

PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is a postprint version which may differ from the publisher's version.

For additional information about this publication click this link.

<http://hdl.handle.net/2066/170043>

Please be advised that this information was generated on 2019-09-20 and may be subject to change.



Visual context enhanced: The joint contribution of iconic gestures and visible speech to degraded speech comprehension.

Journal:	<i>Journal of Speech, Language, and Hearing Research</i>
Manuscript ID	JSLHR-H-16-0101.R1
Manuscript Type:	Research Article
Date Submitted by the Author:	22-Jun-2016
Complete List of Authors:	Drijvers, Linda; Radboud Universiteit Nijmegen, Centre for Language Studies ; Radboud Universiteit Nijmegen Donders Institute for Brain Cognition and Behaviour, Donders Centre for Cognitive Neuroimaging Ozyurek, Asli; Radboud University , Centre for Language Studies ; Max-Planck-Institut für Psycholinguistik, CLS-M; Radboud Universiteit Nijmegen Donders Institute for Brain Cognition and Behaviour, Donders Centre for Cognitive Neuroimaging
Keywords:	Noise, Language, Cognition, Speech perception, Speech

1
2
3
4
5
6
7
8
9
10
11
12 **Visual context enhanced: The joint contribution of iconic gestures and visible speech to**
13 **degraded speech comprehension.**
14

15
16
17 Linda Drijvers^{a,b} & Asli Özyürek^{a,b,c}
18
19

20
21
22
23
24 ^a Radboud University, Centre for Language Studies, Erasmusplein 1, 6525 HT, Nijmegen,
25 The Netherlands
26

27
28 ^b Radboud University, Donders Institute for Brain, Cognition, and Behaviour, Montessorilaan
29 3, 6525 HR, Nijmegen, The Netherlands
30

31
32 ^c Max Planck Institute for Psycholinguistics, Wundtlaan 1, 6525 XD, Nijmegen, The
33 Netherlands
34
35

36
37
38 *** Correspondence to:**

39
40 Linda Drijvers, Radboud University, Centre for Language Studies, Donders Institute for
41 Brain, Cognition and Behaviour, Wundtlaan 1, 6525 XD, Nijmegen, The Netherlands. E-
42 mail: linda.drijvers@mpi.nl, telephone: +31 (0) 24 3521591
43
44

45
46
47 **Conflict of interest:**

48
49 The authors declare no conflict of interest.
50
51
52
53
54
55
56
57
58
59
60

Abstract

Purpose: This study investigated whether and to what extent iconic co-speech gestures contribute to information from visible speech to enhance degraded speech comprehension at different levels of noise-vocoding. Previously, the contributions of these two visual articulators to speech comprehension have only been studied separately.

Method: Twenty participants watched videos of an actress uttering an action verb and completed a free-recall task. The videos were presented in three speech (2-band; 6-band noise-vocoding; clear), three multimodal (Speech+Lips blurred; Speech+VisibleSpeech; Speech+VisibleSpeech+Gesture) and two visual only conditions (VisibleSpeech; VisibleSpeech+Gesture).

Results: Accuracy levels were higher when both visual articulators were present compared to one or none. The enhancement effects of a) visible speech, b) gestural information on top of visible speech and c) both visible speech and iconic gestures were larger in 6-band than 2-band noise-vocoding or visual only conditions. Gestural enhancement in 2-band noise-vocoding did not differ from gestural enhancement in visual only conditions.

Conclusions: When perceiving degraded speech in a visual context, listeners benefit more from having both visual articulators present compared to one. This benefit was larger at 6-band than 2-band noise-vocoding, where listeners can benefit from both phonological cues from visible speech, and semantic cues from iconic gestures to disambiguate speech.

Introduction

Natural, face-to-face communication often involves an audiovisual binding that integrates information from multiple inputs such as speech, visible speech, and iconic co-speech gestures. Notably, the relationship between these two visual articulators and the speech signal seems to differ: Iconic gestures, which can be described as hand movements that illustrate object attributes, actions and space (e.g., Clark, 1996; Goldin-Meadow, 2005; McNeill, 1992), are to be related to speech on a semantic level, due to the similarities to the objects, events and spatial relations they represent. In contrast, the relation between visible speech, consisting of lip movements, tongue movements and teeth, and speech consists of a form-to-form mapping between syllables and visible speech on a phonological level. Previous research has argued that both iconic gestures and visible speech can enhance speech comprehension, especially in adverse listening conditions, such as degraded speech (Holle, Obleser, Rueschemeyer, & Gunter, 2010; Obermeier, Dolk, & Gunter, 2012; Ross, Saint-Amour, Leavitt, Javitt, & Foxe, 2007; Sumbly & Pollock, 1954). However, the contribution of iconic gestures and visual speech to audiovisual enhancement of speech in adverse listening conditions has been mostly studied separately. Since natural, face-to-face communication involves gestures and visual speech as possible visual articulators, this raises the question of whether, to what extent and how the co-occurrence of these two visual articulators influence speech comprehension in adverse listening conditions. To this end, the current study aims to investigate the contribution of both types of visual information to degraded speech comprehension in a joint context.

Iconic gestures are frequently prevalent in natural, face-to-face communication and have both a temporal and semantic relation with the speech they occur with, causing them to be hard to disambiguate without speech. It has been theorized that iconic gestures are an integral part of language (Kendon, 2004; McNeill, 1992): Speech and iconic gestures are

1
2
3 integrated continuously during comprehension, and target linguistic processing on semantic,
4 syntactic and pragmatic levels (Holle et al., 2012; Kelly, Özyürek, & Maris, 2010; McNeill,
5 1992, see for a review and meta-analysis: Hostetter, 2011). Previous research has shown that
6 semantic information from iconic gestures is indeed processed by listeners and that iconic
7 gestures can impact language comprehension, at behavioral and neural levels (e.g. Beattie &
8 Shovelton, 1999; 2002; Holle & Gunter, 2007; Holler et al., 2014; Holler, Kelly, Hagoort, &
9 Özyurek, 2010; Holler, Shovelton, & Beattie, 2009; Kelly, Barr, Church, & Lynch, 1999;
10 Kelly, Healey, Özyürek, & Holler, 2015; Obermeier, Holle, & Gunter, 2011 see for a review,
11 Özyürek, 2014). For example, in an EEG study, Holle & Gunter (2007) showed participants
12 videos of an actor who uttered a sentence while gesturing. Here, the experimental sentences
13 contained an unbalanced homonym in the first part of the sentence (e.g. 'She controlled the
14 ball'). This homonym was disambiguated in the subsequent clause (e.g. 'which during the
15 game'/'which during the dance'). When the actor uttered the homonym, he would
16 simultaneously produce an iconic gesture that either depicted the dominant ('game') or the
17 subordinate meaning ('dance') of the homonym. When the gesture was congruent, they found
18 a smaller N400 as compared to an incongruent gesture. This suggests that listeners use the
19 semantic information from gestures to disambiguate speech.

20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42 So far, it has been argued that in adverse listening conditions, gestures occur more
43 frequently (Hoskin & Herman, 2001; Kendon, 2004) and that listeners take gestures more
44 into account than in clear speech (Rogers, 1978). This was also found by Obermeier et al.,
45 (2011), who used a similar paradigm as Holle & Gunter (2007), to reveal that when there was
46 no temporal overlap of a word and a gesture and participants were not explicitly asked to
47 attend to the gestures, speech-gesture integration did not occur. However, in a subsequent
48 study where the same stimuli were presented in multi-talker babble noise, listeners did
49 incorporate the gestural information with the speech signal to disambiguate the meaning of
50
51
52
53
54
55
56
57
58
59
60

1
2
3 the sentence. This effect was also found for hearing-impaired individuals (Obermeier et al.,
4 2012). These results underline that speech-gesture integration can be modulated by specific
5 characteristics of the communicative situation.
6
7
8

