

## Abstract

*Objective:* The speech intelligibility benefit of visual speech cues during oral communication is well-established (Sumby and Pollack, 1954). Therefore, an ecologically valid approach of auditory assessment should include the processing of both auditory and visual speech cues. This study describes the development and evaluation of a virtual human speaker designed to present speech auditory-visually. *Design:* A male and female virtual human speaker were created and evaluated in two experiments: a visual-only speech reading test of words and sentences and an auditory-visual speech intelligibility sentence test. *Study sample:* A group of five hearing, skilled speech reading adults participated in the speech reading test whereas a group of young normal hearing participants (N=35) was recruited for the intelligibility test. *Results:* Skilled speech readers correctly identified 57 to 67% of the words and sentences uttered by the virtual speakers. The presence of the virtual speaker improved the speech intelligibility of sentences in noise by 1.5 to 2 dB. *Conclusions:* These results demonstrate the potential applicability of virtual humans in future auditory-visual speech assessment paradigms.

## 1 **Introduction**

2 People communicate continuously, in either a spoken or a written way. During oral  
3 communication, listeners do not solely rely on auditory input. In addition, they process  
4 visual speech cues including mouth movements, which improves both the detection and  
5 discrimination of speech (Bernstein, Auer and Takayanagi, 2004; Grant and Seitz,  
6 2000a; Ma et al., 2009; Schwartz, Berthommier and Savariaux, 2004; Sumby and  
7 Pollack, 1954). This visual speech benefit even sustains when listeners do not directly  
8 face the speaker, as shown by research of Grange and Culling (2016). They presented  
9 sentences together with spatially separated noise and demonstrated that head  
10 orientations up to 30° sideways of the speaker did not influence the speech intelligibility  
11 benefit obtained by speech reading. In addition, speech reading is known to be  
12 automatic and therefore an intrinsic part of speech processing (Rosenblum, 2008;  
13 Woodhouse, Hickson and Dodd, 2009). The multimodal nature of speech processing is  
14 established by brain imaging studies showing activation of the superior temporal sulcus,  
15 a brain area related to multimodal integration (Calvert, Campbell and Brammer, 2000;  
16 Wright et al., 2003; Okada et al., 2013). Also, research shows that facial movements can  
17 alter activity in the sensory auditory cortex, in both humans and primates (Ghazanfar,  
18 2009; Chandrasekaran, Lemus and Ghazanfar, 2013). More evidence for the high  
19 impact of visual cues during speech perception is provided by the fact that even  
20 incongruent visual facial cues influence the perception of speech segments. This is  
21 illustrated behaviourally by the McGurk effect where an auditory /ba/ syllable is often  
22 perceived as /da/ when presented together with an incongruent visual /ga/ (McGurk and  
23 MacDonald, 1976). A neuroimaging study by Arnal, Wyart and Giraud (2011) revealed  
24 deviant neural processes when people were faced with incongruent cues, compared to  
25 congruent auditory-visual speech processing. Mismatches between visual speech signals

1 and auditory information activated fast and local brain regions around the primary  
2 sensory areas, whereas matching cues activated slower, higher order language regions in  
3 the brain to process speech at a semantical and syntactical level. Visual cues do not only  
4 influence speech perception per se, they also affect the cognitive load during listening  
5 by altering the listening effort during speech perception (Gosselin and Gagné, 2011;  
6 Stevens et al., 2013).

7       When the goal is to obtain an ecologically valid measure of daily speech  
8 processing abilities, adding accurate visual speech cues to the speech signal is of major  
9 concern. A straightforward way of implementing these cues into an auditory speech test  
10 is by presenting videomaterials. Unfortunately, recording videos of human speakers  
11 does not allow application of existing speech materials, even though many optimised  
12 and high quality speech intelligibility-measurement materials and procedures have been  
13 developed during the last two decades. New recordings would require a time-consuming  
14 optimisation, evaluation and validation of the speech material (Akeroyd et al., 2015),  
15 including characterisation of speech perception characteristics like the speech reception  
16 threshold (SRT), i.e. the signal-to-noise ratio (SNR) at which listeners understand 50%  
17 of the presented speech signal. In addition, visual cues that could steer the gaze of the  
18 listener, like lighting conditions, clothing or head movements of the speaker, are hard to  
19 control in videos. Hence, the interest in implementing speaking virtual human-like 3D  
20 characters in auditory assessment emerged. In this paper, we will refer to these 3D  
21 characters as “virtual humans”, even though most of the time the character consists of a  
22 dynamic head only.

23       In the last decades, many research projects related to human-computer  
24 interactions, language or auditory-visual speech processing and automated dialogue  
25 systems implemented auditory-visual speech synthesis. Human-computer interactions

1 have for example been integrated in a tutoring system aimed at enhancing science  
2 learning in elementary school (Ward et al., 2013). In this project, children who received  
3 face-to-face tutoring from a virtual human speaker developed by Ma et al. (2006),  
4 achieved larger learning gains compared to students in control classrooms. Virtual  
5 human speakers were also implemented in speech and language tutoring applications  
6 like training systems on pronunciation (Dey, Maddock and Nicolson, 2010; Peng et al.,  
7 2018) or on speech perception and production for hearing impaired individuals  
8 (Massaro and Light, 2004). More recently, Schreitmüller et al. (2017) demonstrated the  
9 applicability of “MASSY”, a virtual human developed by Fagel and Clemens (2004), to  
10 investigate lipreading, visual benefit and auditory-visual integration of speech in both  
11 normal hearing and hearing impaired individuals. Even telephone conversations have  
12 been extended with virtual humans to augment speech intelligibility (Siciliano et al.,  
13 2003; Salvi et al., 2009).

14 Visual modelling of virtual humans can be divided in three categories: terminal-  
15 analog systems, anatomy-based systems or performance-driven animations (Mattheyses  
16 and Verhelst, 2015). The first approach refers to virtual humans in which the visual  
17 speech gestures are simulated by directly controlling the facial parameters that are  
18 linked to visual articulators (Massaro and Cohen, 1990; Salvi et al., 2009), whereas in  
19 anatomy-based systems the deformations are the result of an indirect modelling of the  
20 facial anatomy like skin, muscles and bones (Sifakis, Neverov and Fedkiw, 2005). In  
21 performance driven animations, facial motions are learned by mapping gestures from a  
22 real human speaker onto a 3D facial map by means of e.g. motion capture techniques  
23 (Williams, 1990; Fagel, Bailly and Elisei, 2007). Two approaches can also be  
24 combined, as is the case for “MASSY”, where motion capture techniques are applied to  
25 determine the facial parameters of a terminal-analog system (Fagel and Clemens, 2004).

