

The Development And Evaluation Of A Speech To Sign Translation System To Assist Transactions

Stephen Cox * Michael Lincoln* Melanie Nakisa†
Mark Wells‡ Marcus Tutt‡ Sanja Abbott‡

November 6, 2001

Abstract

We describe the design, development and evaluation of an experimental translation system that aims to aid transactions between a deaf person and a clerk in a Post Office. The system uses a speech recogniser to recognise speech from the Post Office clerk and then synthesises the recognised phrase in British Sign language (BSL) using a specially developed avatar. Our main objective in developing this prototype system was to determine how useful it would be to a customer whose first language was BSL and to discover what areas of the system required more research and development to make it more effective. The system was evaluated by six pre-lingually profoundly deaf people and three Post Office clerks. Deaf users and Post Office clerks were supportive of the system, but the former group required a higher quality of signing from the avatar and the latter a system that was less constrained in the phrases it could recognise: both these areas are being addressed in the next phase of the development.

1 Introduction

There has recently been considerable research activity in developing automatic systems which can understand and output speech to provide information services or to perform transactions with customers [MZ95]. Most of these systems have been developed for use over the telephone network with the goal of replacing completely or assisting a human operator [J⁺97]. A key aspect of them is that they operate in a rather restricted

*School of Information Systems, University of East Anglia, Norwich, U.K.

†Royal National Institute for Deaf People, 19–23 Featherstone Street, London EC1Y 8SL, U.K.

‡Televirtual Ltd., Anglia House, Agricultural Hall Plain, Norwich, U.K.

domain of discourse (e.g. train timetable enquiries [L⁺97], e-mail access [W⁺98], directory enquiries [YMS95]) and this gives them some robustness to the difficult problems of variability and “noise” in the language used by the speakers, the speech signal and the telephony channel. There has also been work on interactive speech-to-speech translation systems e.g. [Wah00], [RABC94], [K⁺95] [Mor93]. These systems are designed to provide translation of conversational speech between languages with a potentially very large vocabulary [Wai96]. We have been developing a system which combines aspects of both kinds of systems mentioned above. It is an interactive translation system but it operates in a very restricted domain and is designed to assist in the completion of a transaction between a Post Office (PO) clerk and a deaf customer. The system translates the clerk’s speech into British Sign Language (BSL) and displays the signs using a specially-developed avatar.

A comprehensive approach to the task of enabling humans who cannot sign to communicate using sign-language would clearly require the development of a general purpose speech to sign-language converter. This in turn requires the solution of the following problems:

1. automatic speech to text conversion (speech recognition);
2. automatic translation of arbitrary English text into a suitable representation of sign language;
3. display of this representation as a sequence of signs using computer graphics techniques.

By choosing to develop a system for use in a PO, where very many transactions are highly predictable in their scope and progress, we can sidestep or simplify many of these difficult problems by defining a limited set of phrases that can be recognised and displayed in sign-language. Although this imposes restrictions on what can be “translated”, it is still likely to form a useful system because the task is a narrow one. Our concern in developing the current system was not to attempt to solve these general problems, but to assess the utility to deaf people who use sign language of a simple translation system, and to learn how such an approach could be optimised.

The system has been developed with the collaboration of the UK Post Office, and research continues as part of the European Union fifth framework project, ViSiCAST [BCL⁺00], which aims to benefit deaf citizens by facilitating access to information and services in sign language.

2 Overview of the system

2.1 Design philosophy

Our goal was to develop a system to enable a Post Office counter-clerk to communicate with a deaf customer using automatically-generated sign-language, and hence to aid completion of a transaction. *A priori*, it might seem that recognising the clerk’s speech and displaying it as text to the deaf customer would be adequate. However, for many people who have been profoundly deaf from a young age, signing is their first language so they learn to read and write English as a second language [Con79]. As a result, many deaf people have below-average reading abilities for English text and prefer to communicate using sign language [WWGH86].

Having previously developed a prototype system (*SignAnim*, described in [BCL⁺00, PMEB99]) that used an avatar to provide signing of sub-titles for television, an avatar system was already available that could be employed to produce signs. A problem with *SignAnim*, and also for developing the system reported in this paper, was translation from text to BSL. Whereas systems to translate text from one spoken language to another are now available and work well within a restricted domain of discourse, translation from text to sign language is still a formidable research problem. BSL is a fully developed language, largely independent of English, with its own signs to express distinct concepts and with its own syntactic and semantic structures [Bri92]. These structures, inherent to sign languages, differ somewhat from those found in spoken languages, and hence translation from text to sign language requires a different approach from the techniques used in automatic translation of spoken languages.

SignAnim circumvented the translation problem by translating sub-titles into Sign Supported English (SSE) rather than BSL. SSE uses the same (or very similar) signs for words as BSL, but uses English language word order. Thus the SSE equivalent of “The

man is standing on the bridge” is MAN + STAND + ON-BRIDGE, and for “The cat jumps on the ball” it is CAT + JUMP + ONTO + BALL. SSE may therefore be regarded as more like a system for ‘encoding’ English. Linguists regard SSE as English translated into signs, and don’t consider it a language *per se*. SignAnim was an important starting point for the system described here: by by-passing many of the difficult problems of translation from English to sign language, it provided an opportunity to develop reliable sign capture methods, to determine how legible a virtual human signer could be and to develop a real-time signing “engine” that integrated the whole system.