9
10 Another fMRI study by Holle, Obleser, Rueschemeyer & Gunter (2010) investigated
11 the integration of iconic gestures and speech by manipulating the signal-to-noise ratio (SNR)
12 of the speech to target areas that were sensitive to bimodal enhancement and inverse
13 effectiveness (i.e. greater bimodal enhancement for unimodally least effective stimuli, i.e. the
14 noisiest speech level). Participants watched videos of an actor with a covered face, who
15 uttered short sentences (e.g. 'And now I grate the cheese') with or without an accompanying
16 iconic co-speech gesture. These videos were presented with speech in a good SNR (+2 dB) or
17 in a moderate SNR (-6 dB), using multitalker babble tracks. Their results revealed that the
18 superior temporal sulcus and superior temporal gyrus in both hemispheres were sensitive to
19 bimodal enhancement, and the neural enhancement for bimodal enhancement was even larger
20 when participants were processing the speech and gestures in the degraded speech conditions.
21 On both a neural as a behavioral level (i.e. response accuracy), this study showed that
22 attending to a gesture under adverse listening conditions can significantly enhance speech
23 comprehension and help in the disambiguation of a speech signal that is difficult to interpret.
24 This gestural enhancement had already been described by Rogers (1978), who manipulated
25 noise levels to show that gestures could only benefit speech comprehension when sufficient
26 noise was added to the speech signal.
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

49 As Holle and colleagues (2010) note however, their study (and other studies, see e.g.
50 Obermeier et al., 2011; Obermeier et al., 2012) have only focused on one visual articulator in
51 speech-related audiovisual integration, namely iconic gestures. Other visual articulators, such
52 as lip movements, were deliberately excluded from the stimuli that were used by blocking the
53 actor's face with a black mask. Yet, these lip movements are inherently part of natural, face-
54
55
56
57
58
59
60

1
2
3 to-face communication: Lip movements can provide temporal information about the speech
4
5 signal (e.g. on the amplitude envelope) and information on the spatial location of a speaker's
6
7 articulators (e.g. place and manner of articulation), which can be specifically useful when
8
9 perceiving speech in adverse listening conditions. Additionally, lip movements can convey
10
11 phonological information, because of the form-form relationship between lip movements and
12
13 syllables or segments that are present in the speech stream (for a recent review see Peelle &
14
15 Sommers, 2015)).
16
17

18
19 The enhancement effect visible speech (consisting of lip movements, tongue
20
21 movements and information from teeth) has on speech in clear and adverse listening
22
23 conditions, has been reported by several studies (e.g. Erber, 1969, 1971; Ma, Zhou, Ross,
24
25 Foxe, & Parra, 2009; Ross et al., 2007; Schwartz et al., 2004; Sumbly & Pollock, 1954).
26
27 Recognizing speech in noise is easier when a visual cue is present than when auditory
28
29 information is presented alone, and has shown to improve recognition accuracy (Tye-Murray,
30
31 Sommers, & Spehar, 2007). Previously, studies have argued that this beneficial effect
32
33 increases as the SNR decreases (Sumbly & Pollack 1954; Erber, 1969; 1975, Callan et al.,
34
35 2003). However, more recent studies have reported that visual enhancement of speech by lip
36
37 movements seems to be largest at "intermediate" SNR's where the auditory input is at a level
38
39 between "perfectly audible" and "completely unintelligible" (Ross et al., 2007; 2008, Ma et
40
41 al., 2009). This has also been reported by Holle et al., (2010), for gestural enhancement of
42
43 speech in noise. Nevertheless, most studies on lip movements as a visual enhancement of
44
45 speech have used stimuli that only showed the lips or lower half of the face (e.g. Callan et al.,
46
47 2003; Ross et al., 2007; Schwartz, Berthommier, & Savariaux, 2004) to eliminate influences
48
49 from the rest of the face or body. This is similar to studies in the domain of gestural
50
51 enhancement of speech in noise, where most studies block the face of the speaker, the mouth,
52
53
54
55
56
57
58
59
60

1
2
3 or just show the torso of the speaker, to eliminate influences from visible speech (e.g.
4
5 Obermeier et al., 2011; Obermeier et al., 2012; Holle et al., 2010).
6
7

8
9 Although there has not been a study that investigated the contribution of visible
10 speech and iconic gestures on speech comprehension in adverse listening conditions, a few
11 studies used both visual articulators in their stimuli. In a fMRI study, Skipper, Goldin-
12 Meadow, Nusbaum, & Small (2009) showed that when clear speech was accompanied by
13 meaningful gestures, there was strong functional connectivity between motor planning and
14 production areas and areas that are thought to mediate semantic aspects of language
15 comprehension. This suggests that the motor system works together with language areas to
16 determine the meaning of those gestures. When just facial information (incl. visual speech)
17 was present, there were strong connectivity patterns between motor planning and production
18 areas and areas that are thought to be involved in phonological processing of speech. These
19 results suggest that information from visible speech is integrated with phonological
20 information, whereas meaningful gestures target semantic aspects of language
21 comprehension. However, it remains unknown how these two articulators interact when both
22 are able to enhance language comprehension in adverse listening conditions.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

40 Two other studies by Kelly et al. (2008) and Hirata & Kelly (2010) examined the
41 effects of lip movements and iconic gestures on auditory learning of second language speech
42 sounds (i.e. prosody and segmental phonology of Japanese). They hypothesized that having
43 both modalities present would benefit learning the most, but found that only lip movements
44 resulted in greater learning. They explain their results by stating that hand gestures might not
45 be suited to learn lower-level acoustic information, such as phoneme contrasts. Again, this
46 study underlines the different relations of visible speech and iconic gestures to speech: visible
47 speech can convey phonological information that can be mapped to the speech signal,
48 whereas gestural information conveys semantic information. It remains unknown how these
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 visual articulators interact when both can enhance language comprehension, such as when
4
5 speech is degraded.
6
7

8 *The present study*
9

10
11 The current study aims to investigate the enhancement effect of iconic gestures and
12 visible speech on degraded speech comprehension, by studying these visual articulators in a
13 joint context. Specifically, we ask what gestural information adds on top of the enhancement
14 of visible speech on degraded speech comprehension, and we test the hypothesis whether the
15 occurrence of two visual articulators (i.e. Speech+VisibleSpeech+Gesture) enhances
16 degraded speech comprehension more than having only visible speech (i.e. Speech +
17 VisibleSpeech) present, or having no visual articulators present (i.e. Speech+Lips blurred).
18 As iconic gestures convey semantic cues that could add to degraded speech comprehension
19 and visible speech conveys phonological cues that could add to degraded speech
20 comprehension, we expect iconic gestures to have an additional enhancement effect on top of
21 the enhancement effect from visible speech.
22
23
24
25
26
27
28
29
30
31
32
33
34
35

36 We hypothesize that the enhancement from visible speech compared to speech alone
37 (i.e. VisualSpeech enhancement: Speech+VisibleSpeech compared to Speech+Lips blurred)
38 will be larger at an intermediate level of degradation compared to a severe level of
39 degradation, allowing a listener to map the phonological information from visible speech to
40 the speech signal. Additionally, we expect the enhancement from iconic gestures on top of
41 visible speech (i.e. Gestural enhancement: Speech+VisibleSpeech+Gesture -
42 Speech+VisibleSpeech) to be largest at an intermediate level of degradation compared to a
43 severe level of degradation, which would indicate that a listener can benefit more from the
44 semantic information from iconic gestures when there are more clear auditory cues to map
45 this information too. Lastly, we predict that the enhancement of both articulators combined
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 (i.e. Double enhancement: Speech+VisibleSpeech+Gesture compared to Speech+Lips
4 blurred), to be largest at an intermediate level of degradation compared to severe degradation.
5
6 Since iconic gestures occur on top of information from visible speech, we expect that that
7
8 should only be possible when enough auditory cues are available to the listeners. This way,
9
10 listeners can benefit from both phonological information that is conveyed by visible speech,
11
12 and from semantic information that is conveyed by iconic gestures.
13
14
15

16
17 Based on previous results on gestural enhancement of degraded speech
18
19 comprehension (Holle et al., 2010, with no information from visible speech present) and
20
21 enhancement of visible speech (e.g. Ross et al., 2007, with no information from iconic
22
23 gestures present), we hypothesize that for double enhancement from both iconic gestures and
24
25 visible speech we find a similar moderate range for optimal integration where our language
26
27 system is weighted to an equal reliance on auditory inputs (speech) and visual inputs (iconic
28
29 gestures and visible speech).
30
31
32

