



Exploring automatic text-to-sign translation in a healthcare setting

Lyke Esselink¹ · Floris Roelofsen¹ · Jakub Dotlačil² · Shani Mende-Gillings³ · Maartje de Meulder^{4,5} · Nienke Sijm⁶ · Anika Smeijers⁷

Accepted: 21 August 2023
© The Author(s) 2023

Abstract

Communication between healthcare professionals and deaf patients has been particularly challenging during the COVID-19 pandemic. We have explored the possibility to automatically translate phrases that are frequently used in the diagnosis and treatment of hospital patients, in particular phrases related to COVID-19, from Dutch or English to Dutch Sign Language (NGT). The prototype system we developed displays translations either by means of pre-recorded videos featuring a deaf human signer (for a limited number of sentences) or by means of animations featuring a computer-generated signing avatar (for a larger, though still restricted number of sentences). We evaluated the comprehensibility of the signing avatar, as compared to the human signer. We found that, while individual signs are recognized correctly when signed by the avatar almost as frequently as when signed by a human, sentence comprehension rates and clarity scores for the avatar are substantially lower than for the human signer. We identify a number of concrete limitations of the JASigning avatar engine that underlies our system. Namely, the engine currently does not offer sufficient control over mouth shapes, the relative speed and intensity of signs in a sentence (prosody), and transitions between signs. These limitations need to be overcome in future work for the engine to become usable in practice.

Keywords Access to healthcare information · Sign language · Avatar technology · User study

1 Introduction

Communication between healthcare professionals and deaf patients is challenging enough under normal circumstances [21], but has been especially difficult during the COVID-19 pandemic [36]. Most healthcare professionals do not know

the national sign language, COVID-19 regulations often did not permit sign language interpreters to enter hospitals and clinics, interpreting via video relay is not always viable, and face masks conceal facial expressions and make lipreading impossible [26].

A survey among 179 deaf people in the Netherlands, carried out by one of the authors of the present article in

L. Esselink and F. Roelofsen contributed equally to this work.

✉ Lyke Esselink
l.d.esselink@uva.nl

Floris Roelofsen
f.roelofsen@uva.nl

Jakub Dotlačil
j.dotlacil@gmail.com

Shani Mende-Gillings
semg98@gmail.com

Maartje de Meulder
maartje.demeulder@hu.nl

Nienke Sijm
nienke.sijm@hu.nl

Anika Smeijers
a.s.smeijers@amsterdamumc.nl

¹ Institute for Logic, Language and Computation, University of Amsterdam, Amsterdam, The Netherlands

² Department of Languages, Literature and Communication, Utrecht University, Utrecht, The Netherlands

³ University of Amsterdam, Amsterdam, The Netherlands

⁴ Faculty of Healthy and Sustainable Living, HU University of Applied Sciences, Utrecht, The Netherlands

⁵ Department of Languages and Intercultural Studies, Heriot-Watt University, Edinburgh, Scotland

⁶ Sign Language & Deaf Studies, HU University of Applied Sciences, Utrecht, The Netherlands

⁷ Amsterdam University Medical Centre, Amsterdam, The Netherlands

January–March 2021, confirmed that the general inability of healthcare professionals to communicate in Dutch Sign Language (Nederlandse Gebarentaal, NGT) was perceived as a very significant threat [48]. For instance, 88% of participants stated that they were worried about communication barriers should they need to be hospitalized with COVID-19; while, for comparison, only 33% stated that they were worried about the fact that friends and relatives would not be allowed to visit them in the hospital.

To address these concerns, we have explored the possibility to automatically translate phrases that are frequently used in the diagnosis and treatment of hospital patients, in particular phrases related to COVID-19, from Dutch or English to NGT. We developed a prototype system which displays translations either by means of pre-recorded videos featuring a deaf human signer (for a limited number of sentences) or by means of animations featuring a computer-generated signing avatar (for a larger, though still restricted number of sentences). We evaluated the comprehensibility of the signing avatar, as compared to the videos of a human signer.

We have concentrated on Dutch and English as the source languages for translation and NGT as the target sign language. The general problem we aim to address, however, is not specific to NGT but manifests itself for other sign languages as well. Therefore, we have aimed to design our prototype system in such a way that it could in principle be extended to include other source and target languages in a relatively straightforward way. In this respect, our system diverges from some existing text-to-sign translation systems, which are tailor-made for a specific target sign language and not easily portable to other languages (see Sect. 3.2 below). In particular, to our knowledge, none of the existing systems allows for translation from Dutch to NGT.

We should emphasize that a qualified human sign language interpreter is, whenever available, always to be preferred over a machine translation system, keeping in mind that even the use of sign language interpreters still has its own limitations [7]. We believe that it is worth investigating the extent to which a machine translation system can be of help in situations in which a human interpreter cannot be employed, including in certain medical settings where effective, instantaneous communication between healthcare professionals and patients can be of critical importance. But the aim of such technology should never be to replace human sign language interpreters across the board, and interpreting services of the highest possible quality remain critical in all domains.

The type of research reported here requires collaboration between researchers from several disciplines, bringing in different kinds of positionalities, knowledge and expertise. Before proceeding, we therefore include a brief note about the members of our research team and their respective contributions to the present project. De Meulder and

Sijm are deaf; Esselink, Roelofsen and Smeijers are hearing new signers, with varying levels of proficiency in NGT; Dotlačil is hearing and not proficient in NGT. De Meulder is a scholar in Deaf Studies and applied linguistics, Sijm has a background in criminology and Deaf Studies. Esselink and Mende-Gillings have a background in Artificial Intelligence, Smeijers is a sign linguist and a medical doctor, Roelofsen has a background in linguistics and Artificial Intelligence, and Dotlačil contributed his expertise in statistical analysis. Esselink, Mende-Gillings, and Roelofsen designed and implemented the prototype system. Smeijers contributed her knowledge of the medical domain, and took main responsibility for the production of video translations used in the prototype system. The evaluation study was designed by Esselink, de Meulder, Roelofsen and Sijm, and was executed by Esselink and Roelofsen. The data from the study was analyzed by Esselink, Roelofsen and Dotlačil.

The article is organized as follows. Section 2 provides relevant background information on sign languages and deaf communities, Sect. 3 discusses the prototype system we developed, Sect. 4 reports on the evaluation study, and Sect. 5 concludes.¹

2 Brief background on sign languages

Evidently, we cannot provide a comprehensive overview here of the linguistic properties of sign languages in general (see, e.g., [1]), nor of NGT in particular (see [33]). We will, however, highlight some important features which any text-to-sign translation system needs to take into account.

First of all, sign languages have naturally evolved in deaf communities around the world [35]. This means that, contrary to a rather common misconception, there is not a single, universal sign language used by all deaf people worldwide, but many different sign languages used on different scales by different deaf and hearing signers [28].

Second, although sign languages exist in language ecologies in close contact with spoken languages, there is generally no direct correspondence between the sign language used in a given country and the spoken language used in that same country. For instance, while English is the mainstream spoken language both in the US and in the UK, American Sign Language (ASL) and British Sign Language (BSL) differ considerably from each other, as well as from spoken English. Such differences do not only pertain to the lexicon, but also to grammatical features such as word order. This means in particular that to translate a sentence from English to ASL or BSL it does not suffice to translate every word in

¹ A preliminary report on the prototype we developed was published as Roelofsen et al. [45].

the sentence into the corresponding sign in ASL/BSL and then put these signs together in the same order as the words in the English sentence.

Third, for healthcare professionals to communicate exclusively through written text would not be satisfactory for most deaf patients. Deaf people have varying levels of access to auditory information, due to variance in hearing loss and differential access to education. Most deaf people, while they have developed skills in visual/tactile communication have no, reduced, or contextual sensory access to spoken languages. Contrary to popular belief, not all deaf people can lipread. In any case, lipreading is always mostly guesswork, where context is paramount. This can be a problem in medical settings with the use of medical jargon and words that are harder to anticipate. Moreover, health literacy has proven to be a barrier for many deaf patients [37, 39, 49]. In a medical setting it is critical to avoid miscommunication, to obtain reliable informed consent for interventions, and to foster an environment in which patients feel maximally safe. Relying exclusively on written text and lipreading will not achieve this.

Fourth, signs are generally not just articulated with the hands, but often also involve facial expressions and/or movements of the head, mouth, shoulders, or upper body. These are referred to as the *non-manual* components of a sign. A text-to-sign translation system has to take both manual and non-manual components of signs into account.

Fifth, related to the previous point, non-manual elements are not only part of the *lexical* make-up of many signs, but are also often used to convey certain *grammatical* information (comparable to intonation in spoken languages). For instance, raised eyebrows may indicate that a given sentence is a question rather than a statement, and a head shake often expresses negation. Such non-manual grammatical markers are typically ‘supra-segmental’, meaning that they do not co-occur with a single lexical sign but rather span across a sequence of signs in a sentence. Sign language linguists use so-called *glosses* to represent sign language utterances. For instance, the gloss in (1) represents the NGT translation of the question *Have you already eaten?*.

- (1) $\frac{\text{brow raise}}{\text{YOU EAT ALREADY}}$

Lexical signs are written in small-caps. They always involve a manual component and often non-manual components as well. The upper tier shows non-manual grammatical markers, and the horizontal line indicates the duration of these non-manual markers. In this case, ‘brow raise’ is used to indicate that the utterance is a question. A text-to-sign translation system should thus be able to integrate non-manual elements that convey grammatical information with manual and non-manual elements that belong to the lexical specification of the signs in a given sentence [52]. This means that

a system which translates sentences word by word, even if it re-orders the corresponding signs in accordance with the word order rules of the target sign language, will not be fully satisfactory. More flexibility is needed: word by word translation can be a first step, but the corresponding signs as specified in the lexicon must generally be adapted when forming part of a sentence to incorporate non-manual markers carrying grammatical information.

3 A modular text-to-sign translation system

We have developed a prototype system which displays sign language translations either by means of pre-recorded videos featuring a deaf human signer, or by means of animations featuring a signing avatar. While video translations are clearly expected to be of higher quality, a translation system solely based on video translations would not scale up. With signing avatars, it may in principle be possible to build a system with much more comprehensive coverage. But how should appropriate avatar-based translations be generated, given the specific domain requirements? And, are such translations comprehensible for the target end users, i.e., a widely varied group of deaf people? These questions have been the focus of our investigation. In what follows, we will therefore not say much about the video-based component of the system but concentrate mainly on the avatar-based component.

3.1 Sign synthesis

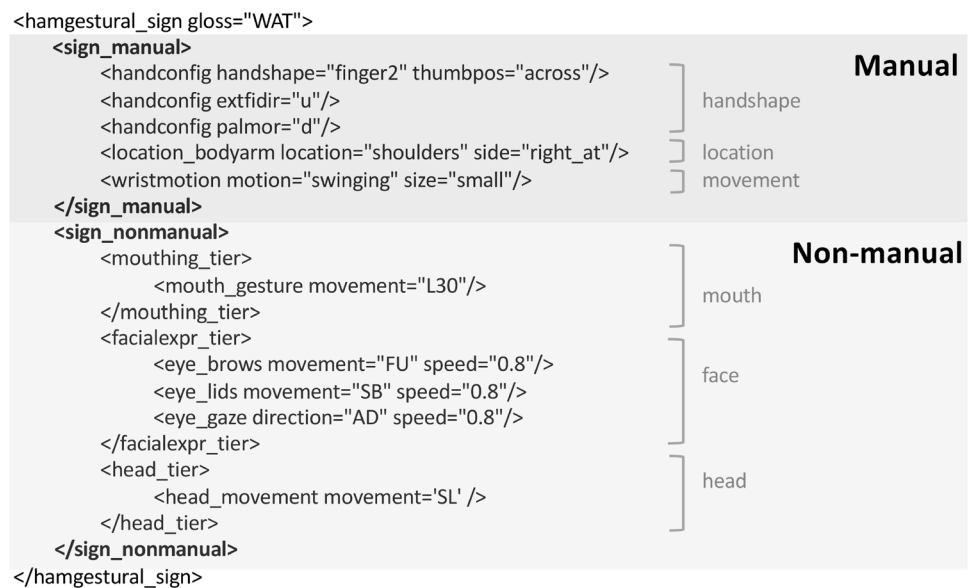
A crucial prerequisite for scalable automated text-to-sign translation is sign *synthesis*: the ability to create animations featuring a signing avatar. Broadly speaking, there are three ways in which this can be achieved: animation based on motion capture, keyframe animation, and scripted animation.

