# Synthesis of virtual reality animation from sign language notation using MPEG-4 body animation parameters

M Papadogiorgaki[1], N Grammalidis[2], N Sarris[3] and M G Strintzis[4]

[1,2,4]Informatics and Telematics Institute
1st Km Thermi-Panorama Road, 57001 Thermi-Thessaloniki, GREECE

[3]Olympic Games Organizing Committee, Athens 2004
Iolkou 8 & Filikis Etairias, 14234 Nea Ionia, GREECE

[1]*mpapad@iti.gr*, [2]*ngramm@iti.gr*, [3]*NESarris@athens2004.com,* [4]*strintzi@iti.gr*

[1,2,4]*www.iti.gr*, [3]*www.athens2004.gr*

## ABSTRACT

This paper presents a novel approach for generating VRML animation sequences from Sign Language notation, based on MPEG-4 Face and Body Animation. Sign Language notation, in the well-known Sign Writing system, is provided as input and is initially converted to SWML (Sign Writing Markup Language), an XML-based format that has recently been developed for the storage, indexing and processing of Sign Writing notation. Each basic sign, namely sign box, is then converted to a sequence of Body Animation Parameters (BAPs) of the MPEG-4 standard, corresponding to the represented gesture. In addition, if a sign contains facial expressions, these are converted to a sequence of MPEG-4 Facial Animation Parameters (FAPs), while exact synchronization between facial and body movements is guaranteed. These sequences, which can also be coded and/or reproduced by MPEG-4 BAP and FAP players, are then used to animate H-anim compliant VRML avatars, reproducing the exact gestures represented in the sign language notation. Envisaged applications include interactive information systems for the persons with hearing disabilities (Web, E-mail, info–kiosks) and automatic translation of written texts to sign language (e.g. for TV newscasts).

## 1. INTRODUCTION

The SignWriting system is a writing system for deaf sign languages developed by Valerie Sutton for the Center of Sutton Movement Writing, in 1974 (SignWriting site, 2004). A basic design concept for this system was to represent movements as they are visually perceived, and not for the eventual meaning that these movements convey. In contrast, most of the other systems that have been proposed for writing deaf sign languages, such as HamNoSys (the Hamburg Notation System) or the Stokoe system employ alphanumeric characters, which represent the linguistic aspects of signs. Almost all international sign languages, including the American Sign Language (ASL) and the Brazilian Sign Language (LIBRAS), can be represented in the SignWriting system. Each sign-box (basic sign) consists of a set of graphical and schematic symbols that are highly intuitive (e.g. denoting specific head, hand or body postures, movements or even facial expressions). The rules for combining symbols are also simple, thus this system provides a simple and effective way for common people with hearing disabilities that have no special training in sign language linguistics, to write in sign languages. Examples of SignWriting symbols are illustrated in Figure 1.



**Figure 1.** *Three examples of representations of American Sign Language in SignWriting system.*

An efficient representation of these graphical symbols in a computer system should facilitate tasks as storage, processing and even indexing of sign language notation. For this purpose, the SignWriting Markup Language (SWML), an XML-based format, has recently been proposed (Costa, 2001). An online converter is currently available, allowing the conversion of sign-boxes in SignWriting format (produced by SignWriter, a popular SignWriting editor) to SWML format.

Another important problem, which is the main focus of this paper, is the visualization of the actual gestures and body movements that correspond to the sign language notation. Grieve-Smith (2001) presented a thorough review of state-of-the art techniques for performing synthetic animation of deaf signing gestures. Traditionally, dictionaries of sign language notation contain videos (or images) describing each sign-box, however the production of these videos is a tedious procedure and has significant storage requirements. On the other hand, recent developments in computer graphics and virtual reality, such as the new Humanoid Animation (H-Anim) (H-anim, 2004) and MPEG-4 SNHC (MPEG-4 document, 1999) standards, allow the fast conversion of sign language notation to Virtual Reality animation sequences, which can be easily visualized using any VRML-enabled Web browser.