1 All visual synthesis methods have their advantages and selection of one of them is  
2 based on an intrinsic trade-off between the desired realistic picture and intelligibility of  
3 the virtual human on the one hand and the ease of creating it on the other hand  
4 (Mattheyses and Verhelst, 2015). According to Beskow (2004), terminal-analog  
5 approaches can outperform performance-driven methods in terms of speech  
6 intelligibility, although the latter tend to provide more natural mouth movements.  
7 However, anatomy-based systems or performance-driven animations typically require  
8 larger computational work.

9       The aim of the present report is to describe the development of a speaking  
10 virtual human and validate the quality of the provided visual facial information to define  
11 its potential applicability in future auditory assessment paradigms. Taking into account  
12 that we have a large amount of speech materials at our disposal, we opted for a modified  
13 terminal-analog method that allowed off-line manual mapping of speech phonemes to  
14 predefined mouth movements. In our case, this method was more efficient compared to  
15 e.g. online terminal-analog mapping or performance-driven auditory-visual mapping  
16 methods such as “MASSY”. Furthermore, we aimed for a realistic looking virtual  
17 human that can be easily integrated in 3D virtual environments to assess, for example,  
18 auditory-visual speech perception in complex daily listening situations. The need for  
19 multimodal, ecologically valid listening situations that can be applied in clinical  
20 practice is high among both researchers and clinicians (Pichora-Fuller et al., 2016).  
21 Therefore, realistic and dynamic 3D listening situations were created at our research  
22 department and the virtual human described in this study was especially designed to be  
23 implemented in these environments. We opted for modern, well-supported 3D rendering  
24 software systems to ensure the virtual humans and listening scenarios can be modified  
25 easily, according to the needs of researchers and clinicians. Furthermore, our virtual

human is innovative in that it represents a full virtual body instead of a talking head only, which is more realistic to implement in a 3D listening scenario.

We predicted that the mouth movements of our virtual human would be comparable to those of a real human speaker. The extent to which the virtual human features are realistic was evaluated in a speech reading experiment. The amount of visual benefit in speech understanding was assessed by comparing an auditory-visual (AV) speech intelligibility (SI) test to an auditory-only (AO) SI test.

## Methods

### *Participants*

Two groups of participants were recruited. First, a total of five hearing adults (one male, four female, aged 24, 25, 26, 29 and 56) participated in the speech reading study. All were native Flemish/Dutch speakers and had minimum one year of experience in speech reading by either following a speech reading course or teaching it to others. Second, a group of 35 young adults (13 male, 22 female) aged 18 to 30 (mean age = 22 years, SD = 2 years) took part in the AV SI experiment. All had bilateral normal audiometric pure-tone thresholds less than or equal to 20 dB HL at 250, 500, 1000, 2000, 4000 and 8000 Hz. Individuals of both groups reported normal or corrected vision. The Medical Ethical Committee of the University Hospitals and University of Leuven approved the experiment and materials (approval number B322201526677) and written consent was obtained from each participant.

### *Stimuli*

#### *Auditory stimuli*

In order to capture a wide range of speech materials in the speech reading experiment,

both meaningful Flemish/Dutch words and sentences were presented to the participants. Monosyllabic (MS) consonant-vocal-consonant (CVC) words were selected from the NVA- and Lilliput-speech material (Wouters, Damman and Bosman, 1994; van Wieringen and Wouters, 2015) in order to cover all Flemish/Dutch phonemes. They were grouped together in three lists of 14 CVC-words. Four original lists of 10 disyllabic (DS) words were selected from the BLU-material (Wouters, Damman and Bosman, 1994) and six lists of 10 sentences from the male and female LIST (Jansen et al., 2014; van Wieringen and Wouters, 2008). The LIST was opted for because of its high standardisation and validation and widespread use in clinical practice in Flanders since 2008. The LIST-sentences are comparable to the American English HINT-sentences (Nilsson, Soli and Sullivan, 1994) and are semantically and syntactically predictable sentences, representative of daily communication, such as “Milk is healthy”. The sentences do not include questions, exclamations or proverbs and exist of two to three keywords (i.e. nouns, adjectives, numerals, adverbs and main verbs) and six to seven other words (i.e. articles, pronouns, prepositions, conjunctions, interjections and linking and auxiliary verbs). The total amount of syllables (89 to 90) and keywords (32 to 33) is balanced between lists. For the AV SI experiment, sentence lists from the female LIST material were randomly selected and presented in both speech weighted noise (SWN) and multitalker babble noise (MTN) (Francart, van Wieringen and Wouters, 2011).

### *Visual models*

Three visual models were used to present the stimuli in the speech reading experiment (Figure 1, A). One 25-year-old male speech therapist was video-recorded uttering the different speech materials. He was instructed to pronounce the stimuli with the same

1 speaking rate as the original material. Recordings were made with a Canon XA10  
2 professional camcorder and a shotgun microphone (Audio-Technica AT875R). The  
3 speaker was seated in a silent room in front of a neutrally colored background and was  
4 captured from shoulders to head.

5 In addition to the video model, two virtual humans were created by means of  
6 open source 3D graphics and animation software packages Blender (Blender, 2015,  
7 Version 2.76b) and Make Human (Make Human, 2015, Version 1.0.2). Blender and  
8 Make Human provide a predefined virtual human model consisting of a mesh, which  
9 simulates the skin of the virtual human, and bones attached to this mesh, which can be  
10 transformed to change the appearance of the mesh. To ensure a natural looking human  
11 appearance when bones are moved, multiple bones are grouped into a single  
12 configuration holding constraints on the transformations of the individual bones. As  
13 such, modification of a single configuration transforms groups of related bones. Starting  
14 from the predefined model, these configurations were modified to create a resting  
15 mouth shape position as well as distinct mouth shapes needed for the articulation of  
16 different speech sounds. A number of 17 mouth shapes was opted for based on the  
17 coding scheme of a semi-automated text to speech synthesis software, Annosoft  
18 Lipsync Tool (Annosoft Lipsync Tool, 2014, Version 4.0260), which was afterwards  
19 used to map the auditory speech signal to the mouth shapes. The mouth shapes were  
20 designed by manually adjusting the position and rotation of 12 facial configurations for  
21 the lower jaw, the lips and the tongue, based on Flemish/Dutch articulatory rules and  
22 expert knowledge. Three additional configurations were adjusted to obtain a realistic  
23 looking face, i.e. the chin appearance, retraction of the nose wings and the thickness of  
24 the upper cheeks. The resulting mouth shapes were static and stored individually in a  
25 facial bone-structure of a virtual human in Blender, containing transformation values for