Using pre-stored SSE “words” enables sentences to be translated into sign language at the expense of using a language that is less acceptable than BSL to deaf people. An alternative approach is to use whole phrase units rather than words. This approach is possible only if a small number of phrases is required, and these phrases can be recorded in BSL rather than SSE. If recording of the signs is done correctly, phrases can be concatenated to a certain extent e.g. amounts of money can be slotted into a carrier phrase such as “The cost is...”. Although this approach imposes considerable restrictions on the meanings that can be conveyed in BSL and hence on the dialogue, we considered that the limited nature of the transactions in the Post Office should mean that most transactions could be completed in this way. Furthermore, the philosophy of using pre-stored phrases enables the speech recognition to be implemented as a finite state network, which as has already been noted, increases the accuracy of the system (see section 2.2). It was important to see how far a BSL system using pre-stored phrases could be taken as the first step towards developing a more general system.

2.2 System components

Figure 1 is a diagram showing the structure of the system. The Post Office clerk wears a headset microphone. In early versions of the system, the clerk operated a “push-to-talk” switch when he wished to communicate with the deaf customer, but in later versions, the speech recogniser was constantly active and would respond when the clerk uttered a “legal” phrase from the grammar. The screen in front of the clerk displays a menu of topics available e.g. “Postage”, “DVLA”, “Bill Payments ”, “Passports”. Speaking any of these words invokes another screen showing a list of phrases relevant to this category

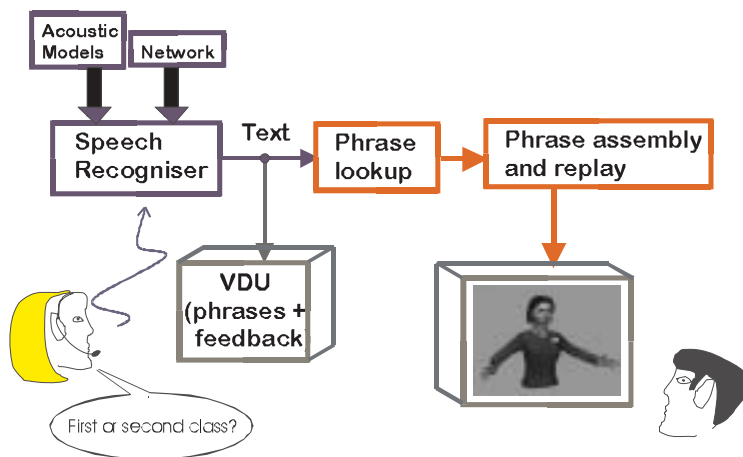


Figure 1: The Post Office translation system

which can be recognised. However, this is only an *aide-memoire* to the clerk; all phrases are active (i.e. can be recognised) at any time, so that switching between categories is seamless. In trials, we found that the clerk could remember many of the most commonly-used phrases without consulting the screen.

Prior to designing the system, we obtained transcripts of recordings of PO transactions at three locations in the UK, in all about sixteen hours of business. Inevitably, much of the dialogue transcribed was in the nature of social interaction and had little to do directly with the transaction in hand. However, analysis of these transcripts was essential for estimating the vocabulary which would be needed by our system to achieve a reasonable coverage of the most popular transactions. At the end of this analysis, a set of 115 phrases was prepared which we estimated should be adequate to cover about 90% of transactions performed. This set of phrases was changed and extended after trials with users (see section 3) and the total number of phrases currently available in the system is about 350.

2.3 Speech recognition

In the first version of the system, the speech recogniser used was the Entropic HAPI (HTK Application Interface) system [OOVW97], which incorporates the HTK (Hidden Markov Model Toolkit) recogniser [JOOW96]. This had the advantage that it allowed us to experiment with using acoustic models that had been prepared in our own laboratory using

the HTK software. The second version used the IBM ViaVoice recogniser, which offered greater range and flexibility in its network definition and in its user interface. Both speech recognisers use the same underlying system: the speech signal is first parameterised into a sequence of vectors, each of which is formed from a 20–30 ms segment of the signal and extracts important information about this segment. The recogniser has stored speech models of several thousand “triphones” (phonemes in left and right context), each model consisting of a hidden Markov model [Cox90] with a multi-variate Gaussian mixture distribution of vectors associated with each state. A network of legal phrases is supplied to the recogniser, which uses a dictionary to rewrite each word within a phrase as a sequence of triphones. Decoding of the speech signal is done using an algorithm that uses the speech models and the network supplied to output the most likely sequence of words given the acoustic input and the network (for a detailed introduction to these topics that is relevant to the operation of the ViaVoice recogniser, see [Jel97]).

An important point about the operation of the recognition system is that both the speech models and the network can be easily changed or adapted. The speech models can be adapted to the voice of each user (“speaker adaptation”), a process which takes about an hour, and the individual’s models are then stored for later use. Speaker adaptation of the models greatly increases the recognition accuracy and hence the usability of the system. The fact that the network can easily be changed means that phrases can be altered or added to the system without the need for any re-compilation.

The network constrains the speech recogniser to a finite number of predefined paths through the available vocabulary. These paths define the set of allowed phrases and consist of a start node (usually denoting silence, or background noise) followed by a number of word nodes or sub-networks, finishing with an end node (again denoting silence). Sub-networks are useful ways of defining phrase segments which can vary. For instance, a sub-network called “one2hundred” represents the legal ways of saying the integers between one and 100, and this can be inserted at any appropriate point into the network. There are other sub-networks called “amounts-of-money”, “days-of-the-week”, “countries” etc. A fragment of the network is shown in Fig 2.