In this paper, we present the design, implementation details and preliminary results of a system for performing such a visualization of sign-boxes, available in SWML. The proposed technique first converts all individual symbols found in each sign box to sequences of MPEG-4 Face and Body Animation Parameters. The resulting sequences can be used to animate any H-anim-compliant VRML avatar using MPEG-4 SNHC BAP and FAP players, provided by EPFL. The system is able to convert all hand symbols as well as the associated movement, contact and movement dynamics symbols contained in any ASL sign-box. Manual (hand) gestures and facial animations are currently supported, while we plan to implement other body movements (e.g. torso) in the near future. The proposed technique has significant advantages:

- Web- (and Internet-) friendly visualization of signs. No special software has to be installed except a VRML plug-in to a Web browser,
- Allows almost real-time visualization of sign language notation, thus enabling interactive applications,
- Avatars can easily be included in any virtual environment created using VRML, which is useful for a number of envisaged applications, such as TV newscasts, automatic translation systems for the deaf, etc.
- Efficient storage and communication of animation sequences, using MPEG-4 coding techniques for BAP/FAP sequences.

Significant similar work for producing VRML animations from signs represented in the HamNoSys transcription system to VRML has been carried out by the EC IST ViSiCAST project (Kennaway, 2001), and its follow-up project "E-Sign"(E-sign site, 2004). Current extensions of HamNoSys are able to transcribe all possible body postures, movements and facial expressions (Hanke, 2002) and significant work towards supporting MPEG-4 BAPs has been made. The main contribution of the proposed approach in this paper is the attempt to work towards the same direction for the most common and popular representation of Sign Languages, which is the SignWriting notation system.

The paper is organized as follows: In Section 2, the proposed technique for converting SWML sign boxes to MPEG-4 Face and Body Animation Parameters is described. The synthesis of animations for H-anim avatars and the design of the experimental "Vsigns" Web page to evaluate the sign synthesis results are outlined in Section 3, while discussion and future work is presented in Section 4.

## 2. CONVERSION OF SWML SIGN BOXES TO MPEG-4 FACE AND BODY ANIMATION PARAMETERS

In this Section, we briefly describe the procedure to convert SignWriting notation in SWML format to MPEG-4 Face and Body Animation Parameters. SWML (SWML Site, 2004) is an XML-based format for the representation of SignWriting notation described by the SWML DTD (currently version 1.0 draft 2) (Costa, 2001).

Each SWML *signbox* consists of a set of symbols, which is specified using the following fields:

a) A shape number (integer) specifying the shape of the symbol,

b) A variation parameter (0 or 1 for hand symbols / 1,2 or 3 for movement and punctuation symbols) specifying possible variations (complementary transformations) of the symbol,

c) A fill parameter (0,1,2 or 3 for hand and punctuation symbols / 0,1 or 2 for movement symbols) specifying the way the shape is filled, generally indicating its facing to the signer,

d) A rotation parameter (0-7) specifying a counter-clockwise rotation applied to symbol, in steps of 45 degrees,

e) A transformation flip parameter (0 or 1) indicating whether the symbol is vertically mirrored or not, relatively to the basic symbol and, finally,

f) The x and y coordinates of the symbol within the sign box.

For sign synthesis, the input for the sign synthesis system consists of the SWML entries of the sign boxes to be visualized. For each sign box, the associated information corresponding to its symbols is parsed.

Currently, symbols from the 1995 version of the Sign Symbol Sequence (SSS-1995) are supported. This sequence comprises an "alphabet" of the SignWriting notation system, while true images (in gif format) of each symbol contained in this sequence are available in (SWML site, 2004). The proposed system is able to convert

- All 106 hand symbols,
- All 95 (hand) movement symbols
- Two punctuation symbols (180,181), which contain synchronization information.
- 27 Facial expression/animation symbols

Other punctuation symbols as well as symbols that represent torso and shoulder movements (12 symbols) are currently not implemented (decoded) by the system. Information from symbols, within each sign-box, that are supported by the sign synthesis application, i.e. hand symbols as well as corresponding movement, contact and movement dynamics symbols, is then used to calculate the MPEG-4 Face and Body Animation Parameters.

The issue body modelling and animation has been addressed by the Synthetic/Natural Hybrid Coding (SNHC) subgroup of the MPEG-4 standardization group (MPEG-4 document, 1999). More specifically, 168 Body Animation Parameters (BAPs) are defined by MPEG-4 SNHC to describe almost any possible body posture. In addition, 68 Face Animation Parameters (FAPs) are used to describe almost any possible facial expression. Most BAPs denote angles of rotation around body joints, while FAPs usually denote movements of specific facial features (Facial Definition Points, FDPs) along a pre-determined axis in 3-D space.