each individual bone. The mouth shapes were first created for a neutral virtual human template (Figure 1, B) and later on copied to a male and female virtual human corresponding more closely to the voice timbre of the auditory stimuli used in our experiments (Figure 1, A). Mouth shapes were afterwards mapped to the audio files of the speech material by means of Annosoft Lipsync Tool. Both the original audio material and the audio signal of the video recordings of the male speech therapist were used for mapping. Plain text was automatically converted into a mouth shape sequence according to a many-to-one phoneme-to-mouth shape mapping scheme. The proposed mouth shapes were revised manually for selection, timing and intensity (Figure 2). The resulting sequences were saved as data strings in an output file, together with information on the interpolation between specific mouth shapes based on Annosoft's internal model of co-articulation. Finally, the output files were loaded in the 3D engine software Unity (Unity, 2016, Version 0.43.4) and parsed into objects holding temporal information on the start and intensity of mouth shapes as well as previously determined transformation values of individual bones for each mouth shape. Co-articulation between two neighbouring mouth shapes was accomplished by an interpolation calculation of the transformation of the individual bones, according to Annosoft's co-articulation model.

## ***Procedure and analysis***

### ***Speech reading experiment***

Participants were seated in a darkened room in front of an acoustically transparent screen (173 x 300 cm) on which both the male and female virtual human and the video recordings were presented alternately in complete silence. The visual models were shown at eye level and the participant sat at a distance of 1.0 m to the screen. The

mouth of the visual models was 3.8 cm wide, comparable to the real mouth size of a speaker during a typical two-side table conversation. Stimuli were presented through software programs Processing (Processing, 2016, Version 3.0.2), Pure Data (Pure Data, 2015, Version 0.43.4) and Unity (Unity, 2016, Version 0.43.4) in three blocks of approximately half an hour each. The order of the lists, words and sentences within a list and the presentation order of the visual models were randomised between participants.

Before the start of each list, the experimenter announced the type of speech material, i.e. MS words, DS words or sentences. The participants were asked to attentively watch the visual presentation of the speech stimuli and to verbally report all the phonemes they could identify. The video or virtual animation was paused by the experimenter after every stimulus and repetitions of the stimuli were provided when participants asked for it. Multiple answers were allowed, but no feedback was given. The experimenter encouraged the participants to always give a response and entered all answers manually into a computer. When multiple answers were given, the most correct answer was taken into account. For the sentences, only the keywords were analysed.

Intelligibility scores of the three visual models were calculated based on the number of mouth movements correctly identified by the participants. Therefore, more general viseme categories, i.e. categories of phonemes that cannot be distinguished on visual features alone (Fisher, 1968), were defined that contained all 17 subtle varying mouth movements (Tables 1 and 2). The viseme categories used for our study were based on previous work of van Son et al. (1994), up to date the most extensive research on categorisation of Dutch phonemes into visemes. Consonants were grouped according to place and manner of articulation (Table 1). Different from the categories of van Son, we created a separate category for the /w/. The pronunciation of this consonant differs

1 between Flanders, i.e. the Flemish/Dutch-speaking part in Belgium used in our research,  
2 and the Netherlands, with the Flemish one being uttered bilabial in contrast to the  
3 labiodental Dutch /w/. Furthermore, we gathered all nonlabial front and back  
4 consonants into one single category, because not all individuals are skilled enough to  
5 observe the differences between these consonants (van Son et al., 1994). In contrast to  
6 van Son et al. (1994), vowels were not categorised according to vowel duration, but to  
7 lip opening and lip rounding only (Table 2). One extra vowel was added as a separate  
8 category, being the mid-central vowel schwa.

9 MS and DS words were analysed according to a CVC or CVCCVC structure  
10 respectively. This means that both presented and answered stimuli were fit into a  
11 structure with three or six possible phoneme chunks. One chunk could consist of a  
12 single consonant or vowel or multiple phonemes such as a consonant cluster. For each  
13 chunk, a viseme category code of the corresponding presented stimuli was assigned, e.g.  
14 V1 for the vowel /e/ (Tables 1 and 2). When a consonant cluster consisted of phonemes  
15 of a different viseme category, a combined code was used for the corresponding chunk,  
16 e.g. C14 for the cluster /pr/. When the code of the presented and answered chunk  
17 aligned, a score of one was assigned to that chunk. Omissions and insertions received a  
18 score of zero if they resulted in a different viseme code. All scores of one list were  
19 summed up and divided by the total amount of chunks in that list, resulting in a %  
20 correct score. The same procedure was applied to keyword scoring of the LIST  
21 sentences. The only difference was that keywords did not always follow a CVC or  
22 CVCCVC structure. Therefore, keywords were first split into syllables and each syllable  
23 was fit into a CVC structure of which the first and initial chunk could be empty. An  
24 example of keyword scoring can be found in Table 3.

### *Auditory-visual speech intelligibility*

In a second experiment, female LIST sentences were uttered by the female virtual human only in a virtual restaurant scenario projected on the screen in front of the participant. The set-up was identical to the speech reading experiment, except for the addition of two loudspeakers. Sentences were presented through a loudspeaker at 0°, standing at eye level 0.11 m behind the screen. A speaker in the zenith, right above the participant, was used for presentation of both SWN and MTN.

SRT's, i.e. the SNR at which 50% of the keywords is correctly repeated, were obtained by presenting lists of 10 sentences in noise, starting at a 60 dB A noise level and 48 dB A speech level, adaptively adjusting the speech level by 2 dB according to the answer of the participant. Only keywords had to be repeated correctly. The SRT was calculated by averaging the SNR's of the last five sentences plus the 11<sup>th</sup> fictive sentence of a single list. Sentences were presented both AO and AV in SWN and MTN. In the AO condition, participants watched a scenario similar to the restaurant but without the virtual human. Each condition was conducted twice and the order was randomly assigned between participants.

### **Results**

All results were analysed in the R software interface (R Core Team, 2017). Additional packages were used when needed, including the nlme package for general linear models (Pinheiro et al., 2017), the ez package for ANOVA-analysis (Lawrence, 2016) and the WRS2 package for robust paired samples t-tests (Mair, Schoenbrodt and Wilcox, 2017). Alpha-levels were set at 0.05 and p-values were corrected for multiple comparisons applying the Bonferroni method.