33 **Methods**

34 *Participants*

35
36
37 Twenty right-handed native speakers of Dutch (11 females, $M_{age} = 23;2$ years, $SD = 4.84$)
38
39 participated in this experiment. All participants reported no neurological or language-related
40
41 disorders, no hearing impairments, and had normal or corrected-to-normal vision. None of
42
43 the participants participated in the pre-test (described below). All participants gave informed
44
45 written consent before the start of the experiment and received a financial compensation for
46
47 participation.
48
49
50

51 *Stimulus materials*

1
2
3 We presented participants with 220 short video clips of a female, native Dutch actress
4 uttering a Dutch action verb. The auditory and visual stimuli consisted of the Dutch high
5 frequent action verbs, to make sure that the verbs could easily be coupled with iconic
6 gestures. All video materials were recorded with a JVC GY-HM100 camcorder. Each
7 recording of an action verb resulted in a video length of 2 seconds with an average speech
8 onset of 680ms after video onset. All videos displayed the female actress from head to knees,
9 appearing in the middle of the screen and wearing neutrally colored clothes (grey and black),
10 in front of a unicolored and neutral background. Upon onset of the recording, the actress'
11 starting position was the same for all videos. She was standing straight, facing the camera,
12 with her arms hanging casually on each side of the body. During recording, she was
13 instructed to utter the action verb while making a hand gesture that she found representative
14 for the verb, without receiving feedback from the experimenter. The gestures she made were
15 not instructed by the experimenter but were created by the actress on the fly. If the actress
16 would have received explicit instructions per gesture, the gestures would have looked
17 unnatural or choreographed, and the conscious effort to make a certain gesture could have
18 drawn the attention to the participants explicitly to the gestures. All gestures that
19 accompanied the action verbs were iconic movements for the actions that the verbs depicted
20 (e.g. a drinking gesture resembling a cup that is raised towards the mouth for the verb 'to
21 drink'). The preparation of all gestures started 120 ms after video onset, and the stroke (the
22 meaning bearing part) of the gestures always coincided with the spoken verb.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

48 The auditory sound files were intensity-scaled to 70 dB and de-noised in *Praat*
49 (Boersma & Weenink, 2015). All sound files were re-combined with their corresponding
50 video files in Adobe Premiere Pro. From each video's clear audio file, we created noise-
51 vocoded degraded versions, using a custom-made script in *Praat*. Noise-vocoding effectively
52 manipulates the spectral or temporal detail while preserving the amplitude envelope of the
53
54
55
56
57
58
59
60

1
2
3 speech signal (Shannon, Zeng, Kamath, Wygonski, & Ekelid, 1995). This way, the speech
4
5 signal remains intelligible to a certain extent, depending on the number of vocoding bands,
6
7 with more bands resulting in a more intelligible speech signal. We bandpass filtered each
8
9 sound file between 50 Hz and 8000 Hz, and divided the signal into logarithmically spaced
10
11 frequency bands between 50 and 8000 Hz. This resulted in cutoff frequencies at 50 Hz, 632.5
12
13 Hz and 8000 Hz for 2-band noise-vocoding and 50 Hz, 116.5 Hz, 271.4 Hz, 632.5 Hz, 1473.6
14
15 Hz, 3433.5 Hz and 8000 Hz for 6-band noise-vocoding. We used the frequencies to filter
16
17 white noise in order to obtain six noise bands. We extracted the amplitude envelope of each
18
19 band by using half-wave rectification. We then multiplied the amplitude envelope with the
20
21 noise bands and recombined the bands to form the distorted signal.
22
23
24
25

26 In addition to clear speech, we included 2-band noise-vocoding and 6-band noise-
27
28 vocoding in our experiment. In total, eleven conditions were created for the experiment (see
29
30 Figure 1 for an overview). First, nine conditions were created in a 3 (Speech+Lips blurred,
31
32 Speech+VisibleSpeech, Speech+VisibleSpeech+Gesture) by 3 (2-band noise-vocoding
33
34 ('severe' degradation), 6-band noise-vocoding ('moderate' degradation), clear speech) design.
35
36 Second, we added two extra conditions without sound (VisibleSpeech only, which is similar
37
38 to lip reading, and VisibleSpeech+Gesture) to test how much information participants can
39
40 resolve from visual input by itself. These conditions did not contain an audio file, so
41
42 participants only could utilize the visual input. The final experimental set contained 220
43
44 videos with 220 distinct verbs that were divided over these eleven conditions (20 per
45
46 condition) to test the different contributions of visible speech and gestures to clear speech
47
48 comprehension and in these two degraded listening conditions.
49
50
51
52

53 [Figure 1]
54

55
56 *Pre-test*
57
58
59
60

1
2
3 To ensure that the verbs that we chose could be disambiguated by the iconic gestures that we
4 recorded we conducted a pre-test to examine whether the gestures that the actress made in the
5 video indeed depicted the verbs we matched them with in our audio files. In this experiment,
6
7
8
9
10 twenty native Dutch speakers (10 female, $M_{age} = 22;2$, $SD = 3,3$) with no motor, neurological,
11
12 visual, hearing or language impairments and who did not participate in the main experiment,
13
14 were presented with 170 video stimuli that contained a gesture (not all 220 videos contained a
15
16 gesture and these videos thus were used in the other conditions), but without the audio file
17
18 that contained the verb. All stimuli were presented on a computer screen using Presentation
19
20 software (Neurobehavioral Systems, Inc.), and presented in a different, randomized order per
21
22 participant. First, participants were presented with a fixation cross for 1000 ms, after which
23
24 the video stimulus started playing. After video offset, participants were asked to type down
25
26 the verbs they associated the movement in the video with. After they filled out the verbs, we
27
28 showed them the verb we originally matched it with in our auditory stimuli, and asked the
29
30 participants to indicate on a 7-point-scale (ranging from "does not fit the movement at all" to
31
32 "fits the movement really well") how iconic they found the movement in the video of the verb
33
34 that was presented on the screen. This way, we could ensure that in the main experiment, the
35
36 spoken verbs matched the gesture and participants could use the information from the
37
38 gestures to disambiguate speech. If the gestures were not a good match with the verb, this
39
40 gestural information would not enhance speech comprehension. All participants completed
41
42 the task in approximately 35 minutes and could take self-paced breaks after every 55 items.
43
44
45
46
47

48 The typed answers on the first question of this pre-test ('Which verb do you associate
49
50 with this video?') were used to determine which verbs had to be renamed to a possibly more
51
52 occurring synonym, or which verbs were not recognizable and had to be discarded. We coded
53
54 the answers either as 'correct', when the correct verb or a synonym was given, or as
55
56 'incorrect', when the input consisted of an unrelated verb. The results revealed a mean
57
58
59
60

1
2
3 recognition rate of 59% over all gesture videos. The percentage reported here indicates that
4
5 the gestures are potentially ambiguous in the absence of speech, which is similar to how they
6
7 are perceived in everyday communication (Krauss et al., 1991). Although this seems like a
8
9 low overall consistency between participants, one must note that co-speech gestures, such as
10
11 the iconic co-speech gestures used in these videos, normally occur in the presence of speech,
12
13 and a higher overall percentage would have indicated that the gestures in our video were
14
15 more like pantomimes, which are often understood and produced without speech. Since our
16
17 study aims to understand the possible effects of iconic co-speech gestures on degraded speech
18
19 comprehension, we did not use pantomimes.
20
21
22
23

24 The second question in this pretest targeted the question whether the video depicted
25
26 the verb we matched it with in our auditory stimuli. Out of all videos, there were six videos
27
28 that did not score above a mean rating of '5' on our 7-point scale (ranging from "does not fit
29
30 the movement at all" (1) to "fits the movement really well" (7), indicating that '5' corresponds
31
32 to "fits the movement"). These videos had a mean score of 4.79, 4.05, 4.15, 4.94, 4.89 and
33
34 4.94) and were not used in this experiment. The mean score on 'iconicity' over the other
35
36 videos was 6.1 ($SD = 0.64$). Interestingly, participants indicated after the experiment that
37
38 when they saw the corresponding verb, they often found that verb (which was often a
39
40 synonym of their own answer) fitting for the gesture in the video as well, even though it did
41
42 not always correspond to their own answer. This shows that the mean recognition rate might
43
44 be negatively biased: even though participants may have filled in a different verb in the first
45
46 task, they still highly agreed that the gesture in the video corresponded to the verb (as
47
48 indicated by the score on the second task).
49
50
51
52