The conversion of the symbols contained in a SWML sign box to BAP sequences starts by first examining the symbols contained within the input sign box. If no symbols describing dynamic information such as hand movements, contact or synchronization exist, the resulting BAP sequence corresponds to just one frame (i.e. a static gesture is reproduced). Information provided by the fields of the (one or two) hand symbols, contained in the sign box, is used to specify the BAPs of the shoulder, arm, wrist and finger joints. On the other hand, if symbols describing dynamic information exist, the resulting BAP sequence contains multiple frames, describing animation key-frames (i.e. a dynamic gesture is reproduced). The first key-frame is generated by decoding the existing hand symbols, as in the case of static gestures. Since the frame rate is constant and explicitly specified within a BAP file, the number of resulting frames may vary, depending on the complexity of the described movement and its dynamics. Synchronization symbols and contact also affect the represented movement and in some cases require special treatment.

When a signbox contains facial expression or animation symbols, the corresponding FAP frame(s) are determined by predefined lookup tables, which provide the FAP values defining one or more FAP frames per facial animation symbol. When two or more facial expression symbols co-exist within the same sign-box, these may either define an animation sequence or have to be combined all together (if each symbol activates different FAPs). The latter case, which is more common, is currently supported by the proposed system.

Smooth and natural-looking transitions between the Face and Body Animation parameters corresponding to each signbox is achieved by generating additional intermediate frames using a FAP/BAP interpolation procedure. A linear interpolation function is used to generate additional FAP/BAP frames to implement:

a) The transition between the neutral face/body position and the first frame of the first sign-box

b) The transition between the end frame of one signbox and the start frame of the next signbox

c) The transition between the end frame of the last sign-box and the neutral body position.

Furthermore, in order to achieve Face/Body synchronization:

a) The frame rates defined for the FAP and BAP sequences should be equal

b) The number of generated FAP frames generated for each sign-box should be always equal to the corresponding number of BAP frames. In order to achieve this goal, the BAP frame sequence is first generated and then specific linear interpolation procedures are used to generate the FAP frame sequence.

