

Evaluating Facial Expressions in American Sign Language Animations for Accessible Online Information

Hernisa Kacorri¹, Pengfei Lu¹, and Matt Huenerfauth²

¹ The City University of New York (CUNY)
Doctoral Program in Computer Science, The Graduate Center,
365 Fifth Ave, New York, NY 10016 USA

{hkacorri@gc, pengfei.lu@gc}.cuny.edu

² The City University of New York (CUNY)
Computer Science Department, CUNY Queens College
Computer Science and Linguistics Programs, CUNY Graduate Center
65-30 Kissena Blvd, Flushing, NY 11367 USA
matt@cs.qc.cuny.edu

Abstract. Facial expressions and head movements communicate essential information during ASL sentences. We aim to improve the facial expressions in ASL animations and make them more understandable, ultimately leading to better accessibility of online information for deaf people with low English literacy. This paper presents how we engineer stimuli and questions to measure whether the viewer has seen and understood the linguistic facial expressions correctly. In two studies, we investigate how changing several parameters (the variety of facial expressions, the language in which the stimuli were invented, and the degree of involvement of a native ASL signer in the stimuli design) affects the results of a user evaluation study of facial expressions in ASL animation.

Keywords: American Sign Language, accessibility technology for people who are deaf, animation, natural language generation, evaluation, user study, stimuli.

1 Accessible Online Information and Documents in ASL

Many people who are deaf in the United States have lower levels of written language literacy [10]; this makes it difficult for them to read English text on TV captioning, websites, or online documents [7][17]. Animations of American Sign Language (ASL) can make online information and services accessible for these individuals. This paper focuses on our research on ASL; however, many of the techniques and methods could be applied to other sign languages. While it is possible to post videos of real human signers on websites, animated avatars are advantageous if the information is frequently updated; it may be prohibitively expensive to continually re-film a human performing ASL for the new information. Assembling video clips of individual signs together into sentences does not produce high-quality results.

One way to produce animations of ASL would be for a skilled animator (fluent in ASL) to create a virtual human that moves in the correct manner using general-purpose 3D animation software. Since this is too time-consuming and depends too much on the skill of the 3D animator, researchers study automated techniques. The most automated approach is to develop “generation” software, to automatically plan the words of a sign language sentence based on some information input. For instance, in an automatic translation system, the input could be an English text, which must be translated into ASL. While some researchers have investigated ASL generation and translation technologies, the state-of-the-art is still rather limited due to the linguistic challenges inherent in planning sign language sentences [10]. A less automated approach for producing ASL animation is to develop “scripting” software, and allow a human to efficiently “word process” a set of ASL sentences, placing individual signs from a dictionary onto a timeline to be performed by an animated character. Such tools (e.g., [19]) make use of pre-built dictionaries of sign animations, and they incorporate software for automating selection of the transition movements between signs, and other detailed (and time-consuming to specify) aspects of the animation.

The ability to efficiently write and revise ASL is a novel development. Currently, there is no standard written form for ASL that has been accepted by the deaf community. While transcription systems have been proposed, e.g., [14], they typically lack details necessary for capturing grammatically correct ASL animation, and those limitations made them not widely used by signers. As scripting technologies improve, this opens the possibility for future “word processing” and “document creation” in ASL: enabling easier creating and sharing of information in accessible manner.

The linguistic complexity of ASL makes developing animation technologies challenging, e.g., various facial expressions and head movements communicate essential information during sentences. State-of-the-art ASL animation systems do not yet handle facial expressions sufficiently to produce clear output, and our lab is studying how to improve this. An important aspect of our work is user-based evaluation, needed to measure the quality of our models. However, designing an experiment in which native ASL signers evaluate facial expressions is not straightforward; careful design is required. In this paper, we design and evaluate experimental stimuli for use in evaluations of facial expressions in ASL animation. Section 2 explains the importance of facial expressions in ASL, section 3 describes related work, section 4 describes our experiments to investigate how changing the parameters in a stimuli design could affect the results of a user study, and section 5 contains conclusions and future work.