The use of a finite-state network may appear to place too much constraint on what can be said by the clerk. However, it is consistent with the philosophy outlined in section

2.1 of using a limited set of pre-stored phrases for signing. Furthermore, once the clerk is familiar with the repertoire of phrases and the recogniser has been adapted to his or her voice, recognition performance is much higher than, for instance, currently available dictation packages. There are essentially two reasons for this:

1. Dictation packages are required to decode a large vocabulary and syntax and therefore use a probabilistic “bigram” language model, in which the decoding of the speech utterance is controlled by the probability of any word in the vocabulary following any other word. Restricting the vocabulary and limiting the syntax to word sequences allowed by a network lowers the number of decoding possibilities very significantly and hence increases accuracy.
2. The recogniser can be operated on a “best-match” basis, so that a phrase which is phonetically “close” but not identical with a phrase in the network will be recognised as the latter. This allows some flexibility for the speech of the clerk. (For instance, the phrase “Put that on the scales, please”, which is not present in the network, would be recognised as “Please put it on the scales”)

High recognition accuracy is very important for our system: the translation process is inherently slow because the avatar signs rather slowly to achieve maximum clarity, and any extra delay due to correcting mistakes made by the recogniser is likely to make the system unusable. Note also that, because there is no separation of speech and language decoding in this system, it does not suffer from inaccuracies in the speech decoding process being forwarded to a language translation process that is also imperfect, an effect that can make more complex systems fail to translate correctly even quite simple phrases. By using pre-stored phrases, we in effect trade flexibility and range for accuracy.

The system described here is the first stage towards a more sophisticated system which will incorporate the techniques used in “speech-understanding” systems to enable a much wider range of transactions to be completed. In our current research system, we are experimenting with using a probabilistic language model recogniser followed by a language processor that attempts to map the output from the recogniser to the correct phrase. This has the benefit of allowing the clerk complete flexibility in what he or she says to the recogniser (as long as the words used are within the 100 000 word vocabulary of the

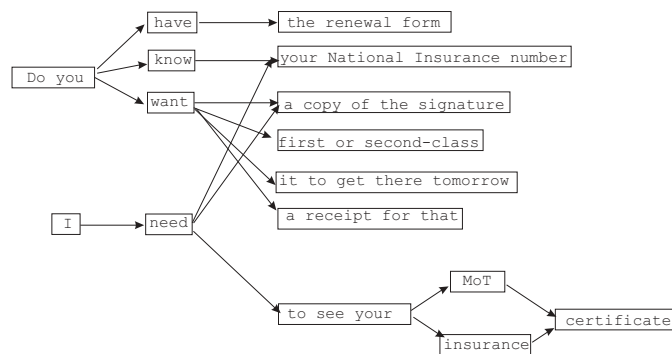


Figure 2: A section of the recognition network

recogniser) at the expense of requiring some language “understanding” to determine the correct sequence of signs to be output. At time of writing, we do not know whether this system will be less accurate than the system that uses a network. In addition, the system can obviously be adapted to translate to another spoken language (using either displayed text or speech output) as well as to sign language, and this possibility is also being explored.

2.4 System software

The system software has the task of enabling communication between the speech recognition module and the avatar module and of controlling the overall progress of a transaction. The sign assembly system is written in TCL and the recognition module incorporated as a TCL extension. The avatar module is written in C++ and communication between this and the other system components is performed using a remote procedure call system via TCP/IP socket connections.

2.5 The signing avatar, TESSA

The simplest way of signing the set of phrases defined for the application would be to store video-recordings of a person signing each phrase and concatenate the appropriate phrases in response to the output from the speech recogniser. However, we have been developing an experimental system that uses a virtual human (avatar) to sign teletext subtitles [WPT⁺99]. In this broadcast application, using an avatar has an important advantage over using video, in that the signing can be transmitted using a very small

bandwidth (only the model positions need to be transmitted at suitable intervals, rather than a full video signal). Although bandwidth is not a consideration for the Post Office system described here, an ultimate aim within the ViSiCAST project is to produce a “text-to-sign” synthesiser that will be capable of synthesising signs from a much less restricted vocabulary; to build such a system using concatenated video clips would not be viable. Another advantage of using an avatar is that different figures can be rendered onto the avatar’s frame, so that a single set of recordings of signs can be used to drive different virtual humans. Conversely, multiple human signers can be used to generate the signed content of the system while using the same avatar for the output signing, making it easy to expand and update the signed content. In addition, concatenation of signing is more fluent and controlled for avatar than for video signing, as the exact positioning of the avatar can be manipulated. For these reasons, we decided to display the signs using an avatar, TESSA, which was based on the avatar used in the SignAnim project.

Research into methods for capturing signing movements directly from video has been reported [SWP98, HH98, LH98, ATLC97]. This approach is highly desirable as it obviates the need to record signs by attaching motion sensors to a human, with the attendant problems of invasiveness, motion restriction, calibration, sensor fusion etc. Unfortunately, capture from video is not yet robust enough to record high quality motion. The alternative is to capture signs using separate sensors for the hands, body and face. This technique appears to capture sufficient movement to generate true and realistic signing from a virtual human.

The motion is captured as follows:

1. Cybergloves with 18 resistive elements for each hand are used to record finger and thumb positions relative to the hand itself.
2. *Polhemus* magnetic sensors record the wrist, upper arm, head and upper torso positions in three-dimensional space relative to a magnetic field source.
3. Facial movements are captured using a helmet mounted camera with infra-red filters and surrounded by infra-red light emitting diodes to illuminate Scotchlight reflectors stuck onto the face. Typically 18 reflectors are placed in regions of interest such as the mouth and eyebrows.

Figure 3 shows this configuration in use.