### *Speech reading experiment*

Figure 3 compares the identification scores obtained in the speech reading experiment for the three visual models, i.e. the video recording of a real person and the two virtual humans, for the three different speech materials. The same participants were tested in all conditions. To account for within-individual variances, we analysed the repeated-measures data using a general linear model. The type of visual model had a significant effect on speech reading scores,  $\chi^2(2) = 16.00$ ,  $p < 0.01$ , as did the type of speech material used in our experiment,  $\chi^2(2) = 53.00$ ,  $p < 0.01$ . No significant interaction effect between visual model and speech material was observed,  $\chi^2(4) = 5.00$ ,  $p = 0.24$ . Contrasts revealed that (1) scores obtained through the video recordings were significantly better compared to the ones from the two virtual humans,  $b = -6.00$ ,  $t(8) = -8.00$ ,  $p < 0.01$  but (2) speech reading scores were comparable for the male and female virtual human,  $b = -2.00$ ,  $t(8) = -1.60$ ,  $p = 0.15$ ; also (3) speech reading scores of the sentences were significantly lower compared to both MS and DS word scores,  $b = 7.00$ ,  $t(24) = 10.20$ ,  $p < 0.01$  albeit (4) scores on MS and DS words were not significantly different from each other,  $b = 2.00$ ,  $t(24) = 1.90$ ,  $p = 0.07$ .

Since the scores of the male and female virtual humans did not differ significantly from each other, they were averaged per participant and subsequently divided by the speech reading score of the video model for the corresponding speech material. In this way, individual proportion correct scores per speech material were obtained (Figure 4) and individual differences in speech reading ability were ruled out. In what follows, these individual proportion correct scores are referred to as ‘normalised proportion correct scores’. Normalised proportion correct scores were analysed through a non-parametrical Friedman’s ANOVA because of non-normal data distributions. The analysis revealed significant differences between the scores  $\chi^2(2) = 7.60$ ,  $p = 0.02$ ,

1 although post-hoc paired samples robust t-tests could not confirm any significant  
2 difference: the normalised proportion correct scores of MS words ( $Mdn = 0.85$ ) were  
3 equal to the scores of DS words ( $Mdn = 0.85$ ),  $p = 1.00$ ,  $r = 0.02$ , DS words ( $Mdn =$   
4  $0.85$ ) to sentences ( $Mdn = 0.62$ ),  $p = 0.37$ ,  $r = 0.75$  and MS words ( $Mdn = 0.85$ ) to  
5 sentences ( $Mdn = 0.62$ ),  $p = 0.52$ ,  $r = 0.79$ .

6       Based on the raw identification scores (Figure 3), scores on the sentences were  
7 significantly lower than the speech reading scores obtained on the MS and DS words.  
8 To investigate this difference in more detail, two extra sentence lists were presented to  
9 four out of five participants for the female virtual human only. In the previous  
10 experiment, sentences had to be speech read completely and no additional information  
11 was given on the possible words of the sentence. Having to remember the phonemes of  
12 an entire sentence could have impaired the speech reading scores, since the cognitive  
13 load on short-term memory during this task is higher compared to the memorisation of  
14 single words. Therefore, for the two extra sentence lists, sheets were handed over to the  
15 participant right before the start of the experiment with the non-keywords of the  
16 sentences printed on it. In this way, the participant had to speech read the keywords  
17 only, achieving a difficulty level more closely related to the MS and DS words task.  
18 Scores on the two extra sentence lists were averaged per participant and compared to  
19 scores on the MS words, DS words and regular sentences (Figure 5). Due to non-normal  
20 distributions, differences between all of the speech materials were investigated by  
21 means of a non-parametrical Friedman's ANOVA, which revealed significant  
22 differences  $\chi^2(3) = 9.30$ ,  $p = 0.03$ . Post-hoc comparisons were conducted via paired  
23 samples robust t-tests. Speech reading scores on the sentences ( $Mdn = 37.5\%$ ) appeared  
24 to be significantly different from the scores on the extra sentences ( $Mdn = 57.2\%$ ),  $p =$   
25  $0.03$ ,  $r = 0.76$ . On the contrary, scores on the MS words ( $Mdn = 67.1\%$ ) and DS words

( $Mdn = 60.8\%$ ) did not differ significantly from the scores on the extra sentences,  $p = 0.55$ ,  $0.42$  and  $r = 0.49$ ,  $0.24$  respectively.

### ***Auditory-visual speech intelligibility***

Figure 6 shows the SRT's of the AV SI experiment for the young adult group ( $N = 35$ ) in the AO and AV condition for both noise materials. The SRT in the AO SWN condition ( $Mean = -12.1$  dB SNR,  $SE = 0.3$  dB) was better than expected based on research of van Wieringen and Wouters (2008), who found a norm SRT of  $-8.0$  dB SNR  $\pm 0.2$  dB. This can be explained by the fact we spatially separated the noise source from the target sound. Listening in free field with two ears instead of one allowed our participants to make use of head diffraction effects and binaural listening phenomena like binaural squelch and binaural redundancy to improve speech understanding. The benefit of adding visual speech cues to the auditory signal was analysed through a factorial repeated-measures ANOVA and revealed a main effect of both mode  $F(1, 34) = 33.58$ ,  $p < 0.01$  and noise  $F(1, 34) = 220.87$ ,  $p < 0.01$ . Addition of visual cues reduced the SRT in SWN on average by  $1.4$  dB ( $SE = 0.4$  dB) and in MTN on average by  $1.8$  dB ( $SE = 0.4$  dB). Overall, SRT's were significantly better in SWN (AV:  $Mean = -13.5$  dB,  $SE = 0.3$  dB; AO:  $Mean = -12.1$  dB,  $SE = 0.3$  dB) compared to SRT's obtained in MTN (AV:  $Mean = -10.3$  dB,  $SE = 0.4$  dB; AO:  $Mean = -8.5$  dB,  $SE = 0.4$  dB), which is expected based on research on the acoustical and informational masking effects of different noise types (Francart, van Wieringen and Wouters, 2011). There was no significant interaction  $F(1, 34) = 0.58$ ,  $p = 0.45$ .