53 *Procedure*

54
55
56
57
58
59
60

1
2
3 In our main experiment, participants were tested in a dimly-lit soundproof booth, and seated
4 in front of a computer with headphones on. Before the experiment started, the experimenter
5 gave a short verbal instruction that prepared the participant for the different videos that were
6 going to be presented. All stimuli were presented full screen on a 1650x1080 monitor using
7 Presentation software (Neurobehavioral Systems, Inc.), at a 70 cm distance in front of the
8 participant. A trial started with a fixation cross of 1000, after which the stimulus was played.
9 Then, in a free-recall task, participants were asked to type which verb they thought the actress
10 tried to convey. After the participants typed in their answers, a new trial began after 500 ms.
11 An answer was coded as 'correct' when a participant wrote down the correct verb, or minor
12 spelling mistakes were made. Synonyms or category-related verbs (e.g. 'to bake' for 'to cook')
13 were counted as incorrect.
14
15
16
17
18
19
20
21
22
23
24
25
26
27

28 All participants were presented with a different pseudo-randomization of the stimuli,
29 with the constraint that a specific condition could not be presented more than twice in a row.
30 The stimuli were presented in blocks of 55 trials, and participants could take a self-paced
31 break in between blocks. All participants completed the tasks within 45 minutes.
32
33
34
35
36
37

38 Results

39 As a first step, we employed a 3 x 3 repeated measures analysis of variance with the factors
40 Visual (Speech+Lips blurred; Speech+VisibleSpeech;
41 Speech+VisibleSpeech+Gesture) and Noise-Vocoding Level (2-band noise-vocoding; 6-band
42 noise-vocoding; clear speech) to subject the percentage of correct answers to. Note that we
43 excluded the Visual Only conditions from this analysis (where we only tested VisibleSpeech
44 and VisibleSpeech+Gesture, and not VisibleSpeech+Lips blurred, as this would result in a
45 silent movie with no movement), since this would make our analysis unbalanced. As
46 hypothesized, we found a significant main effect of Noise-Vocoding ($F(2,38) = 1569.78, p <$
47 $.001, \eta^2 = .96$) indicating that the more the speech signal was noise-vocoded, the less correct
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 answers were given by the participants. We also found a main effect of VisualArticulator
4
5 ($F(2,38) = 504.284, p < .001, \eta^2 = .98$) indicating that the more visual articulators were
6
7 added to the signal, the more correct answers were given. In addition, we found a significant
8
9 interaction between Noise-Vocoding level and VisualArticulator ($F(4,76) = 194.11, p < .001,$
10
11 $\eta^2 = .91$), which seemed to be driven by the relatively higher amount of correct responses in
12
13 the 6-band noise-vocoding condition compared to the other speech conditions (see Figure 2
14
15 for the percentages of correct responses per condition).
16
17

18
19 [Figure 2]
20
21

22 To further investigate this interaction, we compared the differences between and
23
24 within the different noise-vocoding levels and visual articulators in a separate analysis. This
25
26 analysis allowed us to compare the enhancement driven by different visual articulators as
27
28 well as compare those enhancement effects between noise-vocoding levels. In comparing the
29
30 enhancement from the different visual articulators, we recognized that calculating the
31
32 absolute gain in terms of difference scores is limited in appropriately characterizing the
33
34 maximum gain per condition. This is because there is an inverse relationship that exists
35
36 between the performance in the Speech+Lips blurred and Speech+VisibleSpeech conditions
37
38 and the maximum benefit that is derived when calculating the enhancement of the different
39
40 visual articulators (see Grant & Walden, 1996). For example, we found a 2.75% recognition
41
42 rate for Speech+Lips blurred in 2-band noise-vocoding as compared to 11.75% in 6-band
43
44 noise-vocoding. The maximum gain possible on the basis of pure difference scores would
45
46 therefore be 97.75% for 2-band noise-vocoding, and 88.25% for 6-band noise-vocoding,
47
48 which would be hard to compare, since the maximal gain that is possible in 2-band noise-
49
50 vocoding is larger than in 6-band noise-vocoding.
51
52
53
54
55
56
57
58
59
60

Therefore, to avoid possible floor effects and in keeping with previous studies, such as Sumbly & Pollack (1954), we controlled for this by defining three difference scores ((A-B/100-B) (i.e., enhancement types)) for a) VisibleSpeech enhancement: Speech+VisibleSpeech - Speech+Lips blurred; b) Gestural enhancement: Speech+VisibleSpeech+Gesture - Speech+VisibleSpeech; and c) Double enhancement: Speech+VisibleSpeech+Gesture - Speech+Lips blurred, (see Ross et al., 2007 for a discussion of other calculation methods) divided by the maximal possible enhancement (for VisibleSpeech enhancement: 100 -Speech+Lips blurred; for Gestural enhancement: 100 - Speech+VisibleSpeech; for Double enhancement: 100 - Speech+Lips blurred). We subjected these outcomes to a repeated measures ANOVA with the factors Noise-Vocoding (2-band, 6-band, clear) and EnhancementType (VisibleSpeech enhancement, Gestural enhancement, Double enhancement). Our analysis revealed a main effect of Noise-Vocoding ($F(2,38) = 320.23, p < .001, \text{partial } \eta^2 = .94$), indicating that the more degraded the signal was, the less enhancement was present. Moreover, we found a main effect of EnhancementType ($F(1.06, 20.19) = 276.74, p < .001, \text{partial } \eta^2 = .94, \text{Greenhouse-Geisser corrected}$), indicating that the more visual information was present, the more participants answered correctly. Importantly, we found a significant interaction between EnhancementType and Noise-Vocoding ($F(1.97, 37.37) = 102.65, p < .001, \text{partial } \eta^2 = .84, \text{Greenhouse-Geisser corrected}$). Pairwise comparisons (all Bonferroni corrected) showed a significant difference between Gestural enhancement and VisibleSpeech enhancement in both the 2-band noise-vocoding condition ($t(19) = 9.41, p_{\text{bon}} < .001$ and the 6-band noise-vocoding condition ($t(19) = 12.94, p_{\text{bon}} < 0.001$) Furthermore, the difference between Gestural enhancement and VisibleSpeech enhancement was larger for 6-band noise-vocoding than 2-band noise-vocoding ($F(1,19) = 64.48, p_{\text{bon}} < .001, \text{partial } \eta^2 = .77$). Finally, Double enhancement was larger at 6-band noise-vocoding than in 2-band noise-vocoding and ($t(19) = -10.035, p_{\text{bon}} < .001$) (see Figure 3).

1
2
3 Pairwise comparisons showed a significant difference in VisibleSpeech enhancement and
4 Double enhancement in both 2-band noise-vocoding ($t(19) = 12.47, p_{\text{bon}} < .001$) and 6-band
5 noise-vocoding ($t(19) = 20.79, p_{\text{bon}} < .001$). This difference between VisibleSpeech
6 enhancement and Double enhancement was larger in 6-band noise-vocoding than in 2-band
7 noise-vocoding ($F(1,19) = 163.20, p_{\text{bon}} < .001, \text{partial } \eta^2 = .90$). Additionally, pairwise
8 comparisons showed a significant difference in Gestural enhancement and Double
9 enhancement in both 2-band noise-vocoding ($t(19) = 3.36, p_{\text{bon}} < 0.01$) and 6-band noise-
10 vocoding ($t(19) = 7.79, p_{\text{bon}} < .001$), which was again largest in 6-band noise-vocoding
11 ($F(1,19) = 30.44, p_{\text{bon}} < .001, \text{partial } \eta^2 = .62$).

12
13
14 [figure 3]

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Initially, we did not include the two Visual Only conditions (VisibleSpeech only, VisibleSpeech+Gesture) in our main analysis, because they would create an unbalanced design for analyzing all conditions together. However, these conditions were still of interest to determine how much information participants could obtain from visual input alone without speech being present. Therefore, we first tested the difference between the two separate Visual only conditions by means of a paired samples t-test. We found a significant difference between VisibleSpeech only and VisibleSpeech+Gesture ($t(19) = 15.12, p < 0.001$), indicating that response accuracy was higher for trials containing both visible speech and gestures, compared to videos that just contained visible speech (see Figure 2). Subsequently, we compared this difference between VisibleSpeech+Gesture and VisibleSpeech Only (i.e. gestural enhancement, computed as the difference between (VisibleSpeech+Gesture - VisibleSpeech only/100 - VisibleSpeech Only)) to the Gestural enhancement in the context of speech (Speech+VisibleSpeech+Gesture - Speech+VisibleSpeech/100 - Speech+VisibleSpeech) both in the 6-band and 2-band noise-vocoding conditions (see Figure 3). Our analysis revealed a significant difference between Gestural enhancement in the Visual

1
2
3 Only conditions and Gestural enhancement in 6-band noise-vocoding, ($t(19) = -3.23$, $p_{\text{bon}} <$
4 0.05), but not compared to 2-band noise-vocoding condition ($t(19) = 1.1$, $p_{\text{bon}} > .1$). These
5 results confirmed that Gestural enhancement in 6-band noise-vocoding was significantly
6 greater compared to 2-band noise-vocoding and compared to Gestural enhancement in the
7 Visual only conditions. However, Gestural enhancement in the Visual Only conditions was
8 not larger than Gestural enhancement in 2-band noise-vocoding, indicating that if there are no
9 longer reliable auditory cues available (as in 2-band noise-vocoding), comprehension might
10 be comparable to when there is no auditory input at all (as in Visual Only conditions).
11
12
13
14
15
16
17
18
19

20 We have explored the error types per Visual Articulator, per Noise-Vocoding level.
21 However, since the percentage of error type in some conditions was very low, we did not
22 subject these error types to a statistical analysis. To test for possible confounding effects of
23 fatigue or learning, we also compared the amount of correct answers per block. We found no
24 difference between the different blocks in the experiment in correct answers ($p > .1$).
25
26
27
28
29
30
31

32 **Discussion**

33
34
35 The first aim of our study was to reveal whether and to what extent iconic gestures can
36 contribute on top of information from visible speech to enhance degraded speech
37 comprehension, and whether double enhancement from both visual articulators is more
38 beneficial for comprehension than having just visible speech present as a visual articulator, or
39 having no visual articulators present. Whereas previous studies have approached the
40 contribution of these two visual articulators only separately, we investigated the enhancement
41 effects of iconic gestures and visible speech in a joint context. Since iconic gestures can
42 provide information on a semantic level and visible speech can provide information on a
43 phonological level, we expected an additive effect of gestures on top of the enhancement of
44 visible speech during degraded speech comprehension. Our data indeed showed that while
45 perceiving degraded speech in a visual context, listeners benefit most from having both
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 visible speech and iconic gestures present, as compared to having just visible speech present,
4
5 or having only auditory information present. Here, gestures provide an additional benefit on
6
7 top of the enhancement of visible speech.
8
9

10
11 Our second aim was to demarcate the noise conditions under which this double
12
13 enhancement from both visible speech and iconic gestures in the context of visible speech
14
15 add the most to degraded speech comprehension. Our data suggests that at an moderate level
16
17 of noise-vocoding (6-band), there is an optimal range for maximal multimodal integration
18
19 where listeners can benefit most from the visual information. The enhancement effects of
20
21 VisibleSpeech enhancement, Gestural enhancement and Double enhancement were
22
23 significantly larger in 6-band noise-vocoding than in 2-band noise-vocoding or in the Visual
24
25 Only conditions. However, we did not find a difference in Gestural enhancement between 2-
26
27 band noise-vocoding and Visual Only conditions. Taken together, our results showed that at
28
29 this optimal enhancement level of 6-band noise-vocoding, auditory cues were still moderately
30
31 reliable and listeners were able to combine and integrate information from both visible
32
33 speech and iconic co-speech gestures to aid in comprehension, resulting in an additive effect
34
35 of double, multimodal enhancement from visible speech and iconic gestures. Here, semantic
36
37 information from iconic gestures adds to the mapping between the speech signal and
38
39 phonological information that is derived from lip movements in visible speech. Below we
40
41 will discuss these results in more detail.
42
43
44
45

46
47 In line with previous research, we found a significant benefit of adding information
48
49 from visible speech to the speech signal (VisibleSpeech enhancement), in response to stimuli
50
51 from both noise-vocoding levels (e.g. Sumbly & Pollack, 1954). This benefit from solely
52
53 visible speech was significantly larger at a moderate level of noise-vocoding (6-band) than at
54
55 a severe level of noise-vocoding (2-band). Previously, it has been suggested that the benefit
56
57 from visible speech continues to increase as the information that is available from auditory
58
59
60

1
2
3 inputs decreases (Sumbly & Pollack 1954; Erber, 1969; 1975; Meredith & Stein, 1983), as
4
5 would be predicted by the principle of inverse effectiveness. However, recent studies have
6
7 argued that there are minimal levels of auditory information necessary before recognition
8
9 accuracy can be most enhanced by congruent visible input (Ross et al., 2007). Our data
10
11 concurs with this latter idea, by finding an optimal range for multimodal integration and
12
13 enhancement, where auditory cues are moderately reliable, and enhancement from visible
14
15 speech has its maximal effect.
16
17

18
19 Importantly, the current results provide novel evidence by showing that iconic
20
21 gestures can enhance this benefit from visible speech even more: We found a significant
22
23 difference between Gestural enhancement (Speech+VisibleSpeech+Gesture –
24
25 Speech+VisibleSpeech) and VisibleSpeech enhancement (Speech+VisibleSpeech –
26
27 Speech+Lips blurred) at both noise-vocoding levels. In addition, we found significant
28
29 differences between Double enhancement and Gestural enhancement, as well as significant
30
31 differences between Double and VisibleSpeech enhancement at both noise-vocoding levels.
32
33 Our results therefore suggest that although both visual modalities enhance degraded speech
34
35 comprehension, having both iconic gestures and visible speech present (Double
36
37 enhancement) in the input enhances speech comprehension most. This is in line with previous
38
39 literature on the benefits of gestures in language processing and theories of communication
40
41 that postulate that multimodal information combines with speech to aid language
42
43 comprehension (McNeill 1992; Clark 1996; Goldin-Meadow 2003, see for a review Kelly,
44
45 Manning, & Rodak, 2008). Interestingly, the enhancement of both visual articulators (Double
46
47 enhancement) was significantly larger than VisibleSpeech enhancement at both noise-
48
49 vocoding levels. This suggests, in line with previous research, that gestures are actively
50
51 processed and integrated with the speech signal (Kelly et al., 2010; Kendon, 2004), even
52
53 under conditions where visible speech is visible (also see Holler et al., 2014).
54
55
56
57
58
59
60

1
2
3 It is important to note that this double enhancement of both iconic gestures and visible
4 speech is in itself still a product of integrating the auditory (speech) and visual input (iconic
5 gestures and visual speech), and not a result of our participants focusing solely on the visual
6 input. The gain in recognition accuracy in our Visual only (VisibleSpeech+Gesture – Visible
7 Speech only) conditions was significantly smaller than the gain we found in the moderate
8 noise (6-band noise-vocoding) condition. The fact that we did not find a similar difference in
9 enhancement between the Visual only conditions and the severe degradation (2-band noise-
10 vocoding) condition suggests that in 2-band noise-vocoding, visible speech cannot be reliably
11 matched to phonological information in the speech signal, and listeners might have focused
12 more on semantic information from gestures to map to the speech signal for disambiguation.
13 As a result, listeners seem to lose the additive effect of double enhancement from visible
14 speech and gestures for speech comprehension in 2-band noise-vocoding because there are
15 not enough reliable auditory cues present in the speech signal to map visible speech too.
16 Consequently, in 2-band noise-vocoding and Visual Only conditions, Gestural enhancement
17 is solely consisting of what can be picked up semantically from the gesture, in addition to
18 information from visible speech. Taken together, we therefore suggest that listeners are only
19 able to benefit from double enhancement from both gestures and visible speech when
20 auditory information is still moderately reliable, to facilitate a binding that integrates
21 information from visible speech, gestures and speech into one coherent percept that exceeds a
22 certain reliability threshold, forming an optimal range where maximal multimodal integration
23 and enhancement can occur.
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

49
50 In earlier work on the contribution of visible speech and hand gestures to learning
51 non-native speech sounds, Kelly et al. (2008) argued that lip and mouth movements help in
52 auditory encoding of speech, whereas hand gestures only can help to understand the meaning
53 of words in the speech stream when the auditory signal is correctly encoded. Based on their
54
55
56
57
58
59
60

1
2
3 results, Kelly et al. (2008) argue that the benefits of multimodal input target different stages
4 of linguistic processing. Here, mouth movements seem to aid during phonological stages,
5 whereas hand gestures aid during semantic stages, which, according to the authors, fits with
6 McNeill's (1992) interpretation of speech and gesture forming an integrated system during
7 language comprehension.
8
9
10
11
12

13
14
15 The results from the present study indeed concur with the idea that speech and gesture
16 form an integrated system and that the benefits of multimodal input target different stages of
17 linguistic processing. Indeed, visible speech possibly plays a significant role during auditory
18 encoding of speech, but according to our current results, iconic gestures not only benefit
19 comprehension when auditory information can be correctly encoded and understood, but also
20 benefit comprehension under adverse listening conditions (cf. Kelly et al., 2008). Even in 2-
21 band noise-vocoding, when auditory cues are no longer reliable and correct encoding of the
22 auditory input is difficult, gestures significantly enhance comprehension. Instead, our data
23 suggests that when encoding of auditory information is difficult or when auditory cues are
24 largely unreliable, listeners are mostly driven by the semantic information from gestures to
25 guide comprehension, which can be beneficial to disambiguate the auditory cues. However,
26 when auditory cues are moderately reliable and there are enough auditory cues available to
27 map the phonological information of visible speech to, listeners can benefit from a 'double'
28 multimodal enhancement from the two visual articulators, integrating both the phonological
29 information from visible speech and semantic information from gestures with the speech
30 signal. This, in turn, results in an additive effect of the semantic information provided by
31 iconic gestures on top of the phonological information from visible speech. However, in 2-
32 band noise vocoding where phonological information from visible speech can no longer be
33 reliably matched to the speech signal, listeners lose this additive double enhancement effect
34 of visible speech and iconic gestures, and mostly utilize the semantic information from
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 gestures (i.e. Gestural enhancement) to resolve the form of the speech signal. Based on these
4
5 results, we suggest that at least in adverse listening conditions where auditory cues are no
6
7 longer reliable, language processing might be more driven by semantic information that is
8
9 abstracted from iconic co-speech gestures.
10

11
12 Our findings suggest that the use of iconic gestures can play a pivotal role in natural
13
14 face-to-face communication: gestural information can help to access the meaning of a word
15
16 to resolve the form of the speech signal when a listening situation is challenging, such as in
17
18 noise. One limitation of our work can be that our actress uttered the stimuli in a setting with
19
20 optimal listening conditions, without any noise. We edited her auditory input after recording,
21
22 to test the effect of different noise-vocoding bands. In this regard, it is important to note that
23
24 in a natural adverse listening condition, our speaker would have probably adjusted her
25
26 articulatory movements to optimally communicate her message. This effect has been
27
28 previously described as the Lombard effect, which refers to the tendency of speakers to
29
30 increase their vocal effort when speaking in noise to enhance the audibility of their voice
31
32 (which is not limited to loudness, but also to the length of phonemes and syllables, speech
33
34 rate and pitch, amongst others) (Lombard, 1911). Alternatively, this could also have an effect
35
36 on the production of iconic co-speech gestures as well: for example producing a larger iconic
37
38 gesture in an adverse listening condition could have resulted in a larger co-speech gesture
39
40 than in clear speech. Future research could test this possibility by recording stimuli in an
41
42 adverse listening condition and present these videos to participants, to increase ecological
43
44 validity. A second limitation of our study can be that our participants were only presented
45
46 with single action verbs. Future research could investigate whether presenting these verbs in a
47
48 sentence context might have an influence on how much a listener depends on different visual
49
50 articulators. In addition, future endeavors could consider that natural face-to-face
51
52 communication does not only consists of a binding of speech and visual information from
53
54
55
56
57
58
59
60

1
2
3 gestures and visible speech. Instead, research can tap into the influence of other nonverbal
4
5 behavior (such as head and brow movements, see e.g., Kraemer & Swerts, 2007) and their
6
7 co-occurrence with visible speech and gesture to fully understand the optimal conditions for
8
9 visual enhancement of speech in adverse listening conditions. This, in turn, can further
10
11 elucidate the results from the current study, but also inform debates on audiovisual training
12
13 for both clinical populations and educational instruction. Finally, replicating the effects
14
15 found in this study with hearing-impaired populations will provide a better diagnosis of their
16
17 speech comprehension in ecologically valid contexts (i.e., in a multimodal context). This in
18
19 turn could inform debates on audiovisual training for both clinical populations and
20
21 educational instruction.
22
23

24 25 26 **Acknowledgements**

27
28 This research was supported by Gravitation Grant 024.001.006 of the Language in Interaction
29
30 Consortium from Netherlands Organization for Scientific Research. We thank two
31
32 anonymous reviewers for their helpful comments and suggestions that helped to improve the
33
34 paper. We are very grateful to Nick Wood, for helping us in editing the video stimuli and to
35
36 Gina Ginos, for being the actress in the videos.
37
38
39

40 41 **References**

- 42
43 Beattie, G., & Shovelton, H. (1999). Do iconic hand gestures really contribute anything to the
44
45 semantic information conveyed by speech? An experimental investigation. *Semiotica, 1*,
46
47 32-49.
48
49
50 Beattie, G., & Shovelton, H. (1999). Mapping the Range of Information Contained in the
51
52 Iconic Hand Gestures that Accompany Spontaneous Speech. *Journal of Language and*
53
54 *Social Psychology, 18*(4), 438–462. doi:10.1177/0261927X99018004005
55
56
57
58
59
60

- 1
2
3 Beattie, G., & Shovelton, H. (2002). An experimental investigation of some properties of
4
5 individual iconic gestures that mediate their communicative power. *British Journal of*
6
7 *Psychology*, 93(2), 179–192. doi:10.1348/000712602162526
8
9
- 10 Boersma, P., & Weenink, D. (2015). Praat: doing phonetics by computer. [Computer
11
12 program]. Version 6.0.19, retrieved 13 June 2016 from <http://www.praat.org/>
13
14
- 15 Callan, D. E., Jones, J. A., Munhall, K., Callan, A. M., Kroos, C., & Vatikiotis-Bateson, E.
16
17 (2003). Neural processes underlying perceptual enhancement by visual speech gestures.
18
19 *Neuroreport*, 14(17), 2213–8. doi:10.1097/01.wnr.0000095492.38740.8f
20
21
- 22 Clark, H. H. (1996). *Using Language*. Cambridge University Press.
23
- 24 Erber, N. P. (1969). Interaction of audition and vision in the recognition of oral speech
25
26 stimuli. *Journal of Speech and Hearing Research*, 12(2), 423–5.
27
28
- 29 Erber, N. P. (1971). Auditory and audiovisual reception of words in low-frequency noise by
30
31 children with normal hearing and by children with impaired hearing. *Journal of Speech*
32
33 *and Hearing Research*, 14(3), 496–512.
34
35
- 36 Goldin-Meadow, S. (2005). *Hearing Gesture: How Our Hands Help Us Think*. Harvard
37
38 University Press.
39
40
- 41 Grant, K. W., & Walden, B. E. (1996). Evaluating the articulation index for auditory-visual
42
43 consonant recognition. *The Journal of the Acoustical Society of America*, 100(4), 2415–
44
45 24.
46
47
- 48 Hirata, Y., & Kelly, S. D. (2010). Effects of Lips and Hands on Auditory Learning of
49
50 Second-Language Speech Sounds. *Journal of Speech, Language and Hearing Research*,
51
52 53(4), 298–310.
53
54
- 55 Holle, H., & Gunter, T. C. (2007). The role of iconic gestures in speech disambiguation: ERP
56
57 evidence. *Journal of Cognitive Neuroscience*, 19(7), 1175–92.
58
59
60

1
2
3 doi:10.1162/jocn.2007.19.7.1175
4

5 Holle, H., Obermeier, C., Schmidt-Kassow, M., Friederici, A. D., Ward, J., & Gunter, T. C.
6
7 (2012). Gesture facilitates the syntactic analysis of speech. *Frontiers in Psychology*,
8
9 3(3), 74. doi:10.3389/fpsyg.2012.00074
10

11
12 Holle, H., Obleser, J., Rueschemeyer, S.-A., & Gunter, T. C. (2010). Integration of iconic
13
14 gestures and speech in left superior temporal areas boosts speech comprehension under
15
16 adverse listening conditions. *NeuroImage*, 49(1), 875–84.
17
18 doi:10.1016/j.neuroimage.2009.08.058
19
20

21
22 Holler, J., Kelly, S., Hagoort, P., & Ozyurek, A. (2010). When gestures catch the eye : The
23
24 influence of gaze direction on co-speech gesture comprehension in triadic
25
26 communication. In *the 34th Annual Meeting of the Cognitive Science Society (CogSci*
27
28 *2012)* (pp. 467-472). Cognitive Society.
29
30

31
32 Holler, J., Schubotz, L., Kelly, S., Hagoort, P., Schuetze, M., & Özyürek, A. (2014). Social
33
34 eye gaze modulates processing of speech and co-speech gesture. *Cognition*, 133(3),
35
36 692–7. doi:10.1016/j.cognition.2014.08.008
37

38
39 Holler, J., Shovelton, H., & Beattie, G. (2009). Do Iconic Hand Gestures Really Contribute to
40
41 the Communication of Semantic Information in a Face-to-Face Context? *Journal of*
42
43 *Nonverbal Behavior*, 33(2), 73–88. doi:10.1007/s10919-008-0063-9
44

45
46 Hoskin, J., & Herman, R. (2001). The communication, speech and gesture of a group of
47
48 hearing-impaired children. *International Journal of Language & Communication*
49
50 *Disorders / Royal College of Speech & Language Therapists*, 36 Suppl, 206–9.
51

52
53 Hostetter, A. B. (2011). When do gestures communicate? A meta-analysis. *Psychological*
54
55 *Bulletin*, 137(2), 297–315. doi:10.1037/a0022128
56

57
58 Kelly, S. D., Barr, D. J., Church, R. B., & Lynch, K. (1999). Offering a hand to pragmatic
59
60

1
2
3 understanding: The role of speech and gesture in comprehension and memory. *J.Mem.*
4
5 *Lang*, 40, 577–592.
6

7
8 Kelly, S. D., Manning, S. M., & Rodak, S. (2008). Gesture Gives a Hand to Language and
9
10 Learning: Perspectives from Cognitive Neuroscience, Developmental Psychology and
11
12 Education. *Language and Linguistics Compass*, 2(4), 569–588. doi:10.1111/j.1749-
13
14 818X.2008.00067.x
15

16
17 Kelly, S. D., Ozyürek, A., & Maris, E. (2010). Two sides of the same coin: speech and
18
19 gesture mutually interact to enhance comprehension. *Psychological Science*, 21(2), 260–
20
21 7. doi:10.1177/0956797609357327
22

23
24 Kelly, S., Healey, M., Özyürek, A., & Holler, J. (2015). The processing of speech, gesture,
25
26 and action during language comprehension. *Psychonomic Bulletin & Review*, 22(2),
27
28 517–23. doi:10.3758/s13423-014-0681-7
29

30
31 Kelly, S., Hirata, Y., Simester, J., Burch, J., Cullings, E., & Demakakos, J. (2008). Effects of
32
33 hand gesture and lip movements on auditory learning of second language speech sounds.
34
35 *The Journal of the Acoustical Society of America*, 6(10), 2357–2362.
36
37 doi:10.1121/1.2933816
38

39
40
41 Kendon, A. (2004). *Gesture: Visible Action as Utterance*. Cambridge University Press.
42

43
44
45
46 Krahmer, E., & Swerts, M. (2007). The effects of visual beats on prosodic prominence:
47
48 Acoustic analyses, auditory perception and visual perception. *Journal of Memory and*
49
50 *Language*, 57(3), 396-414.
51

52
53 Krauss, R M., Morrel-Samuels, P. & Colsante, C. (1991). Do conversational hand gestures
54
55 communicate? *Journal of Personality and Social Psychology*, 61 (5), 743 - 754.
56

57
58 Lombard, E. (1911). Le signe de l'elevation de la voix. *Annals Maladiers Oreille, Larynx*,
59
60

1
2
3 *Nez, Pharynx*, 37, 101–119.

4
5 Ma, W. J., Zhou, X., Ross, L. A., Foxe, J. J., & Parra, L. C. (2009). Lip-reading aids word
6
7 recognition most in moderate noise: a Bayesian explanation using high-dimensional
8
9 feature space. *PloS One*, 4(3), e4638. doi:10.1371/journal.pone.0004638

10
11
12 McNeill, D. (1992). *Hand and mind: What gestures reveal about thought*. Chicago University
13
14 Press.

15
16
17 Obermeier, C., Dolk, T., & Gunter, T. C. (2012). The benefit of gestures during
18
19 communication: evidence from hearing and hearing-impaired individuals. *Cortex; a*
20
21 *Journal Devoted to the Study of the Nervous System and Behavior*, 48(7), 857–70.
22
23 doi:10.1016/j.cortex.2011.02.007

24
25
26 Obermeier, C., Holle, H., & Gunter, T. C. (2011). What iconic gesture fragments reveal about
27
28 gesture-speech integration: when synchrony is lost, memory can help. *Journal of*
29
30 *Cognitive Neuroscience*, 23(7), 1648–63. doi:10.1162/jocn.2010.21498

31
32
33 Özyürek, A. (2014). Hearing and seeing meaning in speech and gesture: insights from brain
34
35 and behaviour. *Philosophical Transactions of the Royal Society of London. Series B,*
36
37 *Biological Sciences*, 369(1651), 20130296. doi:10.1098/rstb.2013.0296

38
39
40 Peelle, J. E., & Sommers, M. S. (2015). Prediction and constraint in audiovisual speech
41
42 perception. *Cortex; a Journal Devoted to the Study of the Nervous System and Behavior*,
43
44 68, 169–81. doi:10.1016/j.cortex.2015.03.006

45
46
47 Rogers, W. T. (1978). the Contribution of Kinesic Illustrators Toward the Comprehension of
48
49 Verbal Behavior Within Utterances. *Human Communication Research*, 5(1), 54–62.
50
51 doi:10.1111/j.1468-2958.1978.tb00622.x

52
53
54 Ross, L. a, Saint-Amour, D., Leavitt, V. M., Javitt, D. C., & Foxe, J. J. (2007). Do you see
55
56 what I am saying? Exploring visual enhancement of speech comprehension in noisy
57
58
59
60

1
2
3 environments. *Cerebral Cortex (New York, N.Y. : 1991)*, 17(5), 1147–53.

4
5 doi:10.1093/cercor/bhl024

6
7
8 Schwartz, J.-L., Berthommier, F., & Savariaux, C. (2004). Seeing to hear better: evidence for
9
10 early audio-visual interactions in speech identification. *Cognition*, 93(2), B69–78.

11
12 doi:10.1016/j.cognition.2004.01.006

13
14
15 Shannon, R., Zeng, F.-G., Kamath, V., Wygonski, J., & Ekelid, M. (1995). Speech
16
17 Recognition with Primarily Temporal Cues. *Science*, 270(5234), 303–304.

18
19
20 Skipper, J. I., Goldin-Meadow, S., Nusbaum, H. C., & Small, S. L. (2009). Gestures
21
22 orchestrate brain networks for language understanding. *Current Biology : CB*, 19(8),
23
24 661–7. doi:10.1016/j.cub.2009.02.051

25
26
27 Sumbly, W. H., & Pollock, I. (1954). Visual Contribution to Speech Intelligibility in Noise.
28
29 *The Journal of the Acoustical Society of America*, 26(2), 212. doi:10.1121/1.1907309

30
31
32 Tye-Murray, N., Sommers, M. S., & Spehar, B. (2007). Audiovisual integration and
33
34 lipreading abilities of older adults with normal and impaired hearing. *Ear and Hearing*,
35
36 28(5), 656–68. doi:10.1097/AUD.0b013e31812f7185

37
38
39 Wu, Y. C., & Coulson, S. (2007). Iconic gestures prime related concepts: an ERP study.
40
41 *Psychonomic Bulletin & Review*, 14(1), 57–63. Retrieved from
42
43 <http://www.ncbi.nlm.nih.gov/pubmed/17546731>

1
2
3 **Figure captions**
4

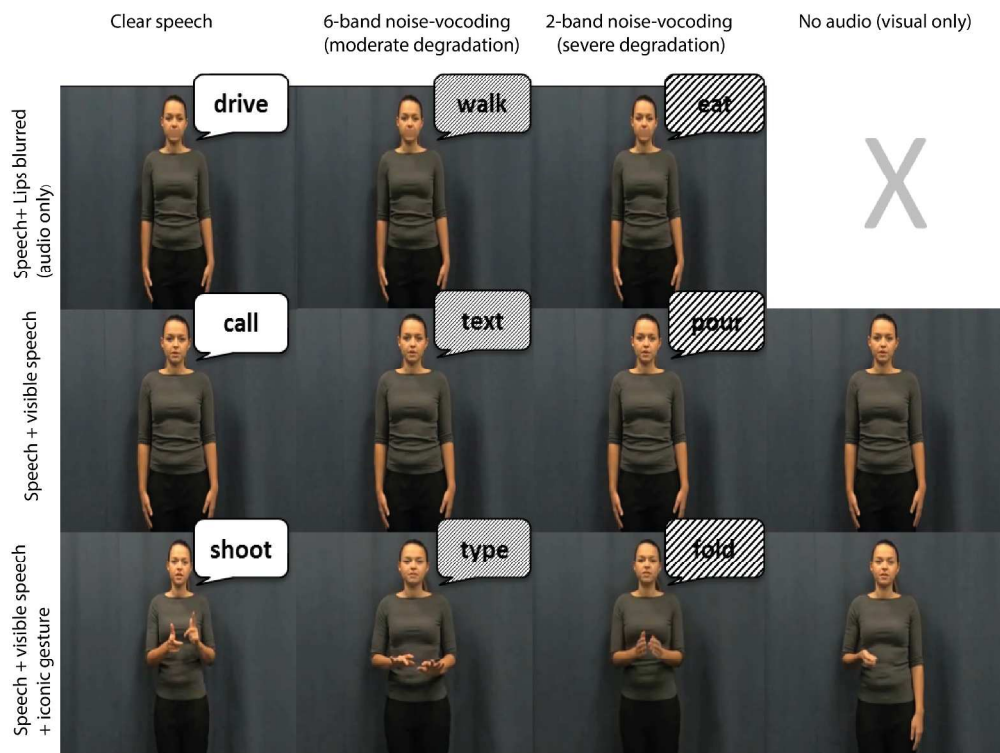
5 **Figure 1:** *Overview of the design and conditions used in the experiment.*
6

7
8 **Figure 2:** *Percentage of correctly identified verbs (% correct) per condition. Error bars*
9 *represent standard deviations.*
10

11
12 **Figure 3:** *Enhancement effect ($A-B/100-B$) corrected for floor effects. Error bars represent*
13 *standard deviations.*
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

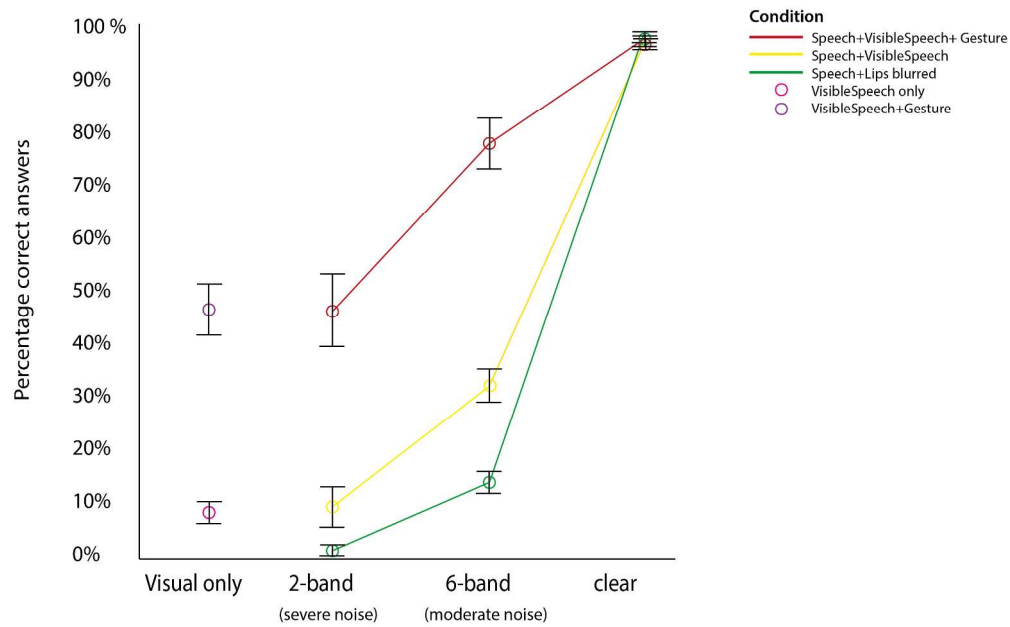
For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



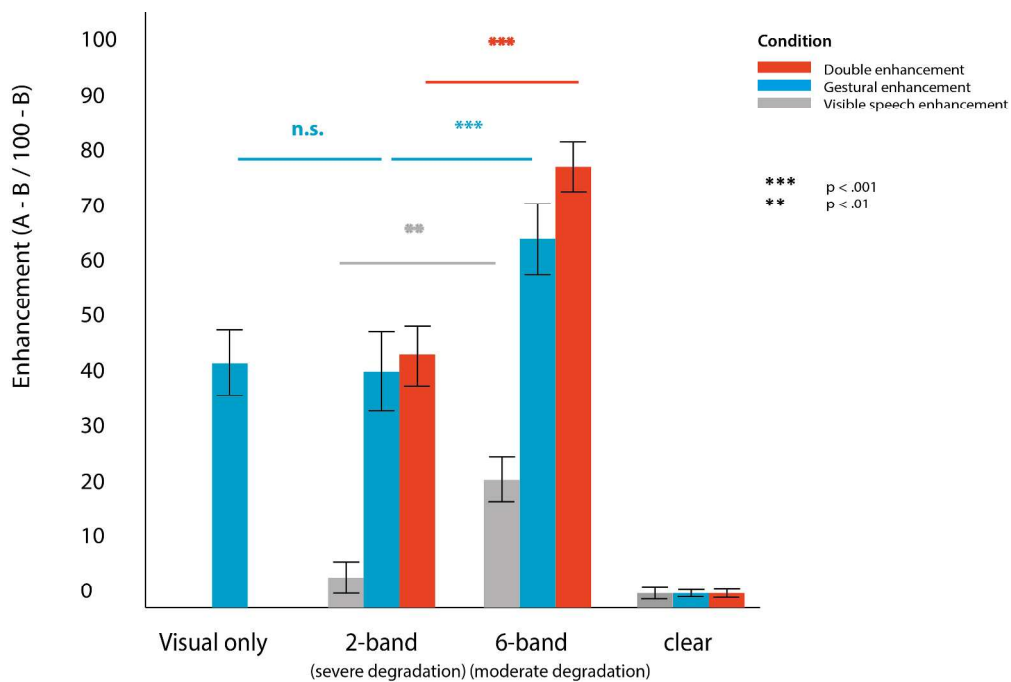
Overview of the design and conditions used in the experiment
365x282mm (300 x 300 DPI)

Review



Percentage of correctly identified verbs (% correct) per condition. Error bars represent standard deviations.
268x220mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Enhancement effect (A-B/100-B) corrected for floor effects. Error bars represent standard deviations. 256x170mm (300 x 300 DPI)

Review