Motion capture makes it possible to obtain a library of high-quality animations for lexical signs, but requires expensive equipment and typically involves a lot of manual post-processing of the original data. Another major challenge under this approach is to modify the animations for lexical signs so as to incorporate non-manual grammatical markers [5]. Although it would in principle be possible to layer manual lexical animations with separate facial animations, exploring this option was not feasible within the time-frame of the present project and must be left for future work.

Keyframe animation results in lower quality lexical sign animations than motion capture. It does not require expensive equipment, but involves a lot of manual labor. Like motion capture, the problem of how to incorporate grammatical non-manual markers also applies to libraries of lexical signs obtained by means of keyframe animation.

The third synthesis method, scripted animation, offers a promising strategy to overcome this problem. On this

Fig. 1 SiGML encoding of the NGT sign WAT ('what')



approach, rather than directly animating each lexical sign, animations of lexical signs are generated *procedurally* based on structured specifications of the phonetic properties of these signs [11]. As in the case of keyframe animation, this also results in lower quality animations than could be obtained with motion capture techniques. However, no expensive equipment is needed, and relatively little manual labor is required. The phonetic properties that make up the required specifications include (but are not limited to) the initial location, shape and orientation of the hands, possibly movements of the hands and other body parts, and facial expressions. Several formalisms have been developed to specify the phonetic properties of signs in a structured, computer-readable fashion (see [5] for an overview). Arguably the most extensively developed and most widely used formalism is the Sign Gesture Markup Language (SiGML) [11, 25], which is based on the HamNoSys notation originally developed for the annotation of sign language corpora [27, 40]. For illustration, our SiGML encoding of the NGT sign WHAT is given in Fig. 1. As can be seen in the figure, both manual features (handshape, location, movement) and non-manual features (mouth, face, head) are encoded.

SiGML specifications can be converted into animations by the JASigning avatar engine [11, 30, 32]. This approach makes it possible, in principle, to integrate non-manual grammatical markers with the lexical signs that make up a sentence, although such functionality has not yet been thoroughly implemented in systems based on SiGML and JASigning to our knowledge.

Given these considerations, we opted to use SiGML and JASigning as a basis for sign language synthesis, and to implement a new functionality to automate the integration of non-manual grammatical markers with lexical signs.

A basic library of SiGML specifications of around 2.000 lexical signs in NGT was already compiled in the course of previous projects ([12], see also [32, 41]). While we have had to extend this library with healthcare-related as well as some general-purpose signs, the availability of an initial repertoire of signs encoded in SiGML was essential for a timely development of the system.

3.2 Translation

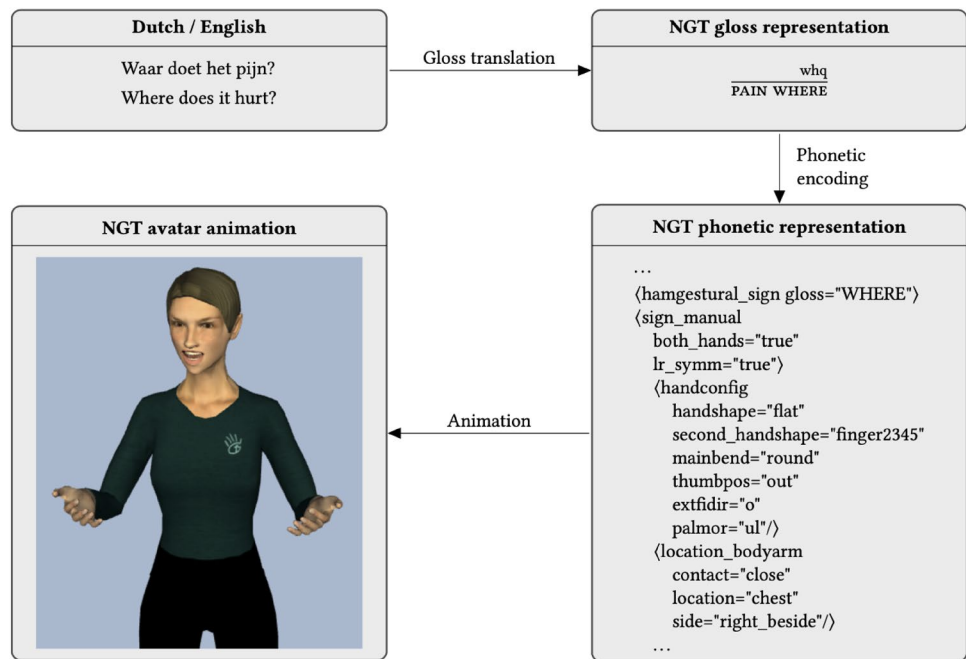
We now turn our attention from sign synthesis to the broader task of text-to-sign translation. Two approaches to this task can be distinguished, differing mainly in the type of intermediate representation that is employed in going from text to sign.

In the first approach, which we will refer to as the *gloss approach*, a given input sentence is transformed into a gloss of the corresponding sign language utterance. Next, based on this gloss representation, an avatar animation is generated.

(2) Gloss approach:

text \Rightarrow gloss \Rightarrow animation

This approach is taken, for instance, by HandTalk, a Brazilian company that provides an automated text-to-sign translation service with Brazilian Portuguese and English as possible source languages, and American Sign Language as well as Brazilian Sign Language (Libras) as possible target languages. HandTalk uses machine learning techniques to map input texts to the corresponding glosses, and a combination of keyframe animation and motion capture techniques to generate animations based on a given gloss.

Fig. 2 Overview of the modular translation pipeline

In the second approach, which we refer to as the *phonetic approach*, the given input sentence is transformed into a sequence of phonetic representations of signs. Next, based on these phonetic representations, an animation is generated.

(3) *Phonetic approach*:

text \Rightarrow phonetic representation \Rightarrow animation

This approach has been taken in work based on SiGML and JASigning (see, e.g., [6, 9, 32, 41, 53]). Unlike in the gloss approach, applying machine learning techniques to carry out the first step—from text to phonetic representations—is not feasible, because it would require the availability of large parallel corpora of texts and the corresponding phonetic sign representations. These do not exist, and would be very costly to create. The process of manually generating phonetic representations requires expert knowledge of SiGML or a similar formalism. Rayner et al. [43] have created a framework to ease this process, which is especially helpful if the sentences that need to be translated are all variations of a limited set of templates. For instance, the framework has been used to develop an application for translating railway announcements [6].

The gloss approach and the phonetic approach have complementary pros and cons. An advantage of the gloss approach is that it enables the use of machine learning technology to carry out the first part of the translation process. Disadvantages are that (i) the animation of each lexical sign involves substantial work, (ii) grammatical non-manual elements cannot be straightforwardly integrated with lexical signs, and (iii) all components of the system

are tailor-made for a particular target sign language, i.e., no part of the system can be re-used when a new target language is considered. In particular, since no gloss-based system currently exists for NGT, this approach was not viable for our purposes.

Advantages of the phonetic approach are that (i) grammatical non-manual features can in principle be integrated with lexical signs (though this possibility remains largely unexplored) and (ii) part of the system, namely the software that generates avatar animations based on phonetic representations (i.e., JASigning or a similar avatar engine) is not language-specific and can in principle be re-used for any target sign language. The main disadvantage is that the initial step from text to phonetic representations involves a lot of manual work.

Given these considerations, we have taken a *modular approach*, which employs *both* a gloss representation *and* a phonetic representation in going from a given input text to an avatar animation of the corresponding sign language utterance. As depicted in Fig. 2, our modular approach breaks the translation process up into three steps:

1. *Gloss translation*

In this step, the given Dutch or English input sentence is mapped to a gloss representation of the corresponding NGT utterance. This can be done with a rule-based grammar or with machine learning, depending on use case requirements and availability of training data;

2. *Phonetic encoding*

In this step, the NGT gloss is transformed into a computer-readable phonetic representation, in our case

formulated in SiGML. This can be fully automated in a rule-based system, which can also integrate grammatical non-manuals;

3. Animation

In this step, an avatar animation is generated based on the given phonetic representation. This procedure is not language specific, thus can be applied universally.

Consider, for instance, the Dutch/English input sentence in (4):

- (4) Waar doet het pijn?
Where does it hurt?

The first step is to convert this sentence into the corresponding NGT gloss in (5), where ‘whq’ stands for the non-manual marking that is characteristic for constituent questions in NGT. While empirical studies have found quite some variation in the actual realization of ‘whq’ in NGT [4, 8], furrowed eyebrows are seen as the most canonical realization [33].

- (5) $\frac{\text{whq}}{\text{PAIN WHERE}}$

The second step is to map this gloss representation to a phonetic representation in SiGML, a fragment of which is displayed in Fig. 2. Finally, this SiGML representation is fed into the JASigning avatar engine, which generates an animation (see Fig. 4 for a snapshot of the user interface of the system).

3.3 Implementation

Implementation choices depend on the specific use case requirements. Is it more important to achieve high precision, which a rule-based system allows, or to achieve broad coverage, which would favor an implementation involving machine learning? If the goal is to have optimal quality of lexical sign animations, one may opt to use motion capture, while scripted animation can be used in a scenario where scalability is of higher importance. The type and amount of resources available inevitably constrain one’s choices as well. Is there enough data for machine learning? Is a rule-based grammar available for the given domain? Is motion capture equipment available? What is the time-frame for development?

3.3.1 Use case requirements and implementation choices

Our main objective was to address the urgent concerns of deaf people in the Netherlands, ensuing from the COVID-19 pandemic, about the general inability of healthcare professionals to communicate in NGT [48]. Two specific

requirements followed from this objective: (i) the system had to be developed within a short time-frame, and (ii) high accuracy of the delivered translations was more important than broad approximate coverage. In addition to these requirements, our implementation choices were also affected by the fact that resources were limited.²

Our aim has therefore *not* been to automate the entire translation process. In particular, automating the process of mapping input sentences to the corresponding NGT glosses using machine learning techniques would not have been feasible within a short time-frame, and would, even in the somewhat longer term, most likely result in an unacceptably low accuracy rate for use in a healthcare setting. We therefore mainly focused on automating the phonetic encoding step, something that significantly reduces the manual labor needed in the overall translation pipeline. Automating the mapping from glosses to phonetic representations has not been done in previous work on NGT [41] and, to the best of our knowledge, not in work on other sign languages either.

3.3.2 Selecting phrases for translation

We selected a set of phrases that are commonly used during the diagnosis and treatment of COVID-19, based on consultation with healthcare professionals at the Amsterdam University Medical Centre (AUMC) as well as direct experience (one of the authors is a medical doctor). We also consulted a list of phrases that was used in the SignTranslate system in the UK [38].³

The resulting corpus was then divided into three categories: video-only, avatar-only, and hybrid. The first category, video-only, consisted mainly of sentences that could be divided into three further categories: emotional, complex, and informed consent. Sentences concerning the patient’s emotional well-being require a high level of empathy to be

² Since the time-frame and available resources for this research project were really quite different than for prototypical academic projects, we provide some details. Initial funding for the project was provided by an ad hoc funding scheme setup by the Netherlands Organization for Innovation in Healthcare (ZonMW) to address urgent COVID-related issues in the healthcare sector. The deadline for proposals in this funding scheme was two weeks after the call for proposals had been announced, funded projects had to start one month later, and had to be completed within six months, with a total budget of 25.000 euros. In this period, we designed and implemented the prototype system. Separate funding was used for the evaluation study.