Figure 3: Data capture: face tracking camera with facial reflectors, Cybergloves for tracking the digits and Polhemus sensors taped onto the back of each hand, upper arm, body and head to track the body.

The sensors are sampled at between 30 and 60 Hz and the separate streams integrated, using interpolation where necessary, into a single raw motion-data stream that can drive the virtual human directly. The system is calibrated at the beginning of each session but, in practice, the main variation lies between signers. For example, the considerable cross-talk between glove sensors depends heavily on how tightly the gloves fit. It is particularly important to ensure good calibration at positions where fingers are supposed to just touch the thumb and where hands touch both each other and the face. These positions are important to clear signing and, to reduce computation times, there is currently no collision detection to prevent body parts sinking into each other. Where individual signs or segments are to be added to the lexicon then signs are altered manually, using a custom editor program, and the beginning and end of each sign marked to aid concatenation. The motion-data stream is displayed using a virtual human. In common with many avatars, a three-dimensional “skeleton” is driven directly from the motion-data. The

skeleton is wrapped in, and elastically attached to, a texture mapped three-dimensional polygon mesh that is controlled by a separate thread (event loop) that tracks the skeleton. We use one of the latest PC-accelerated 3D graphics cards to render the resulting 5000 polygons at 50 frames/s using Direct-X [Cor] on a Pentium class PC. Because TESSA is a full three-dimensional model, her position and pose can be changed by the user during use, an extremely valuable feature that enables users to select the optimal viewing angle and size. In addition, the identity of the virtual human can be changed. TESSA is capable of signing in real time with a refresh rate of approximately 40 frames per second.



Figure 4: Stills from the signs for the four days, Monday, Tuesday, Wednesday and Thursday

3 Evaluation

It is essential that the system conveys useful information in a way that is helpful and acceptable to deaf users. The extent to which TESSA met this aim was assessed using the following evaluation methods:

1. Evaluation of the quality of the signs;
2. Evaluation of the difficulty of performing a transaction with TESSA;
3. Questionnaires to the deaf users and Post Office clerks.

The outcomes of these trials are reported in this section.

3.1 Participants

Six pre-lingually profoundly deaf people whose first language is BSL took part in the evaluations of the system. They were recruited through the deaf-UK e-mail newsgroup or through local UK Royal National Institute for Deaf People (RNID) offices and were paid for their participation. Three clerks were recruited by the Post Office to take part in the evaluations: each had over ten years experience as a clerk and had had experience of serving deaf customers.

3.2 Protocol

The evaluations took place over three sets of two days. Two deaf people and one clerk attended for each pair of days. The first day started with completion of the first part of a questionnaire. Each deaf participant then alternated between identifying a block of signed phrases and attempting a block of staged transactions. At the end of the second day, all participants completed the remainder of the questionnaire and gave any general feedback. BSL/English interpreters were present throughout.

3.3 Quality of signing

The quality of TESSA’s signing was measured in two ways: intelligibility of signs, and acceptability of signs to deaf users. The first of these measurements is an objective one and is clearly important in establishing a baseline for the current system against which future avatars may be evaluated. However, it is well-known that intelligibility on its own is inadequate for assessment of these systems: for instance, synthetic speech can sound fully intelligible but be disliked by users [Joh96]. Hence we also included a subjective measurement of “acceptability” of signing.

3.3.1 Intelligibility

The deaf participants were presented with each signed phrase and asked to write down what they understood. From the 115 distinct phrases, 133 phrases were generated by

incorporating days of the week and numbers to ensure that each day and each number (units and tens) was presented at least once. Signed phrases were presented on the screen without text. The deaf participant controlled presentation of each phrase and was allowed to repeat each phrase up to a maximum of five presentations. Phrases were presented in blocks of between 20 and 24, in groups according to broad categories, for example, postage, bill payment, amounts of money. Accuracy of identification of phrases was assessed in two ways:

1. By the accuracy of identification of complete phrases.
2. By the accuracy of approximate “semantic sign units” within the phrase. For example, the phrase “It should arrive by Tuesday but it’s not guaranteed” requires five sign units, so “should arrive Tuesday not guaranteed” would score 100% and “should arrive Tuesday” 66%.

The 133 phrases gave a total of 444 sign units. While these units were not all distinct (for example, the sign for “pound” was presented several times), identification of each presentation of a unit was scored separately. One experimenter (the third author) judged the accuracy of responses for both measures on the basis of written responses from each deaf participant. Once each phrase had been scored for accuracy of identification, each deaf person was re-presented with each phrase not identified correctly along with the text of the intended phrase. With an interpreter and experimenter, they were asked to indicate whether the signs were considered inappropriate or whether they were just not clear. Any signs considered inappropriate were not necessarily wrong; rather they may have represented different regional variations in sign to those used by the deaf participant¹.