2 Importance of Facial Expressions in ASL Animations

Facial expression conveys grammatical information (questions, negations, etc.) in most sign languages. Eyebrow movements, mouth shape, head tilt/turn, and other facial movements are linguistically required in ASL, and identical hand movements (signs) can have different meanings, depending on the facial expressions performed during the sentence [12]. In a simple case, emotional facial expressions (frustration, sadness, anger) affect the meaning of a sentence. Other facial expressions indicate specific grammatical information about sentences and phrases, e.g.: (1) convert

declarative sentences into Yes-No questions, (2) indicate that a sentence is a WH-word interrogative question (“who, what, where”), (3) invert the logical meaning of a sentence by conveying negation (via head shaking), (4) indicate that some words at the beginning of a sentence are an important “topic” for the upcoming sentence, etc. In this way, a sequence of signs like “BOB LIKE JOHN” could be changed into a Yes-No question (“Does Bob like John?”) or invert its meaning (“Bob doesn’t like John”) by adding a facial expression during an appropriate portion of the sentence. (The timing of the facial expression relative to the signs in the sentence is important).

In a prior study, we experimentally evaluated ASL animations with and without various types of facial expressions, and we found that the inclusion of facial expression led to measurable benefits for the understandability and perceived quality of the animations [11]. However, most prior sign language animation research has not addressed how to synthesize facial expressions [3-5]. To produce an animation with good facial expressions, an animation artist could carefully edit the facial mesh of an animated character to produce beautiful facial expressions, but this is very time-consuming. We want to support automatic synthesis of sign language and scripting of sign language animations. We are studying how to model and generate ASL animations that include facial expressions to convey grammatical syntax information, such as negation; topic; and yes-no, WH-word, and rhetorical questions. Our objective is to determine when signers use these facial expressions, how they perform each, how the timing of these facial expressions occurs in relation to the manual signs, and how the co-occurrence or sequential occurrence of facial expressions affect one another. Thus, we are investigating technologies for automatically planning aspects and timing of face movement. In addition to planning algorithms, we also need a succinct representation of ASL (that can encode a good-quality performance with as few parameters as possible). This makes it practical for a generation system to plan the animation, and it makes it possible for a human using a scripting tool to produce an animation with facial expressions in an efficient manner. We must test both our planning algorithms and our ASL script representation to ensure that they encode sufficient detail for ASL facial expressions that are understandable (and deemed natural) by signers.

3 Related Work on Evaluating ASL and Face Animation

We must evaluate the quality of the facial expressions in an ASL animation to advance research in this field, but it is difficult due to the subtle and complex manner in which facial expressions affect the meaning of sentences. It can be difficult to design experiments that probe whether human participants watching an ASL animation have understood the information that should have been conveyed by facial expressions. The easiest to evaluate is categorical information, e.g., whether (or not) the sentence with a facial expression should be interpreted as a question; it is possible to invent experiments to determine whether a human watching an animation interpreted it as a declarative sentence or as a question, etc. However, some ASL facial expressions convey information in matters of degree, e.g., an emotional facial expression can convey continuous degrees (by intensity of eye-brow movement, etc.). Measuring

whether someone has successfully understood the correct degree is more difficult. In the most challenging case, a facial expression may not affect the superficial meaning of a sentence, but only the implications that can be drawn. For instance, when a signer performs “I didn’t order a soda” with a cluster of behaviors (including frowning and head tilt) during the sign “I,” it can indicate that the signer believes someone else ordered the soda. With facial prominence on the sign “soda,” it could indicate that the signer placed an order, but for something else. In either case, the basic information is the same: the signer did not order a soda, but a different implication can be made.

Researchers studying facial expression of non-signing virtual humans, often evaluate only static faces, e.g., participants must identify the category of the facial expression or assign scores for intensity or sincerity (e.g. [19]). Because sign language facial expressions convey grammatical information and are governed by linguistic rules, additional care is needed to design useful stimuli and questions for evaluations.