³ The SignTranslate system was developed in the UK around 2010 to translate phrases common in a healthcare setting from English to British Sign Language. Translations were displayed by means of videos, not by avatar animations. Evidently, the system was not specifically targeted at COVID-19 healthcare. However, many general-purpose phrases are also relevant in the diagnosis and treatment of COVID-19.

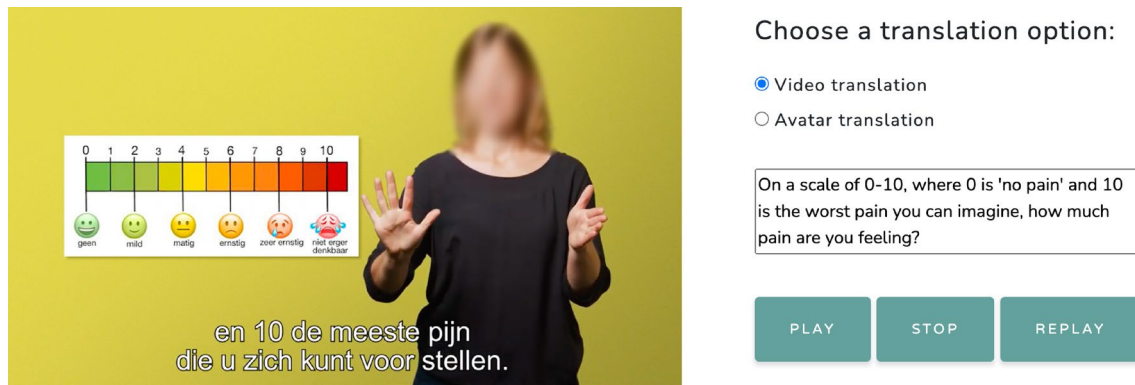


Fig. 3 Example of a video translation of a complex question. It is long and supported by an image

conveyed, which is difficult to achieve in a satisfactory way with an avatar given the current state of the art. We therefore deemed that video translations were necessary for these sentences. Sentences were classified as complex when they involved a combination of several statements and/or questions, or required a demonstration of pictures or diagrams along with an explanation (as shown in Fig. 3). Finally, in the case of questions and statements concerning informed consent, it is especially important to leave no room for potential misunderstandings. To ensure this, we chose to always offer video translations of these sentences.

The second category, avatar-only, consisted of sentences with many variations differing by only one word or phrase, indicating for instance the time of day or a number of weeks. It would not have been feasible to record a video translation for each version of these sentences.

The third category, hybrid, consisted of sentences that did not fall into one of the other two categories. For these, the system offers both a video translation and an avatar translation. In some cases, the avatar translation is slightly simplified compared to the video translation (e.g., some long sentences were broken up into several smaller ones).

After categorizing all of the sentences, those from the first and third category were translated into NGT and recorded by a team consisting of a sign language interpreter and a deaf signer. Translations were checked by one of the authors (Smeijers), who is a sign linguist and a medical doctor. This resulted in a collection of 139 video translations. The sentences from the second and third category (including all variations) together comprised 7.720 sentences for avatar translation.

3.3.3 Constructing phonetic representations

In order for the system to operate fast at run-time, we pre-processed all sentences and stored phonetic SiGML representations of their translations in a database. At run-time, the system only queries this database and does not compute

any translations on the fly. The complete database of SiGML representations can be found at [18, 19].

To construct the SiGML representations of full sentences, we developed a system which, when given the gloss representation of a sentence in NGT, creates the SiGML code for that sentence. It first retrieves the SiGML code for all lexical signs in the given gloss from a lexical database, and then adapts this code to add non-manual grammatical elements. For instance, in the case of questions, the program ensures that the sentence ends with PALMS-UP, a sign that can be used in NGT to mark questions, and adds raised eyebrows, both to PALMS-UP and to the sign preceding it.⁴

3.4 User interface

We developed an online user interface (Fig. 4). Healthcare professionals can choose a translation format (video or avatar) and enter a sequence of search terms. Based on their input they are presented with a list of available sentences from the database. These sentences may differ depending on the translation format chosen (video/avatar). After selecting a sentence, the translation is offered in the chosen format.

As mentioned earlier, some of the possible input sentences differ only in one word or phrase. These sentences can be thought of as involving a general template with a variable that can take several values, such as a day of the week, a time of day, or a number of times/minutes/hours/days/weeks/months. When a user wants to translate such a sentence, they first select the template and then provide the intended value for the variable. For example, they may select the template

⁴ PALMS-UP and raised eyebrows are prototypical question markers in NGT, but questions can be marked in other ways as well [4, 8]. Furrowed eyebrows, for instance, are also sometimes used for this purpose. We always included PALMS-UP and raised eyebrows, but more research is needed to determine exactly under which conditions these question markers are used in NGT and under which conditions they are replaced by or combined with other markers.

Fig. 4 User interface of the system

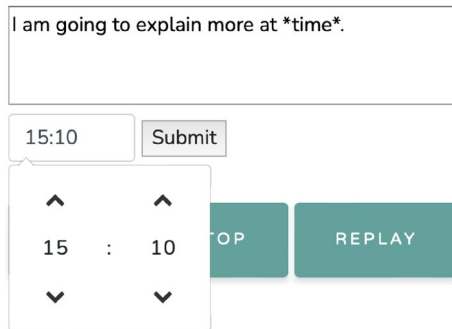
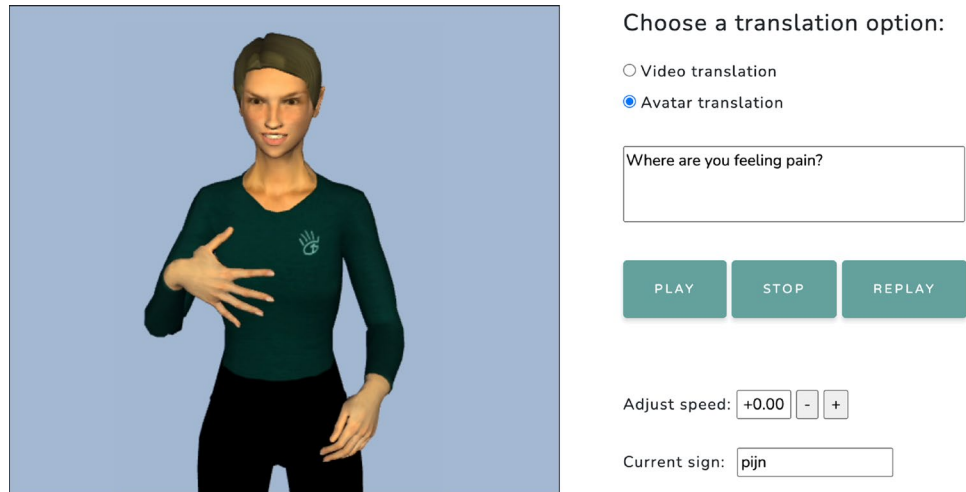


Fig. 5 Selecting a time value for a variable in a sentence

“I am going to explain more at *time*”, and then select a particular time (as illustrated in Fig. 5).

While JASigning in principle offers several avatars for sign language animation, there are differences in execution between these avatars. Our user interface therefore only makes use of one of them, Francoise (see Fig. 4), and does not allow the user to choose different options.

4 Evaluation

To evaluate the implemented prototype system, we conducted an online survey among 22 deaf NGT users. There is no generally accepted methodology for evaluating the comprehensibility of avatars for text-to-sign translation, let alone for doing so *online*. Evaluation procedures designed in previous work generally involve on-site interaction between experimenters and participants ([6, 9, 29, 31, 50], with exceptions noted in [24, 42, 47]). The COVID-19 pandemic made it necessary to turn to online procedures,

which come with additional methodological challenges. On the bright side, such online procedures, if effective, may also have benefits in a post-COVID-19 world.

4.1 Research questions

We focused on the following research questions:

RQ1. *Comprehensibility at the level of individual signs*

To what extent are individual signs understood as intended when performed by the avatar?

RQ2. *Comprehensibility at sentence level*

To what extent is the NGT translation of a given input sentence understood as intended when performed by the avatar?

RQ3. *Influence of interaction between participants and experimenters*

When evaluating the comprehensibility of a signing avatar in an online environment, how does interaction between participants and experimenters, or the lack thereof, affect the results?

The first two research questions concerned the comprehensibility of the implemented system. The third research question on the other hand concerns the methodology of evaluating signing avatars in an online setting. Of course, besides these three research questions there are other pertinent ones as well, concerning for instance the attitude of potential end-users toward signing avatars in general and toward our system in particular. We did include some questions in our survey to probe such attitudes, and will briefly report participants' responses to these questions below, but leave an in-depth investigation for a separate study.

4.2 Methodology

4.2.1 Participants

To recruit participants, we recorded a video in NGT which briefly explained our project and invited people to sign up as a participant of the evaluation study. This video was posted on various social media platforms, including multiple Facebook groups, Instagram, and LinkedIn. We also asked personal contacts in the deaf community to distribute the video and ask their contacts to do the same. To sign up, interested people were asked to fill in a short form collecting their contact information and some demographic information. They were also asked whether they would prefer to participate in the online study with or without guidance of the experimenters. All these questions were asked both in written Dutch and in NGT through videos featuring a deaf signer.

We recruited 23 participants in total, but 1 of them did not complete the survey in the end, so the results of the survey are based on the responses of 22 participants. Table 1 provides an overview of their demographic information (the table contains aggregate data as well as data for both experimental groups, GUIDED and UNGUIDED, see Sect. 4.2.2 below on the experimental design). Participants were spread across age groups relatively evenly. There were many more female (18) than male (4) participants. Most participants were from Central, Northern, or Western provinces.

As for language background, 19 participants identified only NGT as their mother tongue, and 3 participants identified both Dutch and NGT as their mother tongue. No participants identified Dutch as their sole mother tongue. 20 participants indicated that they use NGT *daily*, while 2 participants used NGT *regularly* (a few days a week). 16 participants indicated that they *regularly* make use of an interpreter, 5 indicated that they *occasionally* do (a few days a month), and 1 indicated that they *rarely* do (a few days a year).

When communicating with non-signers in the absence of an interpreter, participants indicate that they use a combination of various communication methods. The most frequent methods are lipreading (20 participants), speaking (18 participants), and writing (15 participants). The majority of participants employ these methods either *daily* (7 participants) or *regularly* (10 participants).

We collected this rather detailed demographic information at the recruitment stage with the intention to create groups of participants that were optimally counter-balanced in terms of age group, region, and language background. However, because the total number of recruited participants was relatively low, we decided to include all of them in the survey.