### **Discussion**

The presented work describes the development and evaluation of two virtual human speakers, designed to be implemented in future auditory assessment methods where

speech will be presented auditory-visually in order to obtain an ecologically valid measure of speech processing. Both virtual humans, a male and female speaker, yielded the same intelligibility results, indicating they can be used interchangeably. Yet, present implementation of the virtual humans resulted in slightly lower speech reading scores than the videotaped speaker. This agrees with many reports showing that natural speakers almost always outperform their synthetic counterparts (Benoît and Le Goff, 1998; Geiger, Ezzat and Poggio, 2003; Siciliano et al., 2003; Fagel and Clemens, 2004; Fagel, Bailly and Elisei, 2007; Dey, Maddock and Nicolson, 2010).

Direct comparison of speech reading scores on our virtual humans with scores on other virtual human speakers is not straightforward, since modelling techniques differ greatly and evaluation methods are not uniform. In addition, the large variability of individuals' speech reading skills hampers direct comparison of raw intelligibility scores across participants. Therefore, we calculated normalised proportion correct scores, which are suitable to compare to other research. Fagel et al. (2007) set up a forced-choice speech reading experiment with German CVC-sequences presented through a performance-driven virtual human and natural video images. Proportion correct scores were about 0.60 (virtual human: 19.4%, video: 32.5%), which is equivalent to our results, i.e. a normalised proportion correct of 0.62 when speech reading sentences. Nevertheless, we found a normalised proportion correct score of 0.85 for speech reading both MS and DS words, outperforming the scores on CVC-sequences in the research of Fagel et al. (2007). Geiger et al. (2003) developed virtual human speakers by pasting synthetic mouth regions onto real faces, preserving the original head and eye movements. The proportion correct score of phonemes for English single words and sentences was about 0.71 (virtual human: 21.2%, video: 30.0%), which is similar to our results. Comparable scores were also found by Siciliano



et al. (2003), who reported phoneme scores correct of 0.58 (virtual human: 13.6%, video: 23.4%) for English CVC sequences presented by their virtual model. Although normalised proportion correct scores were similar between studies, it seems that our virtual humans resulted in higher raw speech reading scores compared to other virtual implementations (Geiger, Ezzat and Poggio, 2003; Siciliano et al., 2003; Fagel, Bailly and Elisei, 2007). In a study of Schreitmüller et al. (2017) for example, normal hearing and cochlear-implanted individuals obtained average scores of 12% and 38% words correct on speech reading nonsense matrix sentences. In our study, individuals scored on average 70%, 64% and 46% visemes correct on MS words, DS words and keywords of meaningful sentences respectively (Figure 3). However, it should be noted that scoring criteria and methods varied across studies.

Both MS and DS words as well as sentences were presented to our participants to investigate potential differences in difficulty. Phonemes and thus visemes were harder to identify when presented in sentences compared to single words. Memory storage of the presented phonemes could have impaired scores on the sentences, since the amount of phonemes per sentence was much higher than in MS and DS words. In other words, task-related difficulties rather than inconclusive mouth movements may have caused these lower results. The notion of increased memory load is in line with participants' spontaneous remarks during the experiment. Overall, participants asked for more repetitions during sentence presentations since they found it hard to remember all presented phonemes. Higher identification scores on two extra sentence lists, relative to scores on the regular sentences, further support this idea. In these lists, the amount of phonemes to be memorised was reduced by revealing all non-keywords of the sentence upfront. One could argue that revealing extra information on the words did not only reduce the needed memory capacity, but also provided more contextual information,

1 facilitating speech reading because of semantic and morphosyntactic redundancy  
2 (Martin et al., 1983; Grant and Seitz, 2000b). Interestingly, the overall benefit of  
3 context was smaller than expected based on research of Boothroyd and Nitttrouer  
4 (1988), who designed a formula to predict word recognition in sentences based on  
5 recognition of isolated words, taking sentence context into account as a free parameter.  
6 According to their formula, assuming a sentence context factor of 2.7, and based on the  
7 median scores of our participants on the MS and DS words, i.e. 67% and 61% correct  
8 respectively (Figure 5), we expected our participants to achieve sentence speech reading  
9 scores of 92% to 95%. In our study, sentence scores were about 57%. This could be due  
10 to the fact we only provided non-keywords, which were mostly articles and were less  
11 informative than keywords, typically nouns and verbs. Work from Martin et al. (1983)  
12 supports this idea, showing that words are easier to speech read when they are highly  
13 predictable. Moreover, other relevant signal cues like word boundaries are hard to  
14 extract in absolute silence which was the case in our experiments. Nevertheless, in an  
15 auditory-visual environment, these cues help to gain access to relevant context  
16 information and thereby improve performance (Grant and Seitz, 2000b).

17 In terms of improvement in speech intelligibility, we found that the presentation  
18 of our female virtual human reduced the SRT's on average by 1.4 dB in SWN and 1.8  
19 dB in MTN, in contrast to an auditory-only listening condition. This visual speech  
20 benefit is in agreement with research demonstrating that visual cues improve speech  
21 understanding in noisy conditions (Sumby and Pollack, 1954; MacLeod and  
22 Summerfield, 1987, 1990; Grant, Walden and Seitz, 1998) and provides support for the  
23 applicability of our virtual humans in future research. Whereas many studies have  
24 focused on improvements in percentage correct speech intelligibility, only a few  
25 measured changes in SRT's, yielding benefits in the range of 2 to 6 dB when video

images are presented (Bernstein and Grant, 2009; Brungart, Sheffield and Kubli, 2014). These benefits are in line with our results, given the fact we used synthetic human speakers instead of real video images. Moreover, our benefit could have been larger when presenting sentences at lower SNR's approximating intelligibility levels below the SRT, since speech reading scores tend to be higher the more noise is presented (Benoît and Le Goff, 1998). This idea is supported by the principle of inverse effectiveness, which assumes that auditory-visual integration becomes more prominent when unimodal processing of either auditory or visual cues declines, as is the case in noisy conditions (Tye-Murray et al., 2010).