Few researchers have explicitly discussed methodological aspects of stimuli design for facial expressions in sign language animation user-studies. Prior research differs as to whether researchers invent their stimuli originally as sign language sentences [6] [16] or as written/spoken language sentences that are translated into sign language stimuli [2][15]. In section 4, we compare both methods, and we study a wider variety of facial expression types; we also investigate the use of comprehension questions, which had not been employed in previous sign language facial expression studies. For inspiration as to how to use comprehension questions, we consider prior research on how humans interpret and understand speech with various prosody [1][13] (speed, loudness, and pitch changes). Researchers designed sets of sentences that, in the absence of prosodic information, contain ambiguity in how they can be interpreted. When prosodic information is added, then one interpretation is clearly correct. Participants in the studies listen to audio performances of these sentences and answer questions about their meaning. These questions are carefully engineered such that someone would answer the question differently – based on which of the alternative possible interpretations of the spoken sentence they had mentally constructed. For example, someone who heard the sentence “I didn’t order a soda” (with prominence on “I”) may be more likely to respond affirmatively to a question asking: “Does the speaker think that someone else ordered a soda?” In designing the studies in section 4, we have used similar experimental design, stimuli, and comprehension questions.

4 Experimental Stimuli Evaluation

The goal of this paper is to identify a methodology for designing stimuli and conducting experiments to measure the quality of facial expressions in an ASL animation. We want to evaluate whether facial expressions in an ASL animation enable participants who view the animations to identify the content of the sentences being performed. In the studies presented in this section, participants look at animations of a virtual human character telling a short story in ASL, and they answer questions about each story. Each story includes one category of facial expression (e.g., Yes-No questions, sadness, etc.). The animations displayed are one of two types: (i) with facial expressions carefully produced by a human animator or (ii) without appropriate facial expressions (i.e., the face doesn’t move). While the experiments in this paper are only pilot studies

used to confirm our methodology, in future work, when we begin to investigate facial expression animation synthesis, our experiments will contain a third type of animation: (iii) with facial expressions planned by our automatic synthesis software.

Deciding on these short stories, creating the animations, and creating the comprehension questions for each story is a process that we refer to as “stimuli design,” and the manner in which this is done can affect the scores collected in a study. The two studies presented in this paper compare two alternative methods for “stimuli design” to determine which is best for conducting ASL facial expression user-studies. There are several variables that we investigate in this paper: (i) whether the stimuli stories originated in English or in ASL, (ii) the amount of involvement of a professional native signer in designing the stimuli, (iii) the categories of facial expressions included in the stimuli, and (iv) the complexity of stimuli (i.e., number of words per story).

In the remainder of this section, we investigate how some of these variables affect users’ opinion by comparing two stimuli sets, which we refer to as our “English-to-ASL” stimuli and our “ASL-originated” stimuli. The primary difference between the sets is the degree of involvement of a native ASL signer in the stimuli-creation process (leading to more fluent ASL sentences in the “ASL-originated” set) and the categories of facial expressions included in each stimuli set. The English-to-ASL stimuli were first evaluated in a study in 2011 [11], and the “ASL-originated” stimuli were evaluated in a new study in 2013.

In both studies, native ASL signers watched animations of a virtual human character, as shown in Fig. 1, telling a short story in ASL. The story was either (i) with facial expressions added by a native ASL signer or (ii) without facial expressions added. Then, participants answered comprehension questions carefully engineered to capture the possible confusion introduced by a misinterpretation of the face. A native signer, who is a professional interpreter, conducted all the instructions and interactions. The animations in both studies were created using identical commercial sign language animation software, Vcom3D Sign Smith Studio [18].



Fig. 1. ASL character and some of the available facial expressions

Our methodology for creating the “English-to-ASL” stimuli was to begin with an English sentence whose meaning would change with/without prosody, and then we attempted to translate the sentence into ASL in a manner that would preserve this reliance on the prosodic information (conveyed by facial expression instead of spoken prosody). Specifically, we asked a native signer (who works as a professional interpreter) to: (i) translate each of the English passages we use into ASL, (ii) use the Vcom3D Sign Smith Studio to produce the ASL animations, and (iii) choose from the available facial expressions repertoire the facial expressions that she thinks linguistically or naturally conveyed the prosodic information for the ASL stimuli. Each animation was produced in two versions: with and without facial expressions added.