The average number of times each phrase was presented before an attempt at identification was made was 1.8. Attempts at identification were made after one presentation for the majority of phrases (51%) and required more than two presentations for 20% of phrases. The average accuracy of identification of complete phrases was 61% and ranged

¹Variation in signs is a more difficult problem to contend with than variations in accent or dialect in spoken languages, as hearing people can use a standard written language as a reference, which is not available to those who communicate using only signs [KW85]

| Acceptability rating | | % of phrases |
|----------------------|---|--------------|
| High | 3 | 20.2 |
| | 2 | 43.2 |
| Low | 1 | 36.6 |

Table 1: The percentage of phrases rated in each category of acceptability.

from 42% to 70% across deaf participants (Figure 5a). For the identification of sign units in phrases, average accuracy was 81% and ranged from 67% to 89% (Figure 5b). Subsequent analysis of the sign units which were wrongly identified indicated that on

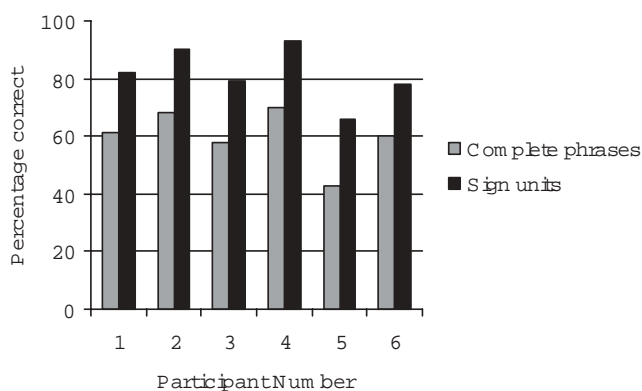


Figure 5: Average percentage recognition scores achieved by each signer for complete phrases and sign units within phrases.

average 30% of errors (6% of all sign units) were due to signs considered inappropriate and the remaining 70% (13% of all sign units) were due to unclear signing.

3.3.2 Acceptability

Participants were asked to rate how acceptable the phrase was as an example of BSL (on a 3-point scale from 1–“Low” to 3–“High”). Table 1 shows the percentage of phrases which were rated in each category of acceptability. The average acceptability rating was 2.2 and ranged from 1.7 to 2.8.

3.3.3 Discussion

Accuracy of identification of the signed phrases was 61% for complete phrases and 81% for sign units, with quite a wide range in accuracy across deaf participants (ranges of 28% and 20%, respectively). This range in accuracy suggests it is important to use many sign-language users for a true assessment of signed content of these systems. In future, it may be more appropriate to use more than six deaf people from a range of UK regions to assess sign quality.

The majority of identification errors (70%) were due to signs being unclear rather than due to inappropriate signs. The percentage of errors for inappropriate signs did not differ greatly between subjects, with personal averages ranging from 4.7% to 6.6%. This pattern might suggest that the same signs were considered inappropriate by all deaf participants. However, inspection of the pattern of errors across deaf participants for each phrase indicated that this was not necessarily the case. Of the 46 phrases where one or more sign was considered inappropriate by any deaf participant, in 34 of these (74%) a sign was considered inappropriate by no more than two of the deaf participants. This result suggests that regional variations or differences in personal signing style may have played a role in phrase intelligibility.

Ratings of acceptability were also given across the scale with 20% of phrases rated as highly acceptable and 63% in one of the top two categories, indicating that there is scope for improving the quality of the avatar’s signing.

3.4 Transactions

Staged Post Office transactions were used to compare completion times and ease and acceptability of communication with and without TESSA. Each deaf participant attempted 30 transactions with a single Post Office clerk. Transactions were selected by the Post Office as those achievable with the phrases available. There were 18 distinct transactions, 6 were denoted “simple”, 6 “average difficulty” and 6 “complex”. The average difficulty and complex transactions were attempted twice by each deaf participant/clerk pair, once with an open counter and once behind a fortified counter where a transparent screen separates clerk and customer. Use of different counter styles did not appear to affect

performance hence results are not reported separately here.

Half of all transactions were attempted with TESSA and half without. The phrases presented with or without TESSA were counter-balanced between deaf participants. Practice transactions were performed with TESSA at the start of each session so that the clerk, deaf participant and interpreter could get used to using TESSA and the format of the evaluation. Transactions were performed in blocks of six, three with TESSA and three without. The approximate time taken to successfully complete each transaction was recorded. On completion of each transaction, both deaf participants and clerks were asked to rate each transaction for acceptability on a 3-point scale from 1–“Low” to 3–“High”.

3.4.1 Timings

On average, transactions took longer to complete with TESSA than without [$F(1,178)=61.2, p < 0.001$] (Figure 6). Average times for transactions were 57s without TESSA and 112s with TESSA. On average, communication in transactions completed

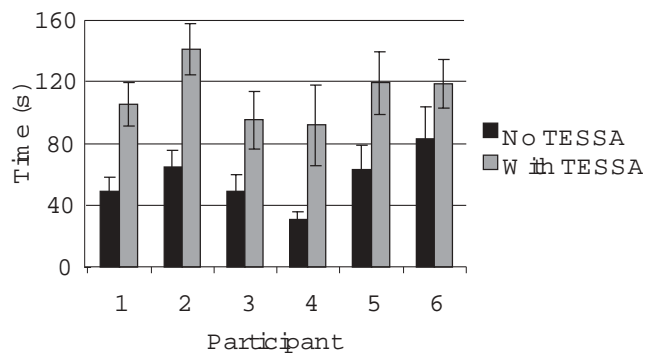


Figure 6: Average times of transactions without TESSA (dark bars) and with TESSA (light bars) for each deaf participant. Error bars show the 95% confidence intervals of the means.

with TESSA was rated as less acceptable than in transactions completed without TESSA [$U(1,178)=6025, p < 0.001$] (Figure 7). On the 3-point scale (from 1–“Low” to 3–“High”) average ratings were 1.9 with TESSA and 2.6 without.