Table 1 Demographic information about the participants

	Guided	Unguided	Total
<i>Age group</i>			
18–30	0	6	6
31–40	5	3	8
41–50	2	2	4
51+	2	2	4
<i>Gender</i>			
Female	8	10	18
Male	1	3	4
<i>Region^a</i>			
Central	2	4	6
Northern	3	5	8
Eastern	0	0	0
Southern	0	1	1
Western	4	3	7
<i>Mother tongue</i>			
NGT	7	12	19
Dutch	0	0	0
NGT and Dutch	2	1	3
<i>Frequency of NGT use^b</i>			
Daily	9	11	20
Regularly	0	2	2
Occasionally	0	0	0
Rarely	0	0	0
<i>Frequency of use of interpreter</i>			
Daily	0	0	0
Regularly	7	9	16
Occasionally	1	4	5
Rarely	1	0	1
<i>Communication methods with non-signers in the absence of an interpreter</i>			
Lipreading	7	13	20
Speech recognition	3	2	5
Writing	8	7	15
Signing	2	8	10
Using voice	6	12	18
Pointing	2	0	2
<i>Frequency of use of other communication methods</i>			
Daily	3	4	7
Regularly	4	6	10
Occasionally	2	2	4
Rarely	0	1	1

^aRegions are divided into the following provinces of the Netherlands: Central (Utrecht, Flevoland); Northern (Groningen, Friesland, Drenthe); Eastern (Gelderland, Overijssel); Southern (Brabant, Zeeland, Limburg); and Western (Noord-Holland, Zuid-Holland)

^bFrequencies are defined as: daily; regularly (a few days a week); occasionally (a few days a month); and rarely (a few days a year)

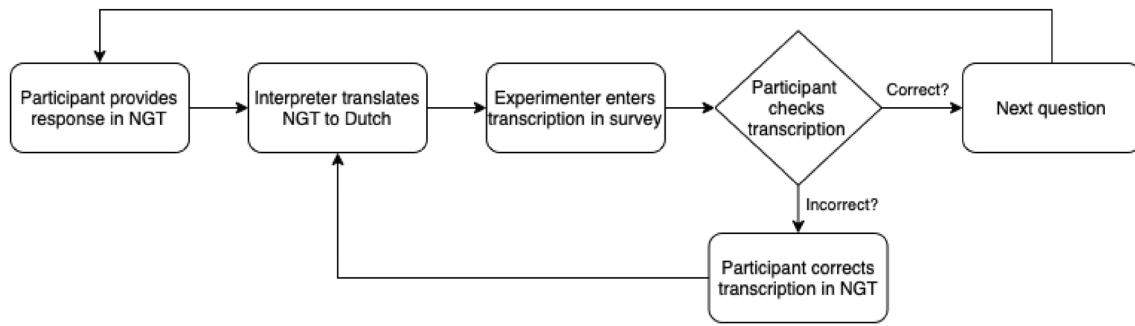


Fig. 7 Feedback loop to ensure faithful textual transcription of responses provided in NGT

Table 2 Encoding key

Code	Individual signs	Sentences
0	No response	No response
1	Wrong response provided	Wrong response provided
2	Manual component recognized correctly, but oral component not	Sentence radical recognized correctly, but sentence type not
3	Two possible interpretations provided, one of which correct	Two possible interpretations provided, one of which correct
4	Correct	Correct

4.2.3 Coding comprehension data

For each sentence in Part 2 and 3, we encoded whether or not the participant correctly recognized each individual sign and interpreted the sentence as intended. In case the response was partly correct and partly incorrect we also encoded the type of error that was made. An overview of the different numeric labels we used to code the responses is provided in Table 2.

The characterizations of codes 0, 1, 3, and 4 in Table 2 are self-evident. As for code 2, for individual signs, this code indicates that the participant correctly recognized the manual component of the sign, but did not correctly recognize the oral component (e.g., *RESULT* and *PASSED* share the same manual component, visualized in Fig. 8, but have different oral components). As for entire sentences, code 2 indicates that the participant correctly identified the meaning of the sentence radical, but confused the sentence type (e.g., the response was formulated as a question when the actual sentence was a statement).

4.3 Results

4.3.1 Comprehension

Table 3a shows the proportion of each response type in percentages, for individual signs and sentences, split across the

Fig. 8 The manual component of the signs *RESULT* and *PASSED* in NGT [46]. The two signs differ only in their oral component



two signers (the avatar and the human signer). Note that ‘partly correct’ responses (coded as 2 or 3, depending on the type of error) occurred very infrequently. Therefore, in Table 3b, the results are presented in a more compact format, taking all incorrect responses (coded as 1, 2, or 3) to form a single category. Finally, in Table 3c, we go one step further and present the results in a binary format, distinguishing only between correct responses on the one hand and incorrect or missing responses on the other.

Comprehension rates for individual signs. We zoom in now on Table 3b and first consider individual signs. Note that the comprehension rate of individual signs performed by the avatar (86.93%) is lower than the comprehension rate of individual signs performed by a human signer (92.28%), but both rates are quite high and the difference between them is rather small. Also note that, in the case of the human signer, almost all cases in which a sign was not correctly recognized are ones in which *no response* was given at all (6.85%). An incorrect response was only given 0.87% of the time. In the case of the avatar on the other hand, an incorrect response was given 6.75% of the time, while no response was given in 6.32% of the cases.

Comprehension rates for full sentences. We now turn to the comprehension rates for full sentences. Here we see a

Table 3 Proportion of each response type in percentages, for individual signs and sentences, split across the two signers (*Avatar vs Human*)

Code	Response type	Individual signs		Sentences	
		Av	Hu	Av	Hu

(a) Fine-grained

0	No response	6.32	6.85	4.17	0.00
1	Wrong response	4.40	0.72	18.56	1.89
2	Partly correct: A	2.27	0.00	1.89	0.76
3	Partly correct: B	0.07	0.14	0.76	0.00
4	Correct response	86.93	92.28	74.62	97.35

(b) All incorrect responses aggregated

0	No response	6.32	6.85	4.17	0.00
1–3	Incorrect response	6.75	0.87	21.21	2.65
4	Correct response	86.93	92.28	74.62	97.35

(c) All incorrect and absent responses aggregated

0–3	No correct response	13.07	7.72	25.38	2.65
4	Correct response	86.93	92.28	74.62	97.35

bigger disparity between the avatar and the human signer. Participants were almost always able to correctly identify the meaning of a sentence when it was signed by a human (97.35%). This was not the case for sentences signed by the avatar: participants found it harder to correctly identify the meaning of these sentences, with a comprehension rate of 74.62%. They provided an incorrect response 21.21% of the time, and no response at all 4.17% of the time.

Most common mistakes. We now take a closer look at which individual signs were not always recognized correctly. Table 4 lists all signs for which either no response was given more than twice or an incorrect response was given more than twice. For each of these signs, the table indicates the number of responses of type 0 (no response at all), 1 (wrong response), and 2 (manual component recognized correctly, but oral component not) as a percentage of the total number of responses that were elicited for that sign performed by the given signer (avatar or human). The signs are divided into four categories: A (grammatical markers), B (pronouns), C (signs whose interpretation crucially relies not just on the manual component of the sign but also on mouthing), and D (miscellaneous).

The signs in category A are glossed as INDEX and as PALMS-UP. INDEX signs are ones that refer back to something that has been introduced earlier in the same sentence. This ‘doubling’ mechanism is common in NGT but generally seems optional: leaving such signs out usually does not change the meaning of the sentence. In this sense, such signs are purely ‘grammatical’, they do not contribute any content. The PALMS-UP sign is used for various purposes in NGT. In the sentences under consideration, it was always used to mark a sentence as a question. When used for this purpose, PALMS-UP is generally optional as well. If it is left out, the signer’s facial expression is generally sufficient to convey that she is asking a question rather than making a statement.

We see that these two signs were often not recognized explicitly (INDEX 27.3% of the time when signed by the avatar, and 41.8% of the time when signed by a human; PALMS-UP 8.2% of the time when signed by the avatar, and 16.7% of the time when signed by a human). This may well be connected to the fact that these signs are optional, and do not contribute any content to the meaning of the sentence as a whole that is not already conveyed by other elements. This may make them less salient, resulting in participants ‘skipping over’ them. They were also recognized incorrectly in some cases, especially when signed by the avatar (INDEX 4.6% of the time; PALMS-UP 9.1% of the time). A possible explanation for this result is that both signs have several other possible interpretations/grammatical functions as well. Indeed we find some of these among the interpretations provided by our participants (e.g., INDEX was sometimes interpreted as YOU OR FOR, and PALMS-UP was sometimes misinterpreted as WHERE).

The signs in category B are the pronouns YOU and I. We see that these signs were quite often not explicitly recognized by our participants when signed by the human signer, though this did not occur with the avatar. A possible explanation for this is that these signs are often made very fast by human signers, and they are often co-articulated with adjacent signs. The avatar on the other hand, produced these signs at a slower pace, and did not co-articulate them with adjacent signs. When looking more closely at the data, we find that there were three sentences in which our participants failed to explicitly recognize these signs. Glosses of these sentences are given below.

- (6) SORRY I FAIL YOU INTRAVENOUS DRIP I COLLEAGUE CALL
not explicitly recognized: YOU (21 times), I (12 times)
- (7) HEARING-AID OR COCHLEAR-IMPLANT HAVE YOU
not explicitly recognized: YOU (3 times)
- (8) PAST SEVEN DAY YOU ALREADY CORONA TEST
not explicitly recognized: YOU (3 times)

The signs in category C are SOMETIMES, FOR, RESULT, and BACK. To recognize these signs correctly it is crucial to recognize not only their manual component but also the accompanying movement of the mouth. This is because for each of these signs there is at least one other sign with the same manual component but different mouthing.

We see that this is particularly problematic for the avatar. For instance, SOMETIMES was misinterpreted 54.5% of the time as MAYBE, which has the same manual component, and similarly, RESULT was misinterpreted 22.7% of the time as PASSED. Note that such misinterpretations did not arise when the signs were performed by a human signer. For the sign FOR, we see that, again only when performed by the avatar, it was misinterpreted 36.4% of the time and not explicitly recognized 22.7% of the time. A possible explanation for the fact that it was so often not explicitly recognized is that,

Table 4 Individual signs which were either not explicitly recognized more than twice or incorrectly recognized more than twice

Category	Sign	Code 0		Code 1		Code 2	
		Av	Hu	Av	Hu	Av	Hu
A	INDEX	27.3	41.8	4.6	1.8	–	–
	<i>Grammatical</i> PALMS-UP	8.2	16.7	9.1	1.5	–	–
B	YOU	–	11.9	–	–	–	–
<i>Pronouns</i>	I	–	18.2	–	3.0	–	–
C	SOMETIMES	4.6	–	–	–	54.5	–
<i>Mouthing</i>	FOR	22.7	–	–	–	36.4	–
	RESULT	–	–	–	–	22.7	–
	BACK	4.6	–	4.6	–	13.6	–
D	PAST	40.9	–	27.3	–	–	–
<i>Miscellaneous</i>	SEVEN	50.0	–	13.6	–	–	–
	DAY	27.3	–	40.9	–	–	–
	INTRAVENOUS DRIP	31.8	–	9.1	–	–	–
	FAIL	9.1	–	27.3	–	–	–
	SORRY	4.6	–	13.6	–	–	–
	MORE	4.6	–	13.6	–	–	–

In each case, we express the number of incorrect/missing responses as a percentage of the total number of responses elicited for the given sign performed by the given signer (Avatar vs Human)

when the mouthing is not correctly perceived, the sign can easily be mistaken for an INDEX sign, and we already saw that INDEX signs are often not explicitly recognized (possibly because they are optional and therefore perhaps less salient).

Category D contains seven miscellaneous signs, i.e., these are not grammatical markers, pointing signs, or signs whose recognition crucially relies on mouthing. Note that all signs in this category were correctly recognized when signed by a human signer, but not always when signed by the avatar. The first three signs, PAST SEVEN DAY were signed sequentially as one phrase (see example (8) above). When signed by the avatar, this sequence of signs was quite fast, with a high degree of co-articulation. This made it especially difficult to identify the sign SEVEN. An additional intricacy is that DAY is very similar to MONTH, differing only slightly in the location of the hand. This contributed to the frequent misinterpretation of the sign.

Turning now to the other three signs in category D, a possible explanation for the low comprehension rate of INTRAVENOUS DRIP when signed by the avatar is that it is, presumably, a rather infrequent sign. When performed by a human signer, however, it was always recognized correctly. The low comprehension rates of FAIL and MORE are possibly due to the fact that there are several other signs which are very similar and would have made sense in the given contexts (in particular, participants misinterpreted FAIL either as LOST or as WHERE, and MORE either as FIRST or as AGAIN). Finally, our participants were sometimes confused about the sign SORRY when performed by the avatar. The interpretations they provided in these cases

were I, HEART, and CARD. Since the sign SORRY is actually quite different from I, HEART, and CARD, the only possible explanation we can offer for this confusion is that the avatar's rendering of the sign SORRY was not quite successful.

4.3.2 Clarity

Figure 9 provides an overview of the clarity scores participants gave to avatar animations generated by the system, compared to videos of a human signer. Recall that all participants saw avatar animations for 12 sentences, followed by videos of a human signer for the same 12 sentences. Clarity ratings were given on a scale from 0 to 10, where 0 was labeled as 'not clear' and 10 as 'very clear'.