<p>1 Original Spoken English Sentence (transcript of audio): I will go to the new restaurant you suggested. It is Chinese?</p> <p>ASL (glosses and facial expressions): I WILL GO NEW RESTAURANT YOU SUGGEST. IT CHINESE.</p>	<p>Comprehension Questions: Q1: Is Charlie asking you a question? Q2: Does Charlie know what kind of restaurant it is? Q3: Did you already tell Charlie that the restaurant is Chinese? Q4: Will Charlie go to the new restaurant?</p>
<p>2 Original Spoken English Sentence (transcript of audio, emphasis on the word "students" to imply that only the students stayed home): It was raining. The students stayed home today.</p> <p>ASL (glosses and facial expressions): TODAY, RAIN. <u> </u>emphasis THEY STUDENT STAY-HOME.</p>	<p>Comprehension Questions: Q1: Did the teachers also stay home? Q2: Is Charlie upset at the students? Q3: Did the students stay home yesterday? Q4: Was it raining today?</p>

Fig. 2. English-to-ASL Stimuli Set examples: (1) yes/no-question and (2) emphasis

There were a total of 28 stimuli with an average of 9 signs in length, and at least one facial expression per story. Fig. 2 shows two examples of stimuli used, as original English and ASL translated transcriptions, and the corresponding comprehension questions. The bars over the script indicate the facial expression to be performed during some of the signs. The stimuli can be divided into 5 categories (Fig. 3), based on the facial expression. The number of stimuli per category is given in parenthesis.

Y/N-Question (4):	The stimuli contained a yes-or-no question. When translated into ASL, a yes/no facial expression was used, without which, it could be interpreted as a declarative statement. See Fig. 2(1).
Wh-Question (4):	The stimuli contained an interrogative (who/what/where) question. The animation included a wh-question facial expression, without which, the sentence may be interpreted as a relative clause: "Last Friday, I saw Metallica. Which is your favorite band?"
Emphasis (8):	The stimuli contained a single word or phrase emphasized, to indicate contrast or incredulity: "It was raining. The <i>students</i> stayed home today." (This suggests the others did not.) "My sister <i>said</i> she ordered coffee, but the waiter brought tea." (This suggests disbelief.) While human signers convey emphasis via pausing, facial movement, and size/speed of hand movements, our animations included facial expression changes only.
Continue (4):	The prosodic cues in these passages convey that the speaker was not yet finished a thought but was only momentarily pausing: "I like to go to the movies and go to plays..." Once again, this information doesn't only correspond to a linguistically meaningful facial expression in ASL, but is communicated through additional signing parameters of speed and eye-gaze direction.
Emotion (8):	The stimuli were performed with a strong emotion (frustration or sadness) that affected their meaning: "Tomorrow is my 30th birthday. I am excited." (A sad face during the second sentence suggests the signer is not really excited.) "Last Friday, my brother drove my car to school." (With an angry facial expression, this suggests that the signer disapproves what her/his brother did.)

Fig. 3. Five categories of facial expression in the English-to-ASL Stimuli Set

Starting with English speech passages when creating stimuli for an ASL animation study seemed like a good approach given: (i) it is true to the goal of ASL animation synthesis, that is converting English text or speech to comprehensible ASL animations; (ii) it makes use of passages that are carefully engineered and successfully applied to collect users interpretation, and (iii) prosodic information in English is often conveyed by facial expressions in ASL. However, it can lead to various problems. First, the English influence might result in ASL stimuli following an English word order, e.g. the ASL sentence in Fig. 2(1) has a rather English-like word order. Second, some of the categories like Emphasis and Continue are communicated by a cluster of behaviors, not a single ASL facial expression, as discussed in section 4.1.

(1) no participant saw the same story twice, (2) the order of presentation was randomized, and (3) each participant saw every story – in either version (a) or (b). Native ASL signers were recruited from ads posted on Deaf community websites in New York. All instructions and interactions were conducted in ASL by a native signer (a professional interpreter). In [8-9], we discussed why it is important to recruit native signers, and we list best-practices to ensure that responses given by participants are as ASL-accurate as possible. Twelve participants evaluated the English-to-ASL stimuli set: 8 participants used ASL since birth, 3 began using ASL prior to age 10 and attended a school using ASL, and 1 participant learned ASL at age 18. This final participant used ASL for over 22 years, attended a university with instruction in ASL, and uses ASL daily to communicate with a spouse. There were 7 men and 5 women of ages 21-46 (median age 32). Sixteen participants evaluated the ASL-originated stimuli set: 10 participants learned ASL prior to age 5, and 6 participants attended residential schools using ASL since early childhood. The remaining 10 participants had used ASL for over 9 years, learned ASL as adolescents, attended a university with classroom instruction in ASL, and used ASL daily to communicate with a significant other or family member. There were 11 men and 5 women of ages 20-41 (median age 31).