In general, the amount of correctly identified visemes and the significant visual benefit obtained by our virtual humans establishes the realism and quality of both the male and female virtual human. They differ from previously developed virtual humans in that the mouth movements are incorporated in full body virtual humans, as opposed to virtual talking heads, and that they are designed with easily accessible 3D rendering software systems, facilitating the implementation in ecologically valid 3D scenarios. Hence, they open avenues for the use of full-body virtual humans in a clinical tool to assess auditory-visual speech processing abilities in complex daily listening situations. Nevertheless, improvements can still be made in terms of intelligibility. A detailed analysis of the answers based on virtual human and video projection is shown in the confusion matrices in Figure 7, revealing which categories of visemes were harder to identify and thus need further investigation. Omitted answers and answers that could not be assigned to one of the stimulus categories (e.g. answering with a vowel to a consonant stimulus) were left out. When it comes to consonants, the most striking result was the difficulty in recognising the /w/ (C3), which tended to be identified as a nonlabial consonant (C4). In terms of vowels, V3 and V4 categories seemed to be

regularly interchanged compared to the video images. This is not surprising, given that both vowel categories differ subtly in terms of the amount of lip opening only. In addition, V1 seemed harder to recognise, possibly because these vowels encompass less pronounced mouth movements. Refinements of all above categories will be made in future research, possibly elevating the intelligibility of our virtual humans to an even higher level.

### **Acknowledgements**

The authors would like to thank all the participants who took part in this study, as well as the students from the Speech Language Pathology and Audiology Sciences educational program for their help with data collection. Sam Denys is gratefully acknowledged for his help with the speech reading video recordings. This work was supported by the Oticon Foundation and a TBM-FWO grant from the Research Foundation-Flanders under Grant [T002216N].

### **Declaration of interest**

The authors report no conflicts of interest.

## 1   **References**

- 2   Akeroyd, M. A., Arlinger, S., Bentler, R. A., Boothroyd, A., Dillier, N., Dreschler, W.  
3   A., Gagné, J. P., Lutman, M., Wouters, J., Wong, L. and Kollmeier, B. (2015).  
4   ‘International Collegium of Rehabilitative Audiology (ICRA) recommendations for the  
5   construction of multilingual speech tests: ICRA Working Group on Multilingual Speech  
6   Tests’, *Int J Audiol*, 54(sup2), pp. 17–22.
- 7   Annosoft Lipsync Tool. (2014). Available at: <http://www.annosoft.com/>.
- 8   Arnal, L. H., Wyart, V. and Giraud, A. L. (2011). ‘Transitions in neural oscillations  
9   reflect prediction errors generated in audiovisual speech’, *Nat Neurosci*, 14(6), pp. 797–  
10   801.
- 11   Benoît, C. and Le Goff, B. (1998). ‘Audio-visual speech synthesis from French text:  
12   Eight years of models, designs and evaluation at the ICP’, *Speech Commun*, 26(1), pp.  
13   117–129.
- 14   Bernstein, J. G. and Grant, K. W. (2009). ‘Auditory and auditory-visual intelligibility of  
15   speech in fluctuating maskers for normal-hearing and hearing-impaired listeners’, *J*  
16   *Acoust Soc Am*, 125(5), pp. 3358–3372.
- 17   Bernstein, L. E., Auer, E. T. and Takayanagi, S. (2004). ‘Auditory speech detection in  
18   noise enhanced by lipreading’, *Speech Commun*, 44(1), pp. 5–18.
- 19   Beskow, J. (2004). ‘Trainable articulatory control models for visual speech synthesis’,  
20   *Int J Speech Technol*, 7(4), pp. 335–349.
- 21   Blender. (2015). Available at: <https://www.blender.org/>.
- 22   Boothroyd, A. and Nittrouer, S. (1988). ‘Mathematical treatment of context effects in  
23   phoneme and word recognition’, *J Acoust Soc Am*, 84(1), pp. 101–114.
- 24   Brungart, D. S., Sheffield, B. M. and Kubli, L. R. (2014). ‘Development of a test battery

- 1 for evaluating speech perception in complex listening environments', *J Acoust Soc Am*  
2 136(2), pp. 777–790.
- 3 Calvert, G. A., Campbell, R. and Brammer, M. J. (2000). 'Evidence from functional  
4 magnetic resonance imaging of crossmodal binding in the human heteromodal cortex',  
5 *Curr Biol*, 10(11), pp. 649–657.
- 6 Chandrasekaran, C., Lemus, L. and Ghazanfar, A. A. (2013). 'Dynamic faces speed up  
7 the onset of auditory cortical spiking responses during vocal detection', *Proc Natl Acad*  
8 *Sci U S A*, 110(48), pp. E4668–E4677.
- 9 Dey, P., Maddock, S. C. and Nicolson, R. (2010). 'Evaluation of A Viseme-Driven  
10 Talking Head', in *Proc. Theory and Practice of Computer Graphic*, pp. 139-142.
- 11 Fagel, S., Bailly, G. and Elisei, F. (2007). 'Intelligibility of natural and 3D-cloned  
12 German speech', in *International Conference on Auditory-Visual Speech Processing*,  
13 *AVSP*. Hilvarenbeek, The Netherlands, pp. 56–61.
- 14 Fagel, S. and Clemens, C. (2004). 'An articulation model for audiovisual speech  
15 synthesis — Determination, adjustment, evaluation', *Speech Commun*, 44(1), pp. 141–  
16 154.
- 17 Fisher, C. G. (1968). 'Confusions among visually perceived consonants', *J Speech Lang*  
18 *Hear Res*, 11(4), pp. 796-804.
- 19 Francart, T., van Wieringen, A. and Wouters, J. (2011). 'Comparison of fluctuating  
20 maskers for speech recognition tests', *Int J Audiol*, 50(1), pp. 2–13.
- 21 Geiger, G., Ezzat, T. and Poggio, T. (2003). Perceptual Evaluation of Video-Realistic  
22 Speech. Tech Report: CBCL Paper 224/AI Memo 2003-003, MIT Artificial Intelligence  
23 Laboratory, Cambridge.
- 24 Ghazanfar, A. A. (2009). 'The multisensory roles for auditory cortex in primate vocal  
25 communication', *Hear Res*, 258(1), pp. 113–120.