A block diagram of the proposed system for processing each sign-box is illustrated in Figure 2, while additional details about the generation of BAPs for static and dynamic gestures as well as the generation of FAPs for gestures containing facial expressions/animations are provided in the following Subsections.
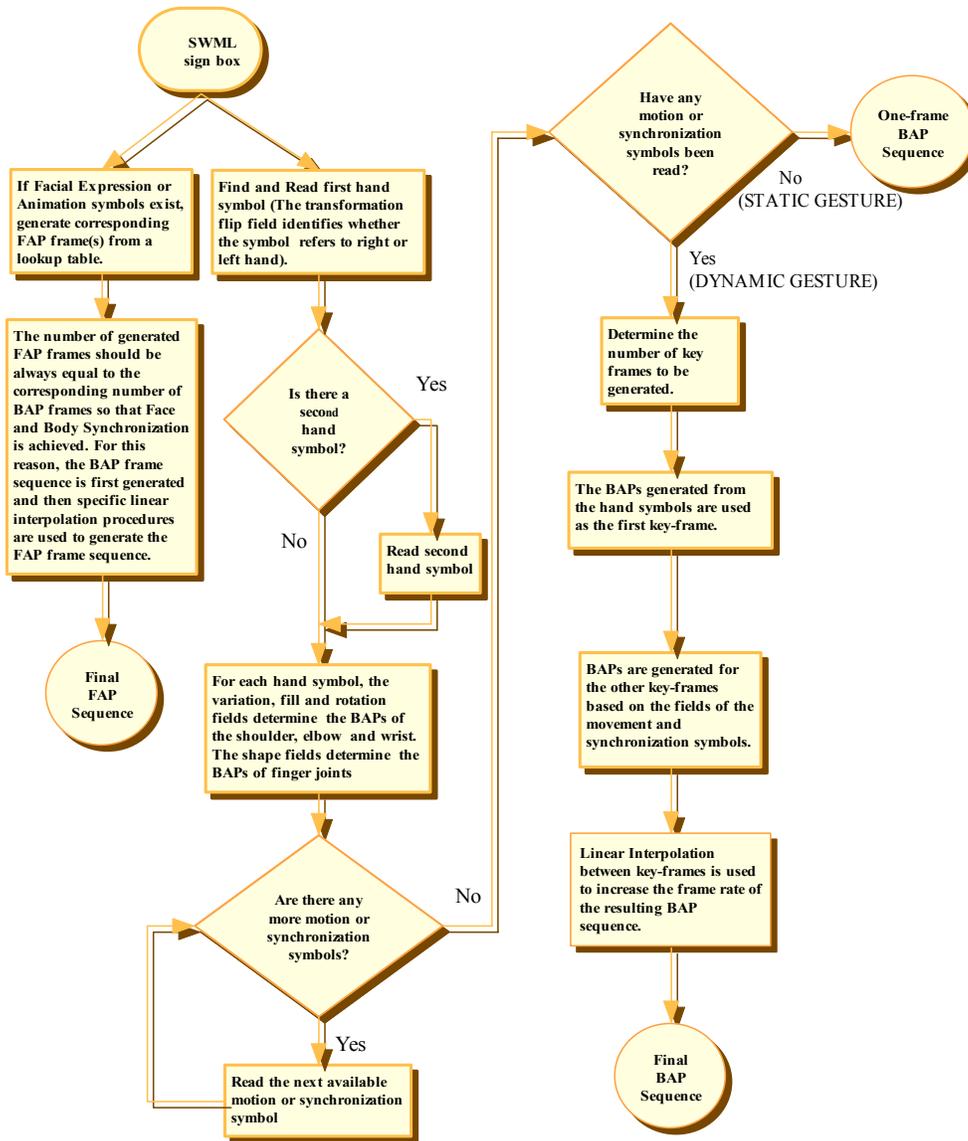


**Figure 2.** *A block diagram of the proposed system.*

## 2.1 Static gestures

In the following, the procedure to extract useful information from the SWML representation of a hand symbol is summarized:

Initially, the binary "transformation flip" parameter is used to identify whether the symbol corresponds to the left or right hand. Then the fill and variation parameters of each symbol are used to determine the animation parameters of the shoulder and elbow joints:

- If (variation, fill)=(0,0),(0,1) or (1,3) then the axis of the arm is parallel to the floor (floor plane).

- If (variation, fill)=(1,0),(1,1) or (1,2) then the axis of the arm is parallel to the human body (wall plane)
- If (variation, fill)=(1,0) or (1,3) then the signer sees his palm
- If (variation, fill)=(1,1) or (0,0) then the signer sees the side of his palm
- If (variation, fill)=(1,2) or (0,1) then the signer sees the back of his palm

In addition, the rotation parameter is used to determine the animation parameters of the wrist joint:

- If the signer sees the side of his palm, the rotation value (multiplied by 45 degrees) is used to define the R_WRIST_FLEXION BAP (for the right hand) or the L_WRIST_FLEXION BAP (for the left hand).
- In the other two cases (signer sees his palm or the back of his palm), the rotation value (multiplied by 45 degrees) is used to define the R_WRIST_PIVOT BAP (for the right hand) or the L_WRIST_PIVOT BAP (for the left hand).

Finally, the symbol shape number is used to specify the animation parameters corresponding to finger joints, using look-up tables of BAP values corresponding to each symbol.

If the sign box contains a second hand symbol, similar procedures are used to extract the Body Animation Parameters of the other hand. After the processing of all existing hand symbols, all Body Animation Parameters corresponding to shoulder, elbow, wrist and finger joints are determined and stored.

## 2.2 Dynamic gestures

MPEG-4 standard allows the description of human body movement using a specific set of Body Animation Parameters corresponding to each time instant. Systems like SignWriting that use a high level animation description define movement by specifying a starting and an ending position, in case of simple motion with constant velocity, or the full trajectory, in case of more complex motion. However, the description of complex motion is also possible by specifying a number of intermediate key-frames. In the following, the procedures for generating these BAP key-frames are briefly described.