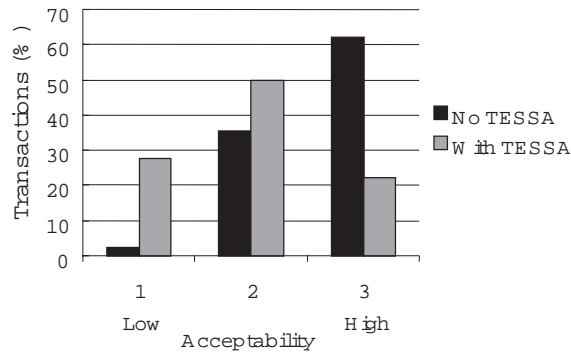


Figure 7: Percentage of transactions rated by the deaf participants in each category of acceptability on a 3-point scale from 1-“Low” to 3- “High” , without TESSA (dark bars) and with TESSA (light bars).

Clerks rated acceptability of transactions completed with TESSA as slightly lower than transactions completed without TESSA. On the 3-point scale average ratings were 2.5 with TESSA and 2.6 without (Figure 8).

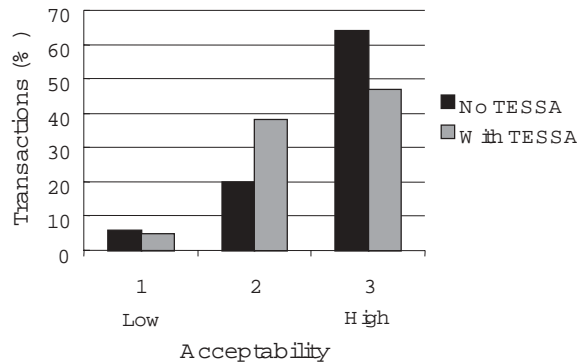


Figure 8: Percentage of transactions rated by the clerks in each category of acceptability on a 3-point scale from 1-“Low” to 3- “High”, without TESSA (dark bars) and with TESSA (light bars)

3.4.2 Discussion

Compared to transactions without TESSA, transactions performed with TESSA took on average nearly twice as long to complete, and the deaf participants, and to a lesser

extent the clerks, rated communication as less acceptable. The main reason most likely to have contributed to these effects was the somewhat disjointed communication with TESSA. As expected, it took the clerks some time to learn which phrases were available and to locate the phrase they wanted so they could read it out word for word. The clerks had only about an hour of practice using the system before the trials. These difficulties should decrease substantially with training and experience on the system. Moreover, the next version of the system, which will incorporate some speech “understanding”, will not require phrases to be repeated verbatim.

Additional factors may have contributed to the longer transaction times and poorer ratings with TESSA:

1. The list of phrases were selected for use in the system as those most commonly used in the PO. These phrases also tended to be those used for the more simple Post Office transactions, for example, buying stamps, cashing a cheque or claiming a pension payment. Hence the transactions used in this evaluation, limited by the phrases available, also tended to be fairly simple or were simplified. This was confirmed by the Post Office staff who selected the transactions and the clerks who said they would usually ask more questions for specific transactions but these were not available in TESSA. The transactions used in the trials therefore tended to represent situations in which communication was fairly easy without TESSA.
2. The deaf participants were all fairly good communicators and all had reasonable written skills. Hence they were able to complete the simple transactions, by lip-reading/speaking and writing notes or asking the clerk to write things down where necessary. This is a consequence of the type of people who would be prepared to attend two days of testing away from their home town, the recruitment process (through e-mail and professional connections) and also the necessary use of textphone, fax and e-mail for the logistics of arranging the trials.
3. The clerks either were “deaf aware” or soon became deaf aware as a result of spending two days with the profoundly-deaf participants. Communication without TESSA was fairly easy as they used good eye contact, spoke clearly and were prepared to write things down if they were not understood.

| Question | Very easy | Fairly easy | Manage-able | Slightly difficult | Very difficult |
|---|-----------|-------------|-------------|--------------------|----------------|
| How easy do you usually find communication in the Post Office? | 2 | | 3 4 6 | | 1 5 |
| How easy did you find communication using TESSA? | | | 1 2 3 4 | 6 | 5 |
| In everyday life, how easy do you think communication would be using TESSA? | | 1 2 | 6 | 4 | 3 5 |

Table 2: Responses made by each deaf participant to the three questions about ease of communication in the PO: previously, in the trials with TESSA and anticipated in everyday life with TESSA. Each number represents the responses from one deaf person.

- There was a delay of a few seconds between recognition of the spoken phrase and the signing of the phrase. Not only did this absolute delay add to overall transaction time but the delay often resulted in loss of attention and the need for the sign to be repeated or the clerk to repeat the phrase.

3.5 Questionnaires

Questionnaires to both deaf participants and clerks were used to obtain subjective views of previous experiences of communication in the PO, and how these experiences differed in the trials and were anticipated to differ in real life using TESSA.

Results from the three questions asking about ease of communication in the PO, previously, in the trials with TESSA and anticipated in everyday life with TESSA are shown for each of the six deaf participants in Table 2. Table 3 shows the responses of deaf participants when asked about how much communication bothers them in the PO, previously, and anticipated with TESSA in everyday life. When asked for a preference, four deaf participants said they would prefer to communicate without TESSA and two preferred with, as an option if needed. All clerks said communication was “Slightly easier” or “Much easier” with TESSA than without, and that in everyday life they anticipated that communication would be “Much easier” with TESSA (Table 5).