Clarity scores given for avatar animations ranged between 0 and 10, with a mean of 6.4. Clarity scores for videos ranged between 6 and 10, with a mean of 9.0. In Fig. 9a we have plotted the relative frequency, in percentages, of each clarity score given to avatar animations and videos, respectively, summed across all participants. This figure shows that scores of 6, 7, and 8 were most frequently given to avatar animations, while scores of 8, 9, and 10 were most frequently given to videos of a human signer. Videos received a maximal score of 10 more than 45% of the time.

In Fig. 9b we have plotted the mean clarity scores for each of the 12 sentences considered in the survey, for avatar animations and videos respectively. Scores for avatar animations were consistently lower than for videos, as expected, but the difference between the two varied considerably across sentences. In particular, the figure shows that while

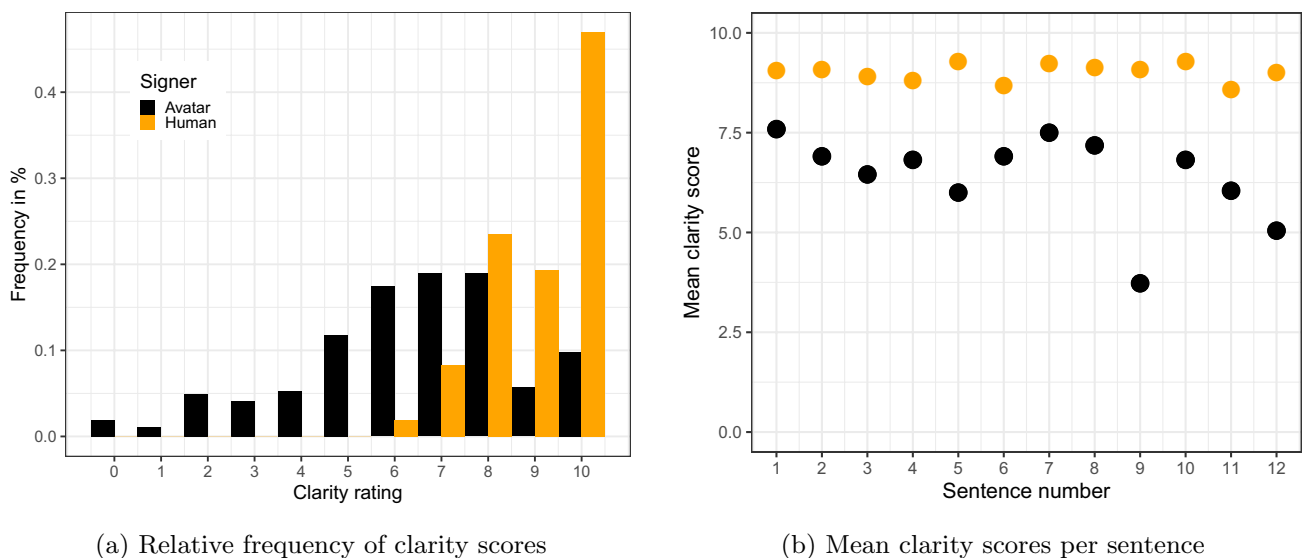


Fig. 9 Clarity scores

all sentences received roughly the same mean score when signed by a human signer (mean 9.0, standard deviation 0.22), the variability among sentences was considerably larger when signed by the avatar (mean 6.4, standard deviation 1.1).

When we zoom in on the scores for avatar animations, we see that sentences 9 and 12 may be considered negative outliers, since their mean scores (3.7 and 5.0, respectively) were more than one standard deviation below the mean across all sentences. Glosses of these two sentences are given in (9) and (10), respectively:

- (9) Sentence 9: PAST SEVEN DAY YOU ALREADY CORONA TEST
 (10) Sentence 12: SORRY I FAIL YOU INTRAVENOUS DRIP I COL-LEAGUE CALL

Sentence 9 starts with the phrase PAST SEVEN DAY, which we already saw was very difficult for participants to recognize (see Table 4). Similarly, sentence 12 also contains three signs which were difficult to recognize when signed by the avatar, SORRY, FAIL, and INTRAVENOUS DRIP. Plausibly, these individual signs contributed to the low clarity scores for these sentences.

If sentences 9 and 12 are disregarded, the mean clarity score for videos is still 9.0 points while the mean score for avatar animations rises to 6.8 points. A possible interpretation of the clarity score data, then, is that there is a ‘root difference’ between videos and avatar animations of about 2.2 points (the difference between the two means if the two outlier sentences are disregarded), and that this difference becomes larger when a sentence contains individual signs or phrases which are particularly difficult to recognize when signed by the avatar in the current

implementation of the system. We hypothesize that the latter effect may in principle be reduced by improving the way in which these individual signs and phrases are rendered by the avatar. On the other hand, we expect that the ‘root difference’ between videos and avatar animations will be more persistent. Closing this gap will not be a matter of ‘quick fixes’ but will require more fundamental improvements of the underlying avatar technology (e.g., by making use of motion capture instead of scripted animation, or a combination of the two).

4.3.3 Guided versus unguided

We now turn to results pertaining to our third research question, comparing results of the GUIDED group with those of the UNGUIDED group.

Comprehension rates for individual signs. We define the comprehension rate of an individual sign as the number of times that the sign was correctly recognized (code 4) divided by the total number of responses for that sign (codes 0–4). So the comprehension rate is a number between 0 and 1, reflecting the proportion of cases in which the sign was correctly recognized. Table 5a shows the mean comprehension rate of all individual signs, split by participant group (GUIDED VS UNGUIDED) and signer (Avatar vs Human). We see that for the GUIDED group, the mean comprehension rate was only 0.05 lower for the avatar (mean 0.82, standard deviation 0.38) than for the human signer (mean 0.87, standard deviation 0.34). For the UNGUIDED group, the mean comprehension rate for the avatar was again close to that of the human signer, with a difference of 0.06 (avatar mean 0.90, standard deviation 0.30; human signer mean 0.96, standard deviation 0.20).

Table 5 Effect of guidance on comprehension rates for individual signs and sentences

	GUIDED		UNGUIDED	
	Av	Hu	Av	Hu
(a) Individual signs				
Mean	0.82	0.87	0.90	0.96
St. dev.	0.38	0.34	0.30	0.20
(b) Sentences				
Mean	0.68	0.95	0.80	0.99
St. dev.	0.47	0.21	0.40	0.11

If we compare the GUIDED and UNGUIDED group per signer, we see that the mean comprehension rate for the avatar was 0.08 higher in the UNGUIDED group than in the GUIDED group, while the mean comprehension rate for the human signer was 0.09 higher in the UNGUIDED group than in the GUIDED group.

Overall, then, results from the GUIDED and the UNGUIDED group were quite similar when it came to the mean comprehension rate of individual signs.

Comprehension rates for full sentences. We define the comprehension rate of a sentence as the number of times that the sentence was interpreted as intended (code 4) divided by the total number of responses elicited for that sentence (codes 0–4). So, just like in the case of individual signs, the comprehension rate for a sentence is a number between 0 and 1, reflecting the proportion of cases in which the sentence was correctly interpreted.

Table 5b shows the mean comprehension rate of all sentences, split by participant group (GUIDED vs UNGUIDED) and signer (Avatar vs Human). We see that for the GUIDED group, the mean comprehension rate was 0.27 lower for the avatar (mean 0.68, standard deviation 0.47) than for the human signer (mean 0.95, standard deviation 0.21). For the UNGUIDED group, the mean comprehension rate for the avatar was 0.19 lower for the avatar (mean 0.80, standard deviation 0.40) than for the human signer (mean 0.99, standard deviation 0.11).

If we compare the GUIDED and UNGUIDED group per signer, we see that the mean comprehension rate for the avatar was 0.12 higher in the UNGUIDED group than in the GUIDED group, while the mean comprehension rate for the human signer was 0.04 higher in the UNGUIDED group than in the GUIDED group.

Overall, then, results from the GUIDED and the UNGUIDED group were quite different when it came to the mean comprehension rate of sentences signed by the *avatar*. In this case, rates were substantially higher in the UNGUIDED group than in the GUIDED group. On the other hand, the comprehension rates of sentences signed by a *human* were more similar across the two groups, closer to what we observed for individual signs.

We further observe that, both in the GUIDED and in the UNGUIDED group, the standard deviation of the comprehension rates for sentences signed by the avatar was much higher than that of the comprehension rates for sentences signed by a human.

To obtain a better understanding of this larger variance in comprehension rates, Fig. 10a plots the mean comprehension rates of all sentences, signed by the avatar and a human signer, respectively. The left pane (for the GUIDED group) shows that nine sentences had a relatively high comprehension rate, while three sentences had a relatively low comprehension rate. Similarly, in the right pane (for the UNGUIDED group) we see that ten sentences had a high comprehension rate, and two a low one. Overall, there is great similarity between the GUIDED and the UNGUIDED group as to which sentences received high rates and which ones received low rates. The only exception is sentence 10, which had a low rate in the GUIDED group but a high rate in the UNGUIDED group.

Clarity scores. Figure 10b compares the GUIDED and the UNGUIDED group with respect to clarity scores. We see that the left pane (for the GUIDED group) is overall very similar to the right pane (for the UNGUIDED group), and both panes are similar to the plot in Fig. 9b above, which displayed clarity scores for all sentences without making a distinction between the GUIDED and the UNGUIDED group. The only salient difference between the left and the right pane in Fig. 10b pertains to sentences 5 and 11. These sentences received lower clarity scores in the GUIDED group than in the UNGUIDED group.

The overall impression that arises from comparing the GUIDED and the UNGUIDED group with respect to comprehension rates and clarity scores is that the two yielded largely similar results. The only case in which we observed a substantial difference between the two groups was in the comprehension rates for sentences signed by the avatar. In Sect. 4.3.4 below we will analyze in which cases the effect of guidance on comprehension rates and clarity scores was statistically significant.

4.3.4 Statistical analysis

We now present statistical models that investigate whether comprehension rates and clarity scores were significantly affected by *guidance* (GUIDED vs UNGUIDED) and the type of *signer* (avatar vs human). Three statistical models are considered: one for the comprehension of individual signs, one for the comprehension of sentences and one for clarity scores. The models are built in R using the lme4 package [3].