When all movement description symbols have been identified, the shape number field identifies their shapes (i.e. the type of movement). First, the total number of key-frames to be produced is specified, based on the number and nature of the available movement, movement dynamics, contact, and synchronization symbols. More specifically, a look-up table is used to define an initial number $k$ of key frames for each movement symbol. Furthermore, the fill parameter specifies whether the motion is slow, normal or fast. In addition, some symbols explicitly specify the movement duration. For this reason, a classification of such symbols into three categories has been defined and a different duration value $D$ is defined for each category:

- Slow motion ($D=3$)
- Normal motion ($D=2$)
- Fast motion ($D=1$)

The total number of frames to be generated when only one motion symbol exists is $N=kDP$, where $P$ is a fixed multiplier (e.g. $P=10$). If the number of such symbols is more than one, the total number of key-frames is the maximum between the numbers of key-frames, corresponding to each symbol. Finally, if the sign box contains a contact symbol, the total number of frames is increased by two (in case of simple contact) or four (in case of double contact).

The initial key-frame is generated by decoding the available hand symbols, exactly as in the case of static gestures. The rotation and transformation flip fields specify the exact direction of movement. Also, the variation field specifies whether the right or the left hand performs the movement. Using information from all available movement, contact and synchronization symbols, the other BAP key-frames of the specific dynamic gesture are then generated using a specific function for each key-frame. Synchronization (Movement Dynamics) symbols (180,181 and 182) are handled in a similar way as movement symbols but an exception exists for the "Un-even alternating" symbol, where first one hand moves, while the other hand is still and then the opposite. To handle this case the total number of key frames is doubled ($N=2kDP$). To produce the first $kDP$ frames, BAPs are generated only for the first hand, so the second hand remains still. Then, BAPs are generated only for the second hand, to produce the next $kDP$ frames, so the first hand remains still.

Finally, when the BAPs for all key-frames have been computed, BAP interpolation is used to increase the frame rate of the resulting BAP sequence. This interpolation procedure results to smoother transitions between key frames.

Interpolation is generally achieved by approximating the motion equation using a mathematical function and then re-sampling this function to obtain the desired intermediate positions at intermediate time instants. Various interpolation functions can be selected in order to improve results. Since Body Animation Parameters represent rotations around specific joints, quaternion interpolation was seen to provide good results [8], but the complexity of the method is increased. For this reason, a linear interpolation technique was applied, which was seen to be very efficient for most signs, since key-frames have been selected so as to simplify the movement description between consecutive key-frames.

*2.3    Gestures containing facial expressions-animations.*

The generation of the FAP frame sequence is performed after the generation of the BAP frame sequence, so that the total number of generated FAP frames is exactly the same as the total number of BAP frames. For each sign-box, the FAP key-frames are determined, based on the existing facial expression/animation symbols, from predefined lookup tables for each symbol. The number of FAP key-frames, $N_{FAP\_keyframes}$, is generally much smaller than the total number of BAP frames $N_{BAP}$ that have been already generated using the procedures described in the previous Subsections. Therefore, if $FAP(k)$, $k = 0, ..., (N_{BAP} - 1)$ denotes the vector of FAPs corresponding to frame $k$, the FAP key frames are first positioned every $step = N_{BAP}/(N_{FAP\_keyframes} - 1)$ frames:

$$FAP(i*step) = FAP\_keyframe(i*step), i = 0, ..., /( N_{FAP\_keyframes} - 1)$$

Then, each of the remaining FAP frames is determined using linear interpolation between the two closest available FAP key frames.

## 3.  SYNTHESIS OF ANIMATIONS USING H-ANIM AVATARS

The "EPFLBody" BAP player (Vergnenegre, 1999), developed by the École Polytechnique Fédérale Lausanne (EPFL) for the Synthetic and Natural Hybrid Coding (SNHC) subgroup of MPEG-4 was used to animate H-anim-compliant avatars using the generated BAP sequences. Since most BAPs represent rotations of body parts around specific body joints, this software calculates and outputs these rotation parameters as animation key-frames to produce a VRML ("animation description") file that can be used for animating any H-anim-compliant VRML avatar. The "Miraface" FAP player, also developed for MPEG-4 SNHC, by MIRALab, University of Geneva and LIG, EPFL was used for Facial Animation. This software had to be modified so that:

a)  VRML animation output is produced using one CoordinateInterpolator node per face model vertex. A problem with the chosen implementation is that the computational demands for the hardware that is reproducing these animations are increased. A possible solution for this problem that should be investigated in the future is to add CoordinateInterpolator nodes only for the points that have actually been moved.

b)  The face model to be animated using the FAP frame sequence was attached to the body to be animated using the BAP frame sequence. Some slight modifications of the VRML face model were also required (e.g. addition of teeth).