Tables 4 and 5, respectively, show responses made by each clerk when asked about communicating with deaf people in the PO, previously, with TESSA in the trials and with

| Question | Very much | Quite a lot | Some | A little | Not at all |
|--|-----------|-------------|------|----------|------------|
| In everyday life, how much does communication in the Post Office upset, annoy or worry you? | 1 | 4 | 3 | 6 | 2 5 |
| In everyday life, how much would communication using TESSA in the Post Office upset, annoy or worry you? | | 5 | | 6 | 1 2 3 4 |

Table 3: Responses made by each deaf participant to the two questions about how much communication in the Post Office bothered them, previously and anticipated with TESSA in everyday life. Each number represents the responses from one deaf person.

| Question | Very easy | Fairly easy | Manageable | Slightly difficult | Very difficult |
|---|-----------|-------------|------------|--------------------|----------------|
| How easy do you usually find communication with deaf customers? | | 1 2 3 | | | |
| How easy did you find communication using TESSA? | 3 | 1 2 | | | |
| In everyday life, how easy do you think communication would be using TESSA? | 2 3 | 1 | | | |

Table 4: Responses made by each clerk to the three questions about ease of communication with deaf customers: previously, in the trials with TESSA and anticipated in everyday life with TESSA. Each number represents the responses from one clerk.

TESSA in everyday life, and the relative ease of communication with and without TESSA, in the trials and anticipated for everyday life. All clerks said that they would prefer to have TESSA available as an option to use when communication became difficult, even though they all thought transactions would take “Slightly longer” with TESSA.

4 Discussion of evaluation results

Of the six deaf participants, one person said that communication would be easier with TESSA and two people said they would prefer to communicate with TESSA in the PO, as an option in case communication became difficult. The three deaf participants who said

| Question | Much easier | Slightly easier | No difference | Slightly worse | Much worse |
|--|-------------|-----------------|---------------|----------------|------------|
| Compared to communication without, do you think TESSA made communication: | 3 | 1 2 | | | |
| In everyday life, do you think that using TESSA in the Post Office would make communication: | 1 2 3 | | | | |

Table 5: Responses made by each clerk to the two questions about comparing communication in the Post Office with and without TESSA, in the trials and anticipated in everyday life. Each number represents the responses from one clerk.

that communication in the PO usually upset or worried them, said they thought using TESSA in the PO would not bother them at all. While this represents positive feedback from some deaf participants, the fact that these responses were not more generally positive does not seem unreasonable at this stage in the life-cycle of the project. These questions were asked about the first version of TESSA to be evaluated by deaf people, and on the basis of use during the trials by clerks with little previous experience of using the system, where communication with TESSA was somewhat lengthy and disjointed. Scores on the visual analogue scales showed a wide range of responses between deaf people for ratings of clarity of signing, acceptability of the avatar as a signer of BSL and the avatar’s appearance. These scales proved easy to use as the deaf participants responded more easily than to the previous questions for which it was often difficult to obtain a categorical response. These scales therefore are likely to be useful outcome measures for evaluating future signing avatars and versions of the system. Scores on the scales were all under 65, hence there is much scope for improvement.

The deaf participants provided much constructive feedback about how TESSA could be improved. Their main points were:

- Facial expressions need to be improved.
- Clearer handshapes, finger configurations and lip patterns are required, especially for numbers and finger-spelling.

- The delay between the end of the spoken phrase and the beginning of signing needs to be reduced.
- The appearance of the avatar needs to improve. In particular, a clearer distinction should be made between the face and hands and the clothing, which should be plain.
- All deaf participants said they would prefer to see both BSL and text rather than just BSL or just text. They also thought that SSE should be available as an option.

When asked to comment on the use of avatars for signing in general, all deaf participants thought that avatars would be most useful for more complex communication needs, e.g. explaining forms to claim social benefits.

All clerks said they would prefer to have the system available as they thought it would make communication with deaf customers easier and more effective. Use of the system for multiple spoken languages and with text sub-titles would ensure more frequent use and hence greater likelihood that the system would be used with deaf people. The clerks also commented that they would like more phrases and an unconstrained speech system, where phrases need not be spoken verbatim.

4.1 General comments and future work

Our goal in developing this trial system was to establish whether the introduction of a limited speech-to-sign translation system for the PO counter clerk would be beneficial to deaf users whose primary means of communication was sign language. Although some of the feedback from the evaluation was critical, we are encouraged by the following points:

- One of the deaf participants said that communication in the Post Office would be easier with TESSA and some said they would prefer to have TESSA available in the Post Office for use when communication became difficult.
- The three deaf participants who said that communication in the Post Office usually upset or worried them said they thought using TESSA in the Post Office would not bother them at all.

- Feedback from the Post Office clerks was generally very positive, despite the very limited time they had to train with the system.

These evaluations, although limited in extent, have indicated that there is much scope for improvement of TESSA, have given some insight into how these improvements could be achieved and provided baseline outcome measures against which improvements can be assessed. The majority of aspects identified for improvement are planned for further development within the ViSiCAST project. Primarily, the development of an unconstrained version, where phrases need not be repeated word for word, will enable much more natural communication and should greatly reduce the time taken for transactions, so is also likely to be more acceptable to both deaf customers and clerks. Other aspects to be explored include research into facial modelling, which should improve avatar facial expressions and lip patterns. New data gloves are also being used to improve recording of finger movements and handshapes. New models of the avatar and clothing will also take account of the comments made by the deaf participants. Less formal evaluations are planned within the deaf community to assess the views of more deaf people and further formal evaluations will continue through the lifetime of the ViSiCAST project.

In tandem with these developments, the ViSiCAST project is also doing basic research into the general problem of converting arbitrary English text into a representation of sign language [SM01], and developing a synthetic avatar that can sign these representations without the need for motion capture [Ken01]. These will feed into the application described here to increase its flexibility and sophistication. The problem of two-way communication is also being addressed by research into sign-language recognition.

Acknowledgments

The development of TESSA was sponsored by the UK Post Office and the evaluation by the European Fifth Framework Programme under the auspices of the ViSiCAST project. We would like to thank the deaf participants, PO clerks and interpreters for their help in the evaluations.