- 1 Gosselin, P. A. and Gagné, J. P. (2011). ‘Older adults expend more listening effort than  
2 young adults recognizing audiovisual speech in noise’, *Int J Audiol*, 50(11), pp. 786–  
3 792.
- 4 Grange, J. A. and Culling, J. F. (2016). ‘Head orientation benefit to speech intelligibility  
5 in noise for cochlear implant users and in realistic listening conditions’, *J Acoust Soc*  
6 *Am* 140(6), pp. 4061–4072.
- 7 Grant, K. W. and Seitz, P. (2000a). ‘The use of visible speech cues for improving  
8 auditory detection of spoken sentences’, *J Acoust Soc Am*, 108(3), pp. 1197–1208.
- 9 Grant, K. W. and Seitz, P. F. (2000b). ‘The recognition of isolated words and words in  
10 sentences: Individual variability in the use of sentence context’, *J Acoust Soc Am*,  
11 107(2), pp. 1000–1011.
- 12 Grant, K. W., Walden, B. E. and Seitz, P. F. (1998). ‘Auditory-visual speech  
13 recognition by hearing-impaired subjects: Consonant recognition, sentence recognition,  
14 and auditory-visual integration’, *J Acoust Soc Am*, 103(5), pp. 2677–2690.
- 15 Jansen, S., Koning, R., Wouters, J. and van Wieringen, A. (2014). ‘Development and  
16 validation of the Leuven intelligibility sentence test with male speaker (LIST-m)’, *Int J*  
17 *Audiol*, 53(1), pp. 55–59.
- 18 Lawrence, M. A. (2016). ‘ez: Easy Analysis and Visualization of Factorial Experiments.  
19 R package version 4.4-0’. Available at: <https://cran.r-project.org/package=ez>.
- 20 Ma, J., Cole, R., Pellom, B., Ward, W. and Wise, B., (2006). ‘Accurate visible speech  
21 synthesis based on concatenating variable length motion capture data’. *IEEE Trans Vis*  
22 *Comput Graph*, 12(2), pp.266-276.
- 23 Ma, W. J., Zhou, X., Ross, L. A., Foxe, J. J. and Parra, L. C. (2009). ‘Lip-Reading Aids  
24 Word Recognition Most in Moderate Noise: A Bayesian Explanation Using High-  
25 Dimensional Feature Space’, *PLoS One*, 4(3), p. e4638.
- 26 MacLeod, A. and Summerfield, Q. (1987). ‘Quantifying the contribution of vision to

- 1 speech perception in noise', *Br J Audiol*, 21(2), pp. 131–141.
- 2 MacLeod, A. and Summerfield, Q. (1990). 'A procedure for measuring auditory and  
3 audiovisual speech-reception thresholds for sentences in noise: Rationale, evaluation,  
4 and recommendations for use', *Br J Audiol*, 24(1), pp. 29–43.
- 5 Mair, P., Schoenbrodt, F. and Wilcox, R. (2017). 'WRS2: Wilcox robust estimation and  
6 testing'.
- 7 Make Human. (2015). Available at: <https://www.makehuman.org/>.
- 8 Martin, L. F. A., Clark, G. M., Seligman, P. M. and Tong, Y. C. (1983). 'A lip-reading  
9 assessment for profoundly deaf patients', *J Laryngol Otol*, 97(4), pp. 343–350.
- 10 Massaro, D. W. and Cohen, M. M. (1990). 'Perception of synthesized audible and  
11 visible speech', *Psychol Sci*, 1(1), pp. 55–63.
- 12 Massaro, D. W. and Light, J. (2004). 'Using Visible Speech to Train Perception and  
13 Production of Speech for Individuals With Hearing Loss', *J Speech Lang Hear Res*,  
14 47(2), pp. 304–320.
- 15 Mattheyses, W. and Verhelst, W. (2015). 'Audiovisual speech synthesis: An overview  
16 of the state-of-the-art', *Speech Commun*, 66, pp. 182–217.
- 17 McGurk, H. and MacDonald, J. (1976). 'Hearing lips and seeing voices', *Nature*,  
18 264(5588), pp. 746–748.
- 19 Nilsson, M., Soli, S.D. & Sullivan, J. A. (1994). 'Development of the hearing in noise  
20 test for the measurement of speech reception thresholds in quiet and in noise', *J Acoust*  
21 *Soc Am*, 95, pp. 1085–1099.
- 22 Okada, K., Venezia, J. H., Matchin, W., Saberi, K. and Hickok, G. (2013). 'An fMRI  
23 Study of Audiovisual Speech Perception Reveals Multisensory Interactions in Auditory  
24 Cortex', *PLoS One*, 8(6), p. e68959.



- 1 Peng, X., Chen, H., Wang, L. and Wang, H. (2018). ‘Evaluating a 3-D virtual talking  
2 head on pronunciation learning’, *Int J Hum Comput Stud*, Elsevier Ltd, 109, pp. 26–40.
- 3 Pichora-Fuller, M. K., Kramer, S. E., Eckert, M. A., Edwards, B., Hornsby, B. W.,  
4 Humes, L. E., Lemke, U., Lunner, T., Matthen, M., Mackersie, C. L. and Naylor, G.,  
5 (2016). ‘Hearing impairment and cognitive energy: The framework for understanding  
6 effortful listening (FUEL)’, *Ear Hear*, 37, pp.5S-27S.
- 7 Pinheiro, J., Bates, D., DebRoy, S., Sarkar, D. and Team, R. C. (2017). ‘nlme: Linear  
8 and Nonlinear Mixed Effects Models. R package version 3.1-131’.
- 9 Processing. (2016). Available at: <https://processing.org/>.
- 10 Pure data. (2015). Available at: <https://puredata.info/>.
- 11 R Core Team (2017). ‘R: A language and environment for statistical computing’.  
12 Vienna, Austria: R Foundation for Statistical Computing. Available at: [https://www.r-](https://www.r-project.org)  
13 [project.org](https://www.r-project.org).
- 14 Rosenblum, L. D. (2008). ‘Speech Perception as a Multimodal Phenomenon’, *Curr Dir*  
15 *Psychol Sci*, 17(6), pp. 405–409.
- 16 Salvi, G., Beskow, J., Al Moubayed, S. and Granström, B. (2009). ‘SynFace — Speech-  
17 Driven Facial Animation for Virtual Speech-Reading Support’,  
18 *EURASIP J Audio Speech Music Process*, 2009(1), p. 191940.
- 19 Schreitmüller, S., Frenken, M., Bentz, L., Ortmann, M., Walger, M. and Meister, H.  
20 (2018). ‘Validating a Method to Assess Lipreading, Audiovisual Gain, and Integration  
21 During Speech Reception With Cochlear-Implanted and Normal-Hearing Subjects  
22 Using a Talking Head’, *Ear Hear*, 39(3), pp. 503-516.
- 23 Schwartz, J. L., Berthommier, F. and Savariaux, C. (2004). ‘Seeing to hear better:  
24 evidence for early audio-visual interactions in speech identification’, *Cognition*, 93(2),  
25 pp. B69–B78.