Two frames from resulting animations are illustrated in Figure 7.



**Figure 7.** *Animation of the "You" sign in ASL using an H-anim avatar*

By including a VRML TouchSensor Node within the VRML file describing the H-anim avatar, the viewer can interactively start and/or replay the animation sequence, by clicking on the avatar. The viewer can also interact by zooming in and out to any specific body region and/or by rotating and translating the model within the 3-D space, in order to fully understand the represented sign.

Furthermore, further evaluation of the proposed sign synthesis system was possible by developing an online system (Vsigns site, 2004) for converting text to Sign Language notation and corresponding VRML animation sequences for H-anim compliant avatars. The application, whose interface is illustrated in Figure 8, is currently based on a 3200-word SWML dictionary file, obtained by the SWML site, which has been parsed and inserted into a relational database. The user is allowed to enter one or more words, which are looked up in this dictionary. If more than one entry is found, all possible interpretations are presented to the user, so that he can choose the desired one. On the other hand, if no entries are found for a specific word, the word is decomposed using its letters (finger-spelling). In any case, the user may choose whether to include a particular term to the selected terms to be used for sign synthesis or not. The user then selects a column corresponding to an H-anim compliant avatar, which is used for sign synthesis of the selected term or terms. A fourth column ("Baxter FBA") allows the user to observe facial animation in addition to body animation, using the modified "Baxter avatar". Furthermore, the user may produce and display the corresponding sign(s) in SignWriting format (in PNG format) and SWML for a specific term or the selected terms.
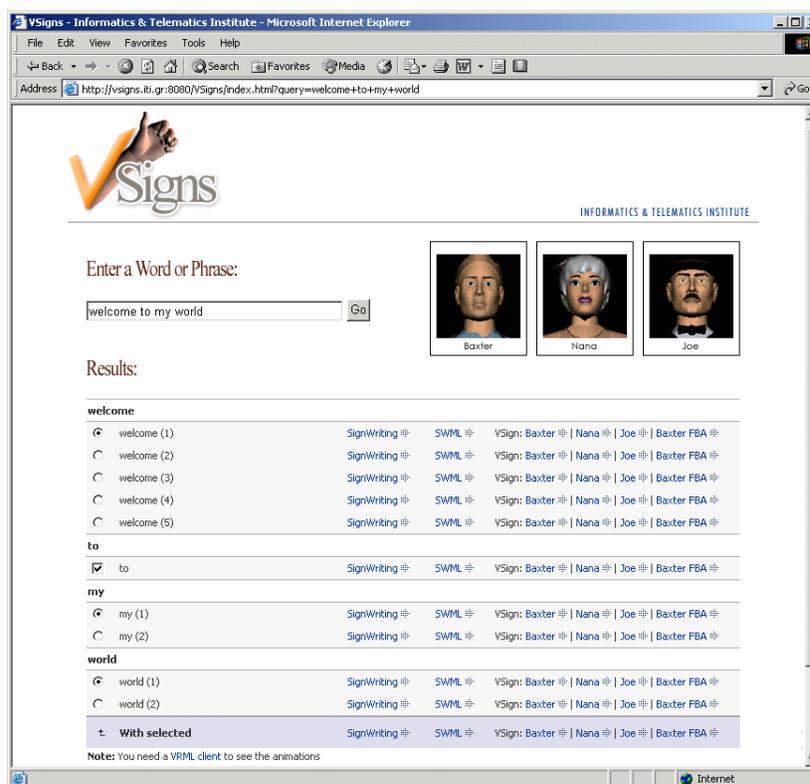


**Figure 8.** *Example query: "Welcome to my world". The user may then select the desired terms and then produce and display sign synthesis results using the selected words or the entire phrase, using any of the available H-anim avatars.*

This experimental Web application has already allowed us to identify significant problems with the synthesis of static and dynamic gestures, which have to be solved in the future, e.g. when contacts and complex movements are involved. A major problem that has to be solved occurs when the sign-box contains contact symbols. In that case the touch between the hands, or the hand and the face is difficult to be achieved. Problems may also occur for complex movements, when the inclinations of the hand joints, which have been estimated in each key frame, are not accurate enough for the exact description of the movement. Both problems can be solved in the future by using inverse kinematics methods, as in (A B. Grieve-Smith, 2001).