Comprehension of individual signs. We concentrate on the binary distinction between correct responses on the one hand (code 4) and incorrect or missing responses on the other hand (code 0–3). We applied a generalized linear mixed effects logistic regression, which models a binary

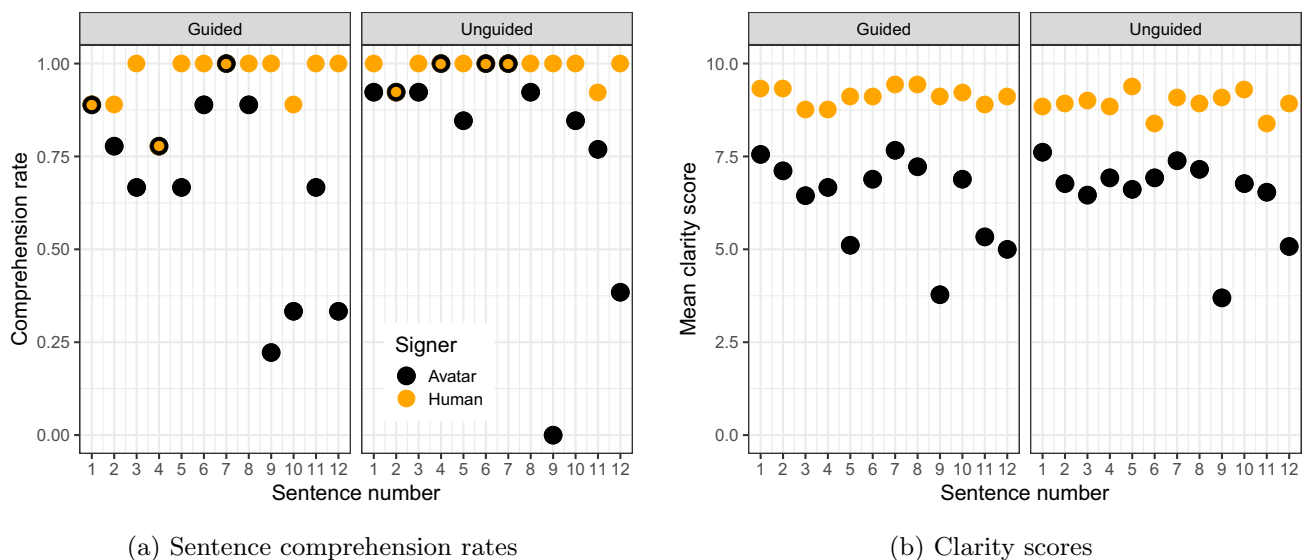


Fig. 10 Effect of guidance on sentence comprehension rates and clarity scores assigned to sentences

outcome as a combination of predictors. These predictors are labeled either as fixed effects or as random effects.⁷ In our case, the outcome is the response (0 = incorrect or missing response, 1 = correct response). Random effects are elements that would vary from one experiment to the next (if the same question was investigated) and which the experimenters do not exert full control over, either because this is impossible or because it is not of interest. In our case, participants and the signs that were used in the experiment are random effects. Fixed effects, on the other hand, are experimental manipulations that are of main interest and that would not vary from one experiment to the next. In our case, the fixed effects are *guidance* (with GUIDED and UNGUIDED as values, coded in the model using sum contrast coding as 0.5 and -0.5 , respectively) and *signer* (with human and avatar as values, coded in the model using sum contrast coding as 0.5 and -0.5 , respectively). We also include the interaction of *guidance* and *signer*, which reveals whether the difference due to one fixed effect (say, *guidance*) differs across different values of the other fixed effect (i.e., *signer*). Following common practice, we use the mixed effects model with the maximal random-effect structure that converges [2], which in our case was a model with a random intercept and random *signer* slopes (but no random *guidance* slopes) for participants and signs.

The obtained model identified both *signer* and *guidance* (but not their interaction) as significant predictors. The effect

of *signer* was positive ($\beta = 4.1, z = 2.5, p = 0.013$), which means that the comprehension rate of signs was significantly higher when performed by the human signer than when performed by the avatar. The effect of *guidance* was negative ($\beta = -1.3, z = -5.3, p < 0.0015$), which means that comprehension rates for individual signs were significantly higher in the UNGUIDED group than in the GUIDED group. Finally, the fact that the interaction between *signer* and *guidance* was not a significant predictor means that we have no evidence that the effect of *signer* differed across the GUIDED and the UNGUIDED group.

Sentence comprehension. We again assumed a binary outcome variable, only distinguishing between cases in which a sentence was interpreted as intended (outcome = 1) and cases in which it was not (outcome = 0). More fine-grained distinctions such as that between missing and incorrect interpretations were disregarded. Just as in the case of individual sign comprehension, we applied mixed effects logistic regression. This time, the random effects were participants and sentences. As before, the fixed effects were *guidance* (with GUIDED and UNGUIDED as values, coded in the model using sum contrast coding as 0.5 and -0.5 , respectively), *signer* (with human and avatar as values, coded in the model using sum contrast coding as 0.5 and -0.5 , respectively), and the interaction between *guidance* and *signer*. The model with the most comprehensive random-effect structure that converged included a random intercept for participants and sentences, and random *signer* slopes for sentences (but not for participants).

The model revealed a significant positive effect of *signer* ($\beta = 2.6, z = 2.9, p = 0.004$), which means that the comprehension rate of sentences was significantly higher when

⁷ For an introduction into mixed effects models in linguistics, see Winter [51]; for a more advanced presentation, see Gelman and Hill [22].

signed by a human than when signed by the avatar. This matches the positive effect of *signer* on the comprehension of individual signs that we found above. The model did not identify the effects of *guidance* and the interaction between *signer* and *guidance* as significant predictors.

Clarity scores. Clarity was judged on a discrete scale from 0 to 10. We applied ordered probit mixed effects regression in this case, which models a discrete and bounded set of outcomes as resulting from a combination of fixed and random effects.⁸ Just like the model for sentence comprehension, this model included participants and sentences as random effects and *guidance* and *signer*, as well as the interaction between *guidance* and *signer*, as fixed effects. The model with the most comprehensive random-effect structure that converged included a random intercept and random *signer* slopes for participants and sentences.

The model identified only one significant effect. Namely, *signer* had a positive effect ($\beta = 4.3, z = 7.6, p < 0.001$), which means that clarity scores for sentences performed by a human signer were significantly higher than those for sentences performed by the avatar. This matches the positive effect of *signer* on the comprehension rates for sentences and individual signs that we found above.

4.3.5 Attitudes toward signing avatars

We now turn to results obtained in Part 4 of the survey. This part did not concern the comprehensibility or clarity of the avatar in our prototype system, but rather probed participants' attitudes toward signing avatars in general, and also asked them for feedback on the setup of our survey.

Participants' attitude toward signing avatars was generally positive: 86.4% found that signing avatars should be further investigated and could potentially be very useful when further developed (see Table 6). Many participants noted explicitly that the technology in its current state is not advanced enough yet to be deployed in real-life settings. Moreover, multiple participants explicitly commented that the looks of the avatar in our current prototype need to be improved. It was perceived as rather stiff, not very friendly, and sometimes even scary. Such comments are reminiscent of the well-known 'uncanny valley' effect.

Those participants who believed that signing avatars should *not* be investigated further (13.6%) felt that it would be impossible for an avatar to ever display human-like facial expressions, and believed that avatars might take away jobs from human sign language interpreters and teachers.

Participants were also asked to reflect on possible use cases for signing avatars. A few example use cases were

Table 6 Participants' attitude toward signing avatars, and possible use cases

		In favor	Against
Attitude	Should be investigated further?	86.4	13.6
Use cases	Travel information	77.2	22.7
	Medical settings	50.0	50.0
	Support for learning	31.8	68.2

given in the survey (see Table 6 and Appendix A.4.1). Most participants indicated that signing avatars could be useful in public places such as train stations and airports, to relay travel information to passengers (77.2% in favor, 22.7% against). Opinions were divided about the use of signing avatars in medical settings (50.0% in favor, 50.0% against) or as part of tools to support people in learning sign language (31.8% in favor, 68.2% against). Concerning both these use cases, some participants were very enthusiastic, but others were strongly opposed. Participants on both sides of the spectrum indicated that the technology would need to be improved significantly before being used in these settings.

In addition to the example use cases listed in the survey, participants suggested other possible use cases as well. Multiple participants indicated that a signing avatar could be useful in waiting rooms, for standardized governmental processes and public services (e.g., renewing a passport), and in supermarkets.

4.3.6 Feedback on the setup of the survey

All participants indicated that the questions in the survey were clearly formulated: mean = 9.22 (standard deviation = 1.09) in the GUIDED group and mean = 8.38 (standard deviation = 1.04) in the UNGUIDED group) and that it was easy to provide answers (mean of 8.78, standard deviation of 1.30 in the GUIDED group and mean of 8.23, standard deviation of 1.01 in the UNGUIDED group).

Moreover, participants from both groups indicated that they felt taken seriously, although GUIDED participants commented on this more often than UNGUIDED ones. Participants listed a number of factors that contributed toward this feeling: (i) all instructions and questions in the survey were given in two formats, Dutch text and NGT videos, (ii) the NGT videos featured a deaf signer, (iii) the transcription of responses involved a feedback loop in the case of the GUIDED group, as described in Sect. 4.2.2.

Some participants from both groups regretted that they were not explicitly asked for suggestions on how to improve the signing quality of the avatar. Systematically collecting such input was intentionally left outside the scope of the survey, to keep sessions manageable in terms of time and cognitive effort. Regardless, some participants from the

⁸ For an introduction to ordered probit mixed effects regression models, see Kruschke [34].

GUIDED group did provide suggestions, and the experimenters made note of these.

Opinions on the added value of being able to respond in NGT, as opposed to entering responses textually, were divided. On the one hand, 92% of participants in the UNGUIDED group, who had to enter their responses textually, indicated that it would *not* have been easier to respond in NGT. On the other hand, while this question was not explicitly posed to participants in the GUIDED group, since they did respond in NGT, 50% of these participants spontaneously mentioned that they appreciated being able to use NGT throughout the survey and not having to enter their responses textually. We note that these findings may well be due, at least in part, to the fact that the participants themselves determined whether they would be part of the GUIDED or the UNGUIDED group.

4.4 Discussion

We now discuss the results of the survey in light of our three main research questions, compare our results to related work, reflect on the limitations of the conclusions that can be drawn from these results, and suggest some avenues for future work.

RQ1: Comprehension of individual signs. The comprehension rate of individual signs performed by the avatar (86.93%) was lower than the comprehension rate of individual signs performed by a human signer (91.34%), but both rates are quite high and the difference between them is rather small. A closer look at the most common mistakes revealed, among other things, that the *mouthings* produced by the avatar can be particularly confusing. For instance, SOMETIMES was misinterpreted 54.5% of the time as MAYBE, which has the same manual component but different mouth-ing, and similarly, RESULT was misinterpreted 22.7% of the time as PASSED. Such misinterpretations did not arise when the signs were performed by a human signer.

The JASigning avatar engine currently offers limited possibilities to produce natural-looking mouthings. More specifically, the engine currently allows for a specification of mouthings in SAMPA notation (Speech Assessment Methods Phonetic Alphabet). This is a phonetic notation system: each SAMPA symbol corresponds to a particular phoneme. There is, however, no one-to-one mapping between phonemes and mouth shapes. For instance, the ‘s’ in ‘sun’ and the ‘s’ in ‘silver’ involve the same phoneme but different mouth shapes because the next vowel is anticipated. This makes it difficult to generate correct mouthings in JASigning. In future work, it would therefore be advisable to reconsider the way in which mouth shapes are handled in the engine. This line of future work may take inspiration from lipsync algorithms for game characters (e.g., [10]).

RQ2: Sentence comprehension and clarity. Sentences signed by the avatar had a comprehension rate of 74.62% and a mean clarity score of 6.4, while ones signed by a human had a comprehension rate of 97.35% and a mean clarity score of 9.0. So here we saw a larger contrast between the avatar and the human signer than in the case of individual signs.

Taking a closer look at the 12 sentences that were used in the survey, we found that 2 of these received particularly low scores when signed by the avatar. These sentences were ones which contained several individual signs with low comprehension rates. Improving the way in which these particular signs are rendered by the avatar may well improve the scores of the sentences that contained them as well.

However, we noted that even if these two negative outlier sentences were to be disregarded, there is still a substantial difference between the comprehension rates and clarity scores for the avatar and those for the human signer. To close this gap, it will not suffice to improve the rendering of some individual signs. Rather, more fundamental improvements of the underlying avatar technology will be necessary. Based on the feedback provided by participants during the survey, we note that sentence prosody (the relative speed and intensity of the signs in the sentence) and the transitions between signs are important elements that strongly influence comprehensibility. The JASigning avatar engine and the SiGML formalism that it makes use of currently offer limited possibilities to control prosody and transitions. To make the engine suitable for practical applications, these functionalities need to be extensively developed in future work. An alternative would be to explore an approach that makes use of motion capture instead of scripted animation (see, e.g., [23]), or a combination of the two.

RQ3: Effect of guidance. In previous work, the evaluation of signing avatars typically involved on-site interaction between experimenters and participants ([6, 9, 29, 31, 50], with exceptions noted in [24, 42, 47]). However, the COVID-19 pandemic made it necessary for us to turn to online procedures, and it is to be expected that in the future researchers may sometimes want to employ online procedures as well. This new experimental setting raises methodological issues. In particular, one basic design choice that needs to be made concerns the online presence of the experimenters while the participants take the survey. Outside the domain of sign language technology, it is most common in online quantitative surveys for the experimenters *not* to be present. For the specific purpose of evaluating signing avatars, however, this has a possible disadvantage, namely that participants have to enter their responses textually rather than in sign language. This may be dispreferred, at least for some participants. To circumvent this potential disadvantage, we offered our participants a choice between a *guided* and an

unguided version of the survey. A comparison between the results from the GUIDED and the UNGUIDED group, as well as the feedback that participants from both groups provided on the setup of the survey, may inform the design of future online evaluation studies.