Fig. 6 shows the results of the studies that compare English-to-ASL stimuli and ASL-originated stimuli, with the results of the “Emotion” category presented separately from the results from all other categories. Error bars indicate standard error of the mean; significant pairwise differences are marked with stars (ANOVA, $p < 0.10$). Our goal is to identify “good” stimuli for use in studies evaluating ASL facial expression animations. Since a human animator has carefully produced the facial expressions for these studies, “good” stimuli should have a big difference in comprehension scores between the without-facial-expression and with-facial-expression versions. It is important to note that the scores across studies can’t be directly compared, since the sentences and questions may have been more difficult in one study.

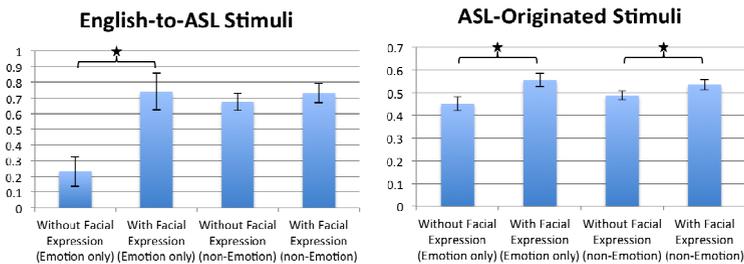


Fig. 6. Comprehension question scores for both types of stimuli, showing results for animations with- and without-facial-expressions, with results for emotion and non-emotion categories

For English-to-ASL stimuli, for the emotion category, adding facial expressions led to significantly higher comprehension scores. However, there was no benefit from adding facial expressions for the non-emotion categories, which didn’t convey the subtle meaning differences that we had intended. Perhaps since the stimuli were first conceived as English stimuli with vocal prosody, something was “lost in translation” when the stimuli were converted into ASL animations with facial expressions.

For the ASL-originated stimuli, adding facial expressions led to significantly higher comprehension scores for both emotion and non-emotion categories. This is a desirable result because it indicates that the stimuli/questions allowed us to distinguish between animations with good or with bad facial expressions (in this case, no facial expressions at all). If we used these stimuli/questions in future studies, we could compare the performance of animations with facial expressions automatically synthesized by our software – to animations with facial expressions produced by a human animator or without any facial expressions. Thus, we could track the performance of our facial-expression synthesis algorithms to guide our research.

5 Discussion and Future Work

This paper has described how we engineered stimuli and questions to measure whether the viewer has understood linguistic facial expressions correctly. Our evaluation methodologies and stimuli will be of interest to other animation researchers studying ASL or other sign languages used internationally. We found that designing stimuli in English and then translating them into ASL was not an effective methodology for designing a sign-language facial expression evaluation study. We have also found that the involvement of native ASL signers in the stimuli design process is important in achieving a high-quality result. In section 4, we preferred the ASL-originated stimuli because we were able to measure a significant benefit from adding facial expressions for both emotion and non-emotion categories. However, the best stimuli for our experiments would show a large difference when facial expressions are added. We note that there was a more dramatic difference in the comprehension scores for the emotion English-to-ASL stimuli. In future work, we may investigate whether we can design ASL-fluent sentences/questions that are analogous to some of the emotion English-to-ASL stimuli, in order to design ASL-originated emotion stimuli with bigger comprehension benefits from good-quality facial expressions. In future work, we also want to investigate how these stimuli would perform during side-by-side comparisons of animations or when evaluated via Likert-scale subjective questions. Guided by the experimental evaluation results we obtain in our studies, we will continue to improve the quality of facial expressions in ASL animations – to increase the naturalness and understandability of those animations – ultimately leading to better accessibility of online information for people who are deaf with low English literacy.

Acknowledgments. This material is based upon work supported in part by the US National Science Foundation under award number 0746556 and 1065009, by the PSC-CUNY Research Award Program, and by Visage Technologies AB through a free academic license. Jonathan Lamberton and Miriam Morrow assisted with the conduct of experimental sessions and provided valuable linguistic insights about ASL.