Further evaluation is planned for the future, using Greek and International SignWriting users, and attempts will be made to solve possible problems in the reproduction of specific signs. Although these

problems indicate that much more work is needed for correct synthesis of all signs, we believe that with this Web tool, a very important step towards automatic Text to Sign synthesis has been made.

## 4. DISCUSSION AND FUTURE WORK

A novel approach for generating VRML animation sequences from Sign Language notation, based on MPEG-4 Body Animation has been presented. The system is able to convert almost all hand symbols as well as the associated movement, contact and movement dynamics symbols contained in any ASL sign-box. Furthermore, most facial expression and animation symbols are also supported, while torso movements will be also supported in the near future. Some facial expressions, e.g. cheek wrinkles, have not been implemented, since no FAPs exist to produce such movements.

Results are satisfactory and are currently being evaluated by SignWriting users and experts, so that problems associated with specific SignWriting symbols are identified and solved. In the future, improved reproduction of difficult movements (e.g. touching) will be made possible using inverse kinematics techniques.

A short-term goal is to design other practical applications of the proposed system, either as a "plug-in" to existing applications (e.g. sign language dictionaries) or as a stand-alone tool for creating animations for TV newscasts (e.g. weather reports). Particular emphasis will be given in applications that can be used and evaluated by the Greek Sign Language community, thus a dictionary of Greek Sign language, in SignWriter notation, is planned to be supported in the near future.

## 5. REFERENCES

Official Sign Writing site, http://www.signwriting.org/

"Generic coding of audio-visual objects - Part 2: Visual", *MPEG Document ISO/IEC JTC 1/SC 29/WG11* N3056, Maui, December 1999

Official SignWriting site, http://www.signwriting.org/

Official Site of SWML, http://swml.ucpel.tche.br/

Moving Pictures Experts Group, (1999). Generic coding of audio-visual objects - Part 2: Visual, MPEG Document ISO/IEC JTC1/SC29/WG11 N3056. Maui.

F Vergnenegre, Tolga K. Capin, and D. Thalmann (1999). Collaborative virtual environments-contributions to MPEG-4 SNHC. ISO/IEC JTC1/SC29/WG11 N2802, http://coven.lancs.ac.uk/mpeg4/

A B. Grieve-Smith (2001). SignSynth: A Sign Language Synthesis Application Using Web3D and Perl. Gesture Workshop, London, pp. 134-145

R. Kennaway (2001). Synthetic Animation of Deaf Signing Gestures. Gesture Workshop, pp. 146-157, London.

Antonio Carlos da Rocha Costa, Cracaliz Pereira Dimuro (2001). Supporting Deaf Sign Languages in Written Form on the Web. The SignWriting Journal, Number 0, Article 1, July. http://gmc.ucpel.tche.br:8081/sw-journal/number0/article1/index.htm.

M. Preda, F. Preteux (2001). Advanced virtual humanoid animation framework based on the MPEG-4 SNHC standard. Proceedings EUROIMAGE International Conference on Augmented, Virtual Environments and Three-Dimensional Imaging (ICAV3D'01), Mykonos, Greece, pp. 311-314.

Humanoid Animation Standard Group. Specification for a Standard Humanoid: H-Anim 1.1. http://h-anim.org/Specifications/H-Anim1.1/

Official site of E-sign (Essential Sign Language Information on Government Networks) project. http://www.visicast.sys.uea.ac.uk/eSIGN/

Th. Hanke (2002). iLex - A tool for sign language lexicography and corpus analysis. In: Proceedings of the Third International Conference on Language Resources and Evaluation, Las Palmas de Gran Canaria, Spain., pp. 923–926.

Vsigns page, http://vsigns.iti.gr