English

# Movement Epenthesis Generation for Video Synthesis of Sign Language

## Ze-Jing Chuang, Chung-Hsien Wu*, Wei-Sheng Chen

**Department of Computer Science and Information Engineering, National Cheng Kung University, Tainan, Taiwan, R.O.C.**
**Email: chwu@csie.ncku.edu.tw**

Individuals with hearing/speech impairment generally have problems in communication skill learning. Since deaf people have difficulty in speaking, many communication-aided approaches, such as sign language, finger spelling, lip-reading, and total communication, have been proposed to enhance their language and communication skills. Because of the impairment of vocal communication, sign-language is still the main communication method between hearing-impaired people. Although the hearing-impaired people can learn sign language with several kinds of assistance, such as books, photographs, and videotapes, none of them can provide a flexible and realistic access of sign language. Accordingly, computerized assistive learning system is proposed for sign language learning. This paper proposes a method that is able to generate a video segment of movement epenthesis using the original video clips of a real signer. Figure 1 shows the diagram of the proposed approach. The basic idea of this approach is to find the best concatenation point between two consecutive sign videos. Based on these two points, this approach can generate the smoothest epenthesis video frames to provide the sign video output.
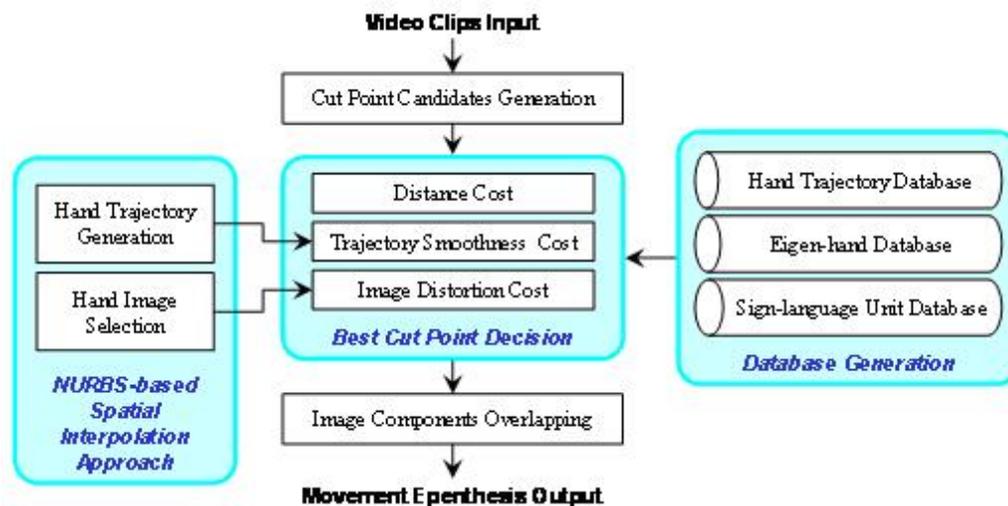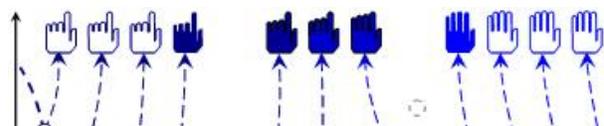


Figure 1. Diagram of video-based sign-language synthesis system.

The first step for generating the movement epenthesis decides where the previous video should stop and where the next video should start, which are called "cut points."

The best cut point pair is selected from the candidates according to the concatenation cost,

which is a linear combination of distance cost, smoothness cost, and image distortion cost. The distance cost is a normalized Euclidian distance between the hand locations at two cut points. The Non-Uniform Rational B-Spline (NURBS) function is used to generate the hand movement trajectory of movement epenthesis. In image distortion cost calculation, the NURBS function is used to generate a hand image change path for the purpose of hand image selection. From the annotation in sign-language database, all hand images in pre-captured video clips are already extracted. Therefore the hand images selection process is done with the NURBS curve. An illustration of hand image selection is shown in Figure 2.
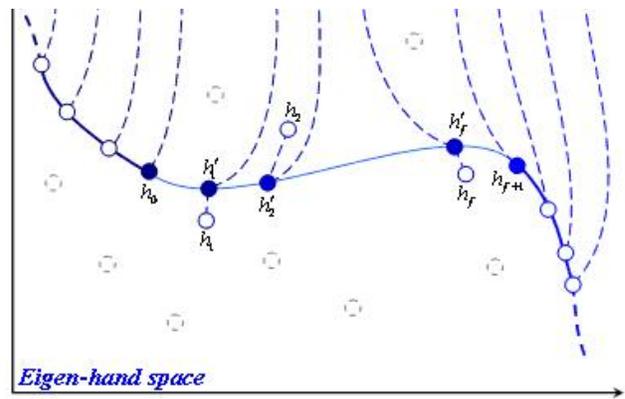


Figure 2. An illustration of the nearest hand image selection.

To support the approach, a set of sign-language databases are collected and preprocessed with image calibration, content annotation, and principle component analysis. Finally, to construct a sign-language synthesis system for demonstration and evaluation, an image component overlapping procedure is also applied as the post-process for generating an output continuous sign-language sentence video.

Using the proposed approach, a video-based Taiwanese sigh-language synthesis system was developed for evaluation on a personal computer. The final integrated system interface is shown in Figure 3. With this interface, user is allowed to input a natural language sentence. All the test video clips were evaluated by 10 normal people and 5 hearing-impaired people using the scores of fluency, understanding, and similarity. The evaluation result demonstrated that the generated movement epenthesis is not only smoothly and naturally, but also similar to the real sign-language sentence video. The result for using the combined concatenation cost also outperforms that using only one of distance cost, smoothness cost, or image distortion cost in both objective and subjective evaluations.

Figure 3. Interface of the constructed sign-language synthesis system.