References

1. Allbritton, D.W., Mckoon, G., Ratcliff, R.: Reliability of prosodic cues for resolving syntactic ambiguity. *J. Exp. Psychol.-Learn. Mem. Cogn.* 22, 714–735 (1996)
2. Boulares, M., Jemni, M.: Toward an example-based machine translation from written text to ASL using virtual agent animation. In: *Proceedings of CoRR* (2012)
3. Elliott, R., Glauert, J., Kennaway, J., Marshall, I., Safar, E.: Linguistic modeling and language-processing technologies for avatar-based sign language presentation. *Univ. Access. Inf. Soc.* 6(4), 375–391 (2008)
4. Filhol, M., Delorme, M., Braffort, A.: Combining constraint-based models for sign language synthesis. In: *Proceedings of 4th Workshop on the Representation and Processing of Sign Languages, Language Resources and Evaluation Conference (LREC)*, Malta (2010)
5. Fotinea, S.E., Efthimiou, E., Caridakis, G., Karpouzis, K.: A knowledge-based sign synthesis architecture. *Univ. Access. Inf. Soc.* 6(4), 405–418 (2008)
6. Gibet, S., Courty, N., Duarte, K., Le Naour, T.: The SignCom system for data-driven animation of interactive virtual signers: methodology and evaluation. *ACM Trans. Interact. Intell. Syst.* 1(1), Article 6 (2011)
7. Holt, J.A.: Stanford achievement test - 8th edn: Reading comprehension subgroup results. *American Annals of the Deaf* 138, 172–175 (1993)
8. Huenerfauth, M.: Evaluation of a psycholinguistically motivated timing model for animations of American Sign Language. In: *The 10th International ACM SIGACCESS Conference on Computers and Accessibility (ASSETS 2008)*, Halifax, Nova Scotia, Canada (2008)
9. Huenerfauth, M., Zhao, L., Gu, E., Allbeck, J.: Evaluation of American Sign Language generation by native ASL signers. *ACM Trans. Access. Comput.* 1(1), 1–27 (2008)
10. Huenerfauth, M., Hanson, V.: Sign Language in the interface: access for deaf signers. In: Stephanidis, C. (ed.) *Universal Access Handbook*, pp. 38.1–38.18. Erlbaum, NJ (2009)
11. Huenerfauth, M., Lu, P., Rosenberg, A.: Evaluating importance of facial expression in American Sign Language and Pidgin Signed English animations. In: *Proceedings of the 13th International ACM SIGACCESS Conference on Computers and Accessibility (ASSETS 2011)*, Dundee, Scotland. ACM Press, New York (2011)
12. Neidle, C., Kegl, J., MacLaughlin, D., Bahan, B., Lee, R.G.: *The syntax of American Sign Language: functional categories & hierarchical structure*. MIT Press, Cambridge (2000)
13. Price, P., Ostendorf, M., Shattuck-Hufnagel, S., Fong, C.: The use of prosody in syntactic disambiguation. *Journal of the Acoustical Society of America* (1991)
14. Prillwitz, S., Leven, R., Zienert, H., Hanke, T., Henning, J.: An introductory guide to HamNoSys Version 2.0: Hamburg notation system for Sign Languages. In: *International Studies on Sign Language and Communication of the Deaf*. Signum, Hamburg (1989)
15. San-Segundo, R., Barra, R., Córdoba, R., D'Haro, L.F., Fernández, F., Ferreiros, J., Lucas, J.M., Macías-Guarasa, J., Montero, J.M., Pardo, J.M.: Speech to sign language translation system for Spanish. *Speech Commun.* 50(11–12), 1009–1020 (2008)
16. Schnepf, J., Wolfe, R., McDonald, J.: Synthetic corpora: A synergy of linguistics and computer animation. In: *4th Workshop on the Representation and Processing of Sign Languages, LREC 2010*, Valetta, Malta (2010)
17. Traxler, C.: The Stanford achievement test, 9th edn: National norming & performance standards for deaf & hard-of-hearing students. *J. Deaf Stud. Deaf Educ.* 5(4), 337–348 (2000)
18. Vcom3D. Homepage (2013), <http://www.vcom3d.com/>
19. Wallraven, C., Breidt, M., Cunningham, D.W., Bülthoff, H.H.: Evaluating perceptual realism of animated facial expressions. *ACM Trans. Appl. Percept.* 4(4), Article 4 (2008)