The main difference between the GUIDED and the UNGUIDED group was that comprehension rates, both for individual signs and for sentences, were generally *lower* in the GUIDED group. This was an unexpected result for us. If anything, we had expected comprehension rates to be higher in the GUIDED group. There is, however, a plausible explanation for why there was in fact a difference in the opposite direction. Namely, it may be that the presence of the experimenters caused a certain amount of social pressure for participants in the GUIDED group. For instance, they may have felt that it would be a burden to ask the experimenters to replay a video, or they may have felt pressure to understand sentences and signs on the first try. They may even have experienced the experiment partly as a memory task rather than a pure comprehension task. Participants from the UNGUIDED group presumably did not experience any such pressures, and may have felt more freedom to replay videos as often as needed. This may be one of the reasons that comprehension rates were higher in the UNGUIDED group.

As for the feedback we received from both groups on the setup of the survey, 92% of participants from the UNGUIDED group indicated that it would *not* have been easier to respond in NGT instead of entering their responses textually. On the other hand, while this question was not explicitly posed to participants in the GUIDED group, since they did respond in NGT, 50% of them spontaneously mentioned that they appreciated being able to use NGT throughout the survey and not having to enter their responses textually.

Overall, then, it is not the case that a GUIDED online procedure is to be strictly preferred over an UNGUIDED procedure for the evaluation of signing avatars, nor vice versa. Both methods have advantages and disadvantages, and which format works best differs from one participant to another. For future work, we can therefore only recommend that, whenever a choice needs to be made between a GUIDED and an UNGUIDED setup, the potential advantages and disadvantages of both options are carefully weighed.

Comparison to results obtained in earlier work. To our knowledge previous studies did not determine comprehension rates of individual signs or sentences in the way that we did (asking participants to provide the interpretation of a given sign or sentence). So as far as comprehension rates are concerned, our results cannot be compared with previous work (for qualitative evaluations of comprehensibility, see, e.g., [6, 9]).

As for clarity scores, our results can be compared to those of Quandt et al. [42], who adopted a similar approach and had similar target items, although the language was different

(ASL) and the study was restricted to individual signs (no full sentences). Concretely, they evaluated eight individual signs produced by a human, a motion capture avatar, and a computer-scripted avatar. For each sign, participants answered the question ‘*The signing was easy for me to understand*’ on a 5-point Likert scale, with ‘Strongly Disagree’ as 1, ‘Neutral’ as 3, and ‘Strongly Agree’ as 5. The clarity scores obtained by Quandt et al. [42] for individual signs are in line with the scores we obtained for sentences. Namely, they report that the human signer scored considerably higher (mean 4.62, standard deviation 0.56) than the computer-scripted avatar (mean 2.62, standard deviation 1.13). In addition, they report that the motion capture avatar scored higher than the computer-scripted avatar, but lower than the human signer (mean 3.79, standard deviation 0.72). The latter finding is not directly comparable with our results, since we did not evaluate a motion capture avatar.

Limitations. There are various factors that limit the generalizability of the results of our survey. First, the design of our survey allowed for two kinds of learning effect to arise. On the one hand, each participant first saw 12 sentences signed by the avatar and then the same 12 sentences signed by a human signer. This may in part explain why the human signer received higher comprehension rates and clarity scores than the avatar. On the other hand, some individual signs appeared in more than one sentence and were therefore seen more often than other signs. This may have positively affected their comprehension rate.

Second, the differences we found between the GUIDED and the UNGUIDED group may in part be due to the fact that participants chose themselves whether to take the GUIDED or the UNGUIDED version of the survey. For instance, this may in part explain why 92% of participants from the UNGUIDED group indicated that it would not have been easier to respond in NGT instead of entering their responses textually.

Finally, some more general limitations apply: our survey involved only a small number of sentences and signs, a rather small participant pool, and was conducted in a controlled environment rather than in a real-life setting. When interpreting our results, these factors should be kept in mind, and future work should investigate how well the results generalize.

5 Conclusion

We have investigated the potential of automated text-to-sign translation to address the challenges that the COVID-19 pandemic implies for the communication between healthcare professionals and deaf patients. We motivated a modular approach to automated text-to-sign translation, and implemented a first prototype system. We conducted a survey among potential end-users to evaluate the comprehensibility and clarity of the avatar. Moreover, we investigated whether the possibility to interact with the experimenters during the

survey and to provide responses in NGT rather than having to enter them textually affected the results when conducting a survey of this sort in an online environment. We have discussed various prospects and limitations of the prototype system we built and of the results of our survey.

For the approach taken here to become viable in practice, the JASigning avatar engine needs to be substantially further developed. At the level of individual signs, the engine should allow for more subtle body movements and facial expressions, and the system for encoding mouth shapes should be revised. At the level of sentences, more control is needed to adapt the relative speed and intensity of the different signs within a sentence (prosody) and to make transitions between signs more natural and smooth. An alternative is to explore an approach based on motion capture instead of scripted animation, or a combination of both.

Finally, we believe that future projects will strongly benefit from a more inclusive and more iterative design process, involving a multi-disciplinary team with a strong representation of deaf researchers and domain experts. The design and implementation phase of the present project was carried out under great time pressure, given the urgency of the issue we aimed to address, and with limited resources. In future work, there should be several design iterations, the design team should include deaf specialists on sign language and Deaf Studies from the start, and focus groups should be organized to receive input from a larger group of potential end-users, aiming for maximal diversity in terms of age, region, and level of education.

Appendix A: Survey

The instruction videos can be found at [14].

Appendix A.1: Part 1—general questions

The videos of the general questions of the survey can be found at [15].

1. What is your age?

- A. 18–30
- B. 31–40
- C. 41–50
- D. 51+

2. What is your gender?

- A. Male
- B. Female

3. Which region are you from?

- A. Central (Utrecht, Flevoland)
- B. North (Groningen, Friesland, Drenthe)
- C. East (Gelderland, Overijssel)
- D. South (Brabant, Zeeland, Limburg)
- E. West (Noord-Holland, Zuid-Holland)

4. What is your mother tongue?

- A. NGT
- B. Dutch
- C. Both
- D. Other, namely...

5. How often do you use NGT?

- A. Daily
- B. Regularly (a few days a week)
- C. Occasionally (a few days a month)
- D. Rarely (a few days a year)

6. How often do you make use of an interpreter?

- A. Daily
- B. Regularly (a few days a week)
- C. Occasionally (a few days a month)
- D. Rarely (a few days a year)

7. How often do you communicate with people that do not use sign language, without making use of an interpreter NGT?

- A. Daily
- B. Regularly (a few days a week)
- C. Occasionally (a few days a month)
- D. Rarely (a few days a year)

8. How do you communicate with people that do not use sign language when an interpreter cannot be present? (multiple options possible)

- A. Lipreading
- B. Speech recognition through the phone
- C. Writing on paper/phone
- D. Signing
- E. Using voice
- F. Other, namely...

Appendix A.2: Part 2—avatar comprehension and clarity

The format of the questions in Part 2 is shown in Fig. 6. The videos of the avatar animations for this part of the survey can be found at [16].

1. Did you sleep well?
YOU GOOD SLEEP PALMS-UP
2. Do you use any medications?
YOU MEDICINE USE PALMS-UP
3. What are you allergic to?
YOU ALLERGIC FOR WHAT PALMS-UP
4. Please stay in bed.
YOU PLEASE BED STAY
5. I will come back later.
I LATER BACK
6. A colleague will come by soon to draw blood.
SOON COLLEAGUE COME BLOOD DRAW
7. Who is your general practitioner?
YOU GENERAL PRACTITIONER WHO INDEX
8. Do you have hearing aids or a cochlear implant?
HEARING-AID OR COCHLEAR-IMPLANT HAVE YOU INDEX PALMS-UP
9. Have you had a Corona test in the past seven days?
PAST SEVEN DAY YOU ALREADY CORONA TEST INDEX PALMS-UP
10. Your Corona test results are negative.
YOU CORONA TEST RESULTS NEGATIVE
11. Sometimes the test is wrong, therefore we have to do more research.
SOMETIMES TEST WRONG THEREFORE WE MORE RESEARCH HAVE TO
12. Sorry, I'm failing to insert the intravenous drip. I'm calling a colleague.
SORRY INTRAVENOUS DRIP FAIL I COLLEAGUE CALL

Appendix A.3: Part 3—human signer comprehension and clarity

The format of the questions in Part 3 is shown in Fig. 6. The videos of the human signer for this part of the survey can be found at [17].

1. Did you sleep well?
YOU GOOD SLEEP INDEX
2. Do you use any medications?
YOU USE MEDICINE INDEX
3. What are you allergic to?
YOU ALLERGIC FOR WHAT PALMS-UP
4. Please stay in bed.
YOU BED STAY
5. I will come back later.
I LATER BACK

6. A colleague will come by to draw blood.
COLLEAGUE COME BLOOD DRAW
7. Who is your general practitioner?
YOU GENERAL PRACTITIONER INDEX WHO INDEX PALMS-UP
8. Do you have hearing aids or a cochlear implant?
HEARING-AID OR COCHLEAR-IMPLANT HAVE YOU
9. Have you had a Corona test in the past seven days?
PAST SEVEN DAY YOU ALREADY CORONA TEST PALMS-UP
10. Your Corona test results are negative.
YOU CORONA TEST RESULTS YOU NEGATIVE
11. Sometimes the test is wrong, therefore we have to do more research.
SOMETIMES TEST INDEX WRONG WE MORE RESEARCH HAVE TO
12. Sorry, I'm failing to insert the intravenous drip. I'm calling a colleague.
SORRY I FAIL YOU INTRAVENOUS DRIP I COLLEAGUE CALL

Appendix A.4: Part 4—final questions

The videos of the final questions of the survey can be found at [15].

Appendix A.4.1: Attitude toward signing avatars

1. There has not been much research on avatar technology for translating text to sign language. Do you think this research should be continued and this technology should be developed further?
 - A. Yes, because...
 - B. No, because...
2. In which situations do you think that avatar technology for translating text to sign language can help? (multiple answers possible)
 - A. For translating travel information in trains and on train stations
 - B. For translating travel information in airplanes and at airports
 - C. As support for people who want to learn sign language
 - D. In hospitals during the COVID-19 crisis
 - E. Other situations, namely...
 - F. I think that this technology is not helpful in any situation

Appendix A.4.2: Feedback on methodology

1. Were the questions in this study clearly formulated? (scale of 0–10)
2. Were the questions in this study easy to answer? (scale of 0–10)

3. Would it have been easier to answer questions in NGT rather than in Dutch text? (only for the UNGUIDED group)
 - A. Yes, because...
 - B. No, because...
4. Which aspects of the setup of the survey were pleasant?
5. Which aspects of the setup of the survey could be improved?

Acknowledgements We are grateful to Tashi Bradford, Richard Cokart, Bastien David, John Glauert, Lisa Hinderks, Richard Kenaway, Lisa van der Mark, Marta Morgado, Joni Oyserman, Marijke Scheffener, Anique Schüller, and Roos Wattel for their help at various stages of this project.

Funding We gratefully acknowledge financial support from the Netherlands Organization for Innovation in Healthcare (ZonMw, Grant No. 10430042010027), the Netherlands Organization for Scientific Research (NWO, Grant No. V1.C.201.014), and the European Research Council (Grant No. 680220).

Data availability The videos used in the survey, as well as the raw survey data, the R analysis script, and all SiGML files pertaining to the prototype translation system that support the findings of this study are available in Figshare [13–20].

Declarations

Conflict of interest On behalf of all authors, the corresponding author states that there is no conflict of interest.

