several HD videos in parallel.

2. Screen Layouts

In our projects, transcribers have screens with native resolutions of either 1920x1200 or 2560x1440. So except for very rare cases, full HD resolution (1920x1080) is not used for transcription as the movie would occupy a good part of the screen. Depending on what they transcribe, we expect users to work more with $\frac{1}{3}$ of full HD (640x360), $\frac{1}{4}$ (480x270) or even $\frac{1}{6}$ (320x180) rather than with $\frac{1}{2}$ (960x540).²

Based on the type of discourse to be described as well as personal preferences, we expect most transcribers to work with one or two movies at a time, optionally with thumbnail-size view (160x90) for the other cameras.

2.1 Focus on one movie at a time

In this layout, clicking on any (movie or still)³ thumbnail zooms the video shown so far out into a thumbnail and the thumbnail video in to the current large size. When needed, a context menu allows to switch to a two-large-movies layout.

¹ In ELAN, switching a video on or off could be easily realised from the transcript if it is the only video in its container. The layout of the videos, however, can only be influenced with respect to a left-to-right order.

² Users can still resize to any in-between value they prefer. iLex uses the next higher available resolution and scales that down.

³ Stills are preferred by some users of moving images in order to reduce visual noise.

2.2 Focus on primary views for both/all informants

With two or more large-size videos shown, thumbnails are bundled to one of the large videos. A click on a thumbnail then exchanges its movie with the bundled one.

2.3 Automatic switching based on tagging

Whenever tagging is available that is a good estimator for what the transcriber will need to focus on, this tagging can be used to switch automatically between different layouts. If for example turns have already been tagged, it makes sense to have the signer in a large view and the addressee in a small view. Good approximations to manual turn tagging can hopefully be in the near future achieved automatically through image processing (cf. e.g. Efthimiou et al. on Dicta-Sign, this volume). Another source of information is knowledge about the tasks informants are currently working on, as logged by Session Director (cf. Hanke et al.: DGS Corpus and Dicta-Sign: The Hamburg Studio Setup, this volume). Of course, thumbnail buttons remain available to either switch to secondary views (such as birds-eye views on a single informant) or to the other informant when needed.

3. Derived Views

In addition to the views available through the films actually shot during the data collection, some derived formats are useful for the transcriber. Top of the list with HD sources certainly is zooming onto particular parts of the video, such as the signer's face. In the beginning, we ask the user to draw a frame around the signer's face. This may have to be repeated for several points in time in the video, whenever the signer moves significantly. In the future, we hope to automate this windowing through image processing (cf. Collet et al., this volume, on interfaces between transcription environments and image processing). Other examples for derived views include results of image processing such as stereo pictures.

Changes in spatial or temporal resolution alone are not considered derived views. We try to give the users the impression that any view can be scaled continuously; therefore resolution pyramids are not immediately visible to the user. As we do not see any need at this point of time to work with reductions in temporal resolution (in fact we would like to have higher resolutions available), such reductions are simply not offered as options.

We are still experimenting how to handle cropping (cutting away border stripes of the image). The idea with cropping is that anything lying outside the marked area is of no interest for transcription, and therefore the cropped movie could replace the original for all further processing. One of the problems is who might be authorised to apply cropping, as all information outside the cropped area would no longer be visible to any transcriber so errors in cropping might pass undetected.

While results of image processing might not immediately become available to the transcriber, zooms are available to the user at the click of a button: iLex just loads a higher-resolution version of the movie and then lets QuickTime crop the image in memory to the part the user is interested in. If such a derived view is used over a longer period of time, iLex marks this view to be produced as a stand-alone movie to save bandwidth and computing power on the client's side.

4. Video Server Infrastructure

Our video server currently consists of three machines with 16 processors each, attached to a SAN with a storage capacity of 100 TBytes. Two thirds of the capacity is reserved for the original footage, one third is available for caching resolution pyramids and other derived video. However, no real caching strategy is in place at this point in time. Instead, cache movies are produced as processing capacity allows. iLex then keeps track of their usage, but purging is currently left to the administrators. Our idea is to observe the system for some time before implementing strategies how to manage cache size. In the current iLex structure which allows the user to copy movies onto the local harddisk in order to work at locations where bandwidth does not allow video server access purging might render local copies useless as iLex would no longer look for them once the database entries are deleted.

Another option for the future is to provide zooming on the server side in real-time. As we currently do this on the client side, we know it can be done in real-time. Implementation on the server side, however, requires much more work, so we will first observe how much this feature will actually be used.

5. Acknowledgements

The research leading to these results has received funding from the European Community's Seventh Framework Programme (FP7/2007-2013) under grant agreement n° 231135 and from the German Academies of Science Programme.

6. References

- Crasborn, O. & Sloetjes, H. (2008). Enhanced ELAN functionality for sign language corpora. In O. Crasborn, T. Hanke, E. Efthimiou, I. Zwitterlood & E. Thoutenhoofd (Eds.), *Construction and Exploitation of Sign Language Corpora. 3rd Workshop on the Representation and Processing of Sign Languages*. ELRA, Paris, pp. 39-43.
- Crasborn, O. & Zwitterlood, I. (2008). The Corpus NGT: an online corpus for professionals and laymen, In O. Crasborn, T. Hanke, E. Efthimiou, I. Zwitterlood & E. Thoutenhoofd (Eds.), *Construction and Exploitation of Sign Language Corpora. 3rd Workshop on the Representation and Processing of Sign Languages*. ELRA, Paris, pp. 44-49.
- Hanke, T. & Prillwitz, S. (1995). SyncWriter. Integrating Video into the Transcription and Analysis of Sign Language. In T. Schermer & H. Bos (Eds.), *Proceedings of the Fourth European Congress on Sign Language Research, Munich, September, 1994*. Signum: Hamburg, pp. 303-312.
- Hanke, T. & Storz, J. (2008). iLex – A Database Tool for Integrating Sign Language Corpus Linguistics and Sign Language Lexicography. In O. Crasborn, T. Hanke, E. Efthimiou, I. Zwitterlood & E. Thoutenhoofd (Eds.), *Construction and Exploitation of Sign Language Corpora. 3rd Workshop on the Representation and Processing of Sign Languages*. ELRA, Paris, pp. 64-67.