References

- [ATLC97] T. Ahmad, C.J. Taylor, A. Lanitis, and T.F. Cootes. Tracking and recognising hand gestures, using statistical shape models. *Image and Vision Computing*, 15(5):345 – 352, May 1997.
- [BCL⁺00] J.A. Bangham, S.J. Cox, M. Lincoln, I. Marshall, M. Tutt, and M. Wells. Signing for the deaf using virtual humans. In *IEE Colloquium on Speech and language processing for disabled and elderly people*, April 2000.
- [Bri92] D. Brien. *Dictionary of British Sign Language / English*. Faber and Faber, 1992.
- [Con79] R. Conrad. *The deaf school child*. Harper and Row, 1979.
- [Cor] Microsoft Corporation. Directx home page. <http://www.microsoft.com/directx>.
- [Cox90] S.J. Cox. Hidden Markov Models for automatic speech recognition: theory and application. In C Wheddon and R Linggard, editors, *Speech and Language Processing*, pages 209–230. Chapman and Hall, 1990.
- [HH98] C.L. Huang and W.Y. Huang. Sign language recognition using model-based tracking and a 3d hopfield neural network. *Machine vision and applications*, 10(5-6):292 – 307, Apr 1998.
- [J⁺97] R.D. Johnston et al. Current and experimental applications of speech technology for Telecom services in Europe. *Speech Communication*, 23(1–2):5–16, 1997.
- [Jel97] F Jelinek. *Statistical Methods for Speech recognition*. The MIT Press, 1997.
- [Joh96] R.D. Johnston. Beyond intelligibility—the performance of text-to-speech synthesisers. *BT Technology Journal*, 14(1):100–110, January 1996.
- [JOOW96] J. Jansen, J. Odell, D. Ollason, and P. Woodland. *The HTK book*. Entropic Research Laboratories Inc., 1996.

- [K⁺95] M.W. Koo et al. KT-STTS: A speech translation system for hotel reservation and a continuous speech recognition system for speech translation. In *Proc. of 4th European Conf. on Speech Communication and Technology*, pages 1227–1230, September 1995.
- [Ken01] J.R. Kennaway. Synthetic animation of deaf signing gesture. In I. Wachsmuth, editor, *4th International Workshop on Gesture and Sign Language Based Human-Computer Interaction*, Springer series Lecture Notes in Artificial Intelligence. Springer Verlag, 2001. To be published.
- [KW85] J. Kyle and B. Woll. *Sign Language: The Study of Deaf People and their Language*. Cambridge University Press, 1985.
- [L⁺97] L.F. Lamel et al. The LIMSI RailTel system: field trial of a telephone service for rail travel information. *Speech Communication*, 23(1–2):67–82, 1997.
- [LH98] C. C. Lien and C. L. Huang. Model-based articulated hand motion tracking for gesture recognition. *Image and Vision Computing*, 16(2):121 – 134, Feb 1998.
- [Mor93] T Morimoto. ATR’s speech translation system: ASURA. In *Proc. 3rd European Conf. on Speech Communication and Technology*, pages 1291–1294, September 1993.
- [MZ95] B. Mazor and B. L. Zeigler. The design of speech-interactive dialogs for transaction- automation systems. *Speech Communication*, 17:313–320, November 1995.
- [OOVW97] J. Odell, D. Ollason, V. Valtchev, and D. Whitehouse. *The HAPI Book*. Entropic Cambridge Research Laboratory, 1997.
- [PMEB99] F. Pezeshkpour, I. Marshall, R. Elliott, and J.A. Bangham. Developing of a legible deaf signing virtual human. In *IEEE Multimedia Systems Conference ’99 (IEEE ICMCS’99)*, pages 333–338, June 1999.

- [RABC94] M Rayner, H Alshawi, I Bretan, and D Carter. A speech to speech translation system built from standard components. In *Proc. ARPA Human Language Technology Workshop '93*, pages 217–222, Princeton, NJ, 1994.
- [SM01] E Safar and I Marshall. The architecture of an English-text-to-Sign-Languages translation system. In G Angleova et al., editors, *Recent Advances in Natural Language Processing (RANLP)*, pages 223–228, September 2001.
- [SWP98] T. Starner, J. Weaver, and A. Pentland. Real-time american sign language recognition using desk and wearable computer based video. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(12):1371 – 1375, Dec 1998.
- [W+98] P.J Wyard et al. Spoken language systems—beyond prompt and response. In F.A. Westall, R.D. Johnston, and A.V. Lewis, editors, *Speech Technology for Communications*, pages 487–520. Chapman and Hall, 1998.
- [Wah00] W Wahlster, editor. *Verbmobil: Foundations of Speech to Speech Translation*. Springer-Verlag, 2000.
- [Wai96] A. Waibel. Interactive translation of conversational speech. *Computer*, 29(7), 1996.
- [WPT+99] M. Wells, F. Pezeshkpour, M. Tutt, J.A. Bangham, and I. Marshall. Simon - an innovative approach to deaf signing on television. In *Proc. International Broadcasting Convention*, pages 477–482, September 1999.
- [WWGH86] D. Wood, H. Wood, A. Griffiths, and I. Howarth. *Teaching and talking with deaf children*. John Wiley and Sons, 1986.
- [YMS95] O. Yoshioka, Y. Minami, and K. Shikano. A speech dialogue system with multimodal interface for telephone directory assistance. *IEICE Transactions on Information and Systems*, E78D(6):616–621, 1995.