Consent to participate Informed consent was obtained from all subjects for being included in the study.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

1. Baker, A., van den Bogaerde, B., Pfau, R., et al.: The Linguistics of Sign Languages: An Introduction. John Benjamins (2016). <https://doi.org/10.1075/z.199>
2. Barr, D.J., Levy, R., Scheepers, C., et al.: Random effects structure for confirmatory hypothesis testing: keep it maximal. *J. Mem. Lang.* **68**(3), 255–278 (2013). <https://doi.org/10.1016/j.jml.2012.11.001>
3. Bates, D., Maechler, M., Bolker, B., et al.: lme4: Linear mixed-effects models using Eigen and S4. <http://CRAN.R-project.org/package=lme4>, r package version 1.1 (2014)
4. Coerts, J.: Nonmanual grammatical markers: an analysis of interrogatives, negations and topicalisations in Sign Language of the Netherlands. PhD thesis, University of Amsterdam (1992)
5. Courty, N., Gibet, S.: Why is the creation of a virtual signer challenging computer animation? In: Motion in Games: Third International Conference (MIG 2010), pp. 290–300. Springer, Berlin (2010). https://doi.org/10.1007/978-3-642-16958-8_27
6. David, B.V.C., Bouillon, P.: Prototype of automatic translation to the Sign Language of French-speaking Belgium. Evaluation by the Deaf community. *Model. Meas. Control C* **79**(4), 162–167 (2018). https://doi.org/10.18280/mmc_c.790402
7. de Meulder, M., Haualand, H.: Sign language interpreting services: a quick fix for inclusion? *Transl. Interpret. Stud. J. Am. Transl. Interpret. Stud. Assoc.* **16**(1), 19–40 (2021). <https://doi.org/10.1075/tis.18008.dem>
8. de Vos, C., van der Kooij, E., Crasborn, O.: Mixed signals: combining linguistic and affective functions of eyebrows in questions in Sign Language of the Netherlands. *Lang. Speech* **52**(2–3), 315–339 (2009). <https://doi.org/10.1177/0023830909103177>
9. Ebling, S., Glauert, J.: Building a Swiss German Sign Language avatar with JASigning and evaluating it among the Deaf community. *Univ. Access Inf. Soc.* **15**(4), 577–587 (2016). <https://doi.org/10.1007/s10209-015-0408-1>
10. Edwards, P., Landreth, C., Fiume, E., et al.: JALI: an animator-centric viseme model for expressive lip synchronization. *ACM Trans. Graph.* **35**(4), 1–11 (2016). <https://doi.org/10.1145/2897824.2925984>
11. Elliott, R., Glauert, J., Jennings, V., et al.: An overview of the SiGML notation and SiGML signing software system. In: Workshop on the Representation and Processing of Sign Languages at the 4th International Conference on Language Resources and Evaluation (LREC 2004). European Language Resources Association, pp. 98–104 (2004)
12. Esselink, L.: Lexical resources for sign language synthesis: the translation of Dutch to Sign Language of the Netherlands, bachelor's thesis. University of Amsterdam. <https://scripties.uba.uva.nl/search?id=715792> (2020)
13. Esselink, L., Roelofsen, F., Dotlačil, J., et al.: Analysis Script. (2023) <https://doi.org/10.21942/uva.22280584>
14. Esselink, L., Roelofsen, F., Dotlačil, J., et al.: Instruction videos (2023). <https://doi.org/10.21942/uva.22276648>
15. Esselink, L., Roelofsen, F., Dotlačil, J., et al.: Part 1 and part 4—questions for participants (2023). <https://doi.org/10.21942/uva.22280545>
16. Esselink, L., Roelofsen, F., Dotlačil, J., et al.: Part 2—Avatar (2023). <https://doi.org/10.21942/uva.22280569>
17. Esselink, L., Roelofsen, F., Dotlačil, J., et al.: Part 3—Signer (2023). <https://doi.org/10.21942/uva.22280575>
18. Esselink, L., Roelofsen, F., Dotlačil, J., et al.: SiGML Codes—variable sentences (2023). <https://doi.org/10.21942/uva.22280590>
19. Esselink, L., Roelofsen, F., Dotlačil, J., et al.: SiGML Codes—words and sentences (2023). <https://doi.org/10.21942/uva.22280587>
20. Esselink, L., Roelofsen, F., Dotlačil, J., et al.: Survey data (2023). <https://doi.org/10.21942/uva.22280581>
21. Feller, J., Holzinger, D., Pollard, R.: Mental health of deaf people. *Lancet* **379**(9820), 1037–1044 (2012). [https://doi.org/10.1016/S0140-6736\(11\)61143-4](https://doi.org/10.1016/S0140-6736(11)61143-4)
22. Gelman, A., Hill, J.: Data Analysis Using Regression and Multilevel/Hierarchical Models. Analytical Methods for Social Research. Cambridge University Press, Cambridge (2006). <https://doi.org/10.1017/CBO9780511790942>
23. Gibet, S.: Building French Sign Language motion capture corpora for signing avatars. In: Workshop on the Representation and Processing of Sign Languages: Involving the Language Community at the 11th International Conference on Language Resources and

- Evaluation (LREC 2018). European Language Resources Association (2018)
24. Gibet, S., Courty, N., Duarte, K., et al.: The SignCom system for data-driven animation of interactive virtual signers: methodology and evaluation. *ACM Trans. Interact. Intell. Syst. (TiIS)* **1**(1), 1–23 (2011). <https://doi.org/10.1145/2030365.2030371>
25. Glauert, J., Elliott, R.: Extending the SiGML notation—a progress report. In: Second International Workshop on Sign Language Translation and Avatar Technology (SLTAT 2011), vol. 23. Association for Computing Machinery (2011)
26. Grote, H., Izagaren, F.: COVID-19: the communication needs of D/deaf healthcare workers and patients are being forgotten. *BMJ* **369** (2020). <https://doi.org/10.1136/bmj.m2372>
27. Hanke, T.: HamNoSys-representing sign language data in language resources and language processing contexts. In: Workshop on the Representation and Processing of Sign Languages at the 4th International Conference on Language Resources and Evaluation (LREC 2004), vol. 4, pp. 1–6. European Language Resources Association (2004)
28. Hou, L., de Vos, C.: Classifications and typologies: labeling sign languages and signing communities. *J. Sociolinguistics* **26**(1), 118–125 (2022). <https://doi.org/10.1111/josl.12490>
29. Huenerfauth, M.: Generating American Sign Language classifier predicates for English-to-ASL machine translation. PhD thesis, University of Pennsylvania (2006)
30. Jennings, V., Elliott, R., Kennaway, R., et al.: Requirements for a signing avatar. In: Workshop on Corpora and Sign Language Technologies at the 7th International Conference on Language Resources and Evaluation (LREC 2010), pp. 33–136. European Language Resources Association (2010)
31. Kacorri, H., Huenerfauth, M., Ebling, S., et al.: Demographic and experiential factors influencing acceptance of sign language animation by deaf users. In: Proceedings of the 17th International ACM SIGACCESS Conference on Computers & Accessibility. Association for Computing Machinery, ASSETS'15, pp. 147–154 (2015). <https://doi.org/10.1145/2700648.2809860>
32. Kennaway, R., Glauert, J., Zwitserlood, I.: Providing signed content on the internet by synthesized animation. *ACM Trans. Comput. Hum. Interact.* **14**(3), 1–29 (2007). <https://doi.org/10.1145/1279700.1279705>
33. Klomp, U.: A descriptive grammar of Sign Language of the Netherlands. PhD thesis, University of Amsterdam (2021)
34. Kruschke, J.K.: Doing Bayesian Data Analysis: A Tutorial with R and BUGS. Academic Press Inc, London (2010)
35. Kusters, A., Lucas, C.: Emergence and evolutions: introducing sign language sociolinguistics. *J. Sociolinguistics* **26**(1), 84–98 (2022). <https://doi.org/10.1111/josl.12522>
36. McKee, M., Moran, C., Zazove, P.: Overcoming additional barriers to care for deaf and hard of hearing patients during COVID-19. *JAMA Otolaryngol. Head Neck Surg.* **146**(9), 781–782 (2020). <https://doi.org/10.1001/jamaoto.2020.1705>
37. McKee, M.M., Paasche-Orlow, M.K., Winters, P.C., et al.: Assessing health literacy in deaf American Sign Language users. *J. Health Commun.* **20**(sup2), 92–100 (2015). <https://doi.org/10.1080/10810730.2015.1066468>
38. Middleton, A., Niruban, A., Girling, G., et al.: Communicating in a healthcare setting with people who have hearing loss. *BMJ* **341** (2010). <https://doi.org/10.1136/bmj.c4672>
39. Napier, J., Kidd, M.R.: English literacy as a barrier to health care information for deaf people who use Auslan. *Aust. Fam. Physician* **42**(12), 896–899 (2013)
40. Prillwitz, S., Leven, R., Zienert, H., et al.: HamNoSys version 2: Hamburg Notation System for Sign Languages—An Introductory Guide. Hamburg Signum Press (1989)
41. Prins, M., Janssen, J.B.: Automated sign language, TNO technical report (2014)
42. Quandt, L.C., Willis, A., Schwenk, M., et al.: Attitudes toward signing avatars vary depending on hearing status, age of signed language acquisition, and avatar type. *Front. Psychol.* **13** (2022). <https://doi.org/10.3389/fpsyg.2022.730917>
43. Rayner, M., Bouillon, P., Ebling, S., et al.: An open web platform for rule-based speech-to-sign translation. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, pp. 162–168. Association for Computational Linguistics (2016). <https://doi.org/10.18653/v1/P16-2027>
44. Roelofsen, F., Esselink, L., Mende-Gillings, S., et al.: Online evaluation of text-to-sign translation by deaf end users: Some methodological recommendations (short paper). In: Proceedings of the 1st International Workshop on Automatic Translation for Signed and Spoken Languages (AT4SSL), pp. 82–87. Association for Machine Translation in the Americas (2021)
45. Roelofsen, F., Esselink, L., Mende-Gillings, S., et al.: Sign language translation in a healthcare setting. In: Proceedings of the Translation and Interpreting Technology Online Conference, pp. 110–124 (2021). https://doi.org/10.26615/978-954-452-071-7_013
46. Schermer, T., Koolhof, C.: Basiswoordenboek Nederlandse Gebarentaal. Van Dale (2009)
47. Schnepf, J., Wolfe, R., Shiver, B., et al.: SignQUOTE: A remote testing facility for eliciting signed qualitative feedback. In: Second International Workshop on Sign Language Translation and Avatar Technology (SLTAT 2011), vol. 23. Association for Computing Machinery (2011)
48. Smeijers, A.S., Roelofsen, F.: Communicatiebehoefte en ervaringen van dove patiënten in Nederland tijdens de COVID-19 pandemie (Communication needs and experiences of deaf patients in the Netherlands during the COVID-19 pandemic) (2021)
49. Smeijers, A.S., Ens-Dokkum, M.H., van den Bogaerde, B., et al.: Clinical practice: the approach to the deaf or hard-of-hearing paediatric patient. *Eur. J. Pediatr.* **170**(11), 1359 (2011). <https://doi.org/10.1007/s00431-011-1530-6>
50. Smith, R.G., Nolan, B.: Emotional facial expressions in synthesised sign language avatars: a manual evaluation. *Univ. Access Inf. Soc.* **15**(4), 567–576 (2016). <https://doi.org/10.1007/s10209-015-0410-7>
51. Winter, B.: Statistics for Linguists: An Introduction Using R. Routledge (2019). <https://doi.org/10.4324/9781315165547>
52. Wolfe, R., Cook, P., McDonald, J.C., et al.: Linguistics as structure in computer animation: toward a more effective synthesis of brow motion in American Sign Language. *Sign Lang. Linguist.* **14**(1), 179–199 (2011). <https://doi.org/10.1075/sll.14.1.09wol>
53. Zwitserlood, I.: Synthetic signing. In: The World of Content Creation, Management, and Delivery (IBC 2005), pp. 352–357 (2005)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.