- 1 Siciliano, C., Williams, G., Beskow, J. and Faulkner, A. (2003). ‘Evaluation of a  
2 Multilingual Synthetic Talking Face as a Communication Aid for the Hearing  
3 Impaired’, in *International Congress of Phonetic Sciences, ICPhS*. Barcelona, Spain,  
4 pp 1-4.
- 5 Sifakis, E., Neverov, I. and Fedkiw, R. (2005). ‘Automatic determination of facial  
6 muscle activations from sparse motion capture marker data’, in *ACM Trans Graph*,  
7 24(3), pp. 417–425.
- 8 Stevens, C. J., Gibert, G., Leung, Y. and Zhang, Z. (2013). ‘Evaluating a synthetic  
9 talking head using a dual task: Modality effects on speech understanding and cognitive  
10 load’, *Int J Hum Comput Stud*, 71(4), pp. 440–454.
- 11 Sumby, W. H. and Pollack, I. (1954). ‘Visual Contribution to Speech Intelligibility in  
12 Noise’, *J Acoust Soc Am*, 26(2), pp. 212–215.
- 13 Tye-Murray, N., Sommers, M., Spehar, B., Myerson, J. and Hale, S. (2010). ‘Aging,  
14 Audiovisual Integration, and the Principle of Inverse Effectiveness’, *Ear Hear*, 31(5),  
15 pp. 636–644.
- 16 Unity. (2016). Available at: <https://unity3d.com/>.
- 17 Ward, W., Cole, R., Bolaños, D., Buchenroth-Martin, C., Svirsky, E. and Weston, T.  
18 (2013). ‘My Science Tutor : A Conversational Multimedia Virtual Tutor’, *J Educ*  
19 *Psychol*, 105(4), pp. 1115–1125.
- 20 van Son, N., Huiskamp, T. M. I., Bosman, A. J. and Smoorenburg, G. F. (1994).  
21 ‘Viseme classifications of Dutch consonants and vowels’, *J Acoust Soc Am*, 96(3), pp.  
22 1341–1355.
- 23 van Wieringen, A. and Wouters, J. (2008). ‘LIST and LINT: sentences and numbers for  
24 quantifying speech understanding in severely impaired listeners for Flanders and the  
25 Netherlands’, *Int J Audiol*, 47(6), pp. 348–355.
- 26 van Wieringen, A. and Wouters, J. (2015). ‘What can we expect of normally-developing

- 1 children implanted at a young age with respect to their auditory, linguistic and cognitive  
2 skills?’, *Hear Res*, 322, pp. 171–179.
- 3 Williams, L. (1990). ‘Performance-Driven Facial Animation’, in *ACM SIGGRAPH*  
4 *Computer Graphics Conf.*, 24(4), pp. 235–242.
- 5 Woodhouse, L., Hickson, L. and Dodd, B. (2009). ‘Review of visual speech perception  
6 by hearing and hearing-impaired people: clinical implications.’, *Int J Lang Commun*  
7 *Disord*, 44(3), pp. 253–270.
- 8 Wouters, J., Damman, W. and Bosman, A. J. (1994). ‘Vlaamse opname van  
9 woordenlijsten voor spraakaudiometrie’, *Logopedie: informatiemedium van de Vlaamse*  
10 *vereniging voor logopedisten*, 7(6), pp. 28–34.
- 11 Wright, T. M., Pelphrey, K. A., Allison, T., McKeown, M. J. and McCarthy, G. (2003).  
12 ‘Polysensory Interactions along Lateral Temporal Regions Evoked by Audiovisual  
13 Speech’, *Cereb Cortex*, 13(10), pp. 1034–1043.
- 14
- 15

1   **Tables**

2   Table 1. Categorisation of Flemish/Dutch consonants.

Consonants	Category	Place and manner of articulation
/p, b, m/	C1	Bilabial plosive/nasal
/f, v/	C2	Labiodental fricative
/w/	C3	Bilabial fricative/semi-vocal
/t, d, n, j, l, s, z, ʃ, k, ʁ, r, x, ŋ, h/	C4	Nonlabial consonants

3

1  
2  
3  
4

Table 2. Categorisation of Flemish/Dutch vowels.

Vowels	Category	Lip opening and lip rounding
/i, ɪ, e, ε/	V1	Closed and half-closed, unrounded
/ɛɪ, ɑ, a/	V2	Half-open and open, unrounded
/u, y, ʌ/	V3	Closed and half-closed, rounded
/ɔ, o, ø/	V4	Half-open, rounded
/œy, ɔu/	V5	Closing, rounding diphthongs
/ə/	V6	Middle, unrounded

1    Table 3. Scoring example of keywords of a LIST-sentence.

Stimulus (S) and answer (A)							
Words and [keywords]	S:	Ze [viel] van de [trap].					
	A:	Ze [vreesde] voor de [trap].					
Phoneme chunk	S:	v	ie	l	tr	a	p
	A:	vr	ee	sde	tr	a	p
Viseme category code	S:	C2	V1	C4	C4	V2	C1
	A:	C24	V1	C4V6	C4	V2	C1
Scoring: 4/6 = 67%		0	1	0	1	1	1

2

## Figure captions

**Figure 1.** Panel A: video-recorded male speaker (left) and male and female virtual human model (middle and right) used in the speech reading experiment (all models) and the AV SI test (female virtual human only). Panel B: neutral virtual human model uttering the vowel /o/.

**Figure 2.** Screenshot of the Annosoft Lipsync Tool for the mapping of the monosyllabic word ‘hang’ with a: auditory waveform, b: time bar, c: selected mouth shapes, d: timing and intensity of selected mouth shapes and e: plain text of speech stimulus.

**Figure 3.** Individual percentage correct identification scores on the speech reading experiment for different speech materials and visual models. The lines connect mean scores per visual model, error bars represent SE of the mean, N = 5.

**Figure 4.** Boxplots representing the normalised proportion of visemes correct when speech reading the virtual humans, N = 5.

**Figure 5.** Individual percentage correct scores and line connecting mean scores for all speech materials, for speech reading the female virtual human. Each symbol represents a different participant, N = 4.

**Figure 6.** Boxplots representing the SRT’s of the speech intelligibility experiment in different noise conditions (MTN = multitalker babble noise, SWN = speech weighted noise) and different listening modes (AO = Auditory-Only, AV = Auditory-Visual) for young normal hearing participants, N = 35.

**Figure 7.** Confusion matrices for the different consonant (top row) and vowel (bottom row) categories for all visual models: A = virtual humans, B = video. Numbers indicate percentage correct answers. Consonant and vowel categories are defined in Table 1 